

# ANALYSIS OF SIGNAL TRANSDUCTION NETWORKS USING ACTIVATION RATIOS

by

Francisco Javier Femenia

S. M. Chemical Engineering Practice, Massachusetts Institute of Technology, June 2000

B. S. Chemical Engineering, University of California, Berkeley, May 1997

Submitted to the Department of Chemical Engineering  
in Partial Fulfillment of the Requirements for the Degree of

DOCTOR OF PHILOSOPHY IN CHEMICAL ENGINEERING

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

February 2004

© Massachusetts Institute of Technology 2004. All rights reserved.

Author .....  
Department of Chemical Engineering  
December 5, 2003

Certified by .....  
Gregory Stephanopoulos  
Professor of Chemical Engineering  
Thesis Supervisor

Accepted by .....  
Daniel Blankschtein  
Professor of Chemical Engineering  
Chairman, Committee for Graduate Students



# ANALYSIS OF SIGNAL TRANSDUCTION NETWORKS USING ACTIVATION RATIOS

by

Francisco Javier Femenia

Submitted to the Department of Chemical Engineering  
on December 5, 2003 in Partial Fulfillment of the Requirements  
for the Degree of Doctor of Philosophy in Chemical Engineering

## ***Abstract***

The molecular processes by which information is incorporated and distributed within a cell are termed signal transduction. These pathways allow cells to interact with each other and with their environments and are critical to the proper cellular function in a variety of contexts. Previously developed methods for analyzing signaling networks have been largely ignored, most likely due to their mathematical complexity and difficulty in application. A novel analysis framework was developed to assist in the examination of signaling networks, both to facilitate the reconstruction of previously undetermined pathways as well as to quantitatively characterize interactions between components.

This approach, termed activation ratio analysis, involves the ratio between active and inactive forms of signaling intermediates at steady state. The activation ratio of an intermediate is shown to depend linearly upon the concentration of the activating enzyme. The slope of the line is defined as the activation factor, and is determined by the kinetic parameters of activation and inactivation. The mathematical functionality of the activation ratio changes for other signaling network arrangements. It is therefore possible to extract the original network structure from a set of measured activation ratios, with activation factors yielding a measure of activation potential between intermediates.

This framework was tested using computational simulations of a small-scale interconnected network, cascades with feedback, and in the presence of experimental noise. In the process, additional tools were developed to automate and evaluate the analysis. The theoretical concepts were also applied to examine the modification cycles of the protein Erk2 by the kinase Mek1 and phosphatases PTP1B and PP2A *in vitro*. Limitations in the accuracy of measurements and experimental setup resulted in high uncertainty in the calculated distribution of Erk states, preventing a quantitative analysis of this system. Nevertheless, qualitative predictions from activation ratio analysis were verified, in particular, the influence of component concentrations on the Erk distribution at steady state. While these issues indicate considerations for future experimental studies, they do not override the ability of activation ratio analysis to investigate signaling networks, where description of interactions in the whole system is more important than detailed examination of the individual steps.

Thesis Supervisor: Gregory Stephanopoulos  
Title: Professor of Chemical Engineering

For Dan and Lia  
In loving memory



## ***Acknowledgements***

The process of the Ph. D. is a long one, but fortunately, not a lonely one. Although the work described here was performed independently, I never would have been able to complete it if not for the support of many wonderful people. I owe a great many my sincere gratitude:

My advisor, Gregory Stephanopoulos, without which this work could never have been done. Thanks for providing a fantastic work environment, and for support, pressure, and patience when I needed them. My thesis committee members, Doug Lauffenburger, Peter Sorger, and (for a time) Martin Yarmush, for open ears, lively discussions, and advice to help me along my way.

My dear friend, colleague, former group- and office-mate, “academic sister” and perennial mentor, Maria Klapa, who has been teaching me since my first day at MIT and continues even now. I have strived to live up to her standard and example, as an exemplary scientist and a wonderful individual. I simply can’t imagine graduate school without her presence or guidance.

My current and former colleagues of the Stephanopoulos group, for providing advice, intellectual debate, a sympathetic ear, or a few laughs over years of coffee, lunch and bad pizza. Special thanks to Saliya Silva, Matt Wong, Bill Schmitt, Jatin Misra, Jose Aleman, Stelios Kouvroukoglou, Gary Jung, and Stefan Wildt.

Brett Roth, Ilda Moura, Janet Fisher, Elaine Aufiero, Suzanne Easterly, Anne Fowler, Jenn Shedd, Mary Keith, and Susan Lanza, for helping with so many little administrative things along the way (I’m sure, many I don’t even know about) that made life at MIT move much more smoothly. Joanne Sorrento and Frances Meale, for helping me juggle the schedules and keep track of three very busy professors.

The National Science Foundation (Graduate Research Fellowship) and Singapore-MIT Alliance for financial support, helping to give me the academic freedom to pursue a project a bit outside the normal realm of my advisor.

Tom Wang, Brian Bucher, Arvind Mallik, Wendy Prud'homme, Gwang-Soo Kim, and Sarah Spurgeon, for being roommates, friends, and family, all rolled up together. You helped make homes exactly that—homes, not mere apartments where I kept my things and slept. We didn't just share a place to stay, but shared lives.

Friends from both coasts, so many that I dare not forget any, for standing by me and supporting me when I needed it, for giving me refuge when I needed that, and accepting my help when I could give it. Folks I knew before coming here: Vineet Gossain, Stacy Mar, Mike Lu, Khanh Ngo, Colleen Yeh, Khang Dao, George Tsao, Betty Chan, Brett Kurtin, Nizar Abdalla, Binita Bhattacharjee, Jimmy Lin, and Gopal Sridhar. And those I met and befriended along the way: Lily Koo, Connie Sun, Carmen Patrick, Hang Lu, Yonathan Thio, Casim Sarkar, Betty Yu, Inn Yuk, Klaudyne Hong, Lacey Southerland, and Lisa Wang. Others I may have missed in name, but you're always in my heart. Without you, I never would have had the strength to get through the dark days or felt the joy celebrating the bright ones.

Poh Lim, Stan Hunter, Helena Chia, AJ Liuba, and other members past and present of the MIT Korean Karate Club, for great exercise and better friendship. Tae Kwon Do was not just a way to exercise in the cold winters, or to work out some frustrations, or clear my head from a long day—it was all of that, and so much more.

Last, but never least, my dear family, for bringing me up, nurturing me along the way, and supporting me through the years even without understanding just what I was up to. My father, whose love of science and engineering led me to where I am. My mother, whose compassion and determination helped me become who I am. My sister, whose consideration, intelligence, and courage have always showed me who I wanted to be. My “Uncle” Dan Larson, whose friendship with my father over decades gave me a perfect example of how friends should be. My Uncle Carlos and Aunt Lia, for showing me what family should be, too. And my Grandma Tota, for stories, for laughs, for that optimism that gave me a glimmer of hope when all was dark.

Thank you all, because without you, I never would have made it to and through MIT, never learned so much along the journey, and never gotten to my own “Ithaca”.



## Table of Contents

<i>Abstract</i> .....	3
<i>Acknowledgements</i> .....	7
<i>Table of Contents</i> .....	9
<i>Table of Figures</i> .....	11
<i>Table of Tables</i> .....	14
<b>1. INTRODUCTION</b> .....	<b>15</b>
1.1 <i>Motivation</i> .....	17
1.2 <i>Background</i> .....	19
1.2.1. Construction and application of models.....	19
1.2.2. Mathematical analysis techniques.....	22
1.2.3. Measurement techniques and experimental considerations.....	25
1.3 <i>Objectives</i> .....	27
1.4 <i>Thesis overview</i> .....	29
1.5 <i>References</i> .....	30
<b>2. ACTIVATION RATIO ANALYSIS</b> .....	<b>35</b>
2.1 <i>Methods</i> .....	35
2.2 <i>Isolated interconverting cycle</i> .....	36
2.3 <i>Extension to simple network arrangements</i> .....	43
2.3.1. Linear Cascade.....	44
2.3.2. Converging Pathways.....	48
2.3.3. Diverging Pathways.....	51
2.3.4. Dual activation steps.....	52
2.3.5. Cascades with feedback.....	53
2.4 <i>Saturating conditions and influence of enzyme-substrate complexes</i> .....	57
2.5 <i>Conclusions</i> .....	63
2.6 <i>References</i> .....	65
<b>3. NETWORK RECONSTRUCTION USING ACTIVATION RATIOS</b> .....	<b>67</b>
3.1 <i>Concepts and Algorithm</i> .....	67
3.1.1. Structural and numerical observability.....	67
3.1.2. Stepwise analysis of networks.....	69
3.1.3. Consistency.....	71
3.2 <i>Analysis of a Model Network</i> .....	73
3.2.1. Network structure and features.....	73
3.2.2. Network analysis using free concentrations.....	74
3.2.3. Network analysis using total activation ratios.....	79
3.3 <i>Automated regression of data and pattern assignment</i> .....	84
3.3.1. Methodology.....	84
3.3.2. Model selection and evaluation.....	87
3.3.3. Example: cascade analysis in the presence of experimental noise.....	89
3.4 <i>Conclusions</i> .....	97
3.5 <i>References</i> .....	98

<b>4. ACTIVATION RATIO ANALYSIS OF ERK PHOSPHORYLATION.....</b>	<b>99</b>
4.1. <i>Experimental system selection and design</i> .....	100
4.2. <i>Development and Operation</i> .....	103
4.2.1. Materials and Methods.....	103
4.2.2. Reaction operating conditions.....	105
4.2.3. Enzyme immobilization.....	109
4.3. <i>Model for Interpreting ELISA Data</i> .....	114
4.3.1. Motivation and concepts.....	114
4.3.2. Selection of data.....	121
4.3.2. Error analysis and model validation.....	123
4.4. <i>Activation Ratios in Erk Phosphorylation Cycles</i> .....	129
4.4.1. Comparison of liquid-phase and immobilized reactions.....	129
4.4.2. Variation of total Erk concentration.....	135
4.4.3. Modulation of phosphatases.....	137
4.5. <i>Conclusions</i> .....	140
4.6. <i>References</i> .....	141
<b>5. CONCLUSIONS – FUTURE WORK.....</b>	<b>145</b>
5.1. <i>References</i> .....	147
<b>6. APPENDICES.....</b>	<b>149</b>
<i>Appendix 1. Single cycle model</i> .....	149
A1.1. Simulation details.....	149
A1.2. MATLAB files.....	150
<i>Appendix 2. Parameter values for simple signaling models</i> .....	154
A2.1. Linear Cascade.....	154
A2.2. Converging Cycles.....	154
A2.3. Diverging Pathways.....	154
A2.4. Cascades with Feedback.....	154
<i>Appendix 3. Model signaling network</i> .....	155
<i>Appendix 4. Extended cascade with noise</i> .....	156
<i>Appendix 5. Sensitivity matrix for ELISA measurement model</i> .....	157

## Table of Figures

<b>Figure 2.1.</b> Single cycle diagrams and reaction scheme. ....	36
<b>Figure 2.2.</b> Simulation results for simple binding under saturating conditions. A) Fraction of A complexed ( $[A \cdot B]/A_T$ ) or B) activation ratio ( $[A \cdot B]/A$ ) plotted against free B (circles) or $B_T$ (squares). ....	39
<b>Figure 2.3.</b> Simulation results for an individual covalent-modification cycle as shown in Figure 2.1. A) Fraction of A activated ( $A^*/A_T$ ) and B) Activation ratios $AR_A$ ( $A^*/A$ ) for the simple cycle, plotted against free activating enzyme $E_1$ . Parameter values: $k_1 = 10$ , $k_2 = 10$ , $K_{m2} = 1$ , $E_{2T} = 1$ , $A_T = 10$ , $K_{m1} = 20$ (diamonds), 10 (squares), 4 (triangles), 2 ( $x^2$ s), 1 (stars), 0.4 (circles), 0.2 ( $+^2$ s). ....	43
<b>Figure 2.4.</b> Diagrams of extended signaling structures. A) linear cascade, B) converging pathways, C) diverging pathways, D) dual activation steps, E) cascade with feedback (single-step activation as in A). ....	44
<b>Figure 2.5.</b> Activation ratios for the linear cascade of Figure 2.4A, plotted against free concentrations of $E_1$ (A), $A^*$ (B), $B^*$ (C), or $C^*$ (D). Ratios for A: diamonds, B: squares, C: triangles. ....	47
<b>Figure 2.6.</b> Fractional activation and activation ratios for the converging cycle as shown in Figure 1C. A) and B), fractional activation ( $A^*/A_T$ ), C) and D), activation ratio ( $A^*/A$ ) plotted against free activating enzyme $E_1$ (A) and (C) or $E_2$ (B) and (D). ....	50
<b>Figure 2.7.</b> Results for diverging branches in Figure 2.4C. A) Activation ratios for A (diamonds) and B (squares) plotted against free activating enzyme $E_1$ , B) Activation ratios for A and B plotted against each other, i.e. $AR_A$ against $B^*$ and $AR_B$ against $A^*$ . . .	52
<b>Figure 2.8.</b> Activation ratios for cascade with positive feedback, unsaturated in feedback step ( $K_{FB} \gg A_T$ ), plotted against free $E_1$ (A), $A^*$ (B), $B^*$ (C), $C^*$ (D). $D^*$ (E), or $E^*$ (F). For clarity, markers were omitted for activation ratios of A in pt. A. ....	55
<b>Figure 2.9.</b> Activation ratios for cascade with negative feedback, unsaturated in feedback step ( $K_{FB} \gg A_T$ ), plotted against free $E_1$ (A), $A^*$ (B), $B^*$ (C), $C^*$ (D). $D^*$ (E), or $E^*$ (F). For clarity, markers were omitted for activation ratios of A in pt. A. ....	56
<b>Figure 2.10.</b> Total activation ratios for the isolated covalent modification cycle (using total active $A^*_T$ and total inactive $A_T$ ) plotted against free $E_1$ , with enzyme conditions A) Saturated: parameters as shown in Figure 2.3. B) Unsaturated: with $E_{iT}/K_{mi} < 0.1$ . ....	59
<b>Figure 3.1.</b> Network with parallel pathways. ....	68
<b>Figure 3.2.</b> Diagram of sample model network, showing activation reaction numbering. Model details and parameter values are included in Appendix 3. ....	73
<b>Figure 3.3.</b> Activation ratios for intermediates in model network of Figure 3.2, calculated using free species only, plotted against A) $E_1$ , B) $A^*$ , C) $C^*$ , D) $E^*$ , E) $G^*$ , and F) $H^*$ . In each case $E_2 = 0$ . ....	75
<b>Figure 3.4.</b> Activation ratios for intermediates in model network of Figure 3.2, calculated using free species only, plotted against A) $E_2$ , B) $B^*$ , C) $C^*$ , D) $D^*$ , E) $E^*$ , and F) $F^*$ . In each case $E_1 = 0$ . ....	76

<b>Figure 3.5.</b> Activation ratios for E (A and B) or F (C and D), plotted as contours against $C^*$ (constant $D^*$ , A and C) or $D^*$ (constant $C^*$ , B and D).....	77
<b>Figure 3.6.</b> Total activation ratios for intermediates in model network of Figure 3.2, plotted against A) $E_1$ , B) $A^*$ , C) $C^*$ , D) $E^*$ , E) $G^*$ , and F) $H^*$ . In each case $E_2 = 0$ .....	81
<b>Figure 3.7.</b> Total activation ratios for intermediates in model network of Figure 3.2, plotted against A) $E_2$ , B) $B^*$ , C) $C^*$ , D) $D^*$ , E) $E^*$ , and F) $F^*$ . In each case $E_1 = 0$ .....	82
<b>Figure 3.8.</b> Total activation ratios for E (A and B) or F (C and D), plotted as contours against $C^*$ (constant $D^*$ , A and C) or $D^*$ (constant $C^*$ , B and D).....	83
<b>Figure 3.9.</b> Models used automated regression of activation ratio data, with graphical significance of parameters and half-saturation points shown.....	85
<b>Figure 3.10.</b> A) Cascade structure and color scheme and B) Expected (and observed) output matrix $R_{ji}$ following automated regression analysis. ....	89
<b>Figure 3.11.</b> Activation ratios (symbols) and best-fit model curves (solid lines) for cascade of Figure 3.10A following automated regression. Activation ratios for A: squares, B: circles, C: triangles, D: diamonds, E: x's.....	90
<b>Figure 3.12.</b> Activation ratios (symbols) and best-fit model curves (solid lines) for cascade with noise added to activation ratios. The noise has a standard deviation of 40% of the true value. Symbols as in Figure 3.11.....	93
<b>Figure 3.13.</b> Activation ratios (symbols) and best-fit model curves (lines) for cascade with noise added in active and inactive species. The noise has a standard deviation of 20% of the true value. Symbols as in Figure 3.11. ....	94
<b>Figure 4.1.</b> Immobilization of enzymes enables separation of free species from enzyme-substrate complexes. ....	101
<b>Figure 4.2.</b> Erk covalent modification cycles under action of Mek, PTP1B and PP2A. ....	103
<b>Figure 4.3.</b> Sample results for <i>in vitro</i> reactions using Erk as a substrate. A) Phosphorylation by N4-Mek (5-10 $\mu$ L) or US-Mek (2 $\mu$ L) in different buffers (see text for composition). For N4-Mek in EB, numbers signify $\mu$ L Mek added (5+5 is 5 $\mu$ L N4-Mek+5 $\mu$ L EB). B) ErkPP dephosphorylation by 1 $\mu$ L phosphatases in RB with or without addition of inhibitors. OA: 100 nM Okadaic Acid, NaVO <sub>4</sub> : 1 mM sodium vanadate... ..	107
<b>Figure 4.4.</b> Phosphorylation of Erk in presence of N4-Mek, PTP1B, and PP2A. A) Dynamic time course of Erk phosphorylation using 0.1 $\mu$ L of each enzyme. B) Distribution of Erk forms after reaction for 2 hr.....	108
<b>Figure 4.5.</b> Enzymatic activities following direct adsorption to polystyrene. A) Phosphorylation of Erk using N4-Mek. B) Dephosphorylation of ErkPP using PTP1B (closed symbols) or PP2A (open symbols). Diamonds represent immobilized reaction conditions, while squares represent data taken for liquid-phase reactions. ....	110
<b>Figure 4.6.</b> Enzymatic activities following immobilization on Protein G-coated plates using capture antibodies without (A, C, E) or with (B, D, F) an extra blocking step. A and B) Phosphorylation of Erk by N4-Mek. C and D) Dephosphorylation of	

phosphotyrosine- (PTP1B) or phosphoserine-containing (PP2A) peptides. E and F) Dephosphorylation of ErkPP by phosphatases. ....	111
<b>Figure 4.7.</b> Enzymatic activities following immobilization using capture antibodies directly adsorbed to plates. A) Phosphorylation of Erk by N4-Mek. B) Dephosphorylation of ErkPP by phosphatases. C) Saturation of immobilization of PP2A. D) Phosphate release from phosphopeptides. Diamonds: immobilized enzymes, Squares: liquid-phase reactions, Triangles: enzymes remaining in supernatant liquid following immobilization. ....	113
<b>Figure 4.8.</b> Expectations for ELISA results with potential issues. A) Cross-reactivity of anti-Erk antibodies (taken from Yao et al [12]), B) Sigmoidal overall response containing an intermediate linear region. ....	114
<b>Figure 4.9.</b> ELISA data modeling and regression as two steps: first, rescaling of data using composition parameters, second, linear fit of resulting rescaled data. Circles: “D” (original ErkPP), squares: “T” (ErkPP + PTP1B), diamonds: “Y” (ErkPP + PP2A), triangles: “N” (ErkPP + both phosphatases). ....	118
<b>Figure 4.10.</b> Sample results from ELISA standards regression. Symbols as in Figure 4.9. Solid lines represent model predictions using best-fit parameters. ....	119
<b>Figure 4.11.</b> Collection of data points from different samples and best-fit regression results. Data is the same as used in generation of Figure 4.10, with $x_T$ rescaled using model parameters. ....	120
<b>Figure 4.12.</b> Sample ELISA data, showing dependence of experimental variance on absolute signal. Diamonds: average signal (triplicate measurements). Squares: ratio of standard deviation to average value at each point. ....	122
<b>Figure 4.13.</b> A) Activation profiles and B) Activation Ratios for Erk cycles vs. normalized volume of N4-Mek added, in liquid phase (closed symbols) or immobilized using capture antibodies (open symbols). “Low” (circles) signifies a maximum of 0.1 $\mu\text{L}$ N4-Mek added, “Norm” (squares): 1 $\mu\text{L}$ , “High” (triangles): 10 $\mu\text{L}$ Mek. ....	131
<b>Figure 4.14.</b> A) Activation profiles (fraction Erk in each form) and B) Activation Ratios vs. volume Mek added, for Erk cycles in liquid phase, varying total concentration of Erk. ....	136
<b>Figure 4.15.</b> A) Activation profiles and B) Activation Ratios vs. volume Mek added, for Erk cycles in liquid phase, varying total amounts of both phosphatases. ....	138
<b>Figure 4.16.</b> A) Activation profiles and B) Activation Ratios vs. volume Mek added, for Erk cycles in liquid phase, varying total amounts of each phosphatase independently. ....	139

## Table of Tables

<b>Table 2.1.</b> Regression results for converging pathway, using parameter values from Figure 2.6.....	49
<b>Table 2.2.</b> Summary of results for simple signaling systems: expressions for activation ratios.....	63
<b>Table 3.1.</b> Expectations for activation ratios when inverting the relationship between species I and J.....	72
<b>Table 3.2.</b> Activation ratio combinations for three-species structures.....	72
<b>Table 3.3.</b> Regression results for simple network, using free activation ratios.....	78
<b>Table 3.4.</b> Models and parameter bounds for automated regression of activation ratio data. Here $\Delta y = y_{\max} - y_{\min}$ , $\Delta x = x_{\max} - x_{\min}$ , $y_{\text{avg}} = (y_{\max} + y_{\min})/2$ , and $x_{\text{avg}} = (x_{\max} + x_{\min})/2$ . $x(y=y_{\text{avg}})$ signifies the value of $x$ nearest to where $y$ equals $y_{\text{avg}}$ , and vice versa for $y(x=x_{\text{avg}})$ .....	86
<b>Table 3.5.</b> Results of automated regression/decision analysis for extended cascade in Figure 3.9A. Results shown are mean value from 10 replicate calculations in selecting an optimal model. 0 represents linear fit, 1 is hyperbolic fit, and -1 is inverse hyperbolic fit. See Figures 3.12-13 for example data.....	91
<b>Table 3.6.</b> Worst-case results of automated regression/decision analysis for extended cascade in Figure 3.10A, where activation ratios have an added noise term of 40% of true value. Data is shown in Figure 3.12. Tables represent model with lowest AIC and weights for each model.....	95
<b>Table 3.7.</b> Worst-case results of automated regression/decision analysis for extended cascade in Figure 3.10A, where active and inactive concentrations have an added noise term of 20% of true value. Data is shown in Figure 3.13. Tables represent model with lowest AIC and weights for each model.....	95
<b>Table 4.1.</b> Buffers used in experimental studies (final values).....	105
<b>Table 4.2.</b> Conditions tested during optimization of enzyme immobilization, and details of final procedure.....	113
<b>Table 4.3.</b> Parameter values from ELISA standards regression (data shown in Figures 4.9-4.11).....	126
<b>Table 4.4.</b> Comparison of predicted (from standards composition) and observed (estimated from measurements) Erk amounts using mixtures of standards (D, T, Y, N) samples. 10 ng of each standard (20 ng total) was added to wells, and ELISA measurements used with Equation 4.20 to calculate $x$ . Predictions were calculated according to Equation 4.2. All entries are shown in ng Erk of each form.....	128

# 1 INTRODUCTION

The normal operation of any cell can be imagined to consist of three interacting systems. The *metabolism* is comprised of the enzymes and intermediates involved in production of energy and synthesis, processing, and recycling of essential building blocks like amino acids, lipids and nucleotides from any materials available to the cell. The *genetic system* contains the DNA and RNA as well as the polymerases, splicing apparatus, ribosomes, and posttranslational machinery that together act to produce structural proteins and enzymes for all cellular reactions. The *signaling system* acts to recognize extracellular cues, relay them across the membrane, and transmit the information accordingly to help direct both the metabolic network, through changes in enzymatic or transporter activities, and the genetic network, through regulation of transcription factors. Of course, this simplified perspective neglects to explicitly consider physical components of the cell like the cytoskeleton or some of the systems involved in cell division such as DNA replication and chromosomal separation. But these aspects of cellular behavior also depend upon and interact with the systems described above, and contain components that in some cases can also be considered members of the metabolic, signaling or genetic networks.

Traditional efforts to discover and characterize individual components of each system have recently begun to give way to more encompassing “systems biology” initiatives. This is in part due to the level of detail to which many of the relevant pathways have already been described. Perhaps more importantly, it is a reflection of the desire to study molecules not in isolation but rather within the contexts of the cellular networks in which they interact. Detailed investigations of specific components provide much important information about the physical and chemical interactions of these molecules, but they cannot address the question of what the actual *in vivo* activity is in a particular cell type or set of environmental conditions. Instead, a set of measurements about the system must be combined with an analytical framework capable of processing the data.

There are already a variety of methods in place to reconstruct and characterize metabolic and genetic networks. In the case of metabolism, analysis of the *flux*, or throughput, between metabolic intermediates is used to provide a measure for the relative engagement of particular reactions as matter is transferred through pathways in the network, based on a set of steady-state mass balances around each intermediate [1]. This methodology has been extended to provide additional insight, such as identification of reversible reactions and cycles, by utilizing substrates labeled with radioactive or uncommon stable isotopes (e.g.  $^{13}\text{C}$ ), measurement by nuclear magnetic resonance (NMR) or mass spectrometry (MS), and adding the appropriate isotopomer balances to the analysis [2]. The correlation between expression patterns for different genes, measured using DNA microarrays, across many sample types or over time has been used to identify connectivity between the genes and assign parameters reflecting the strength of interaction [3, 4]. More recently, analysis of time-lagged correlations has been used to order genes in sequence and therefore add directional information, which may be used to suggest causality [5].

At this time, no analogous method exists to translate measurements of the activation of particular signaling intermediates directly into putative network structures. The connectivity of signaling networks has been predicted by examination of databases describing protein-protein interactions based on two-hybrid experiments or by searching protein sequences for potentially interacting domains [6-8]. However, such approaches are not designed to incorporate information regarding activation of intermediates, and therefore cannot be applied to evaluate the activation of different signaling pathways under a particular set of experimental conditions. Structural identification of the signaling network is separated from a quantitative description, whereas ideally both can be accomplished with the same technique simultaneously.

The methods described above for study of metabolic or genetic systems unfortunately cannot be applied to examine signaling. Metabolic analysis methods depend upon mass flow between intermediates, which does not occur in signal transduction pathways. Genetic analysis depends in general upon the separation in time between activation of each intermediate: production of the corresponding protein is required to activate the next gene in sequence. In signaling, however, activation of one



step can begin before completion of the previous step. While some groups have previously attempted to modify existing tools for the study of signaling networks, in general their results have been mathematically intractable or experimentally infeasible and thus applied only sparingly. Therefore, it was apparent that a novel analysis framework was necessary for the study of signaling systems.

## **1.1 Motivation**

Signal transduction processes allow cells to interact with each other and with their environments and are critical to the proper function of both unicellular eukaryotes and cells in a multicellular organism. Bacteria and yeasts utilize signaling pathways to sense and respond to environmental cues including food, osmotic pressure, or mating factors [9, 10]. In higher-order species including worms, fruit flies, mice and humans, signaling is used to direct proliferation and differentiation of embryonic stem cells into neural, muscle, bone and other tissues, to activate the immune response in macrophages and B and T lymphocytes, and to coordinate organs together, such as the action of insulin to simultaneously regulate liver, pancreatic, fat and muscle tissue following a meal [11-13]. These pathways are therefore centrally involved in the regulation of cellular behavior, from induction of growth and division, death (by apoptosis), movement and shape to even changes in phenotype, either in the metabolic state or patterns of gene expression [14]. Defects in signaling processes have been linked to a variety of diseases, most notably diabetes and cancer [15, 16].

Developing a more thorough understanding of signal transduction phenomena depends on the ability to analyze the simultaneous action of the elements involved. Thus, the same reasons that prompt research in signal transduction in the first place are also pertinent to developing analytical methods. A framework for analysis of signal transduction would give researchers a tool for experimental design, interpretation of results, and a way to visualize the breadth of effects that a signal can have. Together with improved methods in making measurements, this framework would enable the creation and evaluation of models of signal transduction. These could be used in understanding the processes of cellular development, stem cell differentiation, and progression of

disease. Furthermore, methods for analysis of signaling systems could be used for selection of drug targets, evaluation of efficacy, and observation of side effects.

Signal transduction pathways are, by nature, involved in the transfer of information; in many cases this takes the form of intermediates cycling between two or more states; the “information” is contained within the relative amounts of these states and how they influence the states of other intermediates. Individual steps typically involve a translocation of proteins to specific regions of the cell, activation of enzymes, covalent modification of proteins (particularly by phosphorylation), or production and translocation of small “second messenger” molecules, such as  $\text{Ca}^{2+}$ ,  $\text{IP}_3$ , and cAMP [17, 18]. In each case, the state of one species will regulate the activity of an enzyme, either directly or indirectly. This enzyme then acts to modify a new target, and thus the state of one species can be used to direct a change in the state of one (or more) other species. A set of cascades and branching networks is created that ultimately ends at the level of transcription factors or metabolic enzymes. Depending on the specifics of each step, the cascades offer a method for both signal amplification and attenuation [19-23].

The past few decades have witnessed an explosion of research in biology, with a large amount of energy devoted towards the discovery and description of signal transduction phenomena. Completion of sequencing projects for genomes for a variety of organisms, from bacteria to humans, has yielded databases full of putative genes with unknown function, many of which may be involved in signal transduction. Technological advances in sequencing and identification of proteins, characterization and detection of posttranslational modifications, observation of protein-protein complexes, and subcellular microscopy have offered new tools for studying these processes. Powerful new computational techniques for protein and gene sequence alignment and homology identification have been used to classify newly discovered genes, suggest potential interaction partners, and predict catalytic activities. This has combined to give researchers a great number of possible signaling components, but not the ability to quickly identify where these molecules actually reside in the overall signaling network.

Frequently, signal transduction research attempts to investigate in detail particular molecules or short pathways. This approach has proven valuable in learning about

mechanisms of signaling, identification of characteristic domains, and provided information critical for prediction of putative signaling components from genomic and proteomic databases. However, the variety of cellular and effector systems used to study one molecule have yielded results that are sometimes contradictory, and often confusing, when combined with data for other molecules. Exogenous expression, use of knockout or constitutively active forms of intermediates, and investigation in a variety of cell types are all examples of cases where the study of one component, and how it affects other species, may not be reflective of its behavior in the unperturbed system. It is partly for this reason that signaling research is beginning to focus more on the response of networks, not individual components. Without methods in place to analyze data for entire signaling networks, however, it is difficult to integrate data from multiple sources, or more importantly, direct and interpret these systemic studies.

One of the most confusing areas in signaling research so far has been that of signaling specificity. An astounding number of different ligands, receptors, and intracellular signaling intermediates have been discovered and described. However, any particular cell type may respond very differently to the same ligand than another cell type, and the same cell type may respond differently to different ligands [24, 25]. Nevertheless many of the same pathways appear to be activated under a variety of different situations, although perhaps by different upstream mechanisms [26]. Understanding signaling specificity depends upon the ability to accurately describe how different cellular and environmental conditions influence the activation of the entire network.

## **1.2. Background**

### **1.2.1. Construction and application of models**

Models for signaling networks have been used as a tool to help investigate properties of these systems for nearly as long as research in signal transduction has been performed. In the lack of experiments that are easy to develop or without appropriate measurements, models can provide researchers a way to test theories on potential mechanisms. The level of detail can vary from quite abstract to highly detailed, depending upon the amount of information that is available and the particular focus of the

investigation. And while models have an inherent capacity to predict system behavior, it should be noted that they are always constructed so as to fit some sort of training data set. Therefore, models are typically limited in their ability to check their own consistency—i.e., they may not be able to identify inaccuracies in the presumed structure of the network, which would be hidden by the flexibility in values for adjustable parameters.

One way to avoid this problem is to utilize a relatively abstract model, which focuses on representing the network structure rather than mechanistic details. In that case the interaction between system nodes (signaling molecules or modules) is described as probabilities of signal transfer. The resulting models may be constructed in the form of Boolean, logical, neural, or stochastic Petri networks [27]. Structure identification is possible by first enumerating all possible connections between nodes, and nonzero probabilities after fitting to data indicate a structural connection, although an extremely large amount of data is required for this process [28]. On the other hand, models with predefined structure can be used to explore qualitative features of the system, including requirements for different modes of operation (as in T-cell activation vs. relaxation, or metabolic vs. mitogenic signaling by insulin) [29-31].

However, it is difficult to place a physical interpretation on the probabilities or connect them to a particular mechanism. Furthermore, it has been questioned whether signaling can even be thought of as a digital process [32]. To more accurately represent the signal-response characteristics of individual steps, Omholt used “switchlike” sigmoid functions to describe signal transfer during iron homeostasis, although at the expense of requiring additional parameters to describe each reaction [33]. As further modifications to the model are added, to incorporate additional detail regarding the mechanism of each step, the distinction from kinetic models is lost. It is likely that for this reason, such abstract approaches have only been sparingly used in examining signaling systems.

Detailed kinetic models have been used extensively to describe a variety of signaling systems [34-37]. Several groups focused on analyzing the behavior of a single intermediate cycling between two forms [38-40]. It was thus shown that such a signaling intermediate could show a sharp, “ultrasensitive” response to the amount of activating enzyme, reminiscent of the cooperativity seen in the binding of oxygen by hemoglobin

[38]. This behavior was dependent upon saturation of one or both enzymes operating to drive the cycle (zero-order ultrasensitivity), and was observed experimentally for activation of isocitrate dehydrogenase and glycogen phosphorylase [41-44]. Examination of models for short cascades showed that sequences of signaling steps could also yield sharp overall responses (multistep ultrasensitivity) [21, 23, 45-47]. Additional system features, including oscillations or bistability resulting from different feedback modes, interplay with scaffolding proteins, and the effects of multiple activation steps or limitation of diffusion of components across spatial gradients have also been considered [48-51]. Combination with models for receptor-ligand interactions and trafficking has yielded expansive descriptions of signaling induced by growth factors such as EGF [52-55].

Such approaches allow researchers to test theories about details of different mechanisms. However the analysis of model results is essentially qualitative; kinetic parameters could be varied over several orders of magnitude without significantly altering the overall behavior [49, 52, 55]. As the model complexity increases, so does the number of parameters and concentrations of intermediates that need to be included. Some are taken from independent experiments using enzymes purified *in vitro*, which calls into question the validity of the values for *in vivo*, as well as limiting model definition to only a few experimental systems. Others are fit to coincide with data but are rarely validated by additional experiments later, and again can be varied significantly without influencing the fit to experimental data, suggesting that the models may not be completely describing the experimental system.

While increasingly complex systems can be successfully simulated, detailed models nevertheless possess some serious limitations. The model complexity means that results become almost as difficult to interpret as the experimental results that they are trying to emulate. In general, errors in the structure (missing or incorrectly placed reactions or components) cannot be recognized. So models provide little insight on how newly discovered components or entire pathways could be incorporated into the analysis. Furthermore, they provide no easy method by which to gain a general perspective on what components play key roles in the signaling process. With the uncertainty present in the kinetic parameters, little can be said about the relative importance of different

pathways to yield the overall observed behavior. Such characterization generally requires the application of mathematical manipulations to define and extract descriptive parameters.

### 1.2.2. Mathematical analysis techniques

The basic concept behind an analysis framework is that structural and quantitative information about a system can be somehow inferred directly from the data, without having to construct a model (of any level of detail) ahead of time. A model may be used as a starting point to help develop insight into the behavior of the system, but the technique is developed by determining what the proper transformation of data should be, based on some sort of mathematical analysis. The final form of the analysis is independent of the original model, and thus truly describes an alternate method of examining the system.

Perhaps unsurprisingly, the majority of attempts to develop such a technique for the examination of signaling pathways have been based in extending methods originally developed for the analysis of regulation in linear metabolic pathways. Metabolic control analysis (MCA) was created to quantify the effects of changes in the enzyme activities (E) upon the steady-state flux (J) of mass through metabolic networks or the concentrations (X) of intermediates [56, 57]. These effects can be described in terms of the flux control coefficient (FCC)  $C_E^J$  and concentration control coefficient (CCC)  $C_E^X$ :

$$C_E^J = \frac{dJ/J}{dE/E} = \frac{d \ln J}{d \ln E} \quad (1.1)$$

$$C_E^X = \frac{dX/X}{dE/E} = \frac{d \ln X}{d \ln E} \quad (1.2)$$

MCA therefore amounts essentially to a sensitivity analysis, using control coefficients to describe the distribution of regulation that particular enzymes (and therefore particular steps in a pathway) have upon the overall flux. Similar expressions can be written in terms of changes to other system parameters, such as allosteric regulators or concentrations of other metabolites; these are usually called response coefficients ( $R_p^J$  and  $R_p^X$ ) to accentuate that these molecules act indirectly on the system. It should be noted that these coefficients are determined as total differentials, therefore arise from both the direct effects (upon a particular reaction) as well as indirect effects

(by altering concentrations of intermediates that regulate other reactions in the network).

This can be seen by application of the so-called summation and connectivity theorems:

$$\text{FCC summation:} \quad \sum_i C_{E_i}^J = 1 \quad (1.3)$$

$$\text{FCC connectivity:} \quad \sum_i C_{E_i}^J \varepsilon_{X_j}^{E_i} = 0 \quad \forall X_j \quad (1.4)$$

$$\text{CCC summation:} \quad \sum_i C_{E_i}^{X_j} = 0 \quad (1.5)$$

$$\text{CCC connectivity:} \quad \sum_i C_{E_i}^{X_j} \varepsilon_{X_k}^{E_i} = \begin{cases} 0 & (k \neq j) \\ -1 & (k = j) \end{cases} \quad \forall X_j, X_k \quad (1.6)$$

The local effects are described by the elasticities  $\varepsilon_{X_i}^E$ , defined as:

$$\varepsilon_{X_i}^E = \left. \frac{\partial \ln E}{\partial \ln X_i} \right)_{X_j} = \left. \frac{\partial \ln v}{\partial \ln X_i} \right)_{X_j} \quad (1.7)$$

The second equality in Equation 1.7 is based on the fact that the rate of a particular reaction is in general directly proportional to the enzymatic activity. If the kinetics of a particular reaction is known, then the elasticities can be obtained by differentiation of the rate expression. The elasticities and control coefficients can also be determined experimentally, using perturbations to each reaction step, by what is known as the double-modulation method and its extensions [58].

This approach was originally developed for the study of linear metabolic pathways, but through a series of steps has been extended to include cycles and pathways without mass transfer between intermediates [59-63]. Nevertheless, the basis of analysis was focused on measurement of changes in fluxes following changes to enzymatic activities. In signal transduction, the fluxes in question would correspond to the rate of interconversion between forms of each intermediate at steady state—to date, impossible to measure. Furthermore, each reaction step for each interconverting cycle must be included, leading to a prohibitively large number of different permutations required to fully examine the system. Thus these extensions of MCA have only rarely been applied directly to examine signaling systems [64-66].

A novel approach has been recently proposed that is instead based on examining changes to the steady-state concentrations of the intermediates themselves, and also reduces the complexity by separating the system into interacting modules, and focusing attention only on representative molecules each module [67, 68]. This method, called modular response analysis (MRA), utilizes connectivity theorems to translate how the intermodular response coefficients (describing the effect of a molecule from one module onto another module) will result in overall response coefficients for the system. This process can be inverted to determine the intermodular coefficients, which indicate structural connectivity as well as give a quantitative value for module interactions. However, the method suffers from the major limitation that to determine the coefficients, a perturbation must be applied that is specific to each module. For known signaling components, perturbations of the form of enzymatic inhibitors may be available, but that may not be the case for newfound species. And indirect connections arising from missing components (not measured or perturbed directly) will not be recognized with this methodology. Within the species that are being measured, the approach is able to reconstruct the network, but unable to determine exactly where missing steps might be.

Another analytical technique for studying reaction pathways was developed by using time-lagged correlations to infer connectivity between components [69]. The premise here is that the time-dependent behavior of two species will be most similar if they are connected in a reaction network. The correlations can be translated into a matrix of Euclidian distances, and through a series of steps designed to reduce the dimension of the data, a projection into 2D space that ultimately reflects the original structure of the system. Unfortunately, the algorithm does not always obtain the correct structure at the end of analysis; as the method is based on correlations in time then two subsystems with different timescales of operation may not be recognized as being connected. Also, this procedure requires a large number of dynamic data points, where input signals are modulated at a frequency on the same timescale as the remainder of reactions in the network, which for signal transduction would be seconds to minutes. It may be for this reason that this method has not been applied to examine signaling systems, although it has been tested on a segment of glycolysis constructed *in vitro* with purified enzymes [70].



### 1.2.3. Measurement techniques and experimental considerations

Any analysis of signaling pathways depends upon the ability to make quantitative measurements of how an external signal is influencing intracellular components. The complexity inherent in the analytical techniques described above may partially explain the relative lack of experimental applications thus far. But other complications either in design of experiments or methods of sampling may be playing a role. It is therefore worthwhile to consider the various techniques available for measuring signaling intermediates, and the limitations currently placed on conducting experiments.

Signal transduction proceeds through a wide range of different mechanisms, and therefore the ability to measure the amount of “active” intermediate depends on the characteristics of the signaling step. Molecules that are transported to different regions of the cell might be observed through microscopy or by sampling specifically from that region (membrane, cytosol, nucleus, vesicle, etc.). Some intermediates undergo covalent modification, typically by phosphorylation, and thus the modified form must be separated or detected specifically from its original unmodified form. Formation or dissolution of stable noncovalent complexes could be determined by coimmunoprecipitation. As these changes may coincide with induction of enzymatic activity, functional assays are often utilized if a substrate is readily available. Obviously, handling issues related to maintaining the active state become important, whether through appropriate composition of buffers (for example to include inactivating enzyme inhibitors) or sampling conditions (time, temperature, physical separations, etc.)

Furthermore, the type of molecule should be considered. Most signaling intermediates are proteins, therefore are relatively large (10-200 kDa), potentially membrane-bound, and containing complicated surface charges and chemistries. The physical properties of proteins are generally changed only slightly by covalent modification of a few residues. Therefore many protein detection methods involve the use of specific binding reagents such as antibodies, which can preferentially recognize the modified form of the protein. The Western blot is by far the most common method used to detect and quantify signaling proteins, but requires several tedious handling steps and can handle few samples at a time. The multiwell plate version of the assay (ELISA) is gaining popularity, because of the ability to handle more samples, be automated, and is

more readily quantified. Antibodies tagged with fluorophores can be used to track protein localization via microscopy, or protein presence by flow cytometry. In either case, the primary limitation in measurements is the availability of antibodies specific for each protein.

Mass spectrometry (MS) has increasingly been applied for the quantitation and characterization of signaling proteins [71]. One or more separation steps using electrophoresis (gel or capillary) or chromatography are combined with a digestion reaction with specific proteases such as trypsin before the peptides are then applied to the MS. Covalent modifications can be observed as a shift in the mass for a particular peptide in the protein. Quantitation is possible by mixing the test material with a control sample, where one of the two is labeled with a stable isotope to shift the mass slightly [72-74]. While MS-based techniques thus far have primarily been used to identify targets of input stimuli (profiling), it may soon become a dominant technique for protein quantification [74-77].

Unlike proteins, small molecule second messengers often may be measured directly. Phospholipids and their breakdown products can be separated by thin-layer chromatography (TLC) or HPLC and visualized if previously labeled with a fluorophore or radioactivity ( $^{32}\text{P}$ ,  $^3\text{H}$  or  $^{14}\text{C}$ ) [78, 79]. Cyclic nucleotides (cAMP or cGMP) can also be separated from their native forms by TLC or alumina-based chromatography [80, 81]. On the other hand,  $\text{Ca}^{2+}$  concentrations in various regions of the cell are typically measured using secondary reagents, such as fluorescent dyes or enzymes like aequorin and adenylyl cyclase, that show altered activity in the presence of the ion [82-84].

Regardless of the type of molecule or analytical method used, the time and effort involved in preparing samples generally limits the reproducible sampling frequency to the order of minutes, similar in magnitude to the dynamics of most signaling reactions. Cost and time constraints may also reduce the number of measurements. Since most techniques are limited in capacity of measurements at one time, usually less than ten observations are made for any one molecule in an experiment. This is in stark contrast to the capacity of DNA microarrays to measure the expression thousands of genes at one time, and where the timescale of gene expression changes is in the tens of minutes. Thus

analysis of signaling systems must be performed with relatively small numbers of samples and under significant experimental uncertainty. Significant efforts are underway to develop various high-throughput techniques to measure proteins, which will help address the limitation of capacity, but not of sample handling [85-88]. Thus more signaling intermediates may be measured at one time, but for each protein there may still be only a few observations. Analysis techniques such as time-lagged correlations described in Section 1.2.2 would still be infeasible, and thus a new approach appears to be necessary.

### **1.3. Objectives**

In general, when examining signaling networks two questions will arise:

1. What are the pathways involved in response to a particular stimulus?
2. How much are these pathways utilized?

It is by answering these questions that signaling under different conditions can be compared. The first question is essentially qualitative, and is solved by determining the structure of the signaling network downstream of a particular input. This requires not only knowledge of which components are activated, but also how their activation leads to that of other species—in other words, the connectivity of the signaling network. Quantitative descriptions of signaling, which specifically answer the second question, may also be necessary in structural identification. This is because multiple factors may activate some of the same pathways to saturation, and thus the response due to one factor may be concealed by another. With these questions in mind, it is possible to describe the important properties of signaling analysis methods.

First, a method capable of *systemic, network* analysis is needed, rather than an examination of individual pathways or components. No one pathway operates in isolation. It is likely that several unseen factors are simultaneously contributing to activate several pathways to produce the observed effect. Signaling networks are complex, with many possible interactions, and many different sets of external conditions must be compared [25, 26, 36]. Only by examining the full signaling network can interactions between different pathways be seen, and only then can the considerable

complexity of signaling interactions be deconvoluted. Therefore the method must be capable of examining the entire network, and not be dependent upon detailed descriptions of any one step.

Second, a *quantitative* approach is necessary. Although qualitative descriptions of signaling are important to know which pathways are involved in the response, it is only by measuring quantitatively the differences in signaling under different conditions that comparisons can truly be made [89]. Through a quantitative approach, values can be defined to represent signaling under different conditions and relative contributions of different pathways to the response. These values can then be compared between experiments to suggest which factors influence signaling. This necessarily requires that large-scale quantitative measurements of signaling must be available to decipher signal specificity. Furthermore, the measurement methods must be able to provide as much detail as possible about individual states of each component.

Finally, a *practical* method is required that allows visualization of *in vivo* activity under specific sets of conditions. It is of little use to develop analytical methods for study of signaling networks that are mathematically sound but cannot be applied to experimental data for signaling reactions within cells. Previous efforts to analyze signaling networks have been largely ignored, in some cases because of a dependence on unmeasurable values (such as the interconversion rate for a particular intermediate). On the other hand, it may be impossible to characterize signaling networks with measurements that are currently available. In that case, an analytical framework may be useful in directing what types of measurements are necessary for systems analysis.

The objective of this work therefore was to develop a novel analytical approach for the examination of signal transduction networks, with the specific understanding of limitations of experimental methods and lack of *in vivo* kinetic data. Of primary concern was that the framework could be readily applied for the structural analysis of a signaling network yet would still contain quantitative descriptions for the interactions between intermediates. This approach is also useful in experimental design, since it can be used to indicate the types, quantity, and quality of data that will be necessary. The framework should yield simple relationships for simple forms of interactions, and change

appropriately when more complicated interactions are considered, thus enabling the detection of these complicated interactions.

#### **1.4. Thesis overview**

Equations representing the time-dependent behavior of species in signaling systems were combined with simulations written in MATLAB 5.2 (Mathworks, Inc) to investigate how the activation of one component may be described in terms of the other species. The behavior of a single component cycling between two states was studied in detail as a model of the most fundamental unit in signal transduction. This led to the definition of activation ratios as an informative measure of the interaction between target and activator, where all kinetic constants for the reaction system are collapsed into one single factor. Extension to simplified network substructures, such as linear cascades, convergence and divergence points resulted in a set of observations for how the activation ratio for an individual intermediate reflects its position in a larger network. Details for the derivation of activation ratios and sample simulation results for these systems are discussed in Chapter 2.

These observations, however, represent how activation ratios are predicted to behave given a presupposed network structure. In order to invert this process, and thereby reconstruct a network from measurements, the addition of constraints based on self-consistency was required. The resulting algorithm is discussed in Chapter 3, along with an example of application to a small model network, where simulation results from a more detailed model were utilized as theoretical measurements. This process was partially automated using a MATLAB script to regress data against linear and nonlinear models and evaluate fit based on the Akaike Information Criteria. This enabled an examination of how issues with data quality could influence the analysis results.

Activation ratio analysis was applied to a real experimental system by studying the phosphorylation of protein kinase Erk2 *in vitro*, as described in Chapter 4. The experimental setup and tools for processing raw data were developed so as to provide the correct types of measurements needed to calculate activation ratios. Details for development of the system, optimization of reaction conditions, and tools for filtering data are also discussed. Variation of system parameters enabled verification of some

predictions for activation ratios, while highlighting the importance of improvements in measurement capabilities.

## 1.5. References

1. Stephanopoulos, G., A.A. Aristidou, and J. Nielsen, *Metabolic engineering : principles and methodologies*. 1998: Academic Press, San Diego.
2. Klapa, M.I., *High resolution metabolic flux determination using stable isotopes and mass spectrometry*, Ph D Thesis, Dept. of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, MA, 2001.
3. D'Haeseleer, P., S. Liang, and R. Somogyi, *Genetic network inference: from co-expression clustering to reverse engineering*. *Bioinformatics*, 2000. **16**(8): p. 707-26.
4. Ronen, M., et al., *Assigning numbers to the arrows: parameterizing a gene regulation network by using accurate expression kinetics*. *Proc Natl Acad Sci U S A*, 2002. **99**(16): p. 10555-60.
5. Schmitt, W.A., *Extracting Transcriptional Regulatory Information From DNA Microarray Expression Data*, Ph D Thesis, Dept. of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, MA, 2003.
6. Pawson, T. and P. Nash, *Assembly of cell regulatory systems through protein interaction domains*. *Science*, 2003. **300**(5618): p. 445-52.
7. Steffen, M., et al., *Automated modelling of signal transduction networks*. *BMC Bioinformatics*, 2002. **3**(1): p. 34.
8. Yaffe, M.B., et al., *A motif-based profile scanning approach for genome-wide prediction of signaling pathways*. *Nat Biotechnol*, 2001. **19**(4): p. 348-53.
9. Gustin, M.C., et al., *MAP kinase pathways in the yeast Saccharomyces cerevisiae*. *Microbiol Mol Biol Rev*, 1998. **62**(4): p. 1264-300.
10. Hellingwerf, K.J., et al., *Current topics in signal transduction in bacteria*. *Antonie Van Leeuwenhoek International Journal of General and Molecular Microbiology*, 1998. **74**(4): p. 211-227.
11. Whitehead, J.P., et al., *Signalling through the insulin receptor*. *Curr Opin Cell Biol*, 2000. **12**(2): p. 222-8.
12. Tanaka, S., et al., *Signal transduction pathways regulating osteoclast differentiation and function*. *J Bone Miner Metab*, 2003. **21**(3): p. 123-33.
13. Nordin, A.A. and J.J. Proust, *Signal transduction mechanisms in the immune system. Potential implication in immunosenescence*. *Endocrinol Metab Clin North Am*, 1987. **16**(4): p. 919-45.
14. Ferrell, J.E., Jr., *MAP kinases in mitogenesis and development*. *Curr Top Dev Biol*, 1996. **33**: p. 1-60.
15. Marx, J., *Unraveling the causes of diabetes*. *Science*, 2002. **296**(5568): p. 686-9.
16. Spencer, V.A. and J.R. Davie, *Signal transduction pathways and chromatin structure in cancer cells*. *J Cell Biochem Suppl*, 2000. **Suppl 35**: p. 27-35.
17. Berridge, M.J., *The molecular basis of communication within the cell*. *Sci Am*, 1985. **253**(4): p. 142-52.
18. Hunter, T., *Signaling--2000 and beyond*. *Cell*, 2000. **100**(1): p. 113-27.

19. Brown, G.C., J.B. Hoek, and B.N. Kholodenko, *Why do protein kinase cascades have more than one level?* Trends Biochem Sci, 1997. **22**(8): p. 288.
20. Ferrell, J.E., Jr., *How responses get more switch-like as you move down a protein kinase cascade [letter; comment]*. Trends Biochem Sci, 1997. **22**(8): p. 288-9.
21. Ferrell, J.E., Jr., *Tripping the switch fantastic: how a protein kinase cascade can convert graded inputs into switch-like outputs [see comments]*. Trends Biochem Sci, 1996. **21**(12): p. 460-6.
22. Ferrell, J.E., Jr., *Building a cellular switch: more lessons from a good egg*. Bioessays, 1999. **21**(10): p. 866-70.
23. Goldbeter, A. and D.E. Koshland, Jr., *Ultrasensitivity in biochemical systems controlled by covalent modification. Interplay between zero-order and multistep effects*. J Biol Chem, 1984. **259**(23): p. 14441-7.
24. Marshall, C.J., *Specificity of receptor tyrosine kinase signaling: transient versus sustained extracellular signal-regulated kinase activation*. Cell, 1995. **80**(2): p. 179-85.
25. Weng, G., U.S. Bhalla, and R. Iyengar, *Complexity in biological signaling systems*. Science, 1999. **284**(5411): p. 92-6.
26. Schwartz, M.A. and V. Baron, *Interactions between mitogenic stimuli, or, a thousand and one connections*. Curr Opin Cell Biol, 1999. **11**(2): p. 197-202.
27. Bray, D., *Protein molecules as computational elements in living cells*. Nature, 1995. **376**(6538): p. 307-12.
28. Bray, D., *Intracellular signalling as a parallel distributed process*. J Theor Biol, 1990. **143**(2): p. 215-31.
29. Kaufman, M., F. Andris, and O. Leo, *A logical analysis of T cell activation and anergy*. Proc Natl Acad Sci U S A, 1999. **96**(7): p. 3894-9.
30. Shymko, R.M., P. De Meyts, and R. Thomas, *Logical analysis of timing-dependent receptor signalling specificity: application to the insulin receptor metabolic and mitogenic signalling pathways*. Biochem J, 1997. **326**(Pt 2): p. 463-9.
31. Shymko, R.M., et al., *Timing-dependence of insulin-receptor mitogenic versus metabolic signalling: a plausible model based on coincidence of hormone and effector binding*. Biochem J, 1999. **339**(Pt 3): p. 675-83.
32. Agutter, P.S. and D.N. Wheatley, *Information processing and intracellular 'neural' (protein) networks: considerations regarding the diffusion-based hypothesis of Bray*. Biol Cell, 1997. **89**(1): p. 13-18.
33. Omholt, S.W., et al., *Description and analysis of switchlike regulatory networks exemplified by a model of cellular iron homeostasis*. J Theor Biol, 1998. **195**(3): p. 339-50.
34. Asthagiri, A.R. and D.A. Lauffenburger, *Bioengineering models of cell signaling*. Annual Review of Biomedical Engineering, 2000. **2**: p. 31-53.
35. Neves, S.R. and R. Iyengar, *Modeling of signaling networks*. Bioessays, 2002. **24**(12): p. 1110-7.
36. Bhalla, U.S. and R. Iyengar, *Emergent properties of networks of biological signaling pathways*. Science, 1999. **283**(5400): p. 381-7.
37. Tyson, J.J., *Models of cell cycle control in eukaryotes*. J Biotechnol, 1999. **71**(1-3): p. 239-44.

38. Goldbeter, A. and D.E. Koshland, Jr., *An amplified sensitivity arising from covalent modification in biological systems*. Proc Natl Acad Sci U S A, 1981. **78**(11): p. 6840-4.
39. Varon, R. and B.H. Havsteen, *Kinetics of the transient-phase and steady-state of the monocyclic enzyme cascades*. J Theor Biol, 1990. **144**(3): p. 397-413.
40. Stadtman, E.R. and P.B. Chock, *Superiority of interconvertible enzyme cascades in metabolic regulation: analysis of monocyclic systems*. Proc Natl Acad Sci U S A, 1977. **74**(7): p. 2761-5.
41. LaPorte, D.C. and D.E. Koshland, Jr., *Phosphorylation of isocitrate dehydrogenase as a demonstration of enhanced sensitivity in covalent regulation*. Nature, 1983. **305**(5932): p. 286-90.
42. Meinke, M.H., J.S. Bishop, and R.D. Edstrom, *Zero-order ultrasensitivity in the regulation of glycogen phosphorylase*. Proc Natl Acad Sci U S A, 1986. **83**(9): p. 2865-8.
43. Meinke, M.H. and R.D. Edstrom, *Muscle glycogenolysis. Regulation of the cyclic interconversion of phosphorylase a and phosphorylase b*. J Biol Chem, 1991. **266**(4): p. 2259-66.
44. Shacter, E., P.B. Chock, and E.R. Stadtman, *Regulation through phosphorylation/dephosphorylation cascade systems*. J Biol Chem, 1984. **259**(19): p. 12252-9.
45. Chock, P.B. and E.R. Stadtman, *Superiority of interconvertible enzyme cascades in metabolite regulation: analysis of multicyclic systems*. Proc Natl Acad Sci U S A, 1977. **74**(7): p. 2766-70.
46. Varon, R., et al., *Kinetic analysis of reversible closed bicyclic enzyme cascades covering the whole course of the reaction*. Int J Biochem, 1994. **26**(6): p. 787-97.
47. Varon, R., et al., *Kinetic analysis of the opened bicyclic enzyme cascades*. Biol Chem Hoppe Seyler, 1994. **375**(6): p. 365-71.
48. Kholodenko, B.N., *Negative feedback and ultrasensitivity can bring about oscillations in the mitogen-activated protein kinase cascades*. Eur J Biochem, 2000. **267**(6): p. 1583-8.
49. Huang, C.Y. and J.E. Ferrell, Jr., *Ultrasensitivity in the mitogen-activated protein kinase cascade*. Proc Natl Acad Sci U S A, 1996. **93**(19): p. 10078-83.
50. Ferrell, J.E., Jr., *How regulated protein translocation can produce switch-like responses*. Trends Biochem Sci, 1998. **23**(12): p. 461-5.
51. Levchenko, A., J. Bruck, and P.W. Sternberg, *Scaffold proteins may biphasically affect the levels of mitogen-activated protein kinase signaling and reduce its threshold properties*. Proc Natl Acad Sci U S A, 2000. **97**(11): p. 5818-23.
52. Brightman, F.A. and D.A. Fell, *Differential feedback regulation of the MAPK cascade underlies the quantitative differences in EGF and NGF signalling in PC12 cells*. FEBS Lett, 2000. **482**(3): p. 169-74.
53. Haugh, J.M., A. Wells, and D.A. Lauffenburger, *Mathematical modeling of epidermal growth factor receptor signaling through the phospholipase C pathway: mechanistic insights and predictions for molecular interventions*. Biotechnol Bioeng, 2000. **70**(2): p. 225-38.



54. Schoeberl, B., et al., *Computational modeling of the dynamics of the MAP kinase cascade activated by surface and internalized EGF receptors*. Nat Biotechnol, 2002. **20**(4): p. 370-5.
55. Kholodenko, B.N., et al., *Quantification of short term signaling by the epidermal growth factor receptor*. J Biol Chem, 1999. **274**(42): p. 30169-81.
56. Kacser, H. and J.A. Burns, *The control of flux*. Symp Soc Exp Biol, 1973. **27**: p. 65-104.
57. Heinrich, R. and T.A. Rapoport, *A linear steady-state treatment of enzymatic chains. General properties, control and effector strength*. Eur J Biochem, 1974. **42**(1): p. 89-95.
58. Acerenza, L. and A. Cornish-Bowden, *Generalization of the double-modulation method for in situ determination of elasticities*. Biochem J, 1997. **327**(Pt 1): p. 217-24.
59. Fell, D.A. and H.M. Sauro, *Metabolic control and its analysis. Additional relationships between elasticities and control coefficients*. Eur J Biochem, 1985. **148**(3): p. 555-61.
60. Small, J.R. and D.A. Fell, *Covalent modification and metabolic control analysis. Modification to the theorems and their application to metabolic systems containing covalently modifiable enzymes*. Eur J Biochem, 1990. **191**(2): p. 405-11.
61. Small, J.R. and H. Kacser, *Responses of metabolic systems to large changes in enzyme activities and effectors. 1. The linear treatment of unbranched chains*. Eur J Biochem, 1993. **213**(1): p. 613-24.
62. Hofmeyr, J.H., H. Kacser, and K.J. van der Merwe, *Metabolic control analysis of moiety-conserved cycles*. Eur J Biochem, 1986. **155**(3): p. 631-41.
63. Kacser, H., *Recent developments beyond metabolic control analysis*. Biochem Soc Trans, 1995. **23**(2): p. 387-91.
64. Korzeniewski, B. and G.C. Brown, *Quantification of the relative contribution of parallel pathways to signal transfer: application to cellular energy transduction*. Biophys Chem, 1998. **75**(1): p. 73-80.
65. Fell, D.A., *Signal transduction and the control of expression of enzyme activity*. Adv Enzyme Regul, 2000. **40**: p. 35-46.
66. Krauss, S. and M.D. Brand, *Quantitation of signal transduction*. Faseb J, 2000. **14**(15): p. 2581-8.
67. Bruggeman, F.J., et al., *Modular response analysis of cellular regulatory networks*. J Theor Biol, 2002. **218**(4): p. 507-20.
68. Kholodenko, B.N., et al., *Untangling the wires: a strategy to trace functional interactions in signaling and gene networks*. Proc Natl Acad Sci U S A, 2002. **99**(20): p. 12841-6.
69. Arkin, A. and J. Ross, *Statistical Construction of Chemical-Reaction Mechanisms From Measured Time-Series*. Journal of Physical Chemistry, 1995. **99**(3): p. 970-979.
70. Arkin, A., P.D. Shen, and J. Ross, *A test case of correlation metric construction of a reaction pathway from measurements*. Science, 1997. **277**(5330): p. 1275-1279.
71. Resing, K.A. and N.G. Ahn, *Applications of mass spectrometry to signal transduction*. Prog Biophys Mol Biol, 1999. **71**(3-4): p. 501-23.

72. Gygi, S.P., et al., *Quantitative analysis of complex protein mixtures using isotope-coded affinity tags*. Nat Biotechnol, 1999. **17**(10): p. 994-9.
73. Oda, Y., et al., *Accurate quantitation of protein expression and site-specific phosphorylation*. Proc Natl Acad Sci U S A, 1999. **96**(12): p. 6591-6.
74. Pandey, A., et al., *Analysis of receptor signaling pathways by mass spectrometry: identification of vav-2 as a substrate of the epidermal and platelet-derived growth factor receptors*. Proc Natl Acad Sci U S A, 2000. **97**(1): p. 179-84.
75. Zhou, H., J.D. Watts, and R. Aebersold, *A systematic approach to the analysis of protein phosphorylation*. Nat Biotechnol, 2001. **19**(4): p. 375-8.
76. Soskic, V., et al., *Functional proteomics analysis of signal transduction pathways of the platelet-derived growth factor beta receptor*. Biochemistry, 1999. **38**(6): p. 1757-64.
77. Lewis, T.S., et al., *Identification of novel MAP kinase pathway signaling targets by functional proteomics and mass spectrometry*. Mol Cell, 2000. **6**(6): p. 1343-54.
78. Tolia, K. and C.L. Carpenter, *In vitro interaction of phosphoinositide-4-phosphate 5-kinases with Rac*. Methods Enzymol, 2000. **325**: p. 190-200.
79. Casals, I., J.L. Villar, and M. Riera-Codina, *A straightforward method for analysis of highly phosphorylated inositols in blood cells by high-performance liquid chromatography*. Anal Biochem, 2002. **300**(1): p. 69-76.
80. Higashida, H., et al., *Measurement of adenylyl cyclase by separating cyclic AMP on silica gel thin-layer chromatography*. Anal Biochem, 2002. **308**(1): p. 106-11.
81. Johnson, R.A., R. Alvarez, and Y. Salomon, *Determination of adenylyl cyclase catalytic activity using single and double column procedures*. Methods Enzymol, 1994. **238**: p. 31-56.
82. Lee, S.K., et al., *Advantages of calcium green-1 over other fluorescent dyes in measuring cytosolic calcium in platelets*. Thrombosis and Haemostasis, 1999: p. 167-167.
83. George, S.E., et al., *A high-throughput glow-type aequorin assay for measuring receptor-mediated changes in intracellular calcium levels*. Anal Biochem, 2000. **286**(2): p. 231-7.
84. Cooper, D.M., *Calcium-sensitive adenylyl cyclase/aequorin chimeras as sensitive probes for discrete modes of elevation of cytosolic calcium*. Methods Enzymol, 2002. **345**: p. 105-12.
85. Zhu, H., et al., *Analysis of yeast protein kinases using protein chips*. Nat Genet, 2000. **26**(3): p. 283-9.
86. Haab, B.B., M.J. Dunham, and P.O. Brown, *Protein microarrays for highly parallel detection and quantitation of specific proteins and antibodies in complex solutions*. Genome Biol, 2001. **2**(2).
87. Figeys, D. and D. Pinto, *Proteomics on a chip: promising developments*. Electrophoresis, 2001. **22**(2): p. 208-16.
88. Blackstock, W.P. and M.P. Weir, *Proteomics: quantitative and physical mapping of cellular proteins*. Trends Biotechnol, 1999. **17**(3): p. 121-7.
89. Koshland, D.E., Jr., *The era of pathway quantification [comment]*. Science, 1998. **280**(5365): p. 852-3.

# 2 ACTIVATION RATIO ANALYSIS

Any investigation into the description and analysis of a system requires decisions to be made regarding the scope and complexity that will be considered. The scale (in time, space, and concentration) must be selected to sufficiently describe the features of the system while maintaining feasible bounds on intellectual, conceptual and experimental requirements. As the primary objective of this work was to develop methods to characterize and describe the structure of signaling networks, the scope therefore was on examination of interactions between signaling components. This proceeded in two stages: first, the construction and examination of models for several network modules, to determine mathematical expressions based on potentially measurable quantities that reflect the relationship between intermediates while reducing the system complexity; and second, development of an algorithm to reconstruct the modular arrangement based upon those measurements.

In this chapter the first stage is described, where examination of the rate equations describing reactions occurring in signal transduction is combined with practical understanding of measurement limitations to suggest a measurable quantity for a signaling intermediate that reflects its connection to other components in the network. The activation ratio, defined as the ratio between active and inactive forms of an interconverting intermediate, depends quantitatively and qualitatively upon the structural relationship of that intermediate with other network species. This is shown below for a variety of simplified network substructures that can be combined to describe realistic signaling networks.

## **2.1. Methods**

Deterministic, kinetic models were used as simulators of simple signaling arrangements as have been described previously [1-6]. These types of models have been used extensively to examine qualitative features in signaling pathways, which have been later observed experimentally. Enzyme-catalyzed reactions were assumed to follow

simple Michaelis-Menten kinetics, and parameters were varied to explore the patterns of behavior for each network arrangement. Concentrations of non-protein reactants (such as  $Mg^{2+}$ , ATP, and water) were assumed to remain constant. Models were developed as a set of coupled ordinary differential equations in MATLAB 5.2 (Mathworks, Inc) and integrated until steady state using the “ode15s” algorithm. Steady-state concentrations of active and inactive fractions of components at various input stimulus levels were used as “data” in the analytical approach presented here. Further details, parameter values, and sample MATLAB files are included in the Appendices 1 and 2.

## 2.2. Isolated interconverting cycle

The most basic signaling unit can be imagined to consist of a source driving the conversion of a target from one state to another. The input signal is the activity of the source, while the state of the target is the output. Relaxation of the system to its original state requires a competing force to counteract the effect of the source. In biological signaling, this type of interaction is often realized by simple binding between two molecules, such as a ligand to its receptor, or through chemical modification of the target, for example via phosphorylation or acetylation of proteins or cyclization of ATP to cAMP [7]. These cases are shown in Figure 2.1, along with the resulting reaction schemes assuming simple association kinetics for binding or Michaelis-Menten kinetics for an enzymatic modification cycle.

<u>Interaction</u>	<u>Diagram</u>	<u>Reaction scheme</u>
Binding		$A + B \xrightleftharpoons[d]{a} [A \cdot B]$
Covalent modification		$A + E_1 \xrightleftharpoons[d_1]{a_1} [E_1 \cdot A] \xrightarrow{k_1} A^* + E_1$ $A^* + E_2 \xrightleftharpoons[d_2]{a_2} [E_2 \cdot A^*] \xrightarrow{k_2} A + E_2$

Figure 2.1. Single cycle diagrams and reaction scheme.

It may appear peculiar at first to consider binding as an example of an interconverting cycle, but several key features of signaling cycles can be examined very simply in this system. Here, the state (concentration) of molecule B determines the state of the target A, namely, in its free form or bound in the noncovalent complex [A·B]. This interaction can be quantified by the association and dissociation rate constants, a and d. Although we can write expressions for the dynamic interaction between these molecules, it is more common to consider the situation at equilibrium, where the rates of the association and dissociation reactions are equal:

$$a A B = d [A \cdot B] \quad (2.1)$$

Or, by rearranging:

$$AR_A \equiv \frac{[A \cdot B]}{A} = \frac{a}{d} B = K_a B \quad (2.2)$$

According to Equation 2.2, at equilibrium the ratio between the amounts of A complexed with B (“active”) and that of free A (“inactive”) is directly proportional to the amount of free “activator” B. The association constant  $K_a$ , which is the natural parameter used to describe the binding interaction, appears as the coefficient of the activator B. Thus the *activation ratio*  $AR_A$ , defined as the ratio between active and inactive forms of a signaling intermediate A, is linearly dependent upon the concentration of its activator, and the two are related by an *activation factor*, which quantitatively measures the interaction between the source and target. In this case, the activation factor  $\alpha_B^A = K_a$ .

Since the total amount of A,  $A_T = A + [A \cdot B]$ , we can also write:

$$\frac{[A \cdot B]}{A} = \frac{[A \cdot B]}{A_T - [A \cdot B]} = \frac{[A \cdot B]/A_T}{1 - [A \cdot B]/A_T} = K_a B \quad (2.3)$$

$$\frac{[A \cdot B]}{A_T} = \frac{K_a B}{1 + K_a B} \quad (2.4)$$

We can estimate  $K_a$  from a set of measurements for the fraction of A in its “active”, complexed form through either of Equations 2.3-2.4. Normally Equation 2.4 is used and  $K_a$  found via a nonlinear regression procedure, but by rescaling the problem as in Equation 2.3 a simple linear regression approach can be used.

Note that each expression is written in terms of the amount of free activator B, not the total amount  $B_T (= B + [A \cdot B])$ . For most binding studies it is assumed that  $B \approx B_T$ , i.e. that the complexes do not significantly deplete the amount of B available to bind A. More explicitly, this requires that  $K_a A_T \ll 1$ . Although this assumption can be realized during *in vitro* binding experiments, under *in vivo* conditions it is likely to be invalid, since most relevant protein interactions occur with high specificity. Nevertheless, Equations 2.1-2.2 will still hold, and combining them with conservation relationships for A and B we find that:

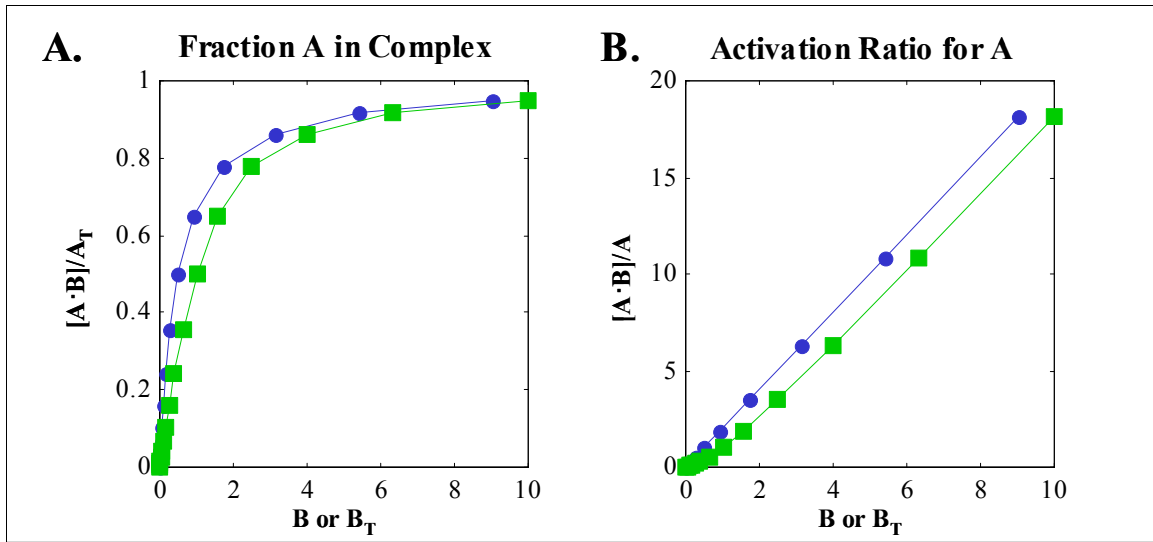
$$\frac{[A \cdot B]}{A_T} = K_a A_T \left( 1 - \frac{[A \cdot B]}{A_T} \right) \left( \frac{B_T}{A_T} - \frac{[A \cdot B]}{A_T} \right) \quad (2.5)$$

This equation can be solved to yield:

$$\frac{[A \cdot B]}{A_T} = \frac{1}{2} \left[ 1 + \frac{B_T}{A_T} + \frac{1}{K_a A_T} - \sqrt{\left( 1 + \frac{B_T}{A_T} + \frac{1}{K_a A_T} \right)^2 - 4 \frac{B_T}{A_T}} \right] \quad (2.6)$$

$$\frac{[A \cdot B]/A_T}{1 - [A \cdot B]/A_T} = \frac{1}{2} K_a A_T \left[ \frac{B_T}{A_T} + \sqrt{\left( 1 + \frac{B_T}{A_T} + \frac{1}{K_a A_T} \right)^2 - 4 \frac{B_T}{A_T}} - \frac{1}{K_a A_T} - 1 \right] \quad (2.7)$$

Under saturating conditions, where  $K_a A_T$  is significant, the expressions for the fraction of A complexed to B (Equation 2.6) as well as the ratio to free A (Equation 2.7) both become significantly more complicated, as well as being much more difficult to extract the relevant quantitative parameter,  $K_a$ . Using a mathematical package such as MATLAB, Equations 2.3-2.4 and 2.6-2.7 can be plotted, as shown in Figure 2.2. The system is saturated ( $K_a = 2$ ,  $A_T = 1$ ), but we can readily observe that the activation ratio shows a generally linear dependence upon  $B_T$ . On the other hand, the fraction of A complexed ( $[A \cdot B]/A_T$ ) appears to be approximately hyperbolic in shape, with some curvature seen for low values of  $B_T$ .



**Figure 2.2.** Simulation results for simple binding under saturating conditions. A) Fraction of A complexed ( $[A \cdot B]/A_T$ ) or B) activation ratio ( $[A \cdot B]/A$ ) plotted against free B (circles) or  $B_T$  (squares).

If we assume for the moment that there are other compounds which may also interact with either A or B, then Equations 2.1-2.2 will still be valid, but new conservation relationships will apply, as will additional binding expressions. The resulting coupled system of equations will have to be solved simultaneously, preventing a closed-form solution for any single component. The concentrations of these other species will play a role in the expressions for the fraction active and thus decoupling the interactions becomes difficult. By focusing only on the free species, the interaction between A and B can be isolated via Equations 2.1-2.2.

Next, consider the case of a covalent modification cycle shown in Figure 2.1, where the target converts between two states A and  $A^*$  by the competing action of enzymes  $E_1$  and  $E_2$ . In protein kinase cascades,  $E_1$  and  $E_2$  are a kinase and phosphatase, respectively, and A and  $A^*$  represent the nonphosphorylated and phosphorylated forms of the intermediate A. Here  $[E_1 \cdot A]$  and  $[E_2 \cdot A^*]$  represent the enzyme-substrate complexes, while  $E_1$  and  $E_2$  are the free concentrations of enzymes, and  $a_i$ ,  $d_i$ , and  $k_i$  are the association, dissociation, and catalytic rate constants for reaction  $i$ . There are a total of six species in this system, and their time-dependent behavior can be described by the following rate equations:

$$\frac{dA}{dt} = -a_1 A E_1 + d_1 [E_1 \cdot A] + k_2 [E_2 \cdot A^*] \quad (2.8)$$

$$\frac{dA^*}{dt} = -a_2 A^* E_2 + d_2 [E_2 \cdot A^*] + k_1 [E_1 \cdot A] \quad (2.9)$$

$$\frac{d[E_1 \cdot A]}{dt} = a_1 A E_1 - (d_1 + k_1) [E_1 \cdot A] \quad (2.10)$$

$$\frac{d[E_2 \cdot A^*]}{dt} = a_2 A^* E_2 - (d_2 + k_2) [E_2 \cdot A^*] \quad (2.11)$$

$$\frac{dE_1}{dt} = -a_1 A E_1 + (d_1 + k_1) [E_1 \cdot A] \quad (2.12)$$

$$\frac{dE_2}{dt} = -a_2 A^* E_2 + (d_2 + k_2) [E_2 \cdot A^*] \quad (2.13)$$

These species are further coupled by conservation equations:

$$E_{1T} = E_1 + [E_1 \cdot A] \quad (2.14)$$

$$E_{2T} = E_2 + [E_2 \cdot A^*] \quad (2.15)$$

$$A_T = A + A^* + [E_1 \cdot A] + [E_2 \cdot A^*] \quad (2.16)$$

If it is possible to measure the time-dependent concentrations of three of the species, for example A,  $[E_1 \cdot A]$ , and  $[E_2 \cdot A^*]$ , along with the total concentrations for each component ( $A_T$ ,  $E_{1T}$ ,  $E_{2T}$ ), then Equations 2.8-2.16 can be used to estimate the values of the parameters of the system, namely the rate constants a, d, and k for both reactions. It is these rate constants that collectively are quantitative measures of the interaction between the components A,  $E_1$ , and  $E_2$  in this simple system. This is a difficult task for even an isolated *in vitro* system, and totally infeasible for a signaling intermediate within a cell. In protein kinase cascades it may be possible to determine the identity and concentration of kinase  $E_1$ , but often the phosphatase  $E_2$  is undefined and may be one of several nonspecific enzymes. Thus, it may be unrealistic to consider  $E_2$  or  $[E_2 \cdot A^*]$  as measurable quantities, and it becomes impossible to solve Equations 2.8-2.16. Furthermore, estimation of the individual rate constants may not be informative, particularly if the behavior of a network of intermediates is being investigated, where a single quantitative parameter describing the interaction between  $E_1$  and A is preferable.

Instead of considering the dynamic behavior of this system, we concentrate on the relationship between intermediates at steady state, analogous to equilibrium for the binding system. In this case Equations 2.8-2.13 are equal to zero. Equations 2.10-2.11 can then be rearranged to yield:



$$[E_1 \cdot A] = \frac{a_1 A E_1}{d_1 + k_1} = \frac{A E_1}{K_{m1}} \quad (2.17)$$

$$[E_2 \cdot A^*] = \frac{a_2 A^* E_2}{d_2 + k_2} = \frac{A^* E_2}{K_{m2}} \quad (2.18)$$

Where  $K_{m1}$  and  $K_{m2}$  are the Michaelis constants for enzymes  $E_1$  and  $E_2$ , respectively. Also, at steady state the two net reaction rates must be equal:

$$k_1 [E_1 \cdot A] = k_2 [E_2 \cdot A^*] \quad (2.19)$$

Substituting Equations 2.17-2.18 into Equation 2.19 and rearranging:

$$AR_A \equiv \frac{A^*}{A} = \frac{k_1 E_1 / K_{m1}}{k_2 E_2 / K_{m2}} \equiv \alpha_1^A E_1 \quad (2.20)$$

From Equation 2.20, it is apparent that the activation ratio  $AR_A$  is once again linearly proportional to the concentration of free activating enzyme  $E_1$  (that is, unbound by substrate A). The unknown kinetic parameters for the enzymes, as well as (most likely unmeasurable)  $E_2$ , are collected together in the activation factor  $\alpha_1^A$ . The activation factor represents the sensitivity of the activation ratio for A with respect to  $E_1$ , and is therefore a quantitative measure of the potential for  $E_1$  to activate A. As  $k_1$  increases or  $K_{m1}$  decreases,  $E_1$  becomes a more powerful activator of A, and  $\alpha_1^A$  increases. Similarly, as  $k_2$  increases or  $K_{m2}$  decreases,  $E_2$  is a more powerful inactivator, and therefore  $E_1$  is a relatively weaker activator of A.

An explicit closed-form solution for  $A^*$  in terms of the other parameters is possible for a few special cases. If the enzyme-substrate complexes can be considered negligible in the conservation relationship for the substrate A (Equation 2.16), such that  $A^* + A \approx A_T$ , it can be shown that [8]:

$$\frac{A^*}{A_T} = \frac{\phi + \sqrt{\phi^2 + 4K_2 \frac{V_1}{V_2} \left( \frac{V_1}{V_2} - 1 \right)}}{2 \left( \frac{V_1}{V_2} - 1 \right)} \quad (2.21)$$

where  $\phi = \left( \frac{V_1}{V_2} - 1 \right) - K_2 \left( \frac{K_1}{K_2} + \frac{V_1}{V_2} \right)$ ,  $V_1 = k_1 E_{1T}$ ,  $V_2 = k_2 E_{2T}$ ,  $K_1 = K_{m1}/A_T$ , and

$K_2 = K_{m2}/A_T$ . To satisfy this assumption,  $E_{1T}/K_{m1}$  and  $E_{2T}/K_{m2}$  must both be much less than one. In signaling networks involving enzymatic cascades, the substrate of one step is the activating enzyme of another step, and concentrations of all species may be significant relative to saturation constants [9]. Thus this may be a poor assumption, but even by making it the expression in Equation 2.21 shows significant complexity as well as the lack of a single quantifiable parameter. Is  $K_1$ ,  $K_2$ , their ratio, or some other parameter a significant representation of this system? Comparing Equations 2.20 and 2.21, it is easy to observe the simplicity of activation ratios and activation factors for describing the interaction between source  $E_1$  and target  $A$ .

In the event that enzyme-substrate complexes constitute a negligible fraction of the total enzyme concentrations as well as substrate, which would occur when both enzymes are far from saturation,  $E_1 \approx E_{1T}$ ,  $E_2 \approx E_{2T}$  and:

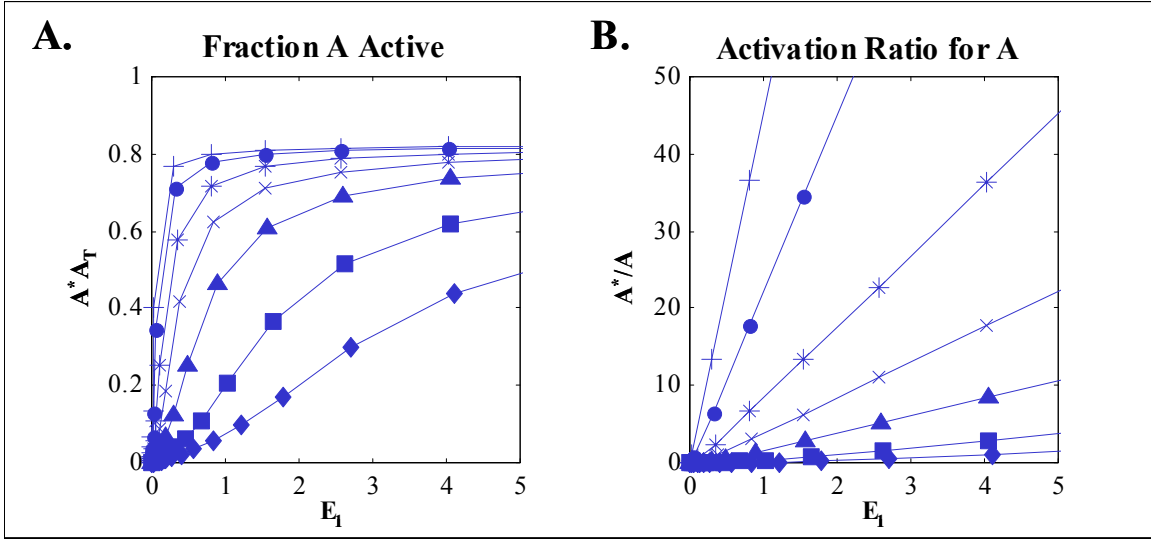
$$\frac{A^*}{A_T} = \frac{k_1 E_{1T} / K_{m1}}{k_1 E_{1T} / K_{m1} + k_2 E_{2T} / K_{m2}} = \frac{\alpha_1^A E_{1T}}{1 + \alpha_1^A E_{1T}} \quad (2.22)$$

where  $\alpha_1^A = \frac{k_1 / K_{m1}}{k_2 E_{2T} / K_{m2}}$ . In this very idealized situation, the obvious

quantitative measure describing this system is again  $\alpha_1^A$ , the same activation factor from Equation 2.20. While Equation 2.22 yields a hyperbolic relationship between  $A^*$  and  $E_{1T}$ , the activation ratio (Equation 2.20) is still linear. In the event that enzyme-substrate complexes cannot be neglected, the fraction active cannot be determined explicitly but is the solution of a third-order equation, and thus mathematical simulation is necessary to determine the steady-state behavior of the system.

Simulation results for the individual cycle are shown in Figure 2.3. The association constant for activation of  $A$  ( $a_1$ ) was varied so as to adjust  $K_{m1}$  50-fold without affecting  $k_1$ , with other parameters set such that enzyme-substrate complexes are significant for balances on substrate as well as enzymes. The fraction of  $A$  in the active form,  $A^*/A_T$ , is shown in Figure 2.3A, while activation ratios  $AR_A$  are plotted in Figure

2.3B. The curvature seen for the fraction active in Figure 2.3A is replaced by a set of lines for the activation ratios in Figure 2.3B. The activation factor  $\alpha_1^A$  as defined in Equation 2.20 is not necessarily constant, because the concentration of  $E_2$  is not necessarily constant and may be indirectly influenced by  $E_1$ . Nevertheless,  $\alpha_1^A$  does approach a limiting value and the curves in Figure 2.3B are well approximated by straight lines, and  $\alpha_1^A$  can be calculated easily as the slope.

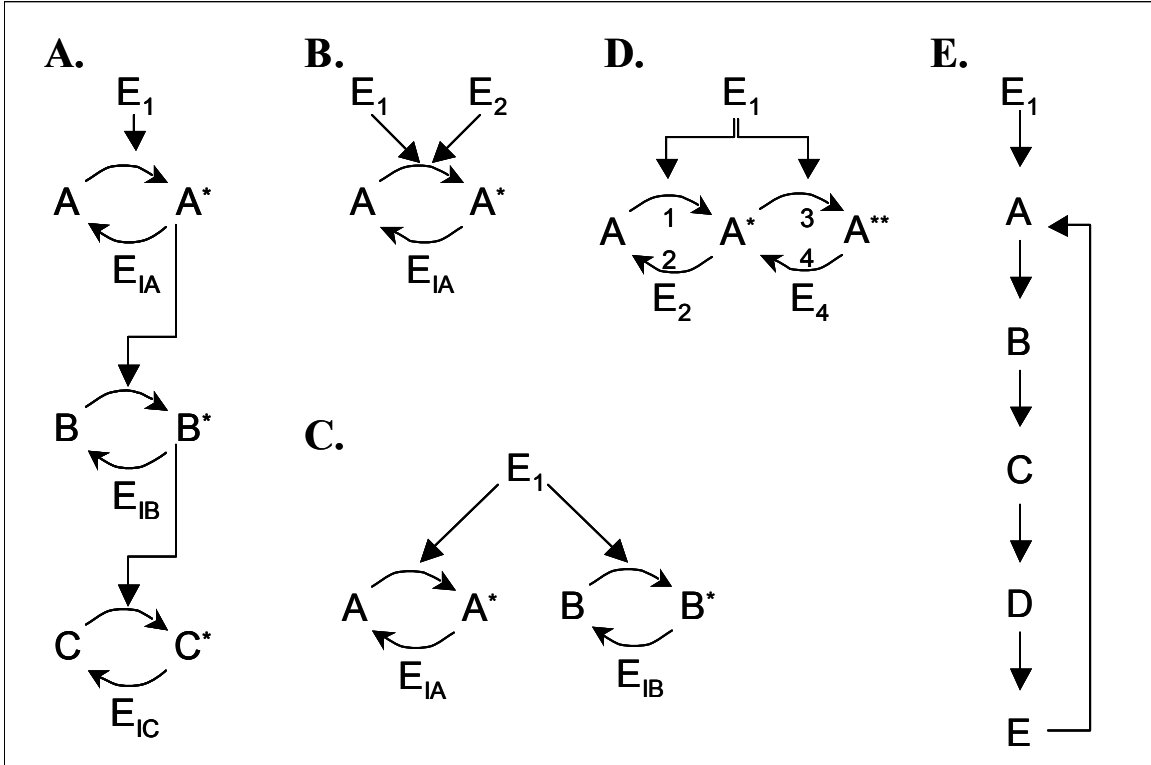


**Figure 2.3.** Simulation results for an individual covalent-modification cycle as shown in Figure 2.1. A) Fraction of A activated ( $A^*/A_T$ ) and B) Activation ratios  $AR_A$  ( $A^*/A$ ) for the simple cycle, plotted against free activating enzyme  $E_1$ . Parameter values:  $k_1 = 10$ ,  $k_2 = 10$ ,  $K_{m2} = 1$ ,  $E_{2T} = 1$ ,  $A_T = 10$ ,  $K_{m1} = 20$  (diamonds), 10 (squares), 4 (triangles), 2 (x's), 1 (stars), 0.4 (circles), 0.2 (+'s).

### 2.3. Extension to simple network arrangements

In the relatively simple case of an isolated cycle, use of activation ratios yields a simple linear relationship between an activator and its target, even under saturating conditions for enzymes and substrates. The power of this approach becomes more apparent as increasingly complicated signaling systems are considered. This is because Equations 2.2 and 2.20 continue to be valid when the cycle is no longer isolated, but rather embedded within a signaling network. Of particular interest are the cases of converging and diverging pathways and linear cascades, as the behavior of these model systems can be combined to examine the effects of multiple inputs on a set of interconnected intermediates, in the absence of feedback. Multiple activation steps and cascades with feedback are further complications that can also be examined using activation ratios.

These arrangements are drawn schematically in Figure 2.4. For each case, expressions governing the steady-state ratio between active and inactive forms of each component are developed. As the systems become complicated, closed-form analytical solutions for the activated fraction of various species become impossible. Therefore, simulations were developed using time-dependent equations analogous to Equations 2.8-2.13 and integrated until steady state to demonstrate the behavior of each system.



**Figure 2.4.** Diagrams of extended signaling structures. A) linear cascade, B) converging pathways, C) diverging pathways, D) dual activation steps, E) cascade with feedback (single-step activation as in A).

### 2.3.1. Linear Cascade

A common arrangement of signaling intermediates is a linear cascade of enzymes, where the activated form of one intermediate catalyzes the activation of the succeeding intermediate, as shown in Figure 2.4A. Following a similar analysis for each cycle in the cascade as done previously for the isolated cycle, using expressions analogous to Equations 2.17-2.20, the activation ratios for the intermediates are:

$$AR_A \equiv \frac{A^*}{A} = \frac{k_1 E_1 / K_{m1}}{k_{IA} E_{IA} / K_{mIA}} = \alpha_1^A E_1 \quad (2.23)$$

$$AR_B \equiv \frac{B^*}{B} = \frac{k_2 A^* / K_{m2}}{k_{IB} E_{IB} / K_{mIB}} = \alpha_A^B A^* \quad (2.24)$$

$$AR_C \equiv \frac{C^*}{C} = \frac{k_3 B^* / K_{m3}}{k_{IC} E_{IC} / K_{mIC}} = \alpha_B^C B^* \quad (2.25)$$

The action of each step upon the next is the same as if the cycle were isolated, as discussed above. However, when considering an indirect effect, for example  $E_1$  upon  $AR_B$ , the results take a quite different form. As described in Section 2.2, the equations are written in terms of free species, unbound by enzymes or targets. If for simplicity enzyme-substrate complexes can be considered negligible, then  $A^* + A \approx A_T$  and  $B^* + B \approx B_T$ , and it can be shown that:

$$AR_B = \frac{\alpha_1^A \alpha_A^B A_T E_1}{1 + \alpha_1^A E_1} \quad (2.26)$$

$$AR_C = \frac{\alpha_A^B \alpha_B^C B_T A^*}{1 + \alpha_A^B A^*} = \frac{\alpha_1^A \alpha_A^B \alpha_B^C B_T A_T E_1}{1 + \alpha_1^A E_1 (1 + \alpha_A^B A_T)} \quad (2.27)$$

The sensitivity of the overall cascade is a product of the sensitivity at each individual level (multistep sensitivity), as has been described previously [9-12]. However, the expression for the activation ratios changes in form, from linear to hyperbolic, depending on which upstream enzyme is being considered. For example, although  $AR_B$  is linear with respect to  $A^*$ , it is hyperbolic with respect to  $E_1$ .  $AR_C$  is hyperbolic with respect both to  $A^*$  and  $E_1$ , but linear with respect to  $B^*$ . This radical change in form can be readily visualized graphically and realized numerically through a linear or hyperbolic regression. We can thus use this approach to suggest if a step is missing between two intermediates of interest. Note, however, it is *not* possible to distinguish between one or more missing steps. In this example, it is possible to know that C is indirectly downstream of  $E_1$  and A, but not by how many steps. By establishing direct links between  $E_1$  to A, A to B, and then B to C, however, the cascade structure can be realized.

Equations 2.26-2.27 explicitly arise from the assumption that enzyme-substrate complexes can be neglected in the conservation relationships for A and B. This assumption specifically requires that  $E_{1T}/K_{m1}$ ,  $E_{IA T}/K_{mIA}$ ,  $A_T/K_{m2}$ , and  $E_{IB T}/K_{mIB}$  are

significantly less than one, essentially meaning that  $K_m$  values for each reaction should be large relative to the concentrations of all species. Prior experience with MAP Kinase cascades indicates this to be highly unlikely [9]. Nevertheless, the *patterns* for activation ratios still hold even if these assumptions are relaxed. Thus, the plots for activation ratios will continue to be linear for direct effects and hyperbolic in shape for indirect effects. Although it is difficult to demonstrate this analytically, it can be readily seen using simulation results for signaling cascades under saturating enzyme conditions.

We have thus far seen that plots of activation ratios of intermediates will be linear or hyperbolic when plotted against intermediates directly and indirectly upstream, respectively. What if we now look at *downstream* intermediates? In other words, what do the plots of activation ratios of intermediates against their direct and indirect targets look like? Again for simplicity, if enzyme-substrate complexes are neglected, then the following results can be obtained:

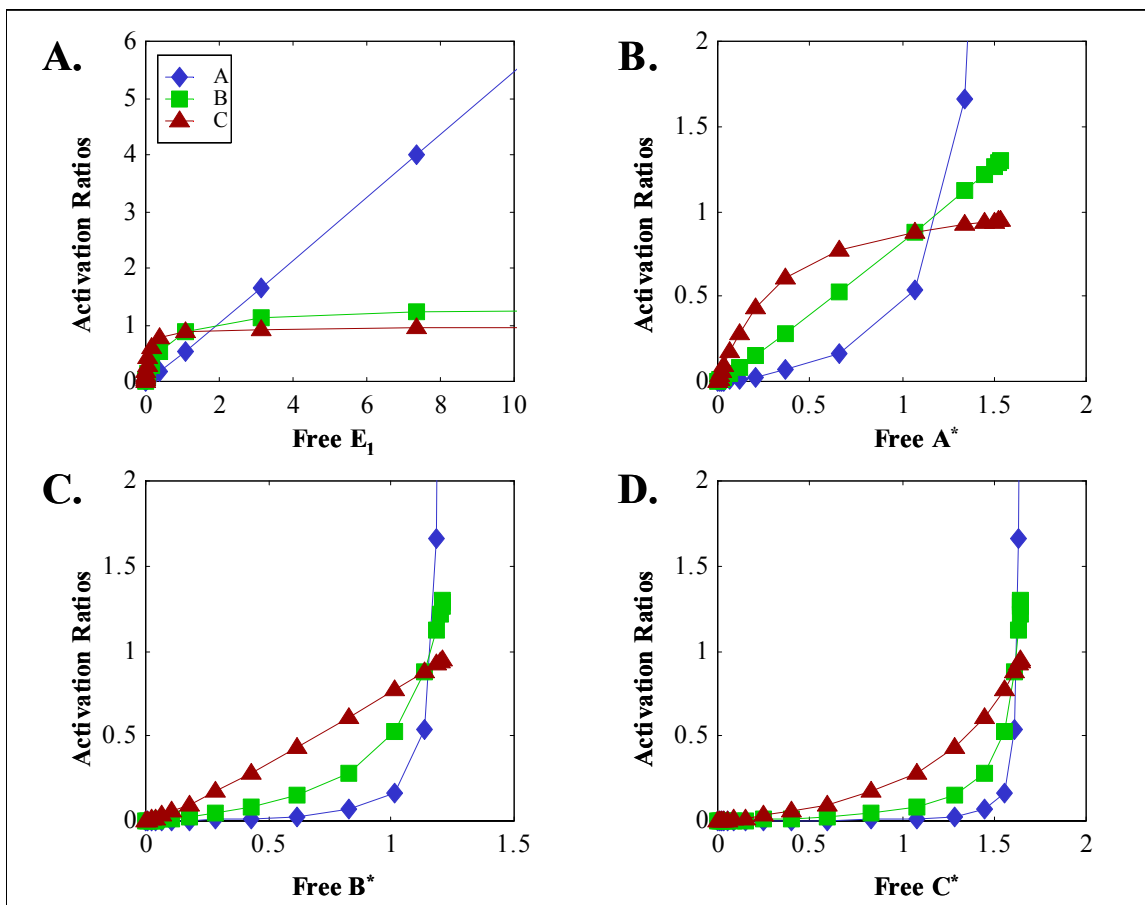
$$AR_B = \frac{C^*}{\alpha_B^C B_T C_T - (1 + \alpha_B^C B_T) C^*} \quad (2.28)$$

$$AR_A = \frac{B^*}{\alpha_A^B A_T B_T - (1 + \alpha_A^B A_T) B^*} \quad (2.29)$$

$$= \frac{C^*}{\alpha_A^B \alpha_B^C A_T B_T C_T - [1 + \alpha_A^B A_T (1 + \alpha_B^C B_T)] C^*}$$

The activation ratios  $AR_A$  and  $AR_B$  in Equations 2.28-2.29 take the form of an *inverse hyperbola* (technically, the upper left quadrant of a hyperbolic section). Thus, the activation ratio for an inverted response has a quite distinct functionality from either type of forward response. It should be further noted that the forms of Equations 2.28-2.29 are the same, albeit with different parameters. Therefore it is impossible to distinguish between a direct and indirect inverted response from the mathematical structure of the relationship; it can only be said that the presumed target and activator are actually in reverse order. As above, Equations 2.28-2.29 hold explicitly only when enzyme-substrate complexes are negligible, but the same patterns as predicted from those equation will still appear if this assumption is relaxed.

Plots of activation ratios for the model cascade of Figure 2.4A are shown in Figure 2.5. Each step in the pathway is saturated, with  $A_T$ ,  $B_T$ , and  $C_T$  all equal to 10 and  $K_m$  values in the range of 0.50-4. Thus, enzyme-substrate complexes will not be negligible compared to the free species in this case. In Figure 2.5A, the activation ratios for each intermediate are plotted against  $E_1$ . As expected from Equations 2.23-2.27, the curve for  $AR_A$  is linear while the curves for  $AR_B$  and  $AR_C$  are hyperbolic. Similarly, in Figure 2.5B the plot of  $AR_B$  against  $A^*$  is linear and  $AR_C$  against  $A^*$  is hyperbolic, and in Figure 2.5C,  $AR_C$  plotted against  $B^*$  is linear. The inverse hyperbolae expected for  $AR_A$  and  $AR_B$  against  $C^*$  are seen in Figure 2.5D, in accordance with Equations 2.28-2.29. Activation ratios for an intermediate plotted against itself, such as  $AR_A$  against  $A^*$ , also appear inverse hyperbolic, which is readily understandable when  $AR_A \approx A^*/(A_T - A^*)$ . The activation ratio data plotted on the ordinate is the same in each graph, and the figures differ only in which component is used for the abscissa.



**Figure 2.5.** Activation ratios for the linear cascade of Figure 2.4A, plotted against free concentrations of  $E_1$  (A),  $A^*$  (B),  $B^*$  (C), or  $C^*$  (D). Ratios for A: diamonds, B: squares, C: triangles.

In summary, comparing Equations 2.24, 2.26, and 2.28 it is possible to see that the activation ratio for an intermediate B ( $AR_B$ ) will be linear, hyperbolic, or inverse hyperbolic when plotted against its direct activator ( $A^*$ ), indirect upstream activator ( $E_1$ ), or downstream target ( $C^*$ ), respectively. The functional form of the activation ratios reflects the relationship between activator and target. Activation ratios can therefore be a powerful tool to arrange intermediates in a cascade based on simultaneous measurements of activation for each component. Furthermore, missing steps can be detected through the lack of any direct steps as demonstrated through linear activation ratio plots.

### 2.3.2. Converging Pathways

In converging pathways, two separate enzymes act independently to activate an intermediate, as shown in Figure 2.4B. One example is the activation of Pbs2p by either Ssk2p/22p isoforms or Ste11p in the yeast high osmolarity (HOG) pathway[13]. In this case, either enzyme  $E_1$  or  $E_2$  can bind and activate A, although they cannot both bind A simultaneously. The activation of A therefore becomes a combination of the effects from the two enzymes, and the expression for the activation ratio is:

$$AR_A = \frac{k_1 E_1 / K_{m1}}{k_{IA} E_{IA} / K_{mIA}} + \frac{k_2 E_2 / K_{m2}}{k_{IA} E_{IA} / K_{mIA}} = \alpha_1^A E_1 + \alpha_2^A E_2 \quad (2.30)$$

Equation 2.30 shows that the activation ratio for an intermediate of converging pathways is a linear combination of terms arising from (and only dependent upon) each activator. The effects of the two activating enzymes are thus completely separated. Moreover, the expression for each enzyme in Equation 2.30 is the same as if  $E_1$  and  $E_2$  were acting upon unrelated substrates. Thus each enzyme is unaffected by the presence of the other. This must be the case since it is conceptually possible to separate one enzyme into two identical pools, and we would expect that the total effect of the two pools would be indistinguishable from the original state.

Since the enzyme effects are separated in Equation 2.30, it is possible to calculate the activation ratio for each enzyme by varying them independently. By keeping  $E_2$  constant at any value and varying  $E_1$ , it is possible to calculate  $\alpha_1^A$ ;  $\alpha_2^A$  can be similarly determined by keeping  $E_1$  constant. These two can be compared to indicate the relative strength of the two branches on the activation of A. Varying both simultaneously will



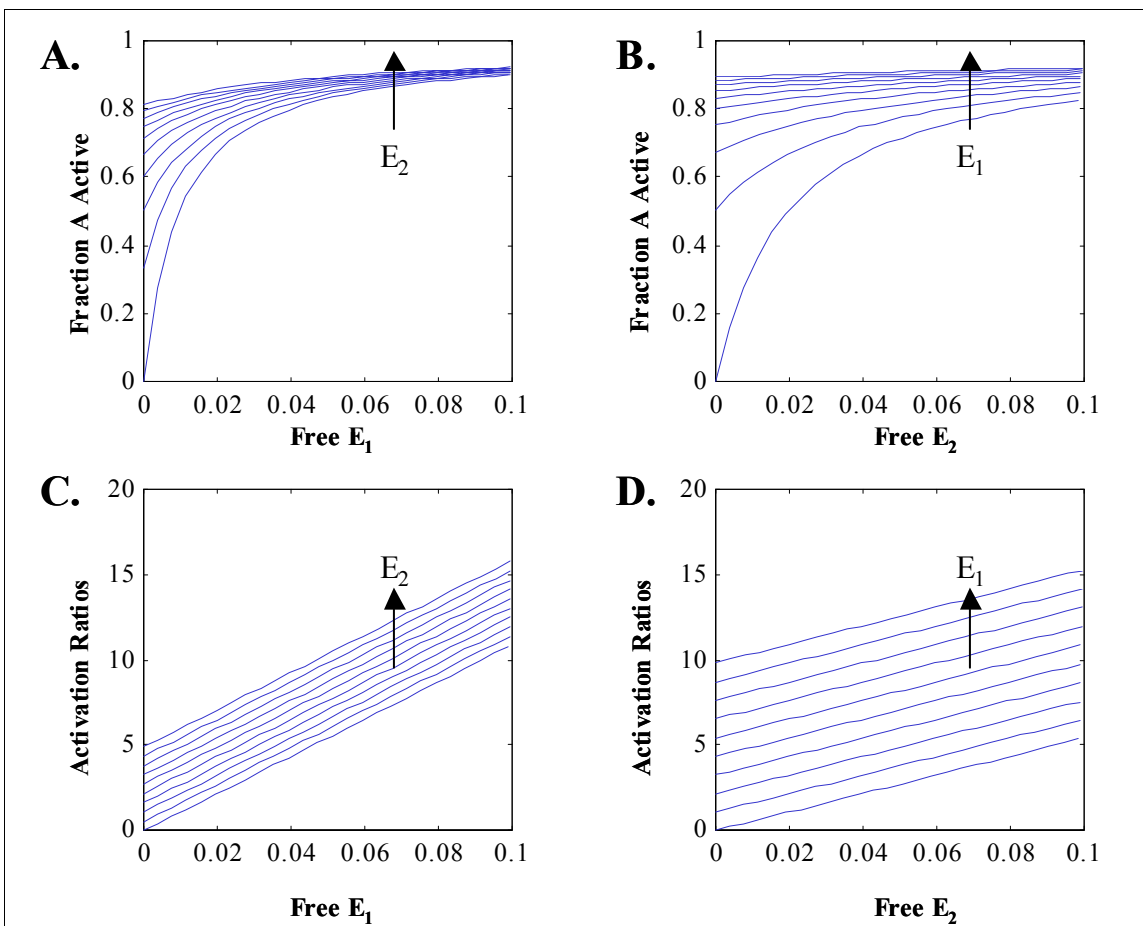
lead to an additive effect that can be predicted using the parameters calculated for each enzyme in isolation, or calculated using multiple linear regression techniques. Also, the presence of a second activating enzyme can be predicted since a plot of  $AR_A$  against  $E_1$ , for example, will not pass through the origin. Graphically, the activation ratio plots of an intermediate at a convergence point will appear as a set of parallel lines when plotted against either enzyme.

Simulation results for the converging pathway are shown in Figure 2.6. Parameter values were chosen such that  $K_{m1} = K_{m2}$  but  $k_1 = 2 k_2$ , thus we expect that  $\alpha_1^A = 2\alpha_2^A$ . In Figures 2.6A and 2.6B, the fraction active ( $A^*/A_T$ ) is shown, while in Figures 2.6C and 2.6D the activation ratios are plotted. From inspection of Figures 2.6A and 2.6B it is difficult to separate the effects of the two enzymes easily, or describe the relative strength of either enzyme in activating A, since the shapes of the curves changes as both enzyme concentrations are varied. On the other hand, plots of activation ratios  $AR_A$  against either enzyme yield the predicted set of parallel lines in Figures 2.6C and 2.6D.

Using simultaneous linear regression for both activators  $E_1$  and  $E_2$ , it is possible to obtain values for the activation factors  $\alpha_1^A$  and  $\alpha_2^A$ . The regression results for the data presented in Figure 2.6 are shown in Table 2.1. The activation factors calculated during regression are slightly lower than the maximum theoretical value, determined by inserting the parameters values into Equation 2.30. (Note that free  $E_{IA} = E_{IAT}/(1 + A^*/K_{mIA}) \approx 0.01/1.1$  for near-complete activation.) This discrepancy is simply due to the fact that some of the substrate A is retained bound to enzymes  $E_1$ ,  $E_2$ , and  $E_{IA}$ , so  $A^* < A_T$ . Nevertheless, the agreement between theoretical and calculated activation factors is excellent, as is the ratio between them (2.002).

**Table 2.1.** Regression results for converging pathway, using parameter values from Figure 2.6.

Enzyme	Average Slope	Range (95% CI)	Theoretical
$E_1$	109.63	109.62-109.64	110
$E_2$	54.76	54.76-54.77	55



**Figure 2.6.** Fractional activation and activation ratios for the converging cycle as shown in Figure 1C. A) and B), fractional activation ( $A^*/A_T$ ), C) and D), activation ratio ( $A^*/A$ ) plotted against free activating enzyme  $E_1$  (A) and (C) or  $E_2$  (B) and (D).

Equation 2.30 shows that the activation ratio for an intermediate A is a linear combination of the effects from the direct activators  $E_1$  and  $E_2$ . Furthermore, the expression for each term is the same as if there were no second activator. What if either  $E_1$  or  $E_2$  (or both) are not direct activators of A? In this case the term for the indirect activator changes from linear to hyperbolic in form, just as was seen in Equations 2.26-2.27 for a linear cascade. Instead of a set of parallel lines, plots of activation ratios for the common target will appear as a set of hyperbolic curves when plotted against the indirect activator. It is still possible to determine that there are two activators for the intermediate, but the plots cannot be used to calculate activation factors for the two direct activators independently.

### 2.3.3. Diverging Pathways

One intermediate may have multiple targets, enabling distribution of the signal. Typically, receptor-linked enzymes, such as the EGFR, insulin (IR), and PDGFR kinases, act on several substrates including other receptor molecules and adapter proteins IRS-1, Src, and Shc [14, 15]. In the simple case of branching pathways as shown in Figure 2.4C, one enzyme  $E_1$  activates two different targets A and B. Since there is no direct interaction between A and B we would not expect one to influence the activation of the other. We can see that this is indeed the case for activation ratios, since:

$$AR_A = \frac{k_1 E_1 / K_{m1}}{k_{IA} E_{IA} / K_{mIA}} = \alpha_1^A E_1 \quad (2.31)$$

$$AR_B = \frac{k_2 E_1 / K_{m2}}{k_{IB} E_{IB} / K_{mIB}} = \alpha_1^B E_1 \quad (2.32)$$

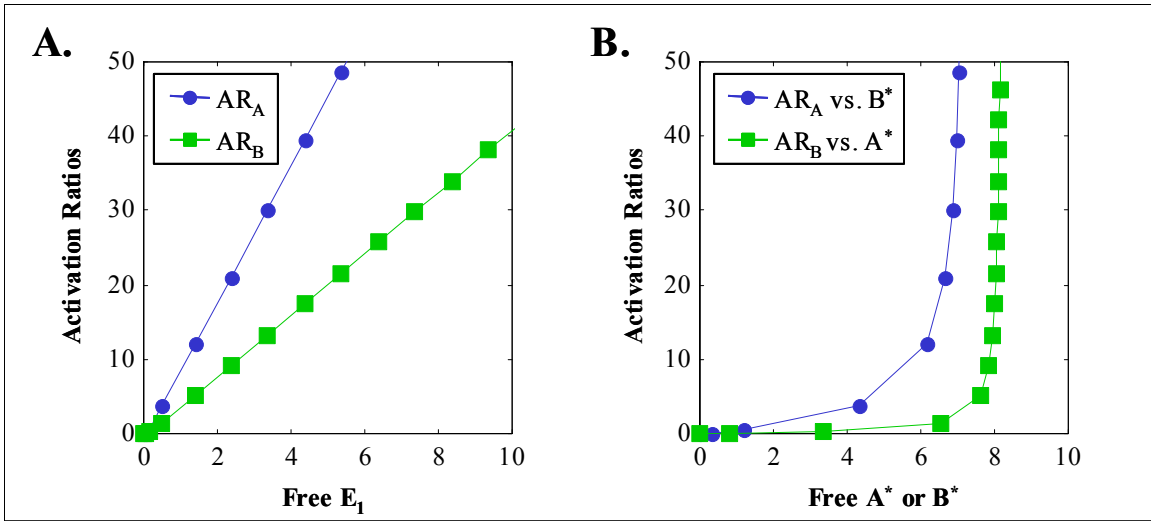
Once again, the expression for the activation ratio of each intermediate is the same as if they were isolated, and the activation factors  $\alpha_1^A$  and  $\alpha_1^B$  are determined only from parameters arising from the interaction between  $E_1$  and A and B, respectively. There is only an indirect interaction between A and B arising from sharing the activating enzyme  $E_1$ . If A and B were actually different pools of the same enzyme, we would rightly expect that Equations 2.31-2.32 have the same form, and the same value for  $\alpha_1^A$  and  $\alpha_1^B$ . Once again by neglecting enzyme-substrate complexes, it can also be shown that:

$$AR_A = \frac{\alpha_1^A}{\alpha_1^B} \frac{B^*}{B_T - B^*} \quad (2.33)$$

$$AR_B = \frac{\alpha_1^B}{\alpha_1^A} \frac{A^*}{A_T - A^*} \quad (2.34)$$

Therefore for branching pathways, we expect that plots of activation ratios for each branch would be linear with respect to their common activator, just as if there were no other branch present. However, plots of activation ratios for the two branch intermediates against each other will be inversely hyperbolic, as if they were the same molecule. Recall from Figure 2.5 that activation ratios for a component plotted against itself also appeared inverse hyperbolic. That this same pattern is observed for different branches should not be surprising, since again we should be able to conceptually divide a

single pool into two identical pools. The diverging case can therefore be resolved from a linear cascade, where one plot would be inversely hyperbolic, and the other would be linear or hyperbolic. These results can be seen readily in Figure 2.7, where  $AR_A$  and  $AR_B$  are plotted against  $E_1$  (Figure 2.7A) or  $B^*$  or  $A^*$ , respectively (Figure 2.7B) for a simulation of diverging pathway.



**Figure 2.7.** Results for diverging branches in Figure 2.4C. A) Activation ratios for A (diamonds) and B (squares) plotted against free activating enzyme  $E_1$ , B) Activation ratios for A and B plotted against each other, i.e.  $AR_A$  against  $B^*$  and  $AR_B$  against  $A^*$ .

### 2.3.4. Dual activation steps

The analysis thus far has considered the case of intermediates interconverting by single-step mechanisms. Many important signaling components, however, require multiple modification steps; for example, activation of MAPKs by MAPKKs require two distinct, nonprocessive phosphorylation events [16, 17]. In such cases complete activation can be considered as overlapping cycles, and activation ratio analysis can still be performed on each cycle. (Multiple processive activation and deactivation steps would be indistinguishable from single-step activation and therefore result in behavior similar to the single-step cycle analyzed in Section 2.2.)

If two steps are required for complete activation of A as shown in Figure 2.4D, each under Michaelis-Menten kinetics, then several activation ratios can be defined:

$$AR_{A1} \equiv \frac{A^*}{A} = \frac{k_1 E_1 / K_{m1}}{k_2 E_2 / K_{m2}} \equiv \alpha_1^A E_1 \quad (2.35)$$

$$AR_{A2} \equiv \frac{A^{**}}{A^*} = \frac{k_3 E_1 / K_{m3}}{k_4 E_4 / K_{m4}} \equiv \alpha_2^A E_1 \quad (2.36)$$

$$AR_{A3} \equiv \frac{A^{**}}{A} = \frac{k_1 k_3 E_1^2}{k_2 k_4 E_2 E_4} = \alpha_1^A \alpha_2^A E_1^2 \quad (2.37)$$

According to Equation 2.37, the “overall” activation ratio  $AR_{A3}$  will be quadratic rather than linear with respect to  $E_1$ . Nevertheless, each of the individual ratios  $AR_{A1}$  and  $AR_{A2}$  will be linear. Each individual factor  $\alpha_1^A$  and  $\alpha_2^A$  reflects the parameters for the individual steps. It can therefore be tested if the second step is faster, slower, or the same as the first, under *in vivo* conditions where the substrate is being deactivated as well as activated.

It is important to note that in cases of multiple activation/deactivation steps, it is necessary to have measurements of the intermediate form (e.g.,  $A^*$ ) as well as the fully active and inactive states ( $A$ ,  $A^{**}$ ). Lacking this information may result in incorrect conclusions regarding the nature of the relationship between target and activator, since a quadratic response may appear similar to an inverse hyperbola, suggesting that target and activator are reversed or parallel targets of some upstream component. Additional information arising from other components in the network would assist in resolving this issue. For example, in a cascade where  $A^{**}$  activates a target B, then plots of  $AR_B$  against  $E_1$  would still be hyperbolic, but linear against  $A^{**}$ , suggesting that  $E_1$  is upstream of A after all. This type of situation would alert the analyst that further investigation is required, and indicates the presence of multiple activation steps if not previously known.

### 2.3.5. Cascades with feedback

The presence of feedback is an important component of signaling pathways, since without feedback there is no possibility for attenuation or adaptation. An increase in the level of an activating stimulus would eventually saturate the signaling machinery, reducing the overall sensitivity of the system to variations in external conditions. Although attenuation is possible simply through downregulation of signaling intermediates, a faster mechanism involves modification of an upstream effector that ultimately leads to deactivation [2, 18, 19]. Positive feedback allows for further sharpening of a signal towards a “step-like” response and has also been observed

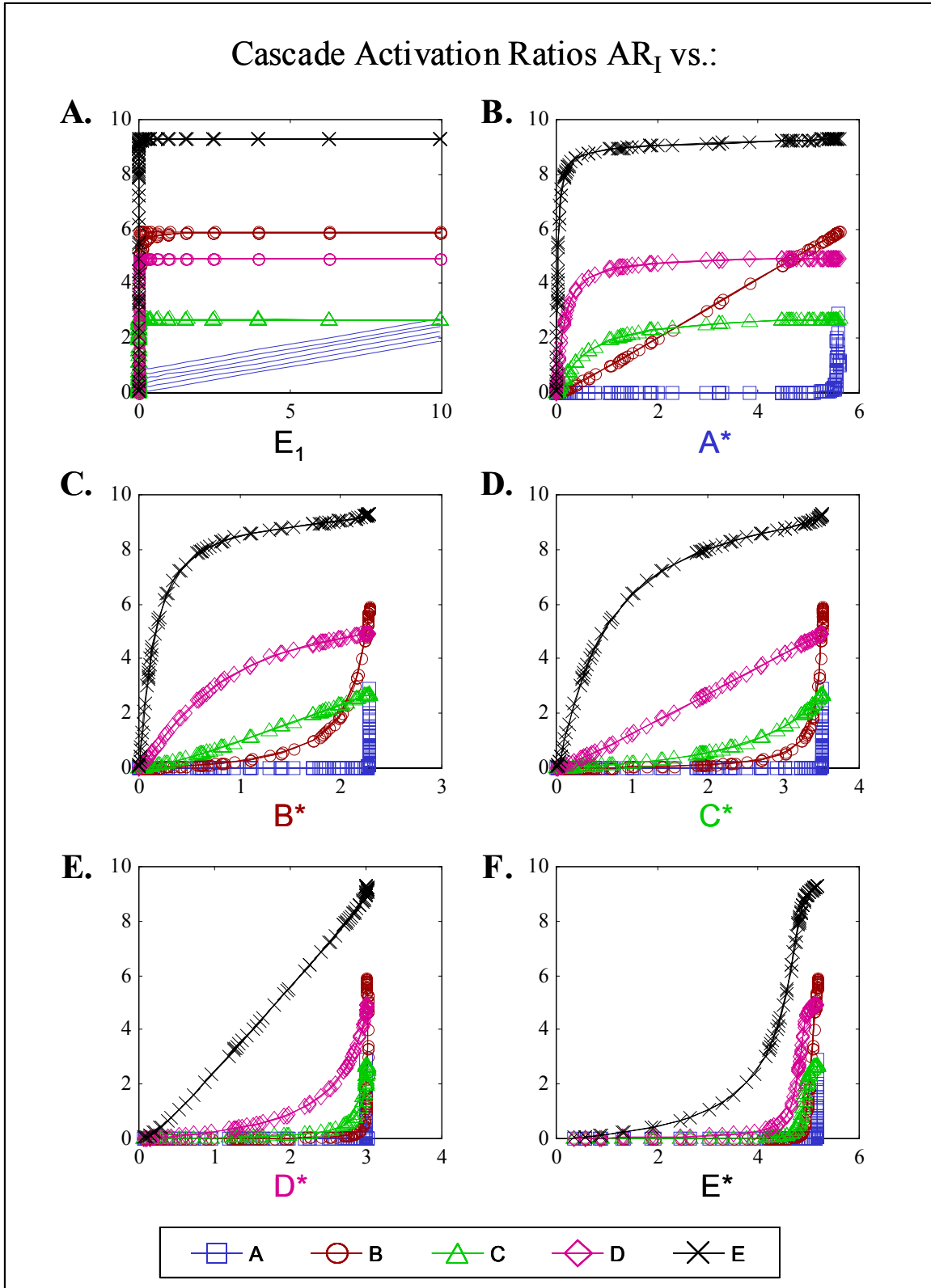
experimentally [20]. It was therefore critical to investigate whether activation ratio analysis would still be valid in feedback systems.

To this end, two cascade models were developed, with the structure shown in Figure 2.4E. The linear cascade of Figure 2.4A was modified to include two additional intermediates (facilitating analysis by extending the number of steps within the feedback loop) and add feedback such that  $E^*$  catalyzes activation of A (positive feedback) or inactivation of  $A^*$  (negative feedback). At steady state, the balance between rates of activation and inactivation of A can be described in one of the following ways:

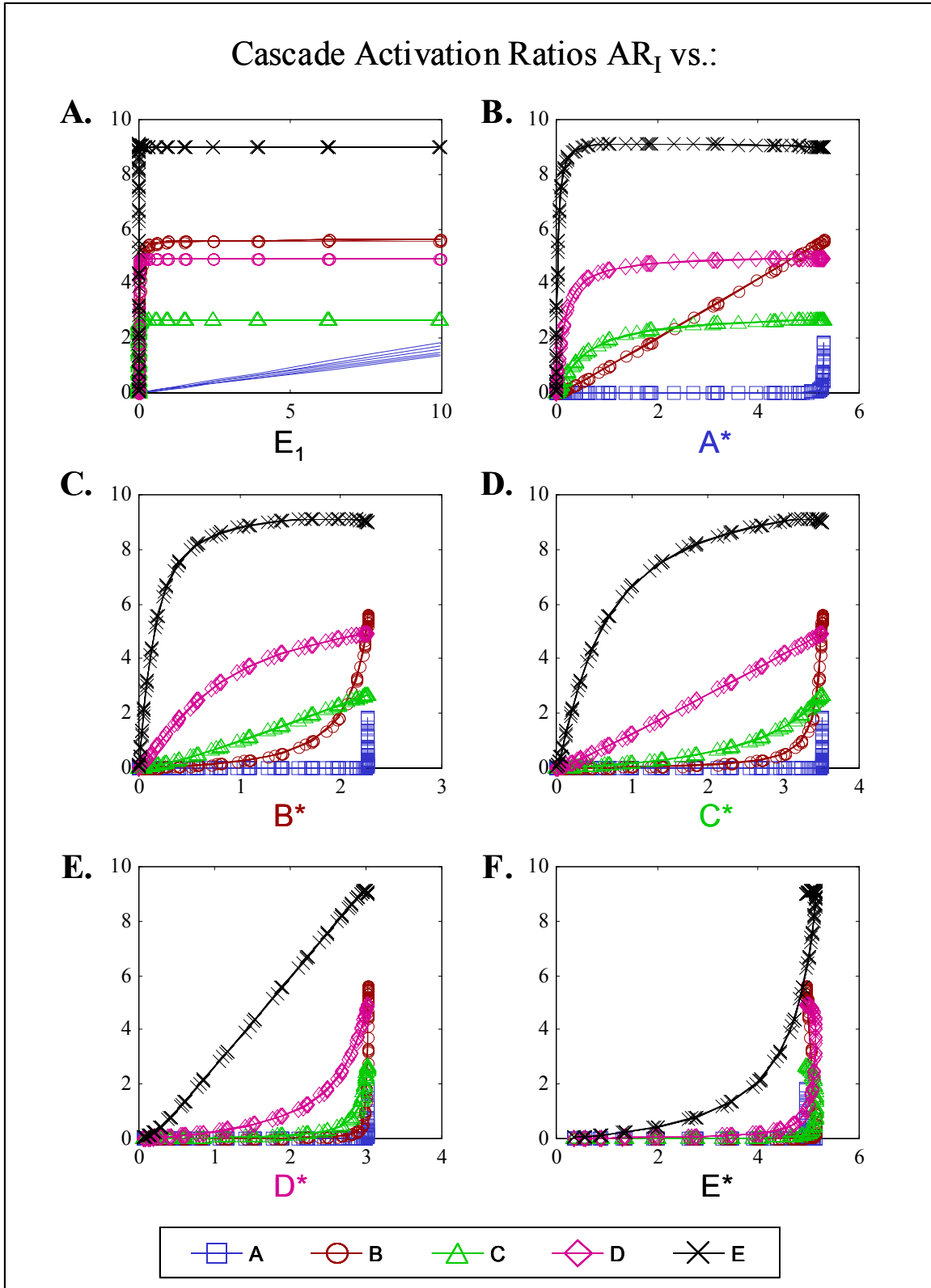
$$\text{Positive feedback: } \frac{k_1 E_1 A}{K_{m1}} + \frac{k_{FB} E^* A}{K_{FB}} = \frac{k_2 E_{AI} A^*}{K_{mAI}} \quad (2.38)$$

$$\text{Negative feedback: } \frac{k_1 E_1 A}{K_{m1}} = \frac{k_2 E_{AI} A^*}{K_{mAI}} + \frac{k_{FB} E^* A^*}{K_{FB}} \quad (2.39)$$

Results for positive and negative feedback models are shown in Figures 2.8 and 2.9, respectively. In these models, the parameter values and total concentrations of each component were identical to results previously shown for an open cascade; each step is saturated such that the total concentrations are greater than the  $K_m$  values. Here the feedback step is unsaturated ( $K_{FB} = 100$ ,  $E_T = 10$ ) and the strength of the feedback response is varied by adjusting  $k_{FB}$  over three orders of magnitude. The activation ratios for each component are plotted against each of the other factors, where the results at each value of  $k_{FB}$  are overlaid. It can be readily seen that feedback has no net effect on the activation ratios of most of the components, and all feedback cases collapse together. Furthermore, the activation ratios continue to show a linear response for one intermediate when plotted against its directly upstream activator (e.g.  $AR_A$  vs.  $E_1$ ,  $AR_B$  vs.  $A^*$ , etc.), a hyperbolic response for indirectly upstream activators ( $AR_B$ ,  $AR_C$ ,  $AR_D$ ,  $AR_E$  vs.  $E_1$ , etc.), and an inverse hyperbolic response when plotted against downstream targets ( $AR_B$  vs.  $C^*$ ,  $D^*$ ,  $E^*$ . etc.). *These results are identical with those observed in the absence of feedback.*



**Figure 2.8.** Activation ratios for cascade with positive feedback, unsaturated in feedback step ( $K_{FB} \gg A_T$ ), plotted against free  $E_1$  (A),  $A^*$  (B),  $B^*$  (C),  $C^*$  (D),  $D^*$  (E), or  $E^*$  (F). For clarity, markers were omitted for activation ratios of A in pt. A.



**Figure 2.9.** Activation ratios for cascade with negative feedback, unsaturated in feedback step ( $K_{FB} \gg A_T$ ), plotted against free  $E_1$  (A),  $A^*$  (B),  $B^*$  (C),  $C^*$  (D),  $D^*$  (E), or  $E^*$  (F). For clarity, markers were omitted for activation ratios of A in pt. A.



The only readily apparent change in the profiles of activation ratios involves the point of feedback, A. In the presence of positive feedback, the results still appear linear, with identical slope but the intercept increasing as  $k_{FB}$  is increased. In the presence of negative feedback, the intercept remains at zero but the slope decreases with increasing  $k_{FB}$ . This can be understood by rearranging Equations 2.38-2.39:

$$\text{Positive feedback: } AR_A = \frac{k_1 E_1 / K_{mI}}{k_{AI} E_{AI} / K_{mAI}} + \frac{k_{FB} E^* / K_{FB}}{k_{AI} E_{AI} / K_{mAI}} = \alpha_1^A E_1 + \alpha_E^A E^* \quad (2.40)$$

$$\text{Negative feedback: } AR_A = \frac{k_1 E_1 / K_{mI}}{k_{AI} E_{AI} / K_{mAI} + k_{FB} E^* / K_{FB}} = (\alpha_1^A)' E_1 \quad (2.41)$$

In order to interpret these results, it is important to realize that, as the endpoint of the cascade, E will become activated quickly, and thus  $E^*$  will remain relatively constant as  $E_1$  is varied (certainly, relative to A and  $A^*$ ). This is a consequence of multistep ultrasensitivity, where the activation of a downstream component saturates before upstream components, even without feedback [10-12]. Therefore, in the positive feedback system, we expect that the activation ratio for A will still be linear with respect to  $E_1$ , with a nonzero intercept that is determined by the feedback parameters  $k_{FB}$  and  $K_{FB}$ . Increasing  $k_{FB}$  therefore increases the intercept by increasing  $\alpha_E^A$ , as seen in Figure 2.8A. With negative feedback, increasing  $k_{FB}$  would increase the constant factor in the denominator of  $\alpha_1^A$ , causing a decreased slope as observed in Figure 2.9A. In either case, analysis of the remainder of the cascade is unaffected by the presence of feedback. This is to be expected as activation ratios are determined by performing local balances around each intermediate, and therefore look at each step as if isolated. In the feedback system, only the feedback point (here, A) is changed.

## **2.4. Saturating conditions and influence of enzyme-substrate complexes**

The method of activation ratios can be quite powerful for the reconstruction of signaling networks, since it enables an isolated examination of the relationship between components regardless of other connections to either species. There is a key consideration, however, which is that *the analysis depends on the measurement of free (unbound) intermediates*. In other words, it is important to resolve between signaling

intermediates in their free active and inactive forms ( $A^*$  and  $A$ ) from when they are bound to other intermediates, including enzyme-substrate complexes (e.g.  $[E_1 \cdot A]$ ,  $[E_2 \cdot A^*]$ ,  $[A^* \cdot B]$ , etc). In the simple case of intermolecular binding shown in Figure 2.1, there must obviously be a separation possible between free from complexed  $A$  to detect “activation” as the concentration of  $[A \cdot B]$ . Moreover, it was shown previously that measurement of free  $B$  would greatly simplify the analysis (compare Equations 2.1-2.4 and 2.5-2.7), although the pattern would be similar (Figure 2.2).

For covalent-modification cycles either in isolation or connected in networks, this issue is even more significant. As an example, consider the isolated cycle of Figure 2.1. In this system there exist six species, i.e.  $A^*$ ,  $A$ , (free active and inactive);  $E_1$ ,  $E_2$ , (free enzyme 1 and 2); and  $[E_1 \cdot A]$ ,  $[E_2 \cdot A^*]$  (enzyme-substrate complexes). At steady state Equations 2.17-2.19 hold, and the activation ratio (as defined in Equation 2.20) is linearly dependent upon  $E_1$ .

If we are unable to resolve the free species from complexes, then we may only be able to measure total active and inactive  $A$  ( $A_{Tot}^*$  and  $A_{Tot}$ , respectively), and total concentrations of enzymes  $E_{1T}$  and  $E_{2T}$ . These total concentrations are related to the components of the system in the following way:

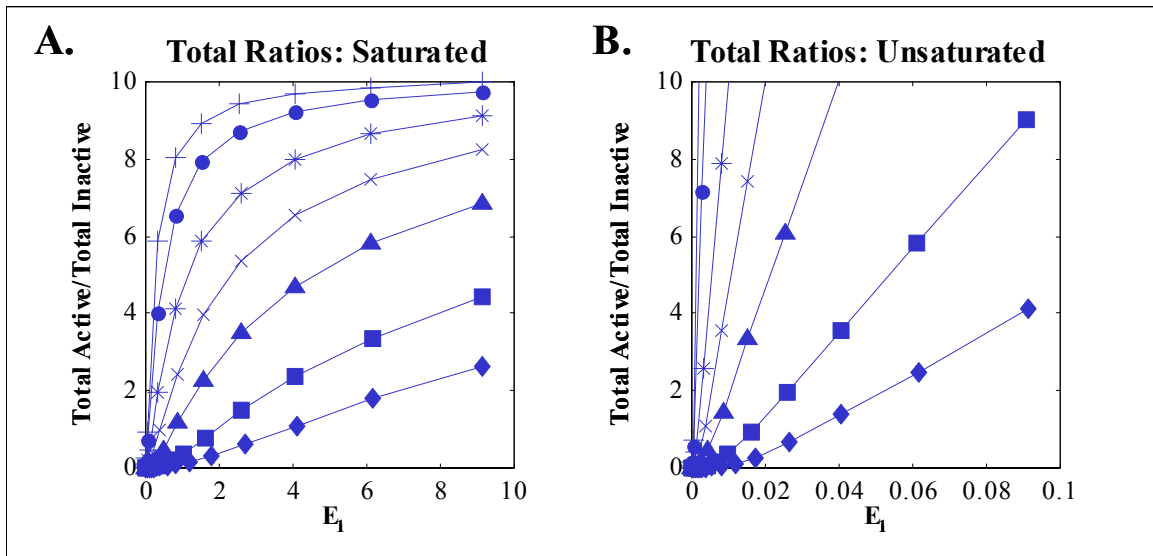
$$A_{Tot} = A + [E_1 \cdot A] \quad (2.42)$$

$$A_{Tot}^* = A^* + [E_2 \cdot A^*] \quad (2.43)$$

Since we probably do not know  $K_{m1}$  and  $K_{m2}$ , we cannot use Equations 2.17-2.18 to help calculate the free species from measurements of the totals  $A_{Tot}$ ,  $A_{Tot}^*$ ,  $E_{1T}$ ,  $E_{2T}$ , and Equations 2.14-2.15, 2.42 and 2.43. Note that  $A_T = A_{Tot} + A_{Tot}^*$ , so only two of the three values need to be measured. We will have six unknowns and four equations, and therefore are unable to calculate or estimate the true activation ratio  $AR_A$  based on total activation measurements. If instead of free species, the total active  $A_{Tot}^*$  and total inactive  $A_{Tot}$  concentrations are used to construct ratios, the results are less useful for network reconstruction. This is because:

$$TAR_A \equiv \frac{A_T^*}{A_T} = \frac{A^* (1 + E_2 / K_{m2})}{A (1 + E_1 / K_{m1})} = \frac{k_1 E_1}{k_2 E_2} \frac{K_{m2} + E_2}{K_{m1} + E_1} \quad (2.44)$$

In Equation 2.44 the “total activation ratio”  $TAR_A$  is no longer linear, but rather hyperbolic, with respect to free  $E_1$ . Thus a direct effect, when examined using total active and inactive species, takes the same form as an indirect effect using only free species. This can be seen in Figure 2.10, where  $TAR_A$  is plotted against free  $E_1$ . In Figure 2.10A, parameter conditions were identical to those used to generate Figure 2.3, where enzyme concentrations were similar in magnitude to the  $K_m$  values. It is readily apparent that the lines of Figure 2.3B have been replaced by hyperbolas in Figure 2.10A. Note that this effect arises only when the enzyme concentrations are significant. If enzyme concentrations are negligible, that is, where  $E_{1T}/K_{m1}$  and  $E_{2T}/K_{m2}$  are much less than one, the total activation ratios should still be linear. This case, shown in Figure 2.10B, was generated by adjusting the parameters to decrease enzyme concentrations 100-fold and increasing kinetic constants accordingly to maintain the same  $K_m$  values. Indeed, the ratio between total active and total inactive A remains linear. However, it should be remembered that this highly idealized case is unlikely to occur in real signaling systems, where enzymes and substrates are proteins often of similar concentrations, and the system is saturated with respect to both.



**Figure 2.10.** Total activation ratios for the isolated covalent modification cycle (using total active  $A_T^*$  and total inactive  $A_T$ ) plotted against free  $E_1$ , with enzyme conditions A) Saturated: parameters as shown in Figure 2.3. B) Unsaturated: with  $E_{iT}/K_{mi} < 0.1$ .

For an individual cycle under saturating conditions, the total activation ratio for the target (substrate) will appear to be hyperbolic with respect to the direct activator.

This can cause confounding with a cascade arrangement where the activator acts indirectly through one or more other species. Furthermore, the activation factor, measured as the slope for the linear plot of the free activation ratio, can no longer be determined nor used as a quantitative description of the kinetic parameters.

Results for the linear cascade are even more complicated, since  $A^*_T$  will now be the sum of free  $A^*$  and that complexed both to the inactivase  $E_{IA}$  and to the target B, and similarly  $B^*_T$  includes complexes for  $[B^* \cdot C]$  and  $[E_{IB} \cdot B^*]$ . Equations 2.23-2.25 are therefore replaced by the following expressions using total forms:

$$TAR_A \equiv \frac{A^*_T}{A_T} = \frac{k_1 E_1}{k_{IA} E_{IA}} \frac{K_{mIA} + E_{IA} + B \frac{K_{mIA}}{K_{m2}}}{K_{m1} + E_1} \quad (2.45)$$

$$TAR_B \equiv \frac{B^*_T}{B_T} = \frac{k_2 A^*}{k_{IB} E_{IB}} \frac{K_{mIB} + E_{IB} + C \frac{K_{mIB}}{K_{m3}}}{K_{m2} + A^*} \quad (2.46)$$

$$TAR_C \equiv \frac{C^*_T}{C_T} = \frac{k_3 B^*}{k_{IC} E_{IC}} \frac{K_{mIB} + E_{IC}}{K_{m2} + B^*} \quad (2.47)$$

The total activation ratios in Equations 2.45-2.47 are hyperbolic with respect to the direct upstream activator, and the concentrations of targets (B, C) also appear, which may further complicate the functional forms of these expressions, since those will also vary as  $E_1$  is increased. The presence of other species in the expressions for total activation ratios is a consequence of the fact that multiple components can form complexes with each intermediate. For example, since B can bind  $A^*$ , increasing the amount of free B will increase the amount of complex  $[A^* \cdot B]$  and therefore total  $A^*_T$ , but the balance of activation around A, determined by free A and  $A^*$ , will remain the same. The theoretical isolation that arises using free species to find activation ratios no longer applies when total activation ratios are calculated.

Simple physical arguments can be used to show that the total activation ratio will still be hyperbolic for indirect effects. If it is assumed that  $A^*$  is approximately a hyperbolic function of  $E_1$ , then substitution into Equation 2.46 would yield an expression for  $TAR_B$  that is also a hyperbolic function of  $E_1$ . A similar argument can be made for  $TAR_C$  as a function of  $E_1$  or  $A^*$ . This can be verified during simulations of cascades

under saturating conditions, where  $K_m$  values are similar in magnitude to the concentrations of all species. Since total activation ratios will appear hyperbolic for both direct and indirect interactions, they can not resolve between these two cases, and thus are more limited in their ability to perform reconstruction of cascade structures, as will be described in Chapter 3.

For converging pathways, the use of total activation ratios has two drawbacks. Since two enzymes can bind A independently, the total amount of A ( $A_T$ ) will be dependent on both enzyme concentrations. Equation 2.30 is therefore replaced by:

$$\text{TAR}_A = \frac{k_1 E_1}{k_{IA} E_{IA}} \frac{K_{mIA} + E_1}{K_{m1} + E_1 + \frac{K_{m1}}{K_{m2}} E_2} + \frac{k_2 E_2}{k_{IA} E_{IA}} \frac{K_{mIA} + E_1}{K_{m2} + \frac{K_{m2}}{K_{m1}} E_1 + E_2} \quad (2.48)$$

The total activation ratio is still a sum of the effects of each input, but each term is hyperbolic with respect to both inputs. This changes the mathematical form of the expression, and the effects of the two inputs are no longer separated into different terms. Plots of total activation ratios under saturating conditions will appear as sets of hyperbolae, but the shape and spacing of plots will change as  $E_1$  and  $E_2$  are varied. This is in stark contrast to free activation ratios (Equation 2.30), where each term contains the effect of only one activator, and thus the activation factors could be determined independently.

It should also be noted that in the analysis thus far the free, not total, amounts of the activating enzyme were used, i.e. activation ratios were plotted against  $E_1$  and not  $E_{1T}$ . The importance of measurement of free species (A and  $A^*$ ) was shown above; the importance of using  $E_1$  instead of  $E_{1T}$  is not so critical. Equation 2.20 shows that the activation ratio for an isolated cycle is linearly dependent upon  $E_1$ , the free (not total) concentration of activating enzyme, by an activation factor  $\alpha_1^A$ . If the total concentrations of enzymes are used then Equation 2.20 becomes:

$$\text{AR}_A = \frac{k_1 E_{1T}}{k_2 E_{2T}} \frac{K_{m2} + A^*}{K_{m1} + A} \equiv \beta_1^A E_{1T} \quad (2.49)$$

In Equation 2.49 as in Equation 2.20, the activation ratio  $\text{AR}_A$  is still linearly proportional to  $E_{1T}$ , but by a new activation factor  $\beta_1^A$ . The difference in these

expressions is the inclusion of  $A^*$  and  $A$ , therefore  $\beta_1^A$  is more explicitly a function of  $A^*$  and  $A$ , and varies more than  $\alpha_1^A$  as  $E_1$  is increased. In fact  $\alpha_1^A$  is also a function of  $A^*$ , since the amount of free  $E_2$  is determined by  $A^*$  according to Equations 2.15 and 2.18, so that:

$$E_2 = \frac{E_{2T}}{(1 + A^*/K_{m2})} \quad (2.50)$$

$$\frac{A^*}{A} = \frac{k_1 E_1}{k_2 E_{2T}} \frac{K_{m2} + A^*}{K_{m1}} \equiv \alpha_1^A E_1 \quad (2.51)$$

It is true that  $\beta_1^A$  is not constant, but approaches a constant value as the fraction of  $A$  in the active form ( $A^*/A_{Tot}$ ) increases. Nevertheless, it represents the local sensitivity of the activation ratio with respect to  $E_{1T}$  and can be a useful measure of the interaction between  $E_1$  and  $A$ . The presence of  $A^*$  and  $A$  in Equations 2.49 and 2.51 result in a slight curvature for plots of activation ratios, in particular at lower values of  $E_{1T}$  where  $A^*$  and  $A$  are varying significantly. As  $E_{1T}$  increases,  $A^*$  and  $A$  approach limiting values, and activation ratio plots appear linear. The severity of this curvature depends upon the degree of saturation of the enzymes; if  $K_{m1}$  and  $K_{m2}$  are much greater than  $A_{Tot}$  the plots will be linear for the full range of  $E_{1T}$ .

A similar result occurs when determining total activation ratios as a function of total enzyme concentrations:

$$TAR_A = \frac{k_1 E_{1T}}{k_2 E_{2T}} \frac{K_{m2} + E_{2T} + A^*}{K_{m1} + E_{1T} + A} \quad (2.52)$$

As before, the total activation ratio is generally hyperbolic with respect to the concentration of the activator. The presence of  $A^*$  and  $A$  in Equation 2.52 once again results in a curvature at lower values of  $E_{1T}$ , yielding an slightly sigmoidal shape overall.

Since it is important to determine concentrations of free active forms of intermediates, to avoid calculation of total activation ratios, then the issue of free or total enzyme is most likely not a major concern. If the free active and free inactive concentrations of each intermediate in the network are measured, then by default the free concentrations of enzymes will be measured. The free activation ratios that will be calculated should yield nearly linear results for direct effects.

## 2.5. Conclusions

A novel method for examination of signaling pathways was developed, based simply on an understanding of the nature of the interconverting cycles that are prominent in these pathways. It was shown that by considering the ratio of active to inactive forms of the intermediate, it is possible to write expressions that focus on isolated interactions between activator and target. Unknown kinetic constants and enzyme concentrations are compressed together into a single quantitative parameter, the activation factor. This approach yields a linear relationship between the activation ratio and the concentration of the enzyme that drives the activating reaction. This simple relationship remains even when the activator and target are embedded in more complicated systems, where converging and diverging pathways, linear cascades, multiple activation steps, or feedback appear. The properties of activation ratios for different systems are summarized in Table 2.2.

**Table 2.2.** Summary of results for simple signaling systems: expressions for activation ratios.

Arrangement	Expression	Mathematical Form of AR <sub>i</sub>
A. Single cycle 1. Binding	$AR_A \equiv \frac{[A \cdot B]}{A} = \frac{a}{d} B = K_a B$	Linear (vs. B or E <sub>1</sub> )
2. Covalent modification	$AR_A \equiv \frac{A^*}{A} = \frac{k_1 E_1 / K_{m1}}{k_2 E_2 / K_{m2}} \equiv \alpha_1^A E_1$	
B. Linear cascade  (Indirect and Downstream: neglect enzyme-substrate complexes)	$AR_A = \frac{A^*}{A} = \alpha_1^A E_1 = \frac{B^*}{\alpha_A^B A_T B_T - (1 + \alpha_A^B A_T) B^*}$ $= \frac{C^*}{\alpha_A^B \alpha_B^C A_T B_T C_T - [1 + \alpha_A^B A_T (1 + \alpha_B^C B_T)] C^*}$ $AR_B = \alpha_A^B A^* = \frac{\alpha_1^A \alpha_A^B A_T E_1}{1 + \alpha_1^A E_1}$ $= \frac{C^*}{\alpha_B^C B_T C_T - (1 + \alpha_B^C B_T) C^*}$ $AR_C = \alpha_B^C B^* = \frac{\alpha_A^B \alpha_B^C B_T A^*}{1 + \alpha_A^B A^*} = \frac{\alpha_1^A \alpha_A^B \alpha_B^C B_T A_T E_1}{1 + \alpha_1^A E_1 (1 + \alpha_A^B A_T)}$	Plotted against (e.g. AR <sub>B</sub> vs.):  Direct activator (A <sup>*</sup> ): Linear  Indirect activator (E <sub>1</sub> ): Hyperbolic  Downstream target (C <sup>*</sup> ): Inverse hyperbolic

Arrangement	Expression	Mathematical Form of AR <sub>i</sub>
C. Converging pathways	$\frac{A^*}{A} = \frac{k_1 E_1 / K_{m1}}{k_{IA} E_{IA} / K_{mIA}} + \frac{k_2 E_2 / K_{m2}}{k_{IA} E_{IA} / K_{mIA}} = \alpha_1^A E_1 + \alpha_2^A E_2$	Linear combination of activators
D. Diverging pathways (Parallel: neglect complexes)	$\frac{A^*}{A} = \alpha_1^A E_1 = \frac{\alpha_1^A}{\alpha_1^B} \frac{B^*}{B_T - B^*}$ $\frac{B^*}{B} = \alpha_1^B E_1 = \frac{\alpha_1^B}{\alpha_1^A} \frac{A^*}{A_T - A^*}$	vs. activator: Linear vs. parallel: Inverse hyperbolic
E. Dual activation steps	$AR_{A1} = \frac{A^*}{A} = \frac{k_1 E_1 / K_{m1}}{k_2 E_2 / K_{m2}} \equiv \alpha_1^A E_1$ $AR_{A2} = \frac{A^{**}}{A^*} = \frac{k_3 E_1 / K_{m3}}{k_4 E_4 / K_{m4}} \equiv \alpha_2^A E_1$ $AR_{A3} \equiv \frac{A^{**}}{A} = \frac{k_1 k_3 E_1^2}{k_2 k_4 E_2 E_4} = \alpha_1^A \alpha_2^A E_1^2$	Linear for each step (AR <sub>A1</sub> , AR <sub>A2</sub> ) Quadratic overall (AR <sub>A3</sub> )
F. Feedback 1. Positive 2. Negative	$AR_A = \frac{k_1 E_1 / K_{m1}}{k_{AI} E_{AI} / K_{mAI}} + \frac{k_{FB} E^* / K_{FB}}{k_{AI} E_{AI} / K_{mAI}} = \alpha_1^A E_1 + \alpha_E^A E^*$ $AR_A = \frac{k_1 E_1 / K_{m1}}{k_{AI} E_{AI} / K_{mAI} + k_{FB} E^* / K_{FB}} = (\alpha_1^A)' E_1$	Effect of FB 1. Nonzero intercept 2. Altered slope
G. Single Cycle, Total Activation ratio	$TAR_A \equiv \frac{A_T^*}{A_T} = \frac{k_1 E_1}{k_{IA} E_{IA}} \frac{K_{mIA} + E_{IA}}{K_{m1} + E_1}$	Hyperbolic
H. Linear cascade, Total Activation	$TAR_A = \frac{k_1 E_1}{k_{IA} E_{IA}} \frac{K_{mIA} + E_{IA} + B \frac{K_{mIA}}{K_{m2}}}{K_{m1} + E_1}$ $TAR_B = \frac{k_2 A^*}{k_{IB} E_{IB}} \frac{K_{mIB} + E_{IB} + C \frac{K_{mIB}}{K_{m3}}}{K_{m2} + A^*}$ $TAR_C = \frac{k_3 B^*}{k_{IC} E_{IC}} \frac{K_{mIB} + E_{IC}}{K_{m2} + B^*}$	Hyperbolic vs. all upstream activators (direct and indirect) Inverse hyperbolic vs. downstream
I. Converging, Total Activation	$TAR_A = \frac{k_1 E_1}{k_{IA} E_{IA}} \frac{K_{mIA} + E_1}{K_{m1} + E_1 + \frac{K_{m1}}{K_{m2}} E_2} + \frac{k_2 E_2}{k_{IA} E_{IA}} \frac{K_{mIA} + E_1}{K_{m2} + \frac{K_{m2}}{K_{m1}} E_1 + E_2}$	Hyperbolic vs. both activators



Activation ratios have great potential in network reconstruction, because the functional form (or graphical pattern, if plotted) changes along with the nature of interaction between two species. Since the network structure is reflected in the form for the activation ratios, then the problem can be inverted; namely, the functionality of activation ratios can be used to determine the structure. This will be discussed in further detail in the next chapter. The ratios can also be used as a descriptive tool during examination of a single step, since the activation factors are reflective of the kinetic parameters, and multi-step, nonprocessive activation kinetics can be resolved from single-step or processive reactions.

The analysis described here was developed assuming steady state in the signaling system. This simplification greatly facilitates the overall analysis development, since concentrations of enzyme-substrate complexes can be written in terms of free species (e.g.  $E_1$  and  $A$ ), as in Equations 2.17-2.18, and individual kinetic parameters ( $k_{\text{cat}}$  and  $K_m$ ) can be collected together as shown in Equation 2.20. In the absence of a steady-state system, such as cells growing in a chemostat, this assumption will not be precisely valid. The transient nature of signaling pathways has been well documented; rather than achieving a steady state often components will peak in activation and relax more slowly [21]. Nevertheless, the method may still be applicable if a *pseudo-steady state* assumption can be made. A pseudo-steady state approximation is often made in examination of enzyme systems, and is based on the assumption that the rate of formation of enzyme-substrate complexes is faster than the rate of decomposition. In this case, the assumption would require that dynamics of component activation (within a step) are faster than transmission (between steps). In such a case it may be possible to apply this pseudo-steady state assumption to the signaling system at each time point over the more global dynamics of the system. Further work in this area is warranted, to investigate the applicability of activation ratios to dynamic signaling systems.

## 2.6. References

1. Bhalla, U.S. and R. Iyengar, *Emergent properties of networks of biological signaling pathways*. Science, 1999. **283**(5400): p. 381-7.
2. Brightman, F.A. and D.A. Fell, *Differential feedback regulation of the MAPK cascade underlies the quantitative differences in EGF and NGF signalling in PC12 cells*. FEBS Lett, 2000. **482**(3): p. 169-74.

3. Huang, C.Y. and J.E. Ferrell, Jr., *Ultrasensitivity in the mitogen-activated protein kinase cascade*. Proc Natl Acad Sci U S A, 1996. **93**(19): p. 10078-83.
4. Kholodenko, B.N., et al., *Quantification of short term signaling by the epidermal growth factor receptor*. J Biol Chem, 1999. **274**(42): p. 30169-81.
5. Kholodenko, B.N., et al., *Quantification of information transfer via cellular signal transduction pathways*. FEBS Lett, 1997. **414**(2): p. 430-4.
6. Levchenko, A., J. Bruck, and P.W. Sternberg, *Scaffold proteins may biphasically affect the levels of mitogen-activated protein kinase signaling and reduce its threshold properties*. Proc Natl Acad Sci U S A, 2000. **97**(11): p. 5818-23.
7. Berridge, M.J., *The molecular basis of communication within the cell*. Sci Am, 1985. **253**(4): p. 142-52.
8. Goldbeter, A. and D.E. Koshland, Jr., *An amplified sensitivity arising from covalent modification in biological systems*. Proc Natl Acad Sci U S A, 1981. **78**(11): p. 6840-4.
9. Ferrell, J.E., Jr., *Tripping the switch fantastic: how a protein kinase cascade can convert graded inputs into switch-like outputs [see comments]*. Trends Biochem Sci, 1996. **21**(12): p. 460-6.
10. Goldbeter, A. and D.E. Koshland, Jr., *Ultrasensitivity in biochemical systems controlled by covalent modification. Interplay between zero-order and multistep effects*. J Biol Chem, 1984. **259**(23): p. 14441-7.
11. Chock, P.B. and E.R. Stadtman, *Superiority of interconvertible enzyme cascades in metabolite regulation: analysis of multicyclic systems*. Proc Natl Acad Sci U S A, 1977. **74**(7): p. 2766-70.
12. Ferrell, J.E., Jr., *How responses get more switch-like as you move down a protein kinase cascade [letter; comment]*. Trends Biochem Sci, 1997. **22**(8): p. 288-9.
13. Posas, F. and H. Saito, *Osmotic activation of the HOG MAPK pathway via Ste11p MAPKKK: scaffold role of Pbs2p MAPKK*. Science, 1997. **276**(5319): p. 1702-5.
14. Miyajima, A., et al., *Cytokine receptors and signal transduction*. Annu Rev Immunol, 1992. **10**: p. 295-331.
15. Roth, R.A., et al., *Substrates and signalling complexes: the tortured path to insulin action*. J Cell Biochem, 1992. **48**(1): p. 12-8.
16. Burack, W.R. and T.W. Sturgill, *The activating dual phosphorylation of MAPK by MEK is nonprocessive*. Biochemistry, 1997. **36**(20): p. 5929-33.
17. Ferrell, J.E., Jr. and R.R. Bhatt, *Mechanistic studies of the dual phosphorylation of mitogen-activated protein kinase*. J Biol Chem, 1997. **272**(30): p. 19008-16.
18. Goldbeter, A. and J.M. Guilmot, *Thresholds and oscillations in enzymatic cascades*. Journal of Physical Chemistry, 1996. **100**(49): p. 19174-19181.
19. Kholodenko, B.N., *Negative feedback and ultrasensitivity can bring about oscillations in the mitogen-activated protein kinase cascades*. Eur J Biochem, 2000. **267**(6): p. 1583-8.
20. Ferrell, J.E., Jr. and E.M. Machleder, *The biochemical basis of an all-or-none cell fate switch in Xenopus oocytes*. Science, 1998. **280**(5365): p. 895-8.
21. Marshall, C.J., *Specificity of receptor tyrosine kinase signaling: transient versus sustained extracellular signal-regulated kinase activation*. Cell, 1995. **80**(2): p. 179-85.

# 3 NETWORK RECONSTRUCTION USING ACTIVATION RATIOS

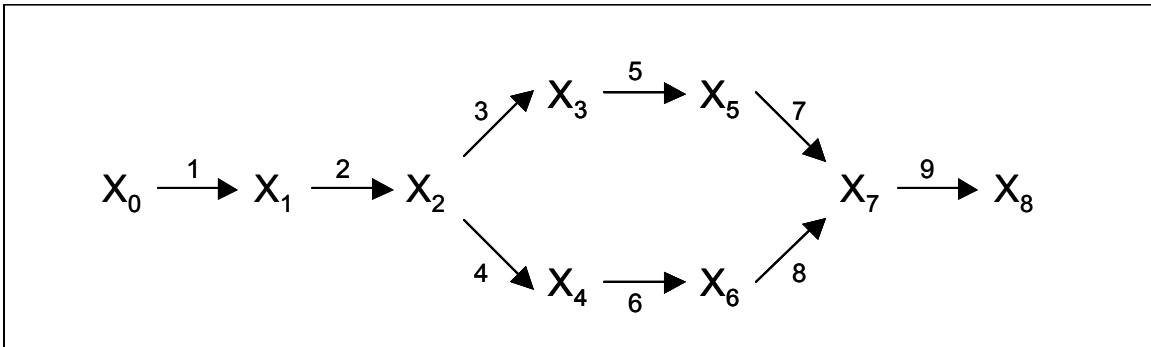
Reconstruction of networks involves the determination of connectivity between components of the network, and can be accomplished through one of two methods. One possibility is to assume *a priori* a network structure, and then utilize a model based upon the assumed structure to predict the behavior of the system. Discrepancies between predictions and experimental observations, after allowing for noise and potential variation of parameter values, can be used to indicate any errors in the assumed structure. The alternative approach is to determine the network structure directly from experimental data, using some sort of inversion method. In the previous chapter, models of simple network structures were used to demonstrate how the structure is reflected in the patterns for activation ratios. This chapter describes an algorithm to deduce the network structure from a set of measured activation ratios, as well as tools useful for automating the process and for evaluating the confidence in the final assignment.

## **3.1. Concepts and Algorithm**

### 3.1.1. Structural and numerical observability

A key issue to consider before delving extensively into procedures for analyzing networks is that of observability, namely, what information can be determined about the network and what features may be invisible. While of course this will depend upon the connectivity of the network, it can also be influenced highly by the type and quality of measurements that are available. Not all measurements may equally contribute information about the network. Furthermore, the numerical sensitivity of some parameters to particular measurements may make the system structurally, but not numerically, observable. In other words, while some measurements may be useful “in theory” for determination of unknown variables, some values of those measurements may lead to numerical singularities whereupon they are rendered essentially meaningless.

As an example, consider the small system shown in Figure 3.1. If this system represents a typical chemical reaction sequence, such as found in cellular metabolism, then the reactions 1-9 may represent uptake, export and conversion between intermediates  $X_0$ - $X_8$ . In that case the *flux*  $v$ , or throughput, of each reaction step is of interest [1]. It is trivial that the fluxes  $v_1$ ,  $v_2$ , and  $v_9$  must be equal at steady state, as well as the total flux between the two parallel paths ( $v_3+v_4$ )—so that measurement of any one can be used to determine all. From the standpoint of structural observability, it should be noted that fluxes  $v_3$  and  $v_4$  (and their respective equivalents) cannot be independently measured if only the total flux is known [2]. In fact, without some additional information about any of the steps in either pathway (such as  $v_3$  or a ratio e.g.  $v_7/v_8$ ) we may as well consider the two paths together as one combined reaction between  $X_2$  and  $X_7$ . In effect, the separation could be ignored since there is no way to independently determine either, and we may not even be aware the split exists. Numerical observability issues appear in this case when there is almost no flux in one branch, or if the capacity to measure the flux is poor. In either situation, the additional information does not actually provide further resolution on the split, simply because of the numerical value (and error) of the measurement.



**Figure 3.1.** Network with parallel pathways.

If the system above is now considered to be a sequencing of signaling reactions, then the physical relationship between intermediates changes drastically. The connections represent transfer of information, which need not be conserved around each node. Thus measurement of “information flux”  $v_1$  (perhaps as an activation factor  $\alpha_0^1$ ) does not give information about  $v_2$ . To resolve the network above, activation of all intermediates  $X_0$ - $X_8$  must be measured so that all activation factors can be determined.

Once again, to correctly identify and quantify the two parallel paths, it is also necessary to obtain some additional information. In this case, measuring the activation of  $X_3$  and  $X_4$  is sufficient to determine the split, since diverging pathways can be observed as described in Section 2.3.3. The convergence at  $X_7$  might be identified as a result of direct connections from both  $X_5$  and  $X_6$ . However, to correctly quantify the convergence, there must be some way to independently regulate the activities of  $X_5$  and  $X_6$ , through additional activation or inhibition with some factor external of  $X_2$ . Otherwise the activities of  $X_5$  and  $X_6$  will be interdependent, and it will be impossible to separately calculate  $\alpha_5^7$  and  $\alpha_6^7$ .

Therefore structural observability limitations in signaling networks arise from the need to measure the activation profile of each intermediate, as well as the capacity to ensure independence of activators at convergence points. For the interactions to be numerically observable, there must be enough samples, with high enough quality in the data, to qualitatively and quantitatively describe lines, hyperbolae and inverse hyperbolae. If too few samples are taken, or only at low activation points, then the conclusions drawn about the network will be incorrect. Hyperbolae will appear as lines, the early curvature for direct interactions may make them appear as inverse hyperbolae, and inverse hyperbolae will appear as very flat lines (low slope). It will be important to increase the spread of data samples, through variation of activators over orders of magnitude, and try to minimize measurement error to reduce scatter that may lead to confounding.

### 3.1.2. Stepwise analysis of networks

Once issues regarding observability have been addressed, by selecting the appropriate number and quality of measurements for a network system, the next step is to combine them an algorithm that utilizes the measurements. As shown in Chapter 2, and summarized in Table 2.2, the arrangement of molecules in simple signaling systems is reflected in the patterns for activation ratios. Reconstruction of signaling networks therefore involves inverting a set of activation ratio measurements to yield one (or potentially several) possible network structures, which may need to be verified by additional experiments. This proceeds by stepwise examination of the activation ratios

for one intermediate I ( $AR_I$ ), plotted against the active concentrations of each other intermediate J ( $J^*$ ), which may yield several possibilities:

1) If the plot is *linear*, then I is immediately downstream of J (most likely, J *directly* activates I). The slope of the line is  $\alpha_J^1$ . Note that there may be some curvature for low values of  $J^*$ , yielding a negative intercept. A positive intercept suggests the possibility of another activator for I.

2) If the plot is *hyperbolic*, then I is further downstream of J (one or more steps exist between J and I). Most likely, J indirectly activates I. The cascade must be determined using direct results for other intermediates between J and I, if available. Also, the hyperbola will be steeper (lower saturation constant) as the distance between J and I increase. Note that if total activation ratios are used, then even direct interactions will appear hyperbolic.

3) If the plot is *inverse hyperbolic*, then either a) I is actually *upstream* of J by one or more steps or b) I and J are on the same level of different branches from an unknown third intermediate. These two possibilities can be resolved by considering plots of  $AR_I$  against  $I^*$  as well as  $AR_I$  and  $AR_J$  against other intermediates  $K^*$ .

4) If the plot is *quadratic* (or higher power), then J directly activates I, through multiple steps. Further study into finding and measuring intermediate forms between I and  $I^*$  may be required.

If there exist multiple inputs to the system, then these studies can be performed for each input individually, as well as two or more together. For intermediates at convergence points, sets of curves appear when both inputs are varied. As described in Section 2.3.2, sets of lines are expected for plots of  $AR_I$  against a direct activator J, while sets of hyperbolic curves are expected for an indirect activator. Regardless, each input must be modulated independently to properly determine the interactions.

Observations regarding the shapes of curves can be determined visually, and often it will be advantageous to examine the activation ratio plots for each intermediate to gain confidence in the method and results. Nevertheless, it is possible to automate this algorithm, by performing both linear and nonlinear regression (using hyperbolic and inverse hyperbolic models) for each pair of intermediates, as will be described in Section 3.3. In either case, the data must be able to allow differentiation between these possibilities by extending as far as possible. Although direct effects can show curvature at early points, the plot will approach a limit of a straight line. An inverse hyperbola will approach an infinite slope at its asymptotic limit, whereas a quadratic plot will remain curved.

### 3.1.3. Consistency

The algorithm described above is based upon examination of each two-way interaction between species, independent of the others in the system. Although the approach seems rather simplistic, its power is expanded by the requirements of consistency. The algorithm can be extended in two ways: first, by looking at the inverse relationship of the two species I and J, and second, by looking at three-way interactions with additional species K. In neither case can the activation ratios take any arbitrary pattern, since they must be determined by the physical relationship between the species. These requirements can be exploited to help resolve any limitations in the data quality (that might lead to misassignment), or point out potential trouble spots where there exists less certainty.

The first approach, looking at inverting two-way interactions, is summarized in Table 3.1. Listed are all the combinations for how activation ratios are expected to change as the target I and effector J are switched. Although  $AR_I$  may take any of the possible patterns when plotted against  $J^*$ , plots of  $AR_J$  against  $I^*$  is restricted to specific possibilities. At least one of the two must be inverse hyperbolic, and no more than one can be linear, hyperbolic or quadratic without violating the rules of consistency. If anything other than these combinations is observed, then there is most likely some problem with the data quality, and it may be important to verify the confidence of both assignments.

**Table 3.1.** Expectations for activation ratios when inverting the relationship between species I and J.

<b>Pattern for AR<sub>I</sub> vs. J*</b>	<b>Relative placement</b>	<b>Likely relationship</b>	<b>Pattern for AR<sub>J</sub> vs. I*</b>
Linear	Direct downstream	$J \rightarrow I$ (direct)	Inverse Hyperbolic
Hyperbolic	Indirect downstream	$J \rightarrow I$ (indirect)	Inverse Hyperbolic
Quadratic	Direct (multistep)	$J \rightarrow I$ (direct)	Inverse Hyperbolic
Inverse Hyperbolic	Upstream	$J \leftarrow I$	Linear, Hyperbolic, Quadratic
	Parallel	$K \begin{matrix} \nearrow I \\ \searrow J \end{matrix}$	Inverse Hyperbolic

Three-way interactions, and further extensions, also can be expected only to have particular combinations of patterns for the activation ratios, which can be assembled from the individual two-way possibilities. Here it becomes much more simple to describe the combinations based upon the possible arrangements for the three species. (Four-way interactions, and beyond, can be sequentially broken down to three-way interactions and so forth.) The possible network structures involving three species, and activation ratio patterns for the three corresponding two-component substructures, are shown in Table 3.2. In each case, the two-way inversion rules as described in Table 3.1 also apply. Again, departure from these combinations is a signal of problems with data quality or pattern assignment.

**Table 3.2.** Activation ratio combinations for three-species structures.

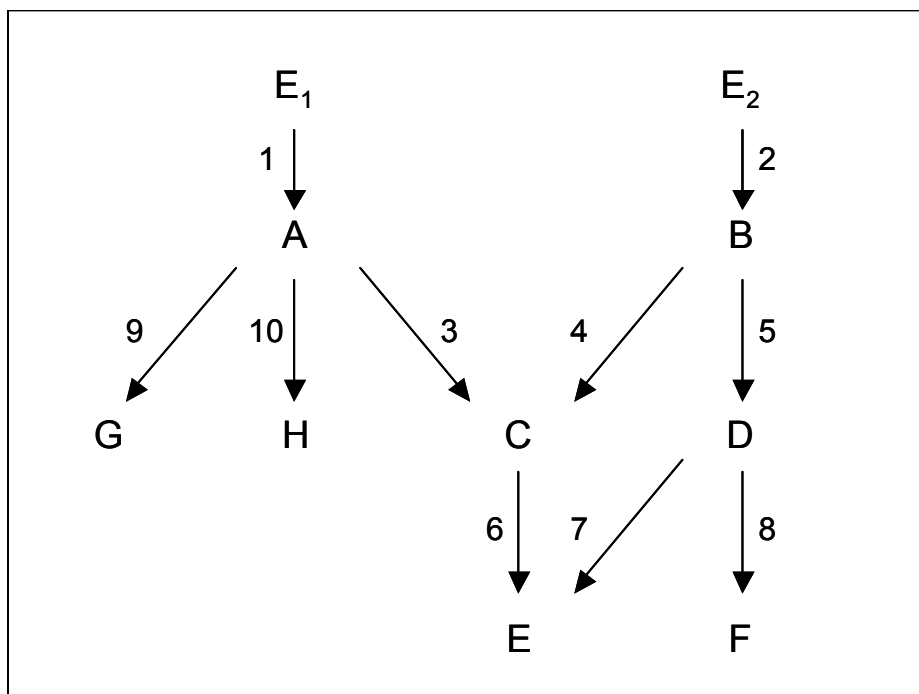
<b>Structure</b>	<b>AR<sub>I</sub> vs. J*</b>	<b>AR<sub>I</sub> vs. K*</b>	<b>AR<sub>J</sub> vs. K*</b>
$K \rightarrow J \rightarrow I$	Linear (direct) or Hyperbolic (indirect)	Hyperbolic	Linear or Hyperbolic
$\begin{matrix} J \\ \searrow \\ K \end{matrix} \nearrow I$	Linear or Hyperbolic	Linear or Hyperbolic	Inverse Hyperbolic
$K \begin{matrix} \nearrow I \\ \searrow J \end{matrix}$	Inverse Hyperbolic	Linear or Hyperbolic	Linear or Hyperbolic



## 3.2. Analysis of a Model Network

### 3.2.1. Network structure and features

To illustrate this algorithm in a somewhat realistic system, a model for a small, interconnected network was constructed, with the structure shown in Figure 3.2. This network contains both distinct and overlapping pathways for the two inputs  $E_1$  and  $E_2$  and cascades with three levels of activation. Parameter values were selected such that in most interactions the  $V_{\max}$  and/or  $K_m$  values for activation or inactivation reactions differ slightly, allowing comparison between the effects of these constants, and using total concentrations such that all reactions are saturated. Further details, including values for all kinetic constants, are included in Appendix 3. Here we will focus on how the analytical framework described above can be used to accurately reconstruct the network structure from a set of simulated measurements. For this purpose, we assume that although all species A-H are known and measurable, we have no *a priori* information about how either  $E_1$  or  $E_2$  may activate these species.



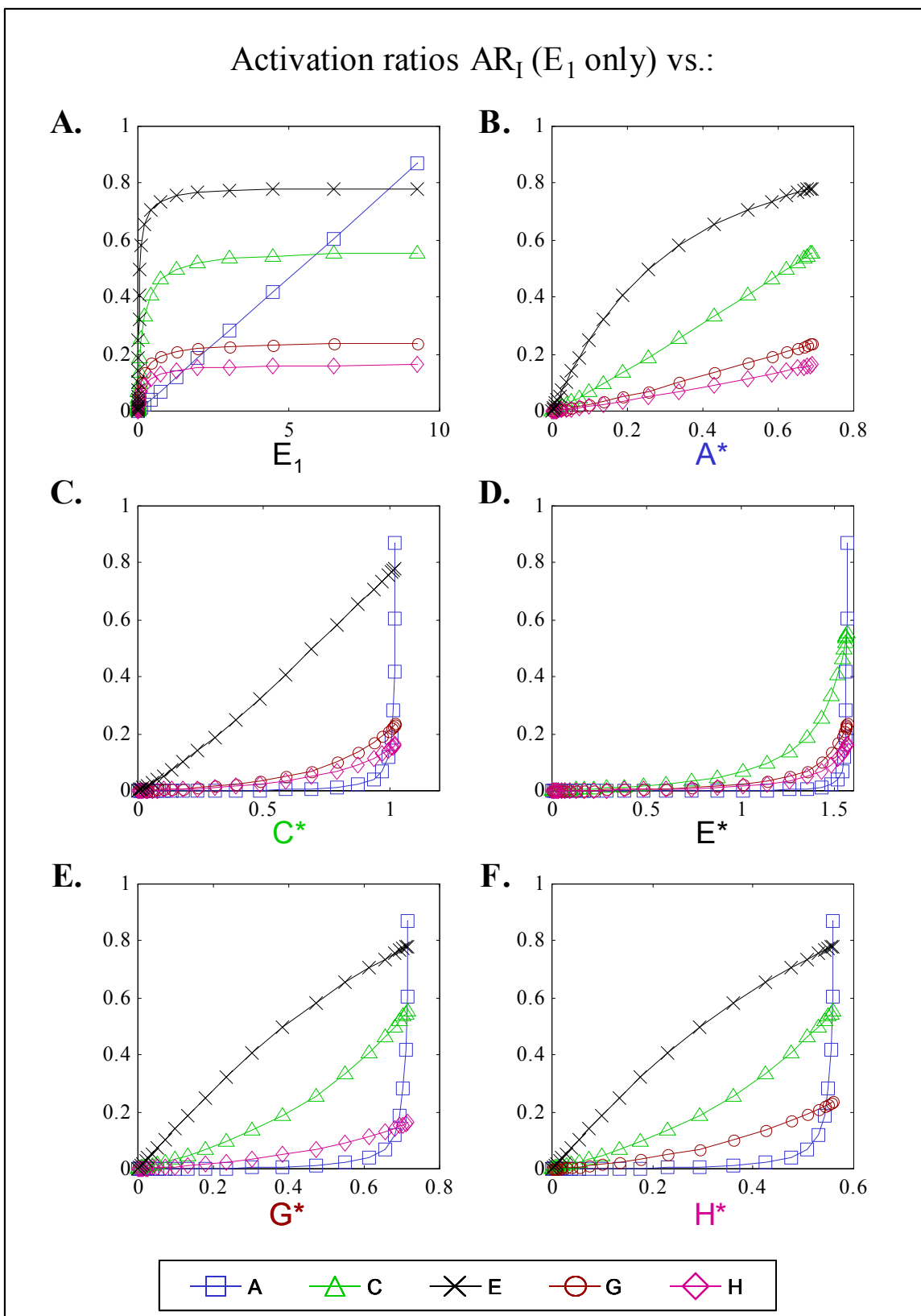
**Figure 3.2.** Diagram of sample model network, showing activation reaction numbering. Model details and parameter values are included in Appendix 3.

### 3.2.2. Network analysis using free concentrations

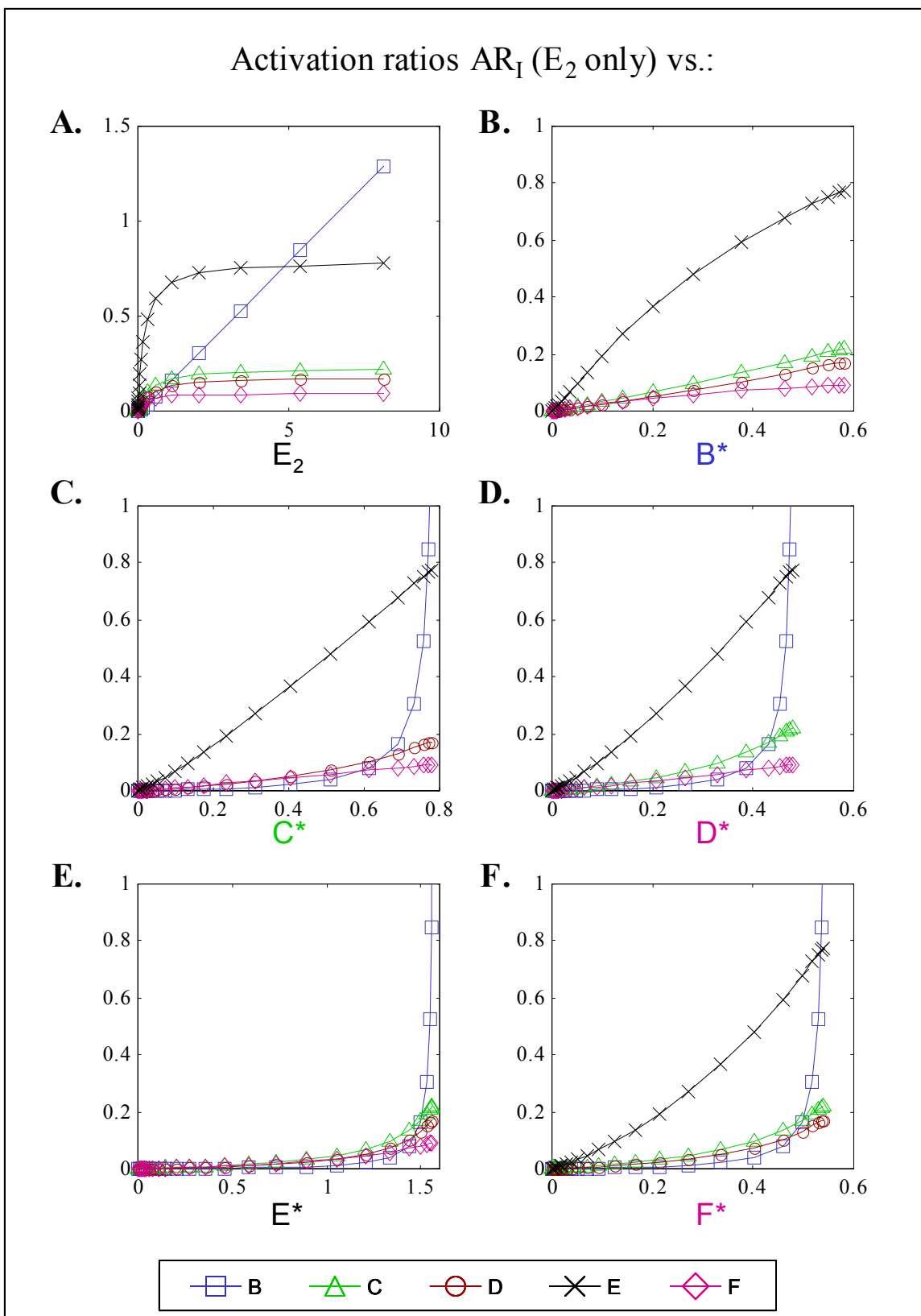
Activation ratios for the network components in response to  $E_1$  and  $E_2$  individually (i.e., for  $E_2$  or  $E_1 = 0$ , respectively) are shown in Figures 3.3-3.4. The color and marker for each curve corresponds with the species being plotted, in accordance with the legends shown in each figure. These activation ratios were calculated using free (unbound) active and inactive concentrations for each component, and are plotted against the free active concentration of each other intermediate. Species that are not included in these plots showed no activation with respect to that input, so B, D, and F were unaffected by the concentration of  $E_1$ , whereas A, G, and H were independent of  $E_2$ . Both  $E_1$  and  $E_2$  activate C and E, indicating an overlap between these two effectors. For the sake of simplicity, first the subnetwork downstream of each input will be examined independently.

Downstream of  $E_1$ , five species (A, C, E, G, H) are activated, and the line in Figure 3.3A indicates that activation of A (blue squares) is direct, whereas the remaining components are indirect, most likely through A. This assumption is strengthened by the plot of activation ratios against  $A^*$  shown in Figure 3.3B, showing lines for C, G, and H (green triangles, red circles and magenta diamonds, respectively). The difference in slopes for  $AR_C$ ,  $AR_G$  and  $AR_H$  shows a preference for the activation of C. The parallel activation of these three species can also be deduced by the fact that activation ratio for each is inverse hyperbolic when plotted against the other two. The activation ratio for E (black x's) is linear when plotted against  $C^*$ , but hyperbolic against  $G^*$  and  $H^*$ , suggesting that E is activated directly by C. All activation ratios are inverse hyperbolic when plotted against  $E^*$ , consistent with the assumption that E is the endpoint of the cascade.

The results in Figure 3.4 indicate that B (blue squares) is directly activated by  $E_2$ , with the remaining species further downstream. Both  $AR_C$  and  $AR_D$  are linear when plotted against  $B^*$ , with almost identical slope. Thus B directly activates C and D, with a slight preference for C. Plots of  $AR_E$  and  $AR_F$  appear linear with respect to both  $C^*$  and  $D^*$ , making it difficult to determine the remaining connections. The results of activation from  $E_1$  only might be used to assume that C activates E alone, while D activates F alone. However, by simultaneously varying both  $C^*$  and  $D^*$  (by varying  $E_1$  and  $E_2$ ), it is possible to clarify the situation.

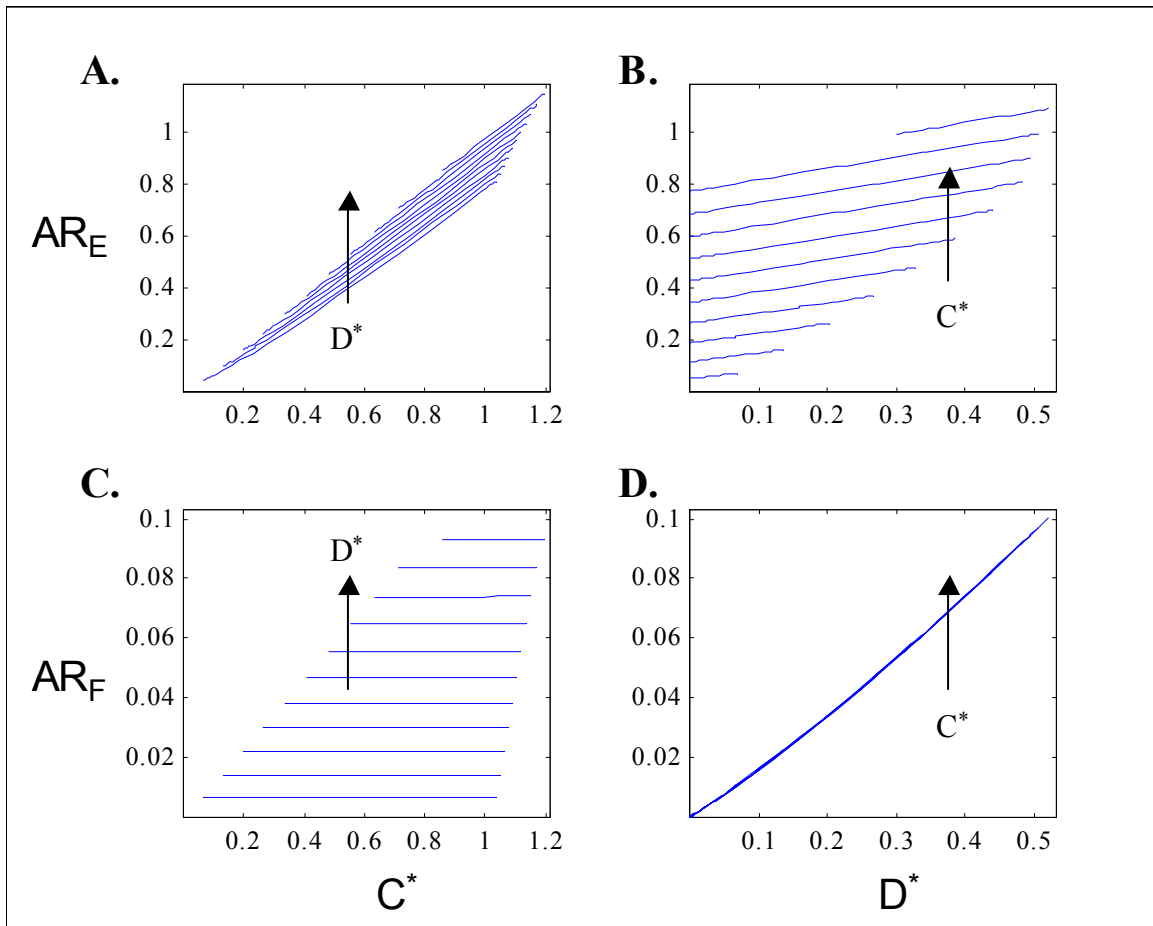


**Figure 3.3.** Activation ratios for intermediates in model network of Figure 3.2, calculated using free species only, plotted against A)  $E_1$ , B)  $A^*$ , C)  $C^*$ , D)  $E^*$ , E)  $G^*$ , and F)  $H^*$ . In each case  $E_2 = 0$ .



**Figure 3.4.** Activation ratios for intermediates in model network of Figure 3.2, calculated using free species only, plotted against A)  $E_2$ , B)  $B^*$ , C)  $C^*$ , D)  $D^*$ , E)  $E^*$ , and F)  $F^*$ . In each case  $E_1 = 0$ .

Activation ratios for the pathway endpoints E and F following activation by both  $E_1$  and  $E_2$  are shown in Figure 3.5. The sets of lines observed in Figures 3.5A-B are typical of converging pathways, indicating that both C and D can independently activate E. However, it is apparent from Figures 3.5C-D that D alone activates F, since all contours of  $C^*$  collapse together. This final information, combined with the results following analysis downstream of  $E_1$  and  $E_2$  individually, allows a complete and accurate reconstruction of the network structure in Figure 3.2.



**Figure 3.5.** Activation ratios for E (A and B) or F (C and D), plotted as contours against  $C^*$  (constant  $D^*$ , A and C) or  $D^*$  (constant  $C^*$ , B and D).

As described in Chapter 2, the activation factors for direct interactions can be estimated through simple and multidimensional linear regression techniques. Each linear plot in Figures 3.3-3.4, as well as the contour plots of Figure 3.5, were therefore fit by linear least squares using the MATLAB “regress” function, which allows for estimation of 95% parameter confidence intervals [3]. Regression results are shown in Table 3.3,

along with the expected values calculated using Equation 2.51, which accounts for variation in the amount of free inactivating enzyme for each reaction. Further data used for calculating these activation ratios is included in Appendix 3. In most cases, there is excellent agreement between the predicted and observed values for the activation factors. One notable exception is for  $\alpha_C^E$ , which is somewhat lower than expected. This is likely due to the curvature seen for low values of  $C^*$  in Figure 3.5A; while it may appear slight, it is enough to decrease the slope for the overall fit. If regression is performed using only data for values of  $C^*$  greater than 0.05, then an accurate estimate of  $\alpha_C^E$  is found ( $0.80 \pm 0.01$ ). Nevertheless, the calculated activation factors vary generally as expected based upon the kinetic constants for each reaction. For example, the  $K_m$  value for activation of E by  $D^*$  is twice as high as for  $C^*$ , so we would expect  $\alpha_C^E \approx 2 \alpha_D^E$ . For the same reason, we expect  $\alpha_A^C \approx 2 \alpha_B^C$ , while  $\alpha_A^G \approx 1.3 \alpha_A^H$  since the  $k_{cat}$  for activation of G by  $A^*$  (2) is 1.3 times higher than for H (1.5).

**Table 3.3.** Regression results for simple network, using free activation ratios.

Reaction No.	$\alpha_J^I$	Linear fit	Regressed $\alpha_J^I$ (95% C.I.)	Theoretical $\alpha_J^I$
1	$\alpha_1^A$	AR <sub>A</sub> vs. E <sub>1</sub>	$0.937 \pm 0.001$	0.938
2	$\alpha_2^B$	AR <sub>B</sub> vs. E <sub>2</sub>	$0.315 \pm 0.001$	0.316
3	$\alpha_A^C$	AR <sub>C</sub> vs. A <sup>*</sup>	$0.799 \pm 0.006$	0.805
4	$\alpha_B^C$	AR <sub>C</sub> vs. B <sup>*</sup>	$0.373 \pm 0.003$	0.378
5	$\alpha_B^D$	AR <sub>D</sub> vs. B <sup>*</sup>	$0.290 \pm 0.004$	0.296
6	$\alpha_C^E$	AR <sub>E</sub> vs. C <sup>*</sup> (E <sub>1</sub> and E <sub>2</sub> )	$0.770 \pm 0.002$	0.800
7	$\alpha_D^E$	AR <sub>E</sub> vs. D <sup>*</sup> (E <sub>1</sub> and E <sub>2</sub> )	$0.399 \pm 0.005$	0.400
8	$\alpha_D^F$	AR <sub>F</sub> vs. D <sup>*</sup>	$0.1878 \pm 0.0004$	0.1904
9	$\alpha_A^G$	AR <sub>G</sub> vs. A <sup>*</sup>	$0.336 \pm 0.006$	0.343
10	$\alpha_A^H$	AR <sub>H</sub> vs. A <sup>*</sup>	$0.230 \pm 0.003$	0.234
6		AR <sub>E</sub> vs. C <sup>*</sup> (E <sub>1</sub> only)	$0.76 \pm 0.01$	0.76
6		AR <sub>E</sub> vs. C <sup>*</sup> (E <sub>2</sub> only)	$0.98 \pm 0.01$	0.76
7		AR <sub>E</sub> vs. D <sup>*</sup> (E <sub>2</sub> only)	$1.58 \pm 0.03$	0.38

It should be noted that the slopes for AR<sub>E</sub> in Figures 3.5A-B, which truly represent the activation of E by C and D, are different than the slopes seen in Figures 3.4C-D, arising when only E<sub>2</sub> is varied. This can also be seen by the difference in activation factors shown at the bottom of Table 3.3. The activation of E is the result of two distinct pathways downstream of B, and AR<sub>E</sub> shown in Figure 3.4 arises from the

activities of both  $C^*$  and  $D^*$ . By varying  $E_2$  only we cannot deconvolute these two pathways, and would calculate incorrect activation factors—in fact, incorrectly predicting that activation is stronger by D. Fortunately, by adding  $E_1$  to modulate C independent of D these two paths can be separated and the accurate activation factors  $\alpha_C^E$  and  $\alpha_D^E$  can be calculated. Since  $E_1$  does not activate D, then  $\alpha_C^E$  can also be estimated using the data for  $E_1$  in the absence of  $E_2$ ; the fact that this value (0.76) is different than using  $E_2$  alone (0.98) hints that another component besides C is involved in the activation of E.

### 3.2.3. Network analysis using total activation ratios

As described in Section 2.4, the presence of enzyme-substrate complexes in measurements of active and inactive species can complicate network reconstruction by activation ratio analysis. These “total activation ratios” behave nonideally in showing dependence upon other intermediates aside from direct activators, and may not be able to distinguish between direct and indirect downstream effects. It is useful to observe how total activation ratios appear in the case of the simple model network, to demonstrate that it is still possible to determine the network structure, although quantitative analysis (through calculation of activation factors) is unavailable.

Total activation ratios for the model network are shown in Figures 3.6-3.8. The same simulation results were used for these plots as for analysis with free species only (Figures 3.3-3.5). However, all enzyme-substrate complexes are included in calculation of total active and inactive forms for each intermediate. As expected, the simple lines for direct effects seen in Figures 3.3-3.5 are replaced by hyperbolae, making network reconstruction somewhat more difficult, as direct and indirect effects cannot be distinguished. Since total activation ratios plotted against upstream activators still appear as an inverse hyperbola, cascades must be deduced by observing which species are above or below a particular intermediate. In this case, no quantitative factors can be calculated, and only a structural analysis is possible.

All total activation ratios are hyperbolic when plotted against  $A^*$  in Figure 3.4B, thus A appears at the top of the network activated by  $E_1$ . Similarly, E must be at the bottom of a cascade since all activation ratios appear inverse hyperbolic. The remainder of the system is more confusing: although C, G, and H all lie in between A and E, it

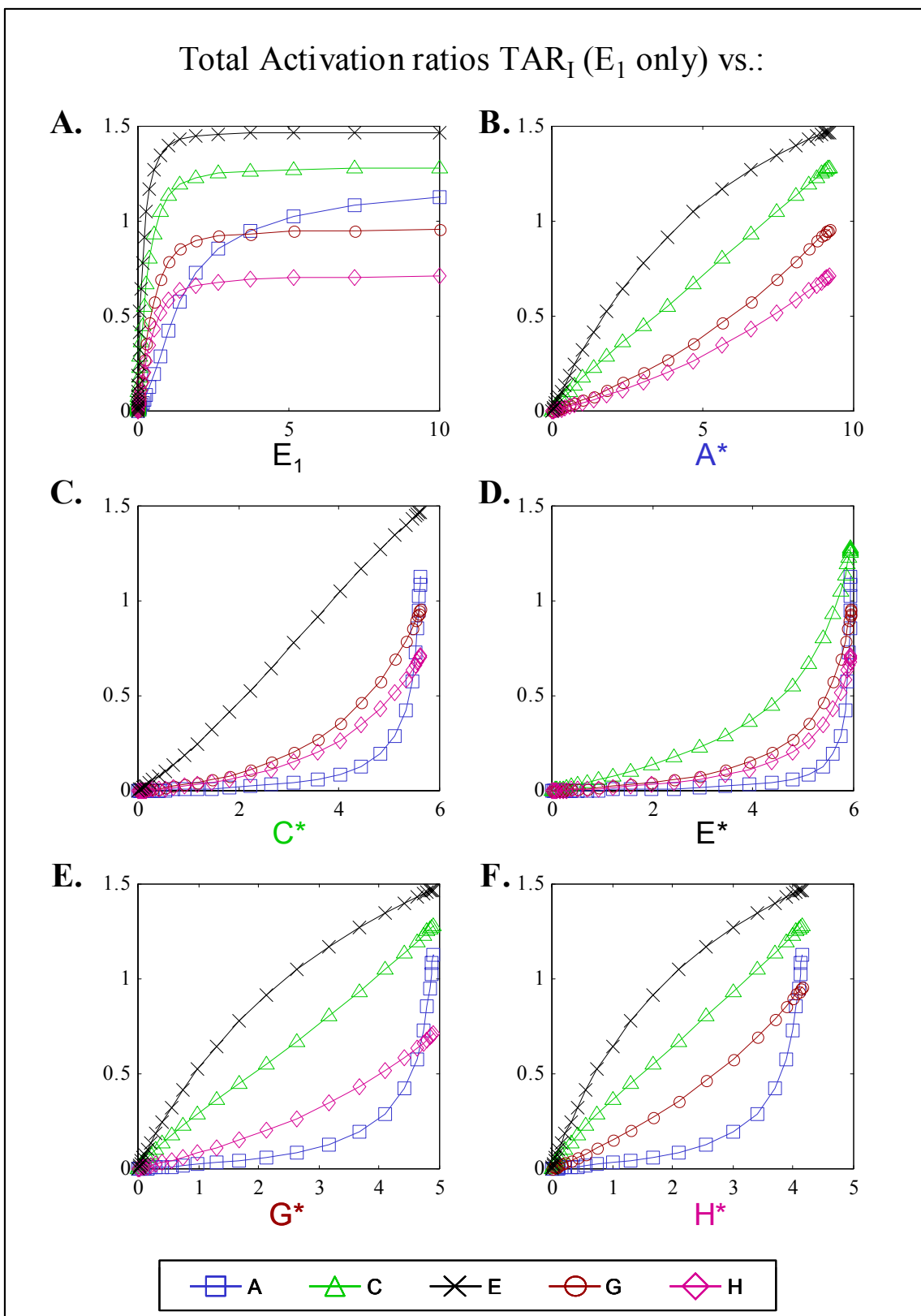
cannot be easily determined which actually activates E. Furthermore, while  $TAR_G$  and  $TAR_H$  are both inverse hyperbolic with respect to  $C^*$ ,  $TAR_C$  is slightly hyperbolic with respect to both  $G^*$  and  $H^*$ , indicating that G and H come before C in a cascade, and may even activate C.  $TAR_G$  appears slightly inverse hyperbolic (almost linear) with respect to  $H^*$ , and vice versa ( $TAR_H$  plotted against  $G^*$ ), which would seem to suggest that G and H are parallel branches downstream of A.

In a similar fashion the results in Figure 3.7 for activation by  $E_2$  alone can be used to conclude that B is the start of the cascade, E is an endpoint, and C and D are parallel branches lying in between, activating E and F with some unknown crosstalk possible. Since  $TAR_E$  is slightly hyperbolic relative to  $F^*$ , at first it might be assumed that F activates E. Unfortunately, not much further detail on either side of the network can be deduced the results from each effector  $E_1$  and  $E_2$  alone.

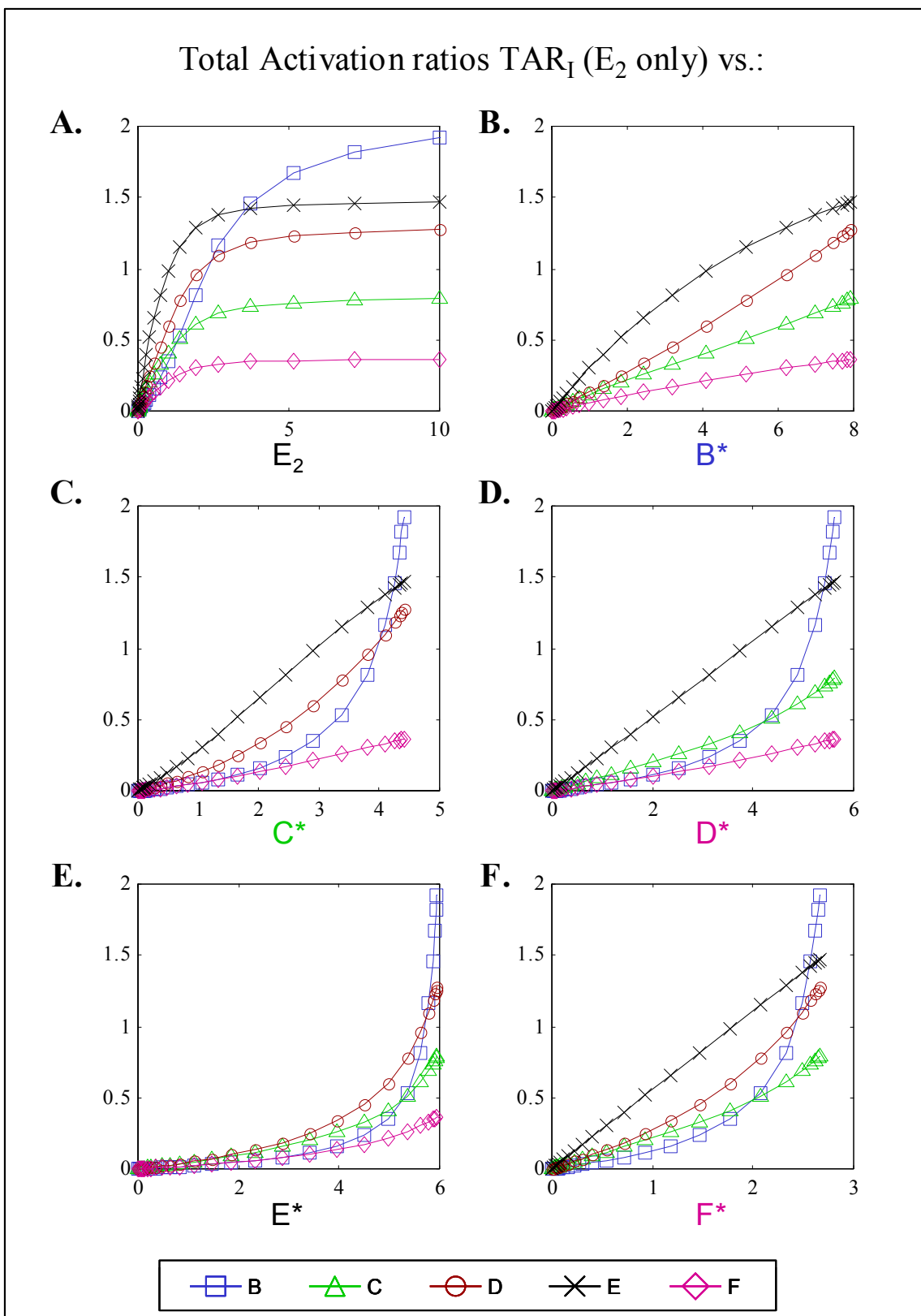
The network can be further clarified by examination of the overlap between  $E_1$  and  $E_2$ , and simultaneous modulation of both. Since  $E_2$  activates C and E, but not G or H, we may conclude that C is parallel to G and H (not downstream of either), and likely activates E. Similarly, F is not required for activation of E downstream of  $E_1$ , so it is unlikely that F activates E. Contour plots of  $TAR_E$  and  $TAR_F$  against  $C^*$  and  $D^*$  again show that D alone activates F, while both activate E. And by plotting contours of  $TAR_E$  against  $F^*$ , which appear as inverse hyperbola, we can again reject the notion that F activates E (data not shown).

We can see that the use of total activation ratios makes structural reconstruction significantly more confusing, even for this relatively simple sample network. Divergence of pathways to form parallel branches can be more difficult to detect, since total activation for each branch may be affected by binding of other components. Similarly, selection between different possible activator candidates may require additional information such as observation of where overlaps appear. For more expansive and interconnected networks, it may be necessary to incorporate knowledge arising from additional sources including *in vitro* reactions, knockouts, or homology to known cascades to help resolve the structure. If free activation ratios can be calculated, however, then the analysis should still be straightforward, as described in Section 3.2.2.

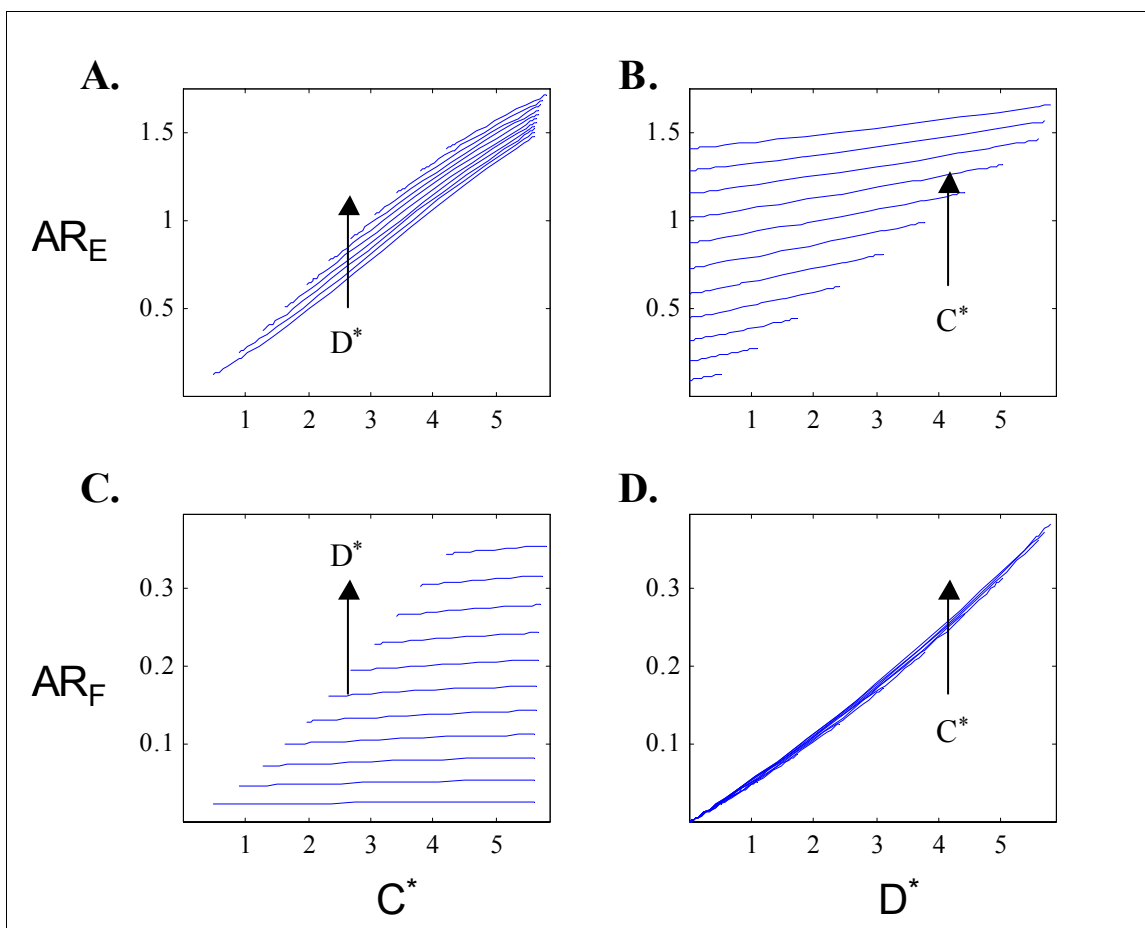




**Figure 3.6.** Total activation ratios for intermediates in model network of Figure 3.2, plotted against A)  $E_1$ , B)  $A^*$ , C)  $C^*$ , D)  $E^*$ , E)  $G^*$ , and F)  $H^*$ . In each case  $E_2 = 0$ .



**Figure 3.7.** Total activation ratios for intermediates in model network of Figure 3.2, plotted against A)  $E_2$ , B)  $B^*$ , C)  $C^*$ , D)  $D^*$ , E)  $E^*$ , and F)  $F^*$ . In each case  $E_1 = 0$ .



**Figure 3.8.** Total activation ratios for E (A and B) or F (C and D), plotted as contours against  $C^*$  (constant  $D^*$ , A and C) or  $D^*$  (constant  $C^*$ , B and D).

A note should be made about network observability. In these examples, it was possible to fully reconstruct the network, and it is also possible to calculate the activation factors when free concentrations are available. This would not have been so simple if data for one or more intermediates were missing. For example, without data on the activation profile for A or B then it would only be possible to say that  $E_1$  and  $E_2$  activate the remaining intermediates indirectly, through some undetermined components. Similarly, without data for C it would not have been possible to determine whether both C and D, or only D, activates F; and whether both C and D activate E. Also, it was critical to examine the activation of E using data obtained from variation of both  $E_1$  and  $E_2$  simultaneously, to rule out activation by F, G, or H and separate effects from C and D. If data were obtained only by varying  $E_2$  it would be possible to determine that B indirectly activates E, but not that there were multiple mechanisms for this activation.

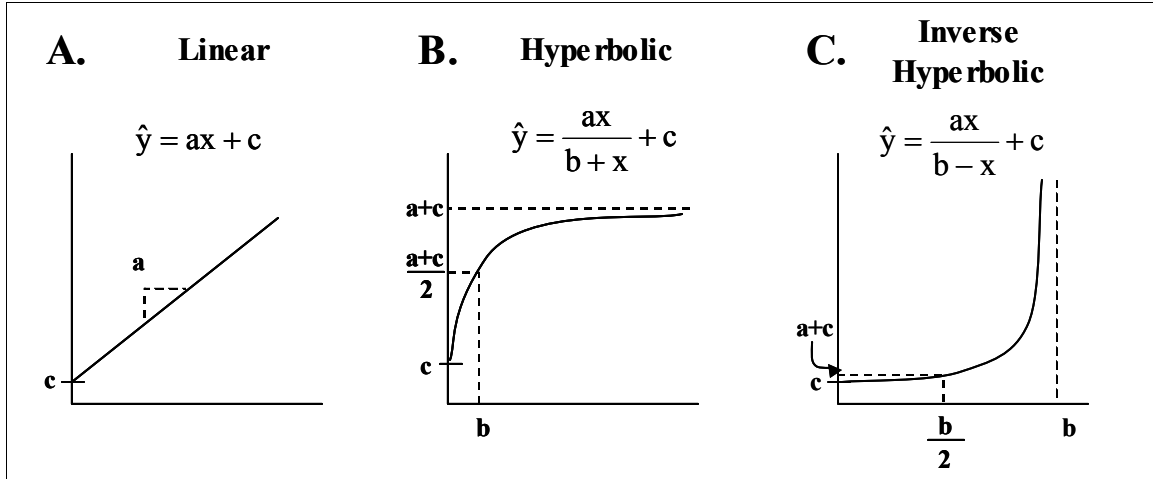
Only by differentially activating the two branches, by the addition of  $E_1$  to further activate C, could the presence of both pathways be discovered.

### **3.3 Automated regression of data and pattern assignment**

#### **3.3.1. Methodology**

Although it is certainly possible to personally examine each plot of activation ratios and develop a network structure by hand, it is quite useful to have a more automatic method, in particular as the number of components and possible interactions increases. Furthermore, an automated procedure would necessarily include some measure of the statistical fit, which can be used in evaluating the accuracy of the assignment. These results could be paired later with a decision-making algorithm to translate a matrix of assignments into one or more putative network structures. This first step involving automated regression and pattern selection steps has been achieved using several functions written in MATLAB.

The pattern assignment methodology has been implemented in three steps. First, matrices X and Y are generated, where the columns of X are vectors containing the activator concentrations ( $E_1$ ,  $A^*$ ,  $B^*$ , etc.) and the columns of Y contain the corresponding activation ratios for targets ( $AR_A$ ,  $AR_B$ , etc.). Second, each column of Y is regressed against each column of X using a linear, hyperbolic, and inverse hyperbolic model, as shown in Figure 3.9. In each case the regression is performed to minimize the total least-squares error between the model fit ( $\hat{y}_i$ ) and data values ( $y_i$ ). Finally, the “best-fitting” model is determined by whichever yields the lowest error, as determined by the Akaike Information Criteria (AIC), as will be described in Section 3.3.2 [4].



**Figure 3.9.** Models used automated regression of activation ratio data, with graphical significance of parameters and half-saturation points shown.

Regression for the linear model parameters is performed by standard linear least squares, using the MATLAB “regress” function [3]. The hyperbolic and inverse hyperbolic fits are achieved by utilizing the “constr” function, which performs constrained optimization of a nonlinear function by Sequential Quadratic Programming (SQP) [5]. These can be written as solution of the optimization problem:

$$\min_{\mathbf{p}=[a,b,c]} \sum_i [\hat{y}_i(x_i, \mathbf{p}) - y_i]^2, \quad \mathbf{p}_{lb} \leq \mathbf{p} \leq \mathbf{p}_{ub} \quad (3.1)$$

Where the model  $\hat{y}$  is defined as shown in Figure 3.9. Thus the objective function for optimization shown in Equation 3.1 is the total sum-of-squares error (SSE, also known as residual sum of squares RSS). Note that in each case the regression is unweighted, with equal importance placed on the fit at each value of  $x$ . An inherent assumption is that the standard deviation  $\sigma_{y,i}$  is approximately equal at each data point  $x_i$ . This procedure may be readily modified to accommodate different weighting for each point, should additional information about the individual errors become available.

Furthermore, the current regression procedure necessarily assumes no error in the  $x$ -values themselves, or at least that the error is small relative to that of  $y$  ( $\sigma_x^2 < \sigma_y^2$ ) [6]. Considering that for activation ratio analysis, the  $x$ -values would represent measured concentrations of free activator, this assumption is unlikely to be strictly valid. (And, since the activation ratios  $y$  are calculated from the free active and inactive

concentrations, the error in y is likely to be correlated to the error in x) The error in y is likely to be greater than that in x, since it involves a ratio. If we express the activation ratio y as the concentration of an active amount z divided by inactive zi, then the error in y can be approximated as [7]:

$$\left(\frac{\sigma_y}{y}\right)^2 = \left(\frac{\sigma_{(z/zi)}}{z/zi}\right)^2 \approx \left(\frac{\sigma_z}{z}\right)^2 + \left(\frac{\sigma_{zi}}{zi}\right)^2 \quad (3.2)$$

Assuming that the variance of measurements for x, z, and zi are all similar, then the variance in y will be approximately twice that of x. Regression in the presence of errors for both x and y is possible, aided significantly by the presence of repeats, but is considerably more complicated [8]. Since the purpose of this work was to demonstrate one potential strategy to automate regression and model selection, however, this extension was not attempted.

The constraints and initial estimates for parameter values when using the hyperbolic and inverse hyperbolic models are shown in Table 3.4. The constraints require that the focal point of curvature (point of half-saturation, as shown in Figure 3.9B-C) occur approximately within the data set. For the hyperbolic model, this point occurs where y is half-maximal and x equals b. If this point occurs at the extreme of the data, then at  $x=x_{\max}=b$ ,  $y=y_{\max}=(a+c)/2$ . Therefore b is constrained to be less than the maximum value observed for x, and a must be no more than twice the range in y. If the data actually appears linear (or is almost indistinguishable from linearity), then the regression solution using the hyperbolic and inverse hyperbolic models will most likely contain some overlap with the parameter constraints.

**Table 3.4.** Models and parameter bounds for automated regression of activation ratio data. Here  $\Delta y = y_{\max} - y_{\min}$ ,  $\Delta x = x_{\max} - x_{\min}$ ,  $y_{\text{avg}} = (y_{\max} + y_{\min})/2$ , and  $x_{\text{avg}} = (x_{\max} + x_{\min})/2$ .  $x(y=y_{\text{avg}})$  signifies the value of x nearest to where y equals  $y_{\text{avg}}$ , and vice versa for  $y(x=x_{\text{avg}})$ .

Model	Functional Form	Parameter (p)	Initial Estimate	Lower Bound (p <sub>lb</sub> )	Upper Bound (p <sub>ub</sub> )
Hyperbolic	$\hat{y} = \frac{ax}{b+x} + c$	a	$\Delta y$	$10^{-4}$	$2 \Delta y$
		b	$x(y=y_{\text{avg}})$	$10^{-4}$	$x_{\max}$
		c	$y_{\min}$	$10^{-4}$	$y_{\max}$
Inverse Hyperbolic	$\hat{y} = \frac{ax}{b-x} + c$	a	$y(x=x_{\text{avg}})$	$10^{-4}$	$y_{\max}$
		b	$1.2 x_{\max}$	$10^{-4}$	$2 x_{\max}$
		c	$y_{\min}$	$10^{-4}$	$y_{\max}$

Initial estimates for the parameter values are calculated assuming that the data fully describes the model. In that case, the estimate of the saturation extreme (a for hyperbolic, b for inverse hyperbolic) is approximately equal to the spread in values for y ( $\Delta y$ ) or maximum value observed for x, respectively. Note that a value slightly higher than  $x_{\max}$  is used for the estimate of b for the inverse hyperbolic model, to avoid a singularity occurring when  $x=b$ . The focal point of curvature for the hyperbolic model (b) is estimated as occurring at the x-value nearest to where y is approximately halfway maximal; and a is determined in an analogous fashion for the inverse hyperbolic model.

### 3.3.2. Model selection and evaluation

The details of different regression models and the general procedures to obtain best-fit parameters were described in the previous section. Once regression has been performed using all three model equations, it becomes necessary to have a criterion that can be used to select which model best represents the data. In general we should not assume any prior knowledge over which of the three models is most likely, and therefore need some sort of quantitative estimate of the likelihood of one the models, given the data, or  $L(\text{model}|\text{data})$ . One option is to simply use the model with smallest residual error, as estimated by the SSE. This approach has two drawbacks. First, for data sets with considerable noise, then all three models may yield high SSE, and small discrepancies between them could incorrectly cause one model to be selected over another. Second, use of SSE alone may bias the selection towards overfitting slightly by preference of the hyperbolic or inverse hyperbolic models over the linear model, in spite of little experimental support for the additional model complexity. One alternative measure is the Akaike Information Criteria (AIC), which is defined as [4]:

$$\text{AIC} = n \ln\left(\frac{\text{SSE}}{n}\right) + 2K + \frac{2K(K+1)}{n-K-1} \quad (3.3)$$

Where n is the number of data points and K is the number of parameters. Equation 3.3 contains a bias-adjustment term, which accounts for relatively low sample sizes ( $n/K < 40$ ), as is likely to occur in signaling experiments. The AIC as defined in Equation 3.3 therefore addresses both issues with using SSE alone: first, the logarithmic transformation helps to compress small differences in large error values, and second, it

contains a penalty for increasing the number of model parameters. Other similar estimators, such as the Mallows  $C_p$  statistic or Bayes Information Criteria (BIC), may also be used, provided that they address these same issues [6]. For the purposes of this work, the AIC was found to be generally satisfactory as a measure of the fit for different models. Selection of the “optimal” model, best representing the data, is achieved by determining which results in the lowest AIC after accounting for variation in number of points and parameters.

One advantage of AIC is that the relative values for the three models can be used to provide an estimate of the relative likelihood of the optimal model compared against the other two. First, the difference between the optimal model and remaining models is calculated [9]:

$$\Delta_i = AIC_i - \min AIC \quad (3.4)$$

If we assume that the optimal model represents the “truth”, then the likelihood that another model also represents the data can be approximated roughly as:

$$L(\text{model}|\text{data}) \propto \exp(-0.5 \Delta_i) \quad (3.5)$$

Comparison between the three models can then be achieved by comparing the likelihood of each, normalized by the total:

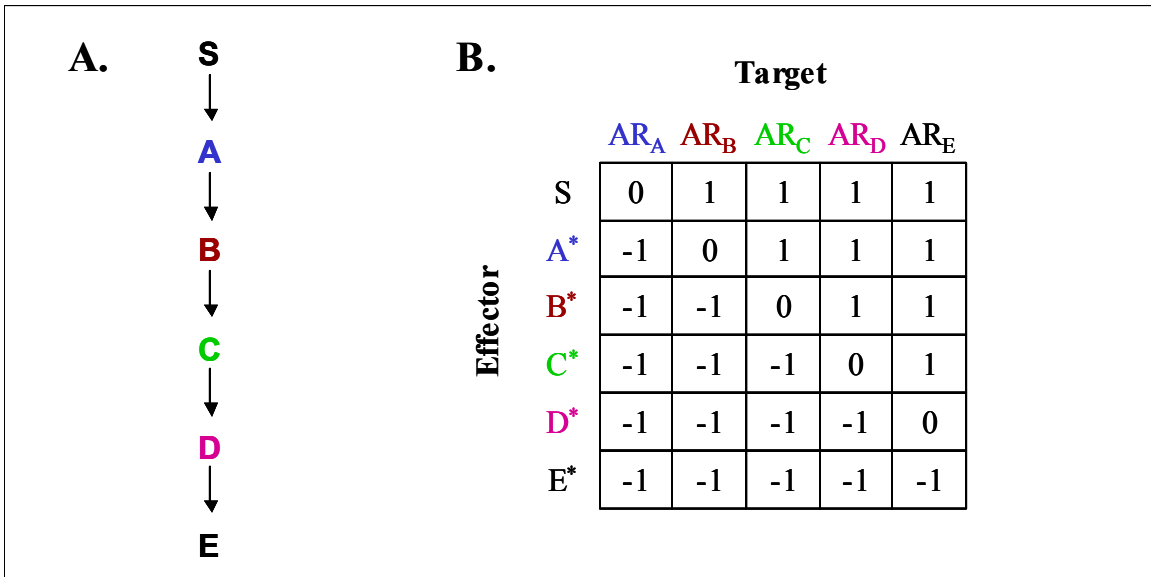
$$w_i = \frac{\exp(-0.5 \Delta_i)}{\sum \exp(-0.5 \Delta_i)} \quad (3.6)$$

The Akaike weights  $w_i$ , calculated in Equation 3.6, range from 0 to 1 and thus provide a readily recognizable value for how each model represents the data. If the weight for the optimal model is greater than 0.9, then the next model can be at most ten times less likely to explain the data. Models with similar weights may be, in essence, nearly indistinguishable in terms of these statistics. In that case additional knowledge may be required to help select between the possible options. This may involve *a priori* information about the interaction between the source and target or application of consistency rules. In any case, the weights can be reviewed after the regression is completed as a measure of confidence in the assignment at each point.

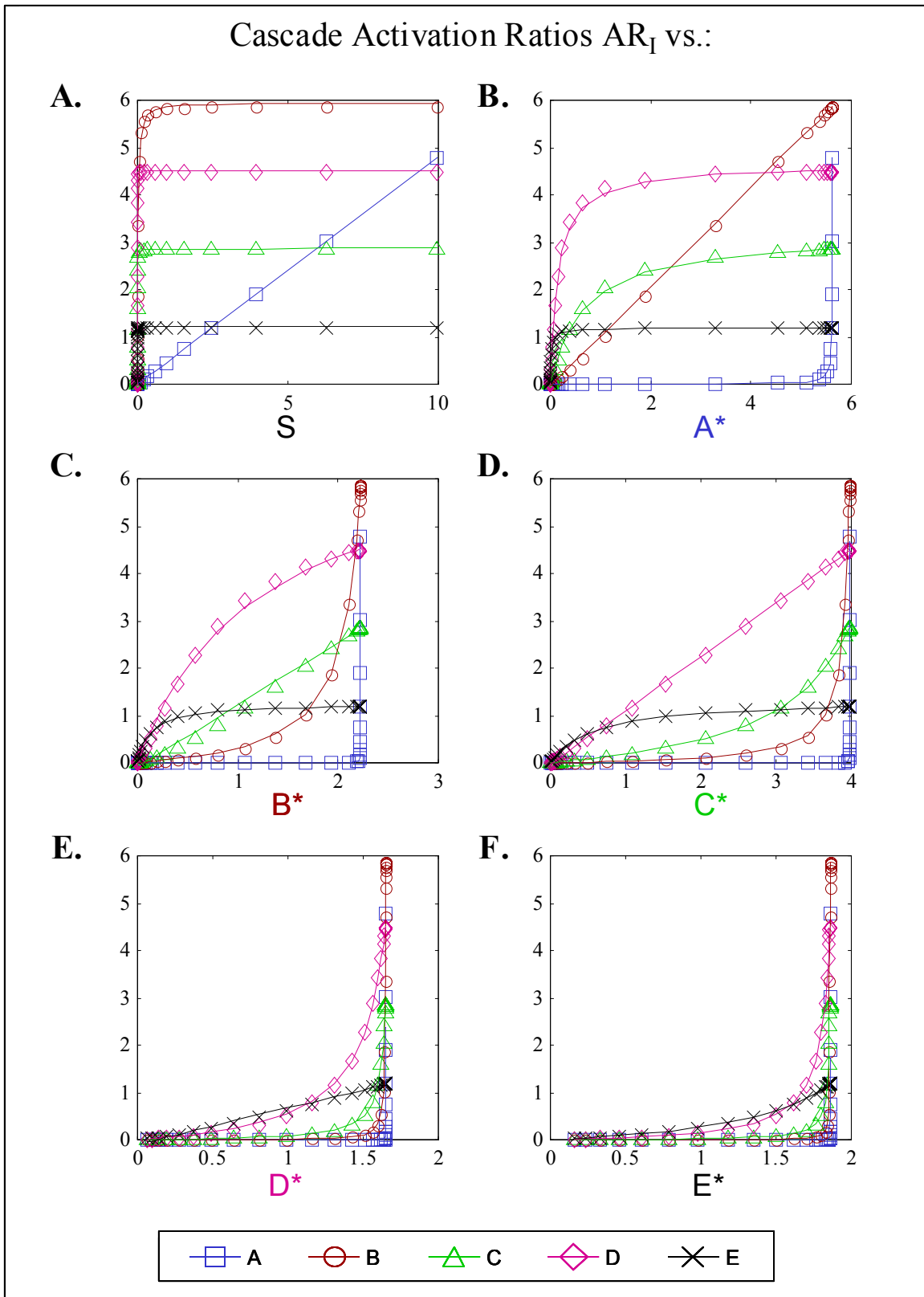


### 3.3.3. Example: cascade analysis in the presence of experimental noise

The algorithms for automated regression and pattern assignment were tested using a model for an extended cascade without feedback, as shown in Figure 3.10A. Each step is saturated and free activation ratios are used for the analysis. Model details can be found in Appendix 4. Initially no error was assumed in any of the variables, therefore the analysis represents the “true” system data. A table with regression and assignment results for the cascade data is shown in Figure 3.10B. Each entry in the table  $R_{ji}$  represents the model selected with the lowest AIC value as calculated using Equation 3.3 following regression of  $AR_i$  against effector concentration  $J^*$ ; 0, 1 and  $-1$  correspond to the linear, hyperbolic, and inverse hyperbolic models, respectively. For this network arrangement and ordering in the table, the contents  $R_{ji}$  are expected to be 0 along the main diagonal (linear for direct targets), 1 above the diagonal (for downstream targets) and  $-1$  below (self-regression and upstream components). This was indeed observed, and in each case the Akaike weights  $w_i$ , calculated using Equation 3.6, for the optimal model was greater than 0.99, indicating that model to be at least 100 times more likely than the next-best fit. Plots of the activation ratio data and best-fit curves using the parameters obtained during the regression are shown in Figure 3.11, exhibiting excellent agreement in all cases.



**Figure 3.10.** A) Cascade structure and color scheme and B) Expected (and observed) output matrix  $R_{ji}$  following automated regression analysis.



**Figure 3.11.** Activation ratios (symbols) and best-fit model curves (solid lines) for cascade of Figure 3.10A following automated regression. Activation ratios for A: squares, B: circles, C: triangles, D: diamonds, E: x's

The robustness of the regression algorithm was tested by addition of noise in two stages. First, error was added only to the y-values, namely the activation ratios themselves. Second, error was added to the individual concentrations of free active and free inactive species, therefore introducing error into both the activation ratios (y) and the predictor values (x). In each case the noise added was normally distributed around the “true” value, with a standard deviation equal to some fraction of the “true” value at each individual point. The modified matrices X' and Y' resulting from addition of error were then used for analysis. This procedure was repeated 10 times for each error type, and the results of each individual regression were compiled. The average regression assignment  $\overline{R_{ji}}$  for the ten replicates for each error method is shown in Table 3.5.

**Table 3.5.** Results of automated regression/decision analysis for extended cascade in Figure 3.9A. Results shown are mean value from 10 replicate calculations in selecting an optimal model. 0 represents linear fit, 1 is hyperbolic fit, and -1 is inverse hyperbolic fit. See Figures 3.12-13 for example data.

	1) $\sigma = 40\%$ , only in AR <sub>i</sub>					2) $\sigma = 20\%$ , I* and I				
	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>
S	-0.3	1.0	1.0	1.0	1.0	-0.1	1.0	0.9	0.9	0.8
A*	-1.0	-0.3	1.0	1.0	0.9	-0.3	-0.5	1.0	1.0	1.0
B*	-1.0	-1.0	-0.2	0.8	1.0	-0.1	-0.7	0.4	0.9	1.0
C*	-1.0	-1.0	-0.9	0.0	1.0	0.0	0.0	-0.3	0.6	1.0
D*	-1.0	-1.0	-1.0	-1.0	0.0	-0.1	0.0	0.0	-0.4	0.6
E*	-1.0	-1.0	-1.0	-1.0	-0.3	-0.4	-0.1	0.0	0.0	-0.1

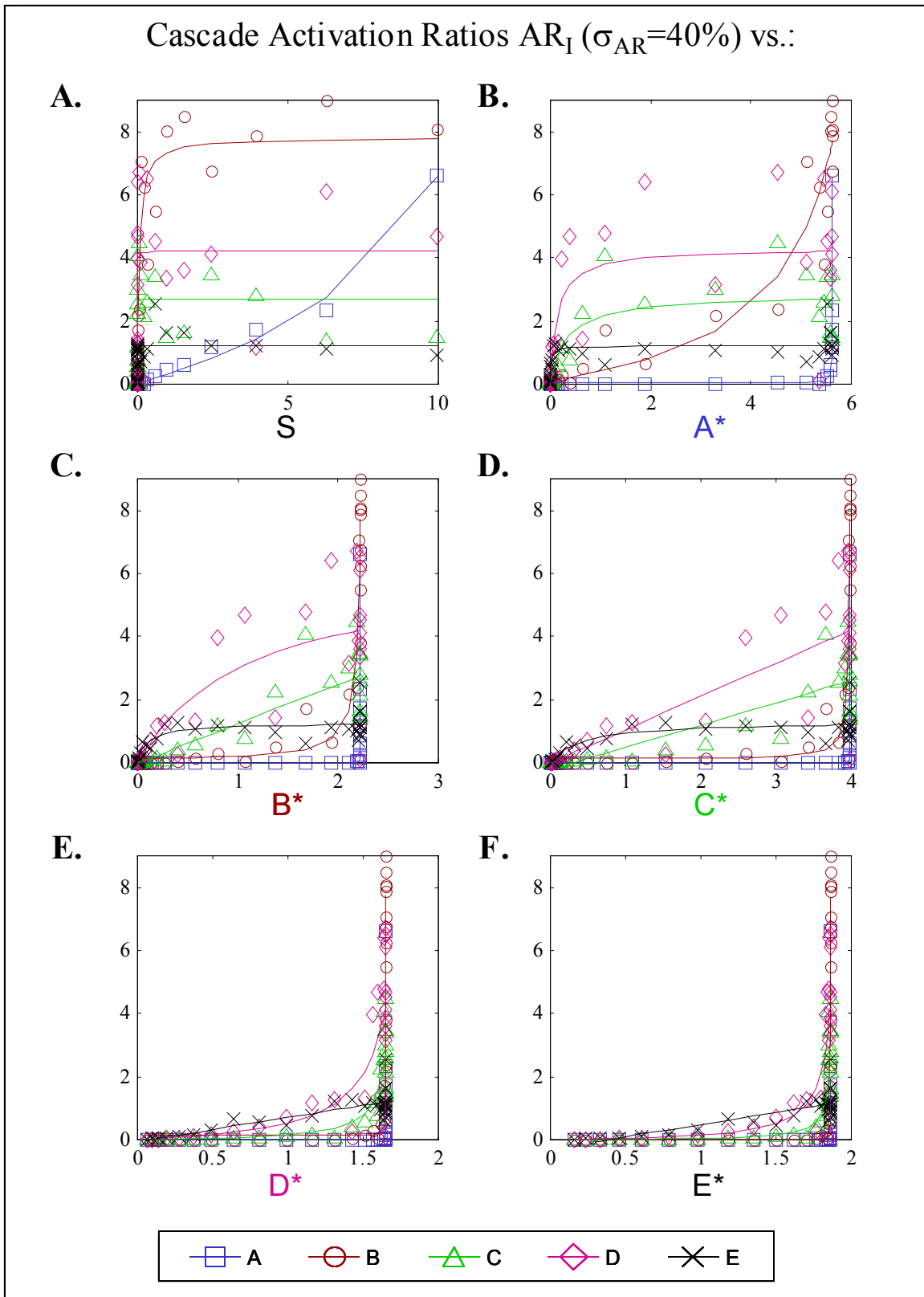
The procedure works well even in the presence of 40% error in the activation ratios AR<sub>i</sub>, selecting the correct model at least 70% of the time (compare against Figure 3.10B). In general, any incorrect assignments are for inverse hyperbola in preference over linear fits (e.g. AR<sub>B</sub> vs. A\*), or vice versa (AR<sub>E</sub> against E\*). Addition of small error (5%) to the predicting activator concentrations J\* also had relatively little detrimental effect on the analysis, but larger error values resulted in a significant deviation from expectations. In particular, the error in x-values consistently resulted in a linear fit improperly being selected over an inverse hyperbola.

This problem, which arises because of the incorrect assumption for regression that the x-values contain no error, seems to be compounded by the nature of the objective function in Equation 3.1. The optimization during regression attempts to minimize the error between model-predicted  $\hat{y}_i$  and observed  $y_i$ . A significant penalty is introduced

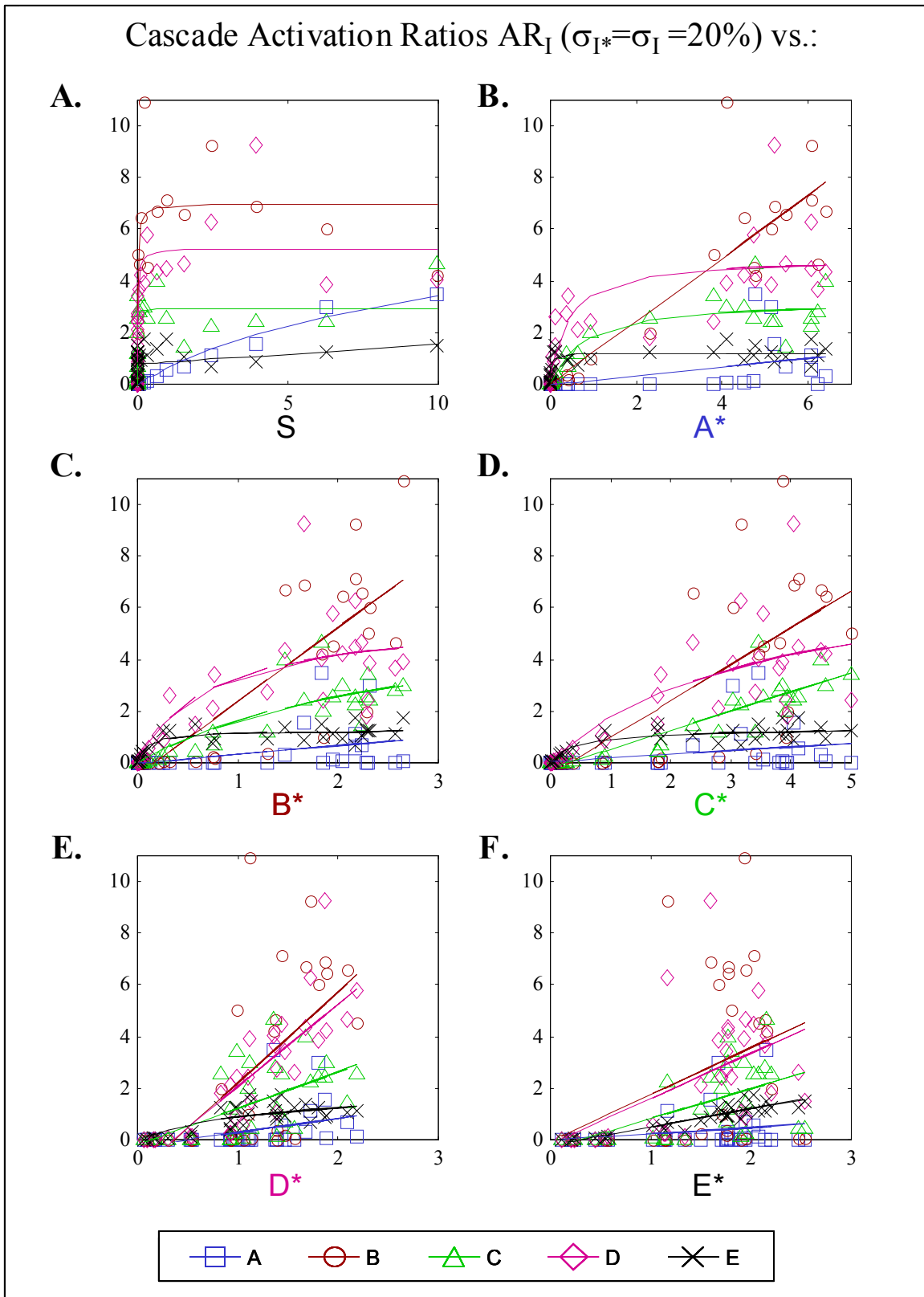
when fitting an inverse hyperbola to sets of data with scattered x-values but similar y values (as observed in the asymptotic region of the inverse hyperbola), since the deviation from the model is actually observed more in the x-direction than in y.

For each error type, the data set that resulted in the worst collection of assignments, as determined by deviation from expectations of Figure 3.10B, was selected for more thorough investigation. The data and best-fit curves (for each optimized model) are shown in Figures 3.12 (for 40% error to  $AR_I$ ) and 3.13 (20% to both  $I^*$  and  $I$ , and therefore also to  $J^*$ ). Although the total variance in the activation ratios is higher in Figure 3.12, the data can still be relatively readily visualized as lines, hyperbolae and inverse hyperbolae. On the other hand, the scatter for active fractions  $I^*$  makes even visual inspection of Figure 3.13 difficult, and unsurprisingly also causes failures during assignment. As described in Section 3.3.2, the Akaike weights can be used as an estimator of the relative probability for different models. The assignment matrices  $\mathbf{R}$  and weights  $w_i$  for the two worst-case sample sets are shown in Tables 3.6-3.7.

Two types of classification errors appear in Table 3.6, when error is added only to activation ratios. The first is when a linear fit is selected over an inverse hyperbola as the activation ratio for a species is plotted against its own active concentration ( $AR_C$  vs.  $C^*$  and  $AR_E$  vs.  $E^*$ ). In both cases, the weight for the linear model ( $w_{lin}$ ) is only slightly higher than for the inverse hyperbolic model ( $w_{inv}$ ), indicating a near-equal probability for the two models. Recall that the activation ratio for a species should always appear inverse hyperbolic when plotted against itself, since  $AR_I \sim I^*/(I_T - I^*)$ . Therefore this error is easy to detect and correct. The second error observed in Table 3.6 is when an inverse hyperbola is a better fit than the line expected for a direct interaction ( $AR_A$  vs.  $S$ ,  $AR_B$  vs.  $A^*$ ). The first of these can be remedied easily, since the external activator  $S$  can never be upstream of any target. Since the activation ratio for  $A$  is also inverse hyperbolic with respect to  $B^*$ , the analysis appears to suggest that these two molecules lie in parallel downstream of  $S$ . Unfortunately, this error cannot be corrected by examination of the weights alone, although it would be detected if repeated experiments were performed, as indicated by the lower value for  $\overline{R_{AB}}$  in Table 3.6.



**Figure 3.12.** Activation ratios (symbols) and best-fit model curves (solid lines) for cascade with noise added to activation ratios. The noise has a standard deviation of 40% of the true value. Symbols as in Figure 3.11.



**Figure 3.13.** Activation ratios (symbols) and best-fit model curves (lines) for cascade with noise added in active and inactive species. The noise has a standard deviation of 20% of the true value. Symbols as in Figure 3.11.

**Table 3.6.** Worst-case results of automated regression/decision analysis for extended cascade in Figure 3.10A, where activation ratios have an added noise term of 40% of true value. Data is shown in Figure 3.12. Tables represent model with lowest AIC and weights for each model.

	Model with lowest AIC					$W_{in}$				
	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>
S	-1	1	1	1	1	0	0	0	0	0
A*	-1	-1	1	1	1	0	0	0	0	0
B*	-1	-1	0	1	1	0	0	0.76	0.46	0
C*	-1	-1	0	0	1	0	0	0.51	0.74	0.04
D*	-1	-1	-1	-1	0	0	0	0	0.06	0.82
E*	-1	-1	-1	-1	0	0.01	0	0	0	0.55

	$W_{hyp}$					$W_{inv}$				
	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>
S	0	1	1	1	1	1	0	0	0	0
A*	0	0	1	1	1	1	1	0	0	0
B*	0	0	0.22	0.52	1	1	1	0.02	0.02	0
C*	0	0	0.01	0.12	0.96	1	1	0.48	0.14	0
D*	0	0	0	0.01	0.06	1	1	1	0.93	0.12
E*	0	0	0	0	0	0.99	1	1	1	0.45

**Table 3.7.** Worst-case results of automated regression/decision analysis for extended cascade in Figure 3.10A, where active and inactive concentrations have an added noise term of 20% of true value. Data is shown in Figure 3.13. Tables represent model with lowest AIC and weights for each model.

	Model with lowest AIC					$W_{in}$				
	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>
S	1	1	1	1	0	0.09	0	0	0	0.57
A*	0	0	1	1	1	0.68	0.59	0	0.35	0
B*	0	0	1	1	1	0.66	0.67	0.48	0.08	0
C*	0	0	0	1	1	0.63	0.78	0.80	0.41	0
D*	0	0	0	0	1	0.73	0.76	0.77	0.74	0.17
E*	0	0	0	0	0	0.64	0.67	0.79	0.70	0.98

	$W_{hyp}$					$W_{inv}$				
	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>	AR <sub>A</sub>	AR <sub>B</sub>	AR <sub>C</sub>	AR <sub>D</sub>	AR <sub>E</sub>
S	0.91	1	1	1	0.28	0	0	0	0	0.15
A*	0.21	0.41	1	0.64	1	0.11	0	0	0.01	0
B*	0.22	0.01	0.51	0.92	1	0.12	0.32	0.01	0	0
C*	0.25	0.06	0	0.58	1	0.12	0.16	0.20	0.01	0
D*	0.09	0.02	0.14	0	0.83	0.18	0.22	0.09	0.26	0
E*	0.20	0.24	0.14	0.24	0	0.16	0.09	0.07	0.06	0.02

If statistical error is added to both inactive and active concentrations of each species, thus affecting both the observed (activation ratios  $y$ ) and predictor (active concentrations  $x$ ) variables, then different problems in assignment appear. Most of the direct interactions, which should yield linear relationships, are instead classified as hyperbolae. In two cases ( $R_{CD}$  and  $R_{DE}$ ), the weight for the hyperbolic model is not much greater than for the linear fit, but in the other two the discrepancy is indeed significant. The introduction of error in this case therefore causes direct interactions to be misinterpreted as indirect effects, but they remain oriented properly in sequence. The misidentification of  $AR_E$  as linear with respect to  $S$  can be recognized partly by the significance of the other weights, and also through a visual inspection of the graph or inspection of the regression parameters, which indicates that the linear fit is nearly flat (slope = 0.08), thus more likely a sharp hyperbola. All upstream interactions and self-regressions should yield inverse hyperbolae, resulting in  $R_{ji}=-1$  underneath the main diagonal. In each case the data is mislabeled as linear. This problem can be recognized by visual inspection of the data. It should also be obvious from application of consistency rules of Section 3.1.3, since examination of any two-way interaction should yield at least one inverse hyperbola. If  $w_{inv}$  values across the main diagonal are compared, then in each case the entry below the main diagonal (corresponding to the inverted sequence, e.g.  $R_{CB}$ ) is lower than its corresponding value above the diagonal ( $R_{BC}$ ), which can be used to correctly place the two components relative to each other.

The relatively simple approach of model fitting, calculation of AIC and verification using consistency rules and Akaike weights can therefore provide a powerful method for performing automated pattern assignment. Introduction of reasonable values for potential experimental error can still be handled, although the algorithm is more susceptible to misclassify interactions. Addition of data replicates, and expansion of the regression procedures to account for errors in predictor variables, may overcome these limitations.



### **3.4. Conclusions**

The characteristics of activation ratios, developed and observed in the previous chapter, were exploited here for the development of a network reconstruction algorithm. A set of experimental measurements for active and inactive concentrations of signaling components can be used to generate activation ratios, which can then be either visually or computationally assigned a functional form. The procedure is subject to similar observability issues seen in examination of other types of networks, such as its incapacity to describe interactions for unmeasured components or to visualize parallel pathways that cannot be independently modulated.

The number and quality of measurements can critically influence pattern assignment. The span of experimental conditions must result in sufficient variation in the measured activation ratios that visible signs of curvature appear. Otherwise, differences between the three model types will not be significant enough to yield high confidence in a particular assignment. Furthermore, excessive scatter in the data can lead to misclassifications as automated regression tools are used. In some cases, application of consistency rules and “common sense” can indicate inaccuracies in assignment. Comparison of Akaike weights to evaluate the relative likelihood of different models can also be effective in proofreading a set of assignments for potential errors.

The algorithms for network reconstruction and automated pattern assignment were applied to two model systems. A small interconnected network was analyzed using both free and total activation ratios. Use of free activation ratios enabled accurate structural determination of the network, and estimates of activation factors that reflected variation of kinetic parameters. Total activation ratios could also be used for reconstruction, but it was significantly more complicated to resolve between several possible structures. In either case, it became critical to combine knowledge resulting from activation by both inputs to properly identify the parallel pathways connecting intermediates B and E. A simple extended cascade was used to demonstrate the pattern assignment algorithm, and demonstrate potential problems arising from several types of experimental error. Although the procedure is relatively robust with respect to variation in activation ratios, scatter in activator concentrations can lead to significant confounding. This may be corrected through manual editing following visual inspection

of the data, or through improvements in the regression and selection methods. Of course, efforts to minimize experimental error will also be powerful in reducing the probability of misclassification. Application of consistency rules will help to resolve these issues, since the regression results are not used independently but rather in the development of a network structure that must yield certain expected pattern sets.

### 3.5. References

1. Stephanopoulos, G., A.A. Aristidou, and J. Nielsen, *Metabolic engineering : principles and methodologies*. 1998, San Diego: Academic Press. xxi, 725.
2. Klapa, M.I. and Massachusetts Institute of Technology. Dept. of Chemical Engineering., *High resolution metabolic flux determination using stable isotopes and mass spectrometry*. 2001. p. 313.
3. Jones, B. and MathWorks Inc., *Statistics toolbox : user's guide*. 1993, Natick, Mass.: MathWorks Inc. 1 v. (various pagings).
4. Akaike, H., et al., *Selected papers of Hirotugu Akaike*. Springer series in statistics. 1998, New York: Springer. viii, 434.
5. Grace, A. and MathWorks Inc., *Optimization toolbox for use with MATLAB*. 1992, Natick, Mass.: MathWorks. 20, 44, 50.
6. Draper, N.R. and H. Smith, *Applied regression analysis*. 2d ed. Wiley series in probability and mathematical statistics. 1981, New York: Wiley. xiv, 709.
7. Taylor, J.R., *An introduction to error analysis : the study of uncertainties in physical measurements*. 2nd ed. 1997, Sausalito, Calif.: University Science Books. xvii, 327.
8. Seber, G.A.F. and C.J. Wild, *Nonlinear regression*. Wiley series in probability and mathematical statistics. Probability and mathematical statistics. 1989, New York: Wiley. xx, 768.
9. Burnham, K.P. and D.R. Anderson, *Model selection and multimodel inference : a practical information-theoretic approach*. 2nd ed. 2002, New York: Springer. xxvi, 488.

# 4 ACTIVATION RATIO ANALYSIS OF ERK PHOSPHORYLATION

The framework of activation ratio analysis was developed with the goal of application to a real experimental system in mind. The hope was to yield an analytical method that could both indicate how experimental data should be obtained and how best to utilize that data. The results described in Chapters 2 and 3 indicate in particular the need to measure amounts of both active and inactive forms of each intermediate, and to separate enzyme-bound from free substrate to enable calculation of free activation ratios. Sufficient activating enzyme must be introduced to drive the cycle forward enough to visualize any curvature in the activation ratios, whether hyperbolic, inverse hyperbolic, or power-law (quadratic or more). Similarly, overall experimental uncertainty for activation ratios must be kept as low as possible, to improve confidence in the assignment of a particular pattern.

The impact of these issues, developed using mathematical analysis and computer model simulations, was therefore addressed using a controlled experimental system. Considering the novelty of the activation ratio framework, it seemed best to investigate its application in as simple and easily controlled system as possible. The experiments would therefore focus in detail on *in vitro* activation of a single species, rather than jump directly into studies using a cascade or larger network.

The covalent modification cycles of the protein kinase Erk2 (mitogen-activated protein (MAP) Kinase 1) were selected as the basis for experimental studies. This enzyme is one of the most highly investigated signaling intermediates, by *in vivo* mutation, expression and characterization experiments as well as detailed *in vitro* mechanistic studies [1-7]. The extensive body of prior research ensured access to materials and information that would be necessary in preparing the experimental system. Additional complexity, to evaluate network-oriented predictions from activation ratios, could later be examined by stepwise addition of other components to the reactive system.

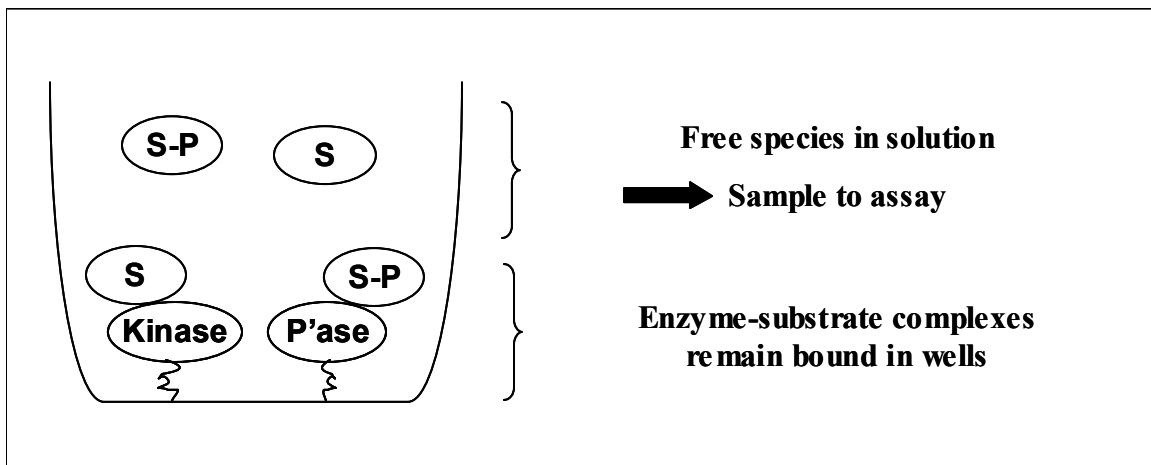
Details regarding the conceptual development of the experimental system are described below in Section 4.1. Each of the system components was evaluated separately, as discussed in Section 4.2, before combining them together for studies of Erk phosphorylation cycles as shown in Section 4.4. Particular limitations in the materials and methods available for quantitation of different Erk phosphorylation forms were addressed by application of a model for the ELISA technique, as described in Section 4.3.

#### **4.1. Experimental system selection and design**

Two key assumptions were made in the development of activation ratio analysis. The first is that signaling systems achieve a steady state. All reactions therefore must be allowed to progress until that steady state is achieved. Second and more critical is that free species can be separated from enzyme-substrate complexes, which inherently requires that the complexes are relatively stable. This must be accomplished without disrupting those complexes, since doing so would yield inaccurate estimates for the activation ratio and potentially lead to incorrect conclusions regarding either the structure or calculation of the activation factor  $\alpha$ . The general argument to assume that this is possible is that often substrates of enzymes can be identified via “pull-down” type co-immunoprecipitation techniques [8, 9]. If the complexes were not stable, these methods would never be able to identify protein-protein interactions in a reproducible manner. Numerous efforts to separate antibody-antigen complexes, as a proxy for enzyme-substrate complexes, by size-exclusion chromatography, nondenaturing polyacrylamide gel electrophoresis (ND-PAGE), or filtration showed extremely poor results. It is likely that these separation techniques all failed because the time required to achieve separation was excessively long compared to the stable lifetime of the complexes. A faster separation method was therefore required to improve the chances of successfully resolving free species from enzyme-substrate complexes.

A simple reaction scheme that addresses these concerns is shown schematically in Figure 4.1. The setup is illustrated for a phosphorylation cycle involving a substrate (S) that is converted to its phosphorylated form (S-P) by a kinase, and back again using a phosphatase. The enzymes are immobilized upon a surface, for example the sides of a well or coated upon a bead. This immobilization does not necessarily need to be covalent

in nature, so long as the enzymes remain bound during the reaction. Substrates and other reagents are incubated with the enzymes until steady state is reached, both in terms of complex formation and net reaction. The liquid supernatant can then be sampled to determine the amounts of the phosphorylated and nonphosphorylated forms, and more importantly, their ratio. The portion of substrate that is abstracted in enzyme-substrate complexes is essentially also immobilized, and therefore readily separated from the free species.



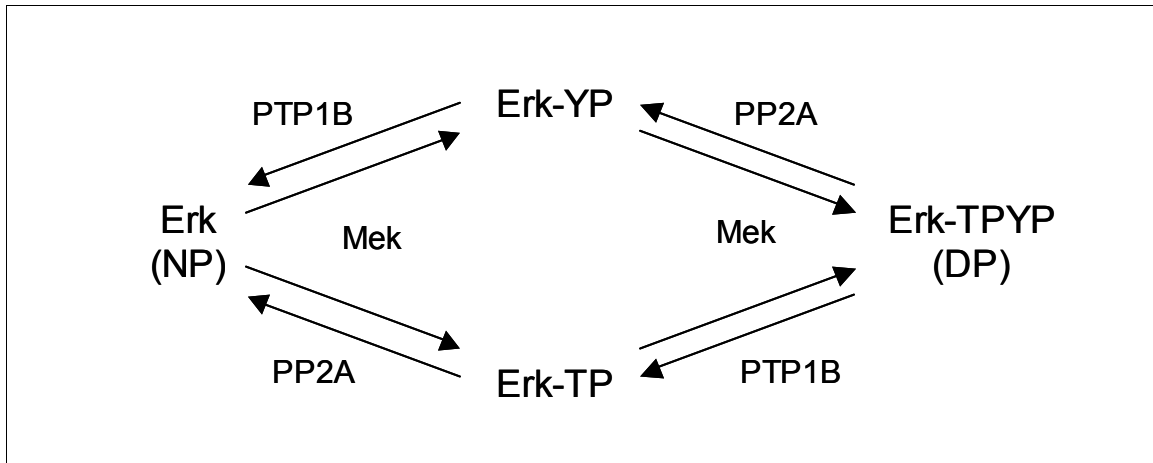
**Figure 4.1.** Immobilization of enzymes enables separation of free species from enzyme-substrate complexes.

The next key issue was selection of a detection method that could be used to independently measure active and inactive forms of a species. Assays that are based upon either enzymatic activity or detection of active-specific label (e.g. radioactive  $^{32}\text{P}$ ) would only enable estimation of the relative or total amount of active form of a species. But in neither case could they be used to determine the inactive fraction that remains inactive. A separate measurement method could be used to determine the total, and the fractions inactive calculated as the difference. However, it is preferable that the same method can be used to measure both inactive and active species, since the potential sources of error such as differences in sample handling, sensitivity, and operating conditions would be identical. In the case of covalent modification cycles such as phosphorylation, it is possible to quantify the amount of nonphosphorylated as well as one or more phosphorylated forms of a protein using isotopic labeling and mass spectrometry (MS) [10, 11]. This discrimination is possible because the modified form

of the protein will have a different mass than the original form, and both can be visualized with the same instrument. Unfortunately, this method was deemed infeasible due to a lack of an available MS and the expertise to perform such measurements.

A more tractable alternative was the use of an enzyme-linked immunoassay (ELISA), since antibodies specific for phosphorylated forms of many proteins are readily available. In the case of Erk1/2 isoforms, monoclonal antibodies specific for the different monophosphorylated and the nonphosphorylated forms have also been developed and are commercially available from Sigma [12]. Erk itself can be purchased readily in purified form from various sources, facilitating development of the reaction system. This helps to make Erk a good choice as the substrate for the reaction system.

The reactions involved in modification of Erk are shown in Figure 4.2. Erk is phosphorylated on two residues (Thr185 and Tyr187 for human Erk2) specifically by Mek1/2 isoforms in a distributed, nonprocessive manner [6, 13, 14]. The tyrosine residue is phosphorylated first in most cases, although threonine-monophosphorylated Erk has been observed during *in vitro* activation experiments. Mek is itself activated by phosphorylation on two serine residues (Ser218 and Ser222 in human Mek1), and can be constitutively activated by site-directed mutagenesis of these residues to aspartate or glutamate [15-19]. Several different phosphatases have been shown to act upon phosphorylated Erk, in particular the dual-specificity phosphatase MKP-3, the tyrosine phosphatase HePTP, and the serine/threonine phosphatase PP2A [20-25]. Since MKP-3 and HePTP are relatively difficult to obtain, the alternate tyrosine phosphatase PTP1B was used together with PP2A to regulate Erk phosphorylation, as illustrated in Figure 4.2.



**Figure 4.2.** Erk covalent modification cycles under action of Mek, PTP1B and PP2A.

## 4.2. Development and Operation

### 4.2.1. Materials and Methods

Reagents were obtained at the highest grade available from Sigma-Aldrich (St. Louis, MO) unless described otherwise. ATP and the Mek inhibitor U0126 were purchased from Cell Signaling Technologies (CST, Beverly, MA). Luria-Bertani broth (LB) was obtained from Becton-Dickinson (Sparks, MD). SuperSignal Pico chemiluminescent detection reagent, Micro BCA and Coomassie Plus protein assay kits, ImmunoPure Normal Rabbit and Goat serum, Reacti-Bind Protein A, Protein G, and Goat anti-Rabbit antibody coated plates were purchased from Pierce (Rockford, IL). Ni<sup>2+</sup>-NTA agarose was obtained from Qiagen (Valencia, CA). Assay kits for Serine/Threonine and Tyrosine phosphatases were purchased from Promega (Madison, WI). Immulon 4HBX high-binding plates were obtained from ThermoLabsystems (Franklin, MA). Costar round-bottom, nonbinding and flat-bottom, high-binding plates were obtained from Corning Life Sciences (Acton, MA). Nunc-Immuno MaxiSorp protein binding plates were purchased from Nalge Nunc International (Rochester, NY).

Monoclonal antibodies against specific Erk phosphorylation forms (aDP: #M8159, aTP: #M3557, aYP: #M3682, aNP: #M3807) were obtained from Sigma. A monoclonal antibody that cross-reacts with all Erk forms (#9107) and a polyclonal anti-Mek antibody (#9122) were obtained from CST. Antibodies against PP2A (#05-421) and PTP1B (#07-088) were obtained from Upstate Biotechnology (Lake Placid, NY). Goat

polyclonal anti-mouse and anti-rabbit antibodies conjugated with horseradish peroxidase were obtained from Pierce.

Purified kinase-deficient mouse Erk2 (hereafter designated as Erk) was obtained from CST. Purified, partially phosphorylated “active” mouse Erk2 (#14-173, ErkPP) was obtained from Upstate and used as a control and in generation of standards for ELISAs. Active human Mek1 (#14-429, designated US-Mek) was also purchased from Upstate as a positive control for Mek purification and activity studies. Human PTP1B produced in *E. coli* and partially purified was obtained from Calbiochem (La Jolla, CA). PP2A purified from human erythrocytes was purchased from Promega.

A sample of *E. coli* transfected with a plasmid encoding a constitutively active form of human Mek1 (hereafter designated N4-Mek) was a generous gift from Natalie Ahn, UC Boulder. This gene product contains a 5kDa addition at the N-terminus including a 6-His sequence useful for purification, deletion of residues 44-51 of the original sequence, and two point mutations (S218E and S222D) that induce activation [16]. This strain was successfully grown, enzyme produced and partially purified in a manner similar to previously described procedures [26]. Briefly, bacteria were grown at 30°C in two 500 mL batches in LB containing 50 µg/mL ampicillin and 34 µg/mL chloramphenicol until an OD<sub>600</sub> of 0.6 was reached. Protein production was induced by addition of IPTG to a final concentration of 100 µM. After 6 hours of growth and protein production, the cells were spun down (2,500g, 25 minutes, 4°C). The cells were resuspended by gentle vortexing in 20 mL water, spun down again and frozen at –80°C. Cells were thawed, resuspended in extraction buffer (EB, see Table 4.1 for composition of buffers) and lysed with two cycles of a French press. Lysates were cleared by centrifugation (15,000g, 30 minutes, 4°C) and adjusted to 1% Tween-20 and 1 M KCl. Ni<sup>2+</sup>-NTA-agarose was added and the mixture rotated end-over-end for 2 hr at 4°C. The resin was loaded into a Poly-Prep gravity-flow column (Bio-Rad, Hercules, CA) and washed with two column volumes of wash buffer (WB). Bound protein was eluted with elution buffer (WB containing 250 mM imidazole) and aliquots stored frozen at –80°C. Final purity was estimated at approximately 50% as judged by SDS-PAGE stained with Coomassie blue.



ELISAs were performed by first adsorbing protein samples to Costar high-binding plates in 50 mM carbonate buffer, pH 9.6, overnight at 4°C. Plates were blocked with 1% BSA in PBS for 2 hours at 25°C and washed once with PBST (PBS with 0.5% Tween-20). Primary antibodies were added at approximately 1:1000 dilution in PBST containing 1% BSA, incubated for 4 hours, and washed three times with PBST. Secondary antibodies were added to 1:2000 dilution in PBST with 1% BSA for 2 hours, the plates were washed three times with PBST, and chemiluminescent reagent was added for 1 minute. Luminescence was read using a Fusion plate reader (Packard Instruments, Meriden, CT).

**Table 4.1.** Buffers used in experimental studies (final values).

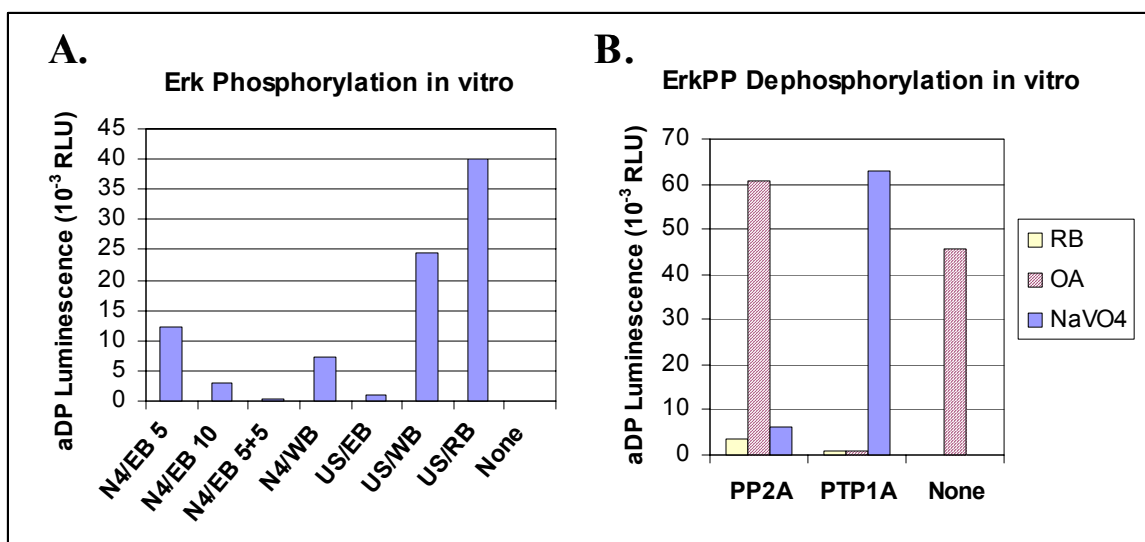
<b>Name</b>	<b>Composition</b>
Extraction Buffer (EB)	50 mM sodium phosphate, pH 8.0, 10% glycerol, 1 mM DTT, 0.25% Tween-20, 1mM PMSF, 1 µg/mL pepstatin A, 2.5 µg/mL leupeptin, and 1 µg/mL aprotinin
Wash Buffer (WB)	50 mM Tris, pH 8.0, 10% glycerol, 0.01% Triton X-100, 1 mM DTT, 20 mM imidazole, 300 mM NaCl
Reaction Buffer (RB)	25 mM Tris, pH 7.5, 2 mM DTT, 15 mM MgCl <sub>2</sub> , 200 µM ATP
RIPA buffer	50 mM Tris, pH 7.5, 150 mM NaCl, 0.1 mM EGTA, 0.1 mM EDTA, 0.1% β-mercaptoethanol
PP2A buffer	50 mM imidazole, pH 7.3, 0.2 mM EGTA, 0.1% BSA, 0.1% β-mercaptoethanol

#### 4.2.2. Reaction operating conditions

Enzymes were tested first in liquid-phase reactions (in microtubes or non-binding plates) to verify individual activities and select reaction conditions. Many previous experiments have been performed with Mek, PTP1B, or PP2A as enzymes and Erk as the substrate, in a variety of different buffers, but never with all components together in one single reaction [6, 7, 13, 14, 19, 21, 24, 27-30]. Almost all contained a buffer base and reducing agent (DTT or β-mercaptoethanol), but other additions varied with application. For example, in many cases phosphatase reactions are performed in the presence of EDTA, while kinase reactions necessarily included a source of Mg<sup>2+</sup>, ATP and on occasion, phosphatase inhibitors. Very few trials have been performed previously with kinase, phosphatase and substrate all together, in these cases, turnover of ATP was enabled by further addition of phosphocreatine and creatine phosphokinase [31-34]. Therefore one key step in system development was selection of buffer conditions that would promote activity of all enzymes together.

Initial trials for Erk phosphorylation using N4-Mek purchased from Sigma in a simple reaction buffer (50 mM Tris, pH 7.5, 1 mM DTT, 50 mM NaCl, 10 mM MgCl<sub>2</sub>, 500 μM ATP) yielded extremely poor results, apparently due to lack of sufficient enzyme. Other trials using activated native human Mek1 purchased from Upstate Biotechnology (US-Mek) under similar buffer conditions were much more successful. Unfortunately this enzyme could not be used for later experiments, since it would be inactivated by the phosphatases, so was useful only as a positive control. Using *E. coli* to produce N4-Mek, it was possible to obtain sufficient constitutively active enzyme to continue further trials.

Sample results using the partially purified N4-Mek are shown in Figure 4.3A. In this early experiment, the N4-Mek was eluted in extraction buffer (EB) containing 250 mM imidazole. Interestingly, less activity is seen when more enzyme is used (N4/EB5 contains 5 μL, N4/EB10 contains 10 μL). These puzzling results appeared to be partially explained by interference of the phosphate contained in the extraction buffer, since simple addition of extra EB to the reaction also reduced Mek activity (N4/EB5+5). Dialysis of the enzyme into wash buffer (WB), which is based on Tris, showed improved activity (N4/WB). This improvement was increased further by directly eluting into WB, thereby avoiding loss of sample during the extra dialysis step. Activity of US-Mek was also poor in EB, better in WB, and better still in reaction buffer (RB: 50 mM Tris, pH 7.5, 1 mM DTT, 15 mM MgCl<sub>2</sub>, 100 μM ATP). It may be that the either salt or detergent, present in the wash buffer but not reaction buffer, inhibits kinase activity by discouraging protein-protein interactions.

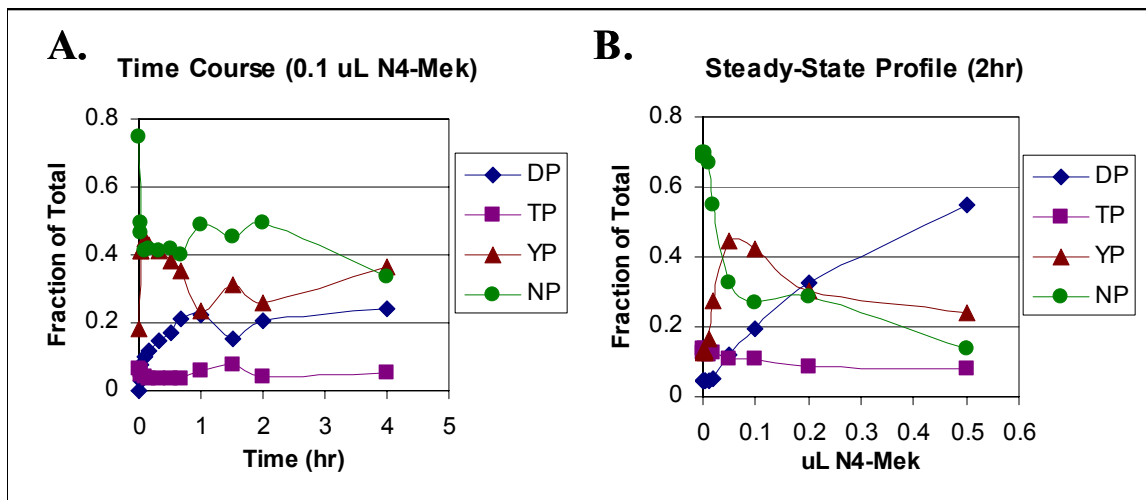


**Figure 4.3.** Sample results for *in vitro* reactions using Erk as a substrate. A) Phosphorylation by N4-Mek (5-10 $\mu$ L) or US-Mek (2 $\mu$ L) in different buffers (see text for composition). For N4-Mek in EB, numbers signify  $\mu$ L Mek added (5+5 is 5 $\mu$ L N4-Mek+5 $\mu$ L EB). B) ErkPP dephosphorylation by 1 $\mu$ L phosphatases in RB with or without addition of inhibitors. OA: 100 nM Okadaic Acid, NaVO<sub>4</sub>: 1 mM sodium vanadate.

As shown in Figure 4.3B, both PTP1B and PP2A also showed activity towards Erk in RB, that could be disrupted using specific inhibitors such okadaic acid (for PP2A) or sodium vanadate (for PTP1B). Addition of up to 150 mM NaCl, 0.05% Tween-20, 1 mM EGTA or 1 mg/mL BSA to base RB had little effect on phosphatase activity towards phosphorylated peptide substrates (data not shown). Therefore, no further modifications were made to the reaction buffer composition, although concentrations of some components were later increased to yield final values of 25 mM Tris, 2 mM DTT, 15 mM MgCl<sub>2</sub>, and 200  $\mu$ M ATP.

With each of the enzymes working well in isolation, the next step was to combine them together to form a complete system. A time course for Erk phosphorylation at 30°C with 0.1  $\mu$ L of each enzyme is shown in Figure 4.4A. The system appears to stabilize within one hour using these conditions, with the Erk distribution unchanging for three additional hours. Formation of tyrosine-monophosphorylated Erk (YP-Erk) peaks within 5 minutes, while subsequent phosphorylation to yield diphosphorylated Erk (DP-Erk) lags behind. Consistent with previously published reports, little threonine-monophosphorylated Erk (TP-Erk) is observed [6, 14]. In fact, the TP-Erk shown in Figure 4.4 may actually be an artifact of cross-reactivity of the aTP antibody with

nonphosphorylated Erk (NP-Erk), since the signal decreases initially. Estimation of each Erk form was accomplished under the assumption of no cross-reactivity of the antibodies, which is unlikely to be valid. Nevertheless, each antibody is reported to have preference for the specific Erk form against which it was raised, and so the trends shown in Figure 4.4 should still reflect the actual system [12]. A more detailed discussion of antibody cross-reactivity, and a model used to account for this effect, is described in Section 4.3.



**Figure 4.4.** Phosphorylation of Erk in presence of N4-Mek, PTP1B, and PP2A. A) Dynamic time course of Erk phosphorylation using 0.1  $\mu$ L of each enzyme. B) Distribution of Erk forms after reaction for 2 hr.

According to Figure 4.4A, two hours reaction should be sufficient to reach steady state. Variation of the total amount of N4-Mek, keeping the phosphatase concentrations constant at 0.1  $\mu$ L stock in 50  $\mu$ L reaction, yielded the activation profiles shown in Figure 4.4B. Again, the primary intermediate observed was YP-Erk, with little if any TP-Erk being formed. Further addition of N4-Mek beyond 0.5  $\mu$ L did little to change the Erk distribution, perhaps because of competition for the substrate, product inhibition, or interference from washing buffer used to store the enzyme. Several groups have noted an inability to completely phosphorylate Erk to DP-Erk during *in vitro* reactions with Mek, in the absence of phosphatases [6, 13, 35]. The activation profiles shown above are consistent with these previous reports, which suggested substrate competition as the most likely source of the deficiency. In any case, the three enzymes appear to work well in concert, and what remained was to obtain a similar activation profile using immobilized enzymes.

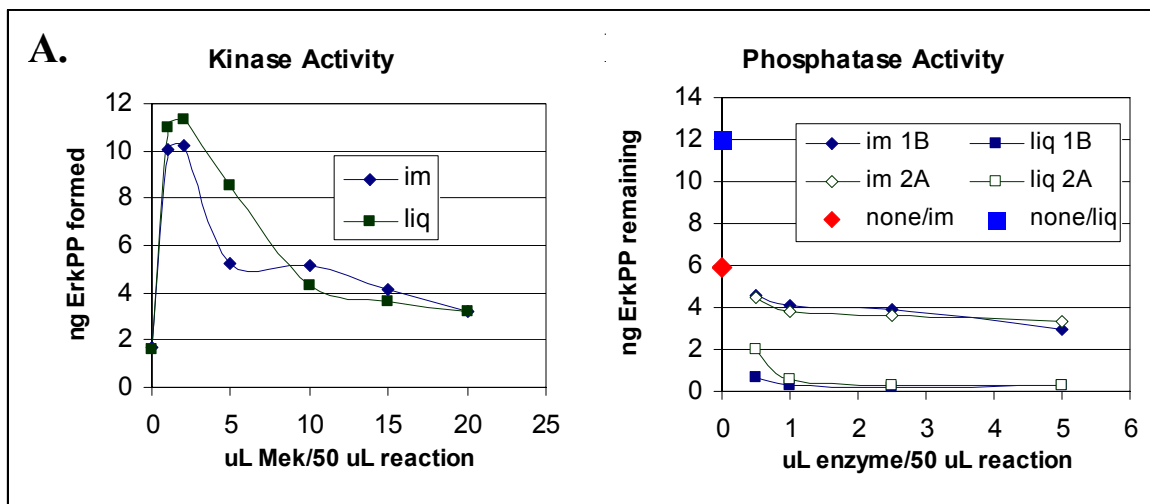
### 4.2.3. Enzyme immobilization

As described in Section 4.1, a key component of the experimental design involves immobilization of the enzymes, to facilitate separation of free substrate from enzyme-substrate complexes. Ideally, as few distinct handling steps and reagents would be involved in the immobilization procedure, since a variety of washes and reactive agents could lead to enzyme inactivation. Immobilization by chemical means, while certainly the most common method, did not appear suitable for this system, since each enzyme contains a different reactive center. A lysine residue is critical for kinase activity in Mek, while tyrosine phosphatases such as PTP1B depend upon a cysteine, and PP2A may utilize a key histidine for phosphate removal [36-38]. Reagents therefore reactive against amine or sulfhydryl groups could potentially inactivate one or more of the enzymes. Two other options for straightforward immobilization were either direct adsorption onto a protein-binding surface or capture by antibodies (immunocapture).

Initial experiments were performed by directly adsorbing enzymes onto treated polystyrene plates (Nunc MaxiSorp), as this would be the simplest procedure. As shown in Figure 4.5A, N4-Mek appeared robust to adsorption, since activity was only slightly decreased during immobilization. Elevated concentrations of N4-Mek actually decreased the activity towards Erk, as described above, even when the enzyme was immobilized. It may be that the washing buffer (WB, as above) composition inhibited adsorption to polystyrene. The plates were washed with PBST before reactions were performed, so any interference of the buffer, potentially the cause of decrease in the liquid-phase reactions, would not occur for immobilized enzyme.

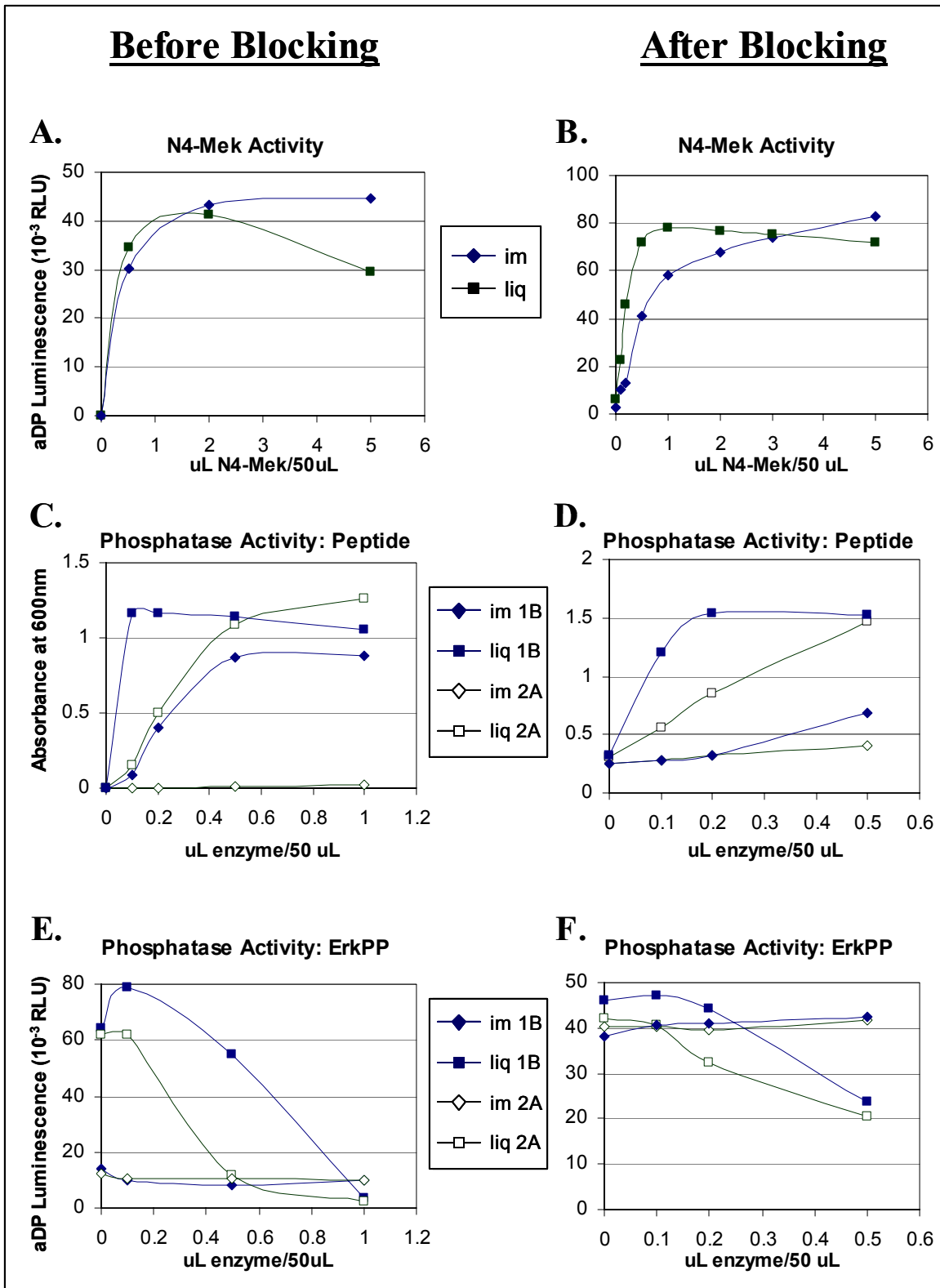
Unfortunately, both phosphatases showed low activity following adsorption, as shown in Figure 4.5B. In fact, most of the depletion of ErkPP was due to nonspecific adsorption to the plate, as evidenced by a decrease of substrate in the absence of enzyme. This nonspecific adsorption was virtually eliminated by optimization of blocking conditions, including decreasing the concentration of BSA in the blocking buffer from 5% to 1% and by using the smaller casein (approximately 24 kDa) rather than BSA (67 kDa). However, the phosphatases remained inactive regardless of adsorption buffer composition, pH, or time (data not shown). It may be that the surface chemistry, which was designed for tight capture of proteins, was overly strong for enzyme adsorption. The

phosphatases may therefore have unfolded and spread out upon the surface, disrupting the structure around the catalytic sites and inactivating the enzyme. As repeated trials failed to exhibit enzymatic activity, it appeared best to switch to immunocapture for enzyme immobilization.



**Figure 4.5.** Enzymatic activities following direct adsorption to polystyrene. A) Phosphorylation of Erk using N4-Mek. B) Dephosphorylation of ErkPP using PTP1B (closed symbols) or PP2A (open symbols). Diamonds represent immobilized reaction conditions, while squares represent data taken for liquid-phase reactions.

Studies using capture antibodies showed improved results, although activity was low when compared against liquid-phase reactions. Antibodies against each enzyme were either captured themselves using plates precoated with Protein A, Protein G, or goat anti-rabbit antibodies or adsorbed directly onto polystyrene plates. The use of antibody-binding surfaces such as Protein G was attempted to increase the efficiency of enzyme capture by ensuring that epitope-binding regions of the antibodies would be presented towards the enzymes. N4-Mek immobilized to pretreated plates showed excellent activity (Figure 4.6A), as did PTP1B when using a phosphotyrosine-containing peptide as test substrate (Figure 4.6C). However, neither phosphatase was active towards Erk; instead, any ErkPP depletion was due to nonspecific adsorption (Figure 4.6E), even though the plates are already blocked. Addition of an extra blocking step inhibited Erk adsorption (Figure 4.6F). However, the phosphatases were still inactive towards ErkPP, and activity towards peptide substrates was diminished (Figure 4.6D). Mek activity was also decreased slightly after blocking (Figure 4.6B).



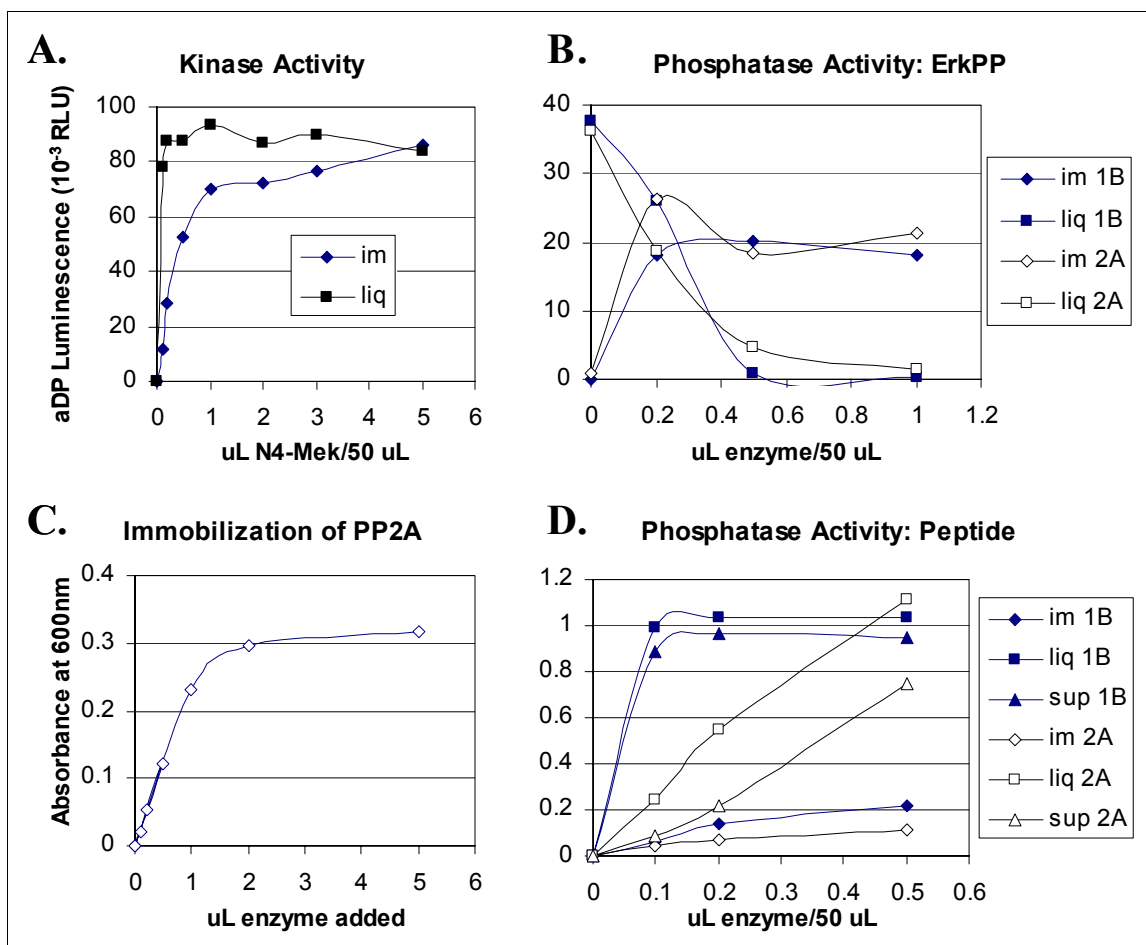
**Figure 4.6.** Enzymatic activities following immobilization on Protein G-coated plates using capture antibodies without (A, C, E) or with (B, D, F) an extra blocking step. A and B) Phosphorylation of Erk by N4-Mek. C and D) Dephosphorylation of phosphotyrosine- (PTP1B) or phosphoserine-containing (PP2A) peptides. E and F) Dephosphorylation of ErkPP by phosphatases.

A benefit of antibody-binding plates should be the ability to orient antibodies such as to utilize completely their capture potential. The Protein G or similar molecules would potentially reduce any steric hindrance from the blocking proteins by providing extra space between the surface and the epitope-binding regions of the antibody. However, it may be that these antibody-binding molecules were showing undesired cross-reactivity with ErkPP. The additional blocking step, which inhibited Erk adsorption but also any enzymatic activity, essentially nullified the benefit of using pretreated plates.

Therefore, the antibodies were next directly adsorbed onto polystyrene plates with various degrees of hydrophobicity (ThermoLabsystems Immulon 1B, 2HB, 4HBX, EB, and UB). As PP2A was generally the least active of the three enzymes following immobilization, the high-binding 4HBX plates, which yielded the highest activity of PP2A, were selected for all further work. Of course, in this case a blocking step was essential following antibody coating, and was optimized both for inhibition of Erk adsorption and promotion of enzyme binding. The coating, blocking, binding, and washing conditions tested during optimization are summarized in Table 4.2 below.

As shown in Figure 4.7, each enzyme was immobilized with low but observable activity. In particular, the phosphatases were indeed active towards ErkPP. (No ErkPP was added to wells for without immobilized enzyme.) In all cases, it appeared that immobilized enzyme activity was approximately 5-10% of that for corresponding concentrations in liquid-phase reactions. Increasing the amount of enzyme added to the wells during capture yielded higher immobilized activity, but only up to a saturation point beyond which no further improvement was seen, as evidenced for PP2A in Figure 4.7C. The supernatant liquid following enzyme binding contained almost all the activity added to the wells (compare triangles and squares in Figure 4.7D). This seems to suggest that the low activity of immobilized enzymes was due simply to poor efficiency during binding, and that what little enzyme was bound was active. Changing the antibody concentration, buffer composition, wash stringency, or binding incubation time did little to increase the capture efficiency. Therefore, the conditions that yielded the highest activity possible for all enzymes were utilized for final studies of Erk activation, and are listed in Table 4.2. All steps were performed at 4°C except for reactions, run at 30°C. Plates were washed once before and after enzyme binding.





**Figure 4.7.** Enzymatic activities following immobilization using capture antibodies directly adsorbed to plates. A) Phosphorylation of Erk by N4-Mek. B) Dephosphorylation of ErkPP by phosphatases. C) Saturation of immobilization of PP2A. D) Phosphate release from phosphopeptides. Diamonds: immobilized enzymes, Squares: liquid-phase reactions, Triangles: enzymes remaining in supernatant liquid following immobilization.

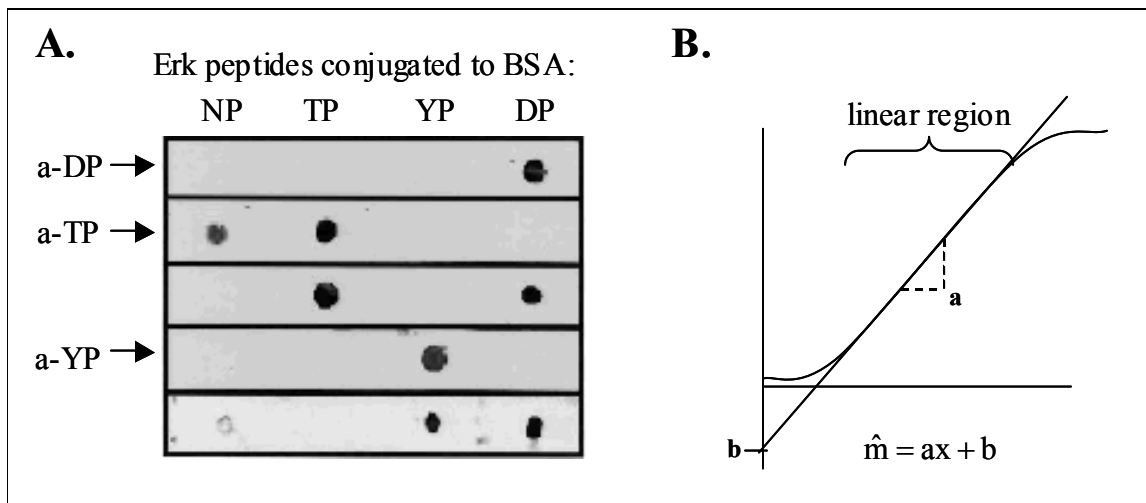
**Table 4.2.** Conditions tested during optimization of enzyme immobilization, and details of final procedure.

Parameter	Variations Tested	Final
Plate	Immulon 1B, 2HB, 4HBX, EB, UB	4HBX
Coating:		
Antibody conc.	1, 2, 5, 10, 20, 50 $\mu\text{g/mL}$	2 $\mu\text{g/mL}$
Coating buffer	TBS, TBST, PBS or carbonate buffer, pH 5.0, 7.4, 9.0	PBS, pH 9.0
Blocking:		
Agent	BSA, casein, bovine, goat, or rabbit sera	casein
Concentration	1%, 2.5%, 5%	1%
Buffer	PBS, TBS, TBST, RB, pH 7.4 or 9.0	PBS, pH 9.0
Enzyme binding:		
Incubation time	4 hr or overnight (~14 hr)	Overnight
Buffer	TBS, RB, RIPA, PP2A	RIPA
Washing buffer	TBS, TBST, RB, RIPA, PP2A	PP2A

### 4.3. Model for Interpreting ELISA Data

#### 4.3.1. Motivation and concepts

The set of monoclonal antibodies available from Sigma (designated aDP, aTP, aYP, aNP) was selected for use in quantitation of Erk in its various forms by ELISA. Unfortunately, this quantitation was complicated by cross-reactivity of the antibodies. The antibody raised against TP-Erk (aTP) also recognizes NP-Erk, as shown in Figure 4.8A. Additionally, aNP showed significant reactivity with TP-Erk and YP-Erk as well as NP-Erk (as reported by Sigma). The remaining antibodies almost complete selectivity for their intended targets. Nevertheless, the failure of a simple one-to-one correspondence required a deconvolution procedure to translate ELISA data (in relative luminescence units, or RLU) into amounts of each form. Furthermore, like many experimental methods ELISA rarely shows a purely linear response to the amount of sample, but may exhibit regions of linear behavior as seen in Figure 4.8B. This nonlinearity arises from the numerous steps involved in ELISAs: equilibrium binding of sample to a solid substrate, primary antibody to sample, secondary antibody to primary antibody, and reaction of a luminescent substrate with secondary antibody-conjugated enzyme. Each of these binding steps may saturate, leading to an overall saturation point of the assay, while sensitivity issues may induce curvature at low sample concentrations.



**Figure 4.8.** Expectations for ELISA results with potential issues. A) Cross-reactivity of anti-Erk antibodies (taken from Yao et al [12]), B) Sigmoidal overall response containing an intermediate linear region.

For the sake of simplicity, we may focus our attention upon the linear response of the assay, with the realization that later this may introduce error if we attempt to

extrapolate beyond this region. If we further assume that the cross-reactivity can be modeled as a simple additive effect of each species, then the following model for the measurements can be generated:

$$\begin{bmatrix} \text{aDP} \\ \text{aTP} \\ \text{aYP} \\ \text{aNP} \end{bmatrix} = \begin{bmatrix} \text{a}_{11} & 0 & 0 & 0 \\ 0 & \text{a}_{22} & 0 & \text{a}_{24} \\ 0 & 0 & \text{a}_{33} & 0 \\ 0 & \text{a}_{42} & \text{a}_{43} & \text{a}_{44} \end{bmatrix} \cdot \begin{bmatrix} \text{DP} \\ \text{TP} \\ \text{YP} \\ \text{NP} \end{bmatrix} + \begin{bmatrix} \text{b}_1 \\ \text{b}_2 \\ \text{b}_3 \\ \text{b}_4 \end{bmatrix} \quad \text{Or} \quad \hat{\mathbf{m}} = \mathbf{Ax} + \mathbf{b} \quad (4.1)$$

In Equation 4.1, the vector  $\mathbf{x}$  contains amounts of each Erk species (DP-Erk, TP-Erk, YP-Erk, and NP-Erk). The model parameters  $\text{a}_{11}$ - $\text{a}_{44}$  reflect the affinities of each antibody for each substrate. Note that the intercepts  $\text{b}_1$ - $\text{b}_4$  will likely be less than zero, so do not represent a physical characteristic of the system such as background signal alone. These parameter values may vary from one experiment to another, since they would depend upon the concentrations and activities of the antibodies and detection reagent. Within one experiment, however, once the parameters are known this same model can be used as a set of calibration curves to estimate the composition in other samples from a known set of measurements  $\mathbf{m}$ , by solving for a now unknown  $\hat{\mathbf{x}}$ .

Unfortunately, pure samples of DP-, TP-, YP-, and NP-Erk were not commercially available, nor could they be readily prepared, making parameter estimation and model verification more difficult. Upstate Biotechnology sells “active” Erk, which is phosphorylated and purified by a proprietary method. It is almost certainly not 100% DP-Erk, and likely contains some monophosphorylated (YP and/or TP) as well as NP. We can assume that it contains some unknown composition  $d$ ,  $t$ , and  $y$  representing the fraction of DP-, TP-, and YP-Erk, respectively. (The fraction of NP-Erk would necessarily equal  $1-d-t-y$ ). By reacting the “active” Erk with PTP1B, PP2A, or both the distribution between Erk forms changes as follows, assuming complete reaction:

$$\begin{array}{ll} \text{Original ErkPP (“D”):} & \mathbf{f}_D = \begin{bmatrix} \text{DP} & \text{TP} & \text{YP} & \text{NP} \\ d & t & y & 1-d-t-y \end{bmatrix}^T \\ \text{Add PTP1B (“T”):} & \mathbf{f}_T = \begin{bmatrix} 0 & d+t & 0 & 1-d-t \end{bmatrix}^T \\ \text{Add PP2A (“Y”):} & \mathbf{f}_Y = \begin{bmatrix} 0 & 0 & d+y & 1-d-y \end{bmatrix}^T \\ \text{Add both (“N”):} & \mathbf{f}_N = \begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix}^T \end{array} \quad (4.2)$$

ELISAs performed on the ErkPP in its original state and after reaction with the phosphatases should yield different results, since although the total amount of Erk remains the same, the amount of each form will change, as indicated in Equation 4.2. Combining the model for ELISAs with the fractional distributions gives a set of equations for the expected data:

$$aDP_D = a_{11} \cdot d \cdot x_T + b_1 \quad (4.3)$$

$$aTP_D = [a_{22} \cdot t + a_{24} \cdot (1-d-t)] \cdot x_T + b_2 \quad (4.4)$$

$$aYP_D = a_{33} \cdot y \cdot x_T + b_3 \quad (4.5)$$

$$aNP_D = [a_{42} \cdot t + a_{43} \cdot y + a_{44} \cdot (1-d-t-y)] \cdot x_T + b_4 \quad (4.6)$$

$$aDP_T = b_1 \quad (4.7)$$

$$aTP_T = [a_{22} \cdot (d+t) + a_{24} \cdot (1-d-t)] \cdot x_T + b_2 \quad (4.8)$$

$$aYP_T = b_3 \quad (4.9)$$

$$aNP_T = [a_{42} \cdot (d+t) + a_{44} \cdot (1-d-t)] \cdot x_T + b_4 \quad (4.10)$$

$$aDP_Y = b_1 \quad (4.11)$$

$$aTP_Y = a_{24} \cdot (1-d-y) \cdot x_T + b_2 \quad (4.12)$$

$$aYP_Y = a_{33} \cdot (d+y) \cdot x_T + b_3 \quad (4.13)$$

$$aNP_Y = [a_{43} \cdot (d+y) + a_{44} \cdot (1-d-y)] \cdot x_T + b_4 \quad (4.14)$$

$$aDP_N = b_1 \quad (4.15)$$

$$aTP_N = a_{24} \cdot x_T + b_2 \quad (4.16)$$

$$aYP_N = b_3 \quad (4.17)$$

$$aNP_N = a_{44} \cdot x_T + b_4 \quad (4.18)$$

This system contains 16 equations in 14 unknowns: the ELISA parameters  $a_{11}$ - $a_{44}$ ,  $b_1$ - $b_4$  as well as the composition  $d$ ,  $t$ , and  $y$ . Further, the system can be written for each of a set of values for total Erk  $x_T$ . The samples of original Erk (D) and following reaction with phosphatases were diluted to yield a total of 25, 20, 15, 10, 5, 2.5, 1, or 0 ng total Erk in each ELISA measurement. The constrained optimization function (“constr”) of MATLAB was then used to find the values for the 14 unknowns that minimize the least-squared error between model and measured values for each antibody. This can be written as the solution of the following optimization problem:

$$\min_{\mathbf{p}=[A,b,d,t,y]} \sum_i \sum_j \frac{(\hat{\mathbf{m}}_{i,j} - \mathbf{m}_{i,j})^2}{\sigma_{i,j}^2} \quad (4.19)$$

subject to:  $d + t + y \leq 1$ ,  $a_{22} \geq a_{24}$ ,  $a_{44} \geq a_{43}$ ,  $a_{44} \geq a_{42}$ ,

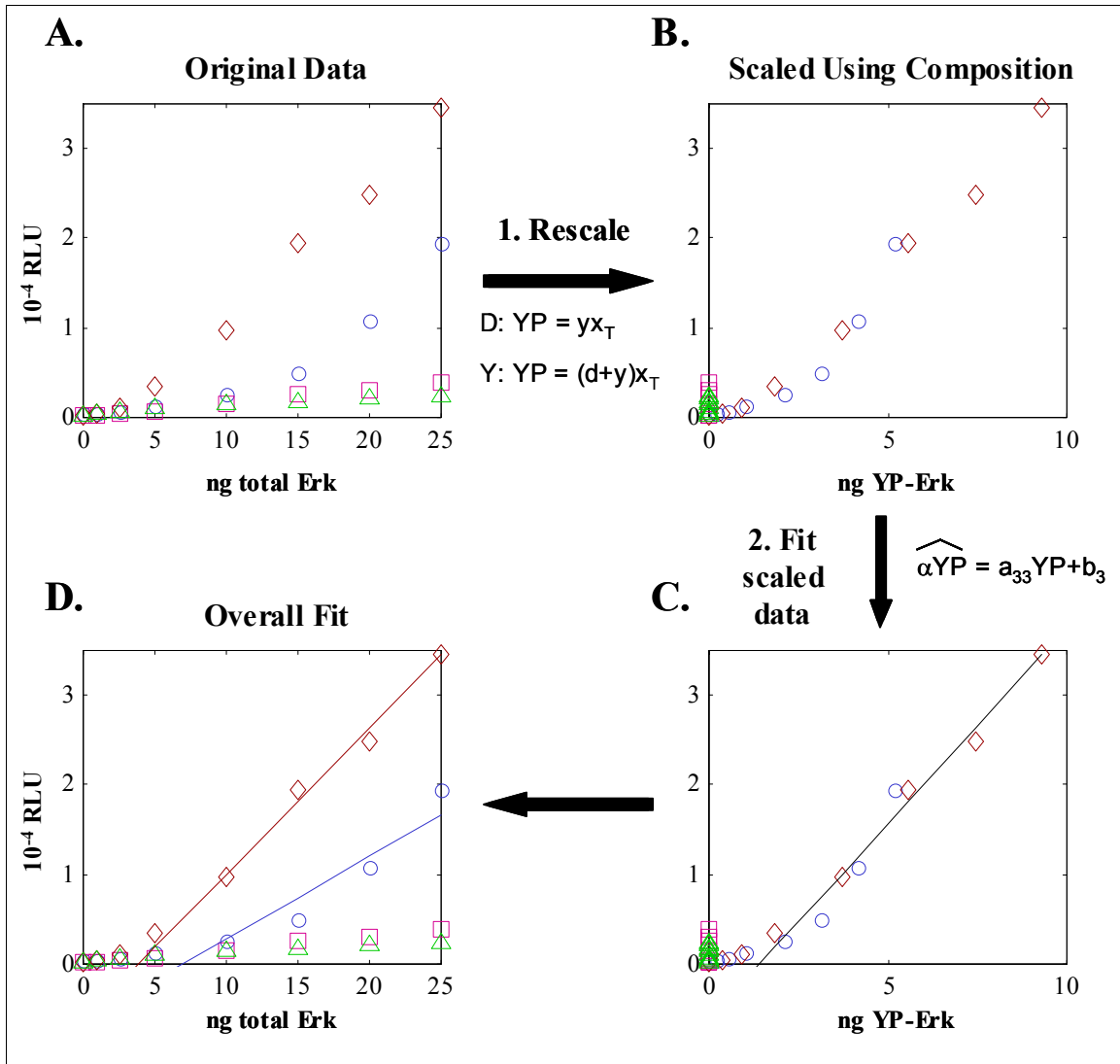
and  $\hat{\mathbf{m}}_{i,j} = \mathbf{A}\mathbf{f}_i \mathbf{x}_{T,j} + \mathbf{b}$  (Eq. 4.3-4.18 above)

where

- $i$  represents the sample  $i \in \{D, T, Y, N\}$ ,
- $j$  designates the particular total amount of Erk  $x_{T,j}$ ,
- $\mathbf{m}_{i,j}$  is the set of ELISA measurements (aDP, aTP, aYP, aNP) observed for sample  $i$  at total Erk  $x_{T,j}$ , and
- $\sigma_{i,j}^2$  is the variance for the measurements  $\mathbf{m}_{i,j}$  (not necessarily constant across  $j$ ).

One way to help understand the regression procedure is to consider the ELISA model as occurring in two interrelated stages. First, the compositional parameters  $d$ ,  $t$ , and  $y$  are selected so as to create the fractional distribution for each Erk species as shown in Equation 4.2. This step acts to rescale the data, written in terms of total Erk  $x_T$ , into amounts of each form DP, TP, YP and NP in the sample. If appropriate parameters are selected, then data for each antibody using all samples should compress together to form a single curve, representing the “real” response towards the Erk form found in all samples. A best-fit line passing through all the compressed data points gives the parameters  $a$  and  $b$ , as in Equation 4.1. These two steps must be performed together, since the scaling by composition can only be considered “appropriate” if they are able to yield a reasonable linear fit during the second stage.

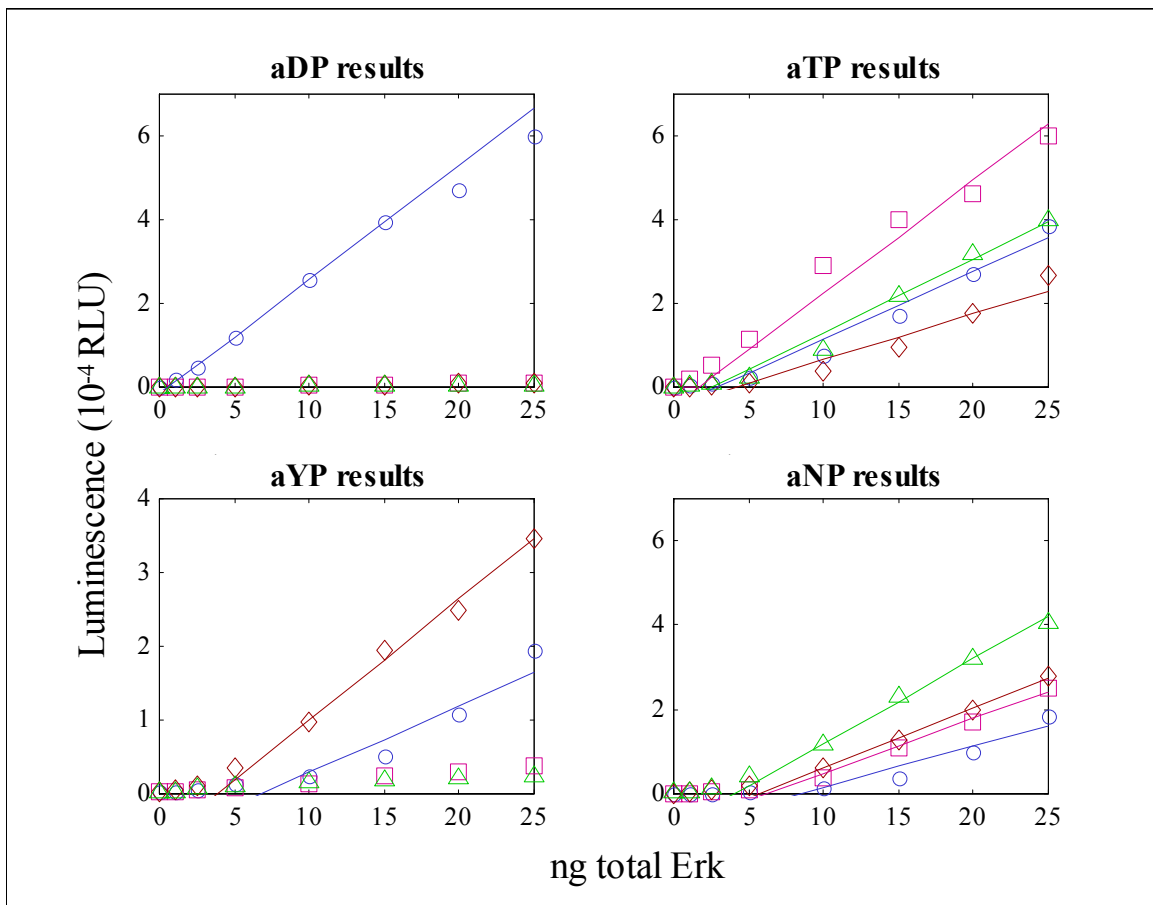
As an example, consider the response of the antibody aYP, which is assumed to be selective for YP-Erk. Sample data for different Erk samples using aYP is shown in Figure 4.9A. As expected from Equation 4.2, samples “D” (blue circles) and “Y” (orange triangles) contain YP-Erk, but at different amounts, so data shown in terms of total Erk appear as two curves. By selecting appropriate values for  $d$  and  $y$ , the two data sets can be collapsed together, as shown in Figure 4.9B. The best-fit line through the data (Figure 4.9C) then yields values for  $a_{33}$  and  $b_3$ , which can be used later to estimate YP-Erk in an unknown sample. Note that different values for  $d$  and  $y$  would not consolidate the data in the same way, and the best-fit line would show significantly more scatter than is observed in Figure 4.9B-C. The regression model can also be drawn using Equations 4.3-4.18 (i.e. with unscaled data), to compare the fits for individual samples, as shown in Figure 4.9D.



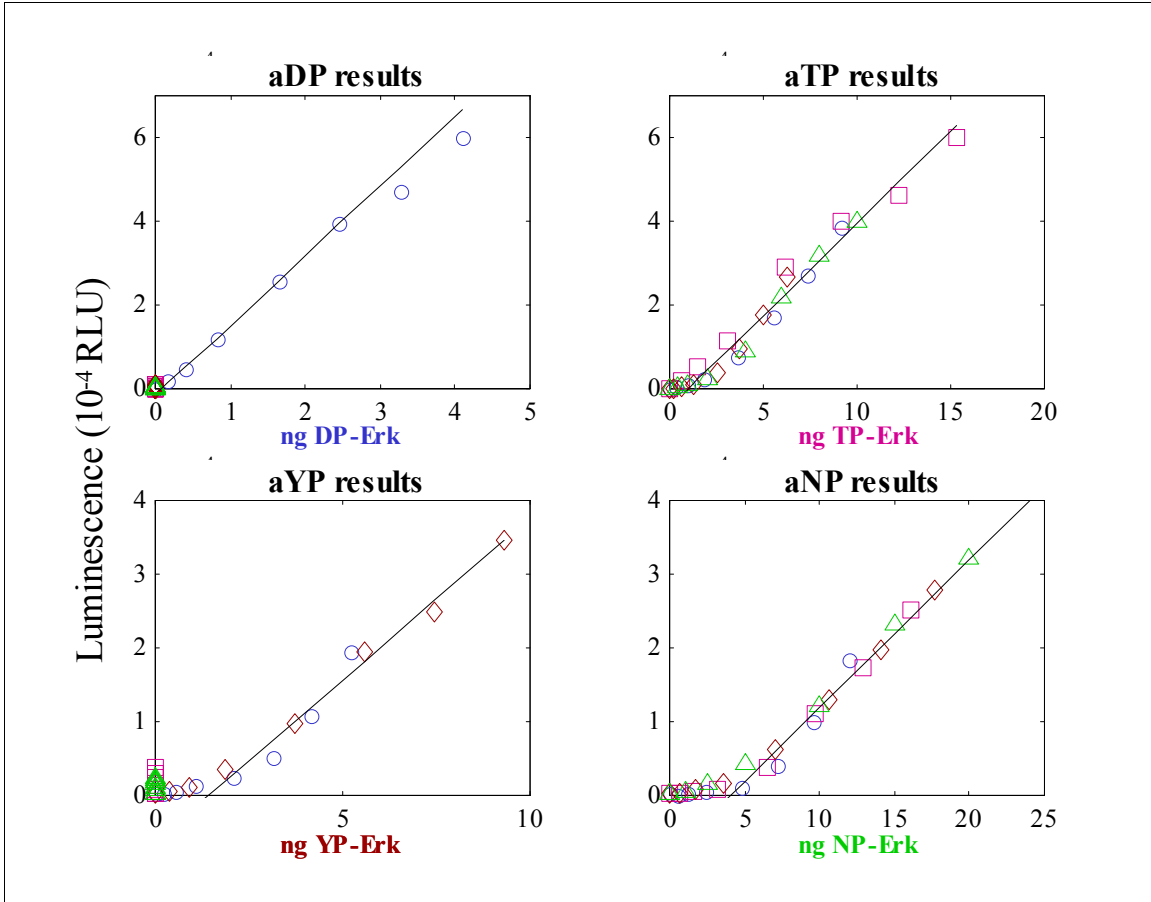
**Figure 4.9.** ELISA data modeling and regression as two steps: first, rescaling of data using composition parameters, second, linear fit of resulting rescaled data. Circles: “D” (original ErkPP), squares: “T” (ErkPP + PTP1B), diamonds: “Y” (ErkPP + PP2A), triangles: “N” (ErkPP + both phosphatases).

Also visible in the figure is curvature for low sample amounts, as described previously and shown schematically in Figure 4.8B. Above 1.5ng YP-Erk, the aYP antibody yields a linear response with increasing Erk, but below this value the signal-response curve is shallow. This helps to justify omission of particular data points as being outside the linear range and to provide an estimate for the lower limit of detection for the assay. When the model parameters are later used to invert a set of measurements to estimate the composition of an unknown sample, calculated values of YP-Erk below 1.5 ng should not be attributed much confidence.

An example set of results for ELISA measurements and regression for standards samples is shown in Figure 4.10-4.11. In Figure 4.10, the data is plotted using total Erk concentrations, whereas in Figure 4.11 the amounts are scaled for particular Erk forms as in Figure 4.9. Open symbols represent the observed values, while solid lines show best-fit model estimates following constrained optimization. This includes the effects of compositional parameters  $d$ ,  $t$ , and  $y$ , so the lines are realizations of Equations 4.3-4.18. In most cases model predictions show excellent agreement with measured values.



**Figure 4.10.** Sample results from ELISA standards regression. Symbols as in Figure 4.9. Solid lines represent model predictions using best-fit parameters.



**Figure 4.11.** Collection of data points from different samples and best-fit regression results. Data is the same as used in generation of Figure 4.10, with  $x_T$  rescaled using model parameters.

The ability to analyze standards, even with unknown initial sample composition, is not the ultimate goal of the ELISA model. Instead, the regression procedure is necessary to separate the parameters  $\mathbf{A}$  and  $\mathbf{b}$ , which correspond to the response of the ELISA measurements to a particular sample, from the standards composition, defined by parameters  $d$ ,  $t$ , and  $y$ . With the response parameters  $\mathbf{A}$  and  $\mathbf{b}$  known, the same model of Equation 4.1 may be used with ELISA measurements for a new test sample to estimate the Erk composition.

In this case, for each sample the model yields a linear system of four equations in four unknowns  $\hat{\mathbf{x}} = [\text{DP TP YP NP}]^T$ . There is one constraint that prevents simply inverting Equation 4.1 and choosing  $\hat{\mathbf{x}} = \mathbf{A}^{-1}(\mathbf{m} - \mathbf{b})$ , namely that the unknowns must all be greater than zero to be physically reasonable. Therefore, analysis of data samples is accomplished using the MATLAB function “nls”, to solve the optimization problem:



$$\begin{aligned} \min_{\mathbf{x}} \quad & \|\mathbf{Ax} - (\mathbf{m} - \mathbf{b})\|^2 \\ \text{subject to} \quad & \hat{\mathbf{x}} \geq 0. \end{aligned} \tag{4.20}$$

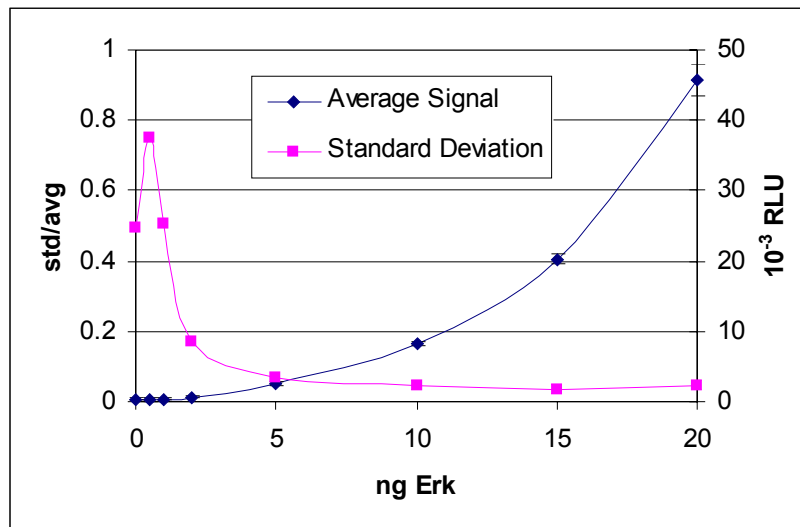
It is extremely important to be able to estimate the expected error in the model predictions for  $\hat{\mathbf{x}}$ , in particular because these estimates will be used later in calculation of activation ratios. Samples that upon analysis yield concentrations of Erk forms with high relative error should necessarily be considered suspect, as they may distort the shape of activation ratio plots. With an estimate for a reasonable confidence interval for the calculated Erk composition, particular trouble points may be identified and discarded if necessary. This estimation procedure is described in Section 4.3.3.

#### 4.3.2. Selection of data

In practice, there was a (small) positive slope seen in the data for aDP for samples T, Y and N, and aYP for T and N, as seen in Figures 4.9-4.11. This arose even in samples reacted with phosphatases for eight hours at 30°C. Either the reactions failed to reach completion, or there existed some nonspecific binding or low cross-reactivity of the antibodies for other Erk forms not incorporated into the model given in Equation 4.1. For the sake of simplicity it appeared easier to ignore this discrepancy, and admit some likely error in the parameter estimates. Therefore the optimization was performed neglecting this data, i.e. omitting the measurements and expressions 4.7, 4.9, 4.11, 4.15 and 4.17 above. This left at best 11 equations \* 8 samples = 88 expressions for 14 unknowns.

In fact, not all measurement points were used, since some appeared to fall outside the linear range for a particular antibody, as shown in Figure 4.9B. Generally this linear range was determined visually, but as a matter of course any signal less than 2000 RLU was discarded. This was justified by observations that often the standard deviation of ELISA measurements was greater than 10% of the mean value for signals less than 2000 RLU, as shown in Figure 4.12. In some cases, measurements above 45,000 RLU appeared to be saturating the detector or assay reagents and were again ignored. Typically, this process excluded one-third of the data points, often in the low signal range (for no or low values of total Erk). In any case more than 30 data points remained, leaving more than enough redundancy to solve the system for the unknown parameters. The omission of data points from consideration in the objective function (Equation 4.19)

was achieved by setting the variance  $\sigma_{i,j}^2$  equal to infinity, or equivalently, its reciprocal equal to zero. The variance was observed experimentally to increase with the magnitude of the measurement signal, but decrease as a proportion of the signal. To avoid bias towards samples with lower absolute variance (and lower signal), the standard deviation was set equal for all remaining points at a value of 2200 RLU. This was the largest value observed during repeated trials with aNP on different samples, used in generation of Figure 4.12.



**Figure 4.12.** Sample ELISA data, showing dependence of experimental variance on absolute signal. Diamonds: average signal (triplicate measurements). Squares: ratio of standard deviation to average value at each point.

Another justification for excluding particular data points arises through consistency analysis for gross errors in measurements, when redundant measurements are available [39, 40]. The objective function in Equation 4.19, representing the total weighted least-squares error at the optimum solution, is also known as the consistency index,  $h$ . Since the consistency index takes the form of a sum of squared values of random variables (the measured ELISA signals  $m$ ), it is expected to follow a  $\chi^2$ -distribution with the same number of degrees of freedom,  $f$ , as redundant measurements (in this case, total number of data points remaining,  $n$ , minus 14 parameters,  $p$ ). At a specified probability  $\alpha$ , if  $h \geq \chi^2_{1-\alpha}(f)$  then it is likely that there are gross errors in the measurements, since the variance from the model cannot be explained by variance in the measurements alone. Individual measurement points may therefore be removed if their

individual contribution,  $h_i$ , to the consistency index is significantly large. Of course, following removal of one measurement point, the regression needs to be redone to determine new best-fit model parameters, and a new consistency index calculated. This process is repeated until the consistency index is less than the  $\chi^2$ -value for the degrees of freedom remaining. For example, during regression for the data shown in Figures 4.10-4.11, 38 data points were omitted, leaving 50 points ( $n$ ) and therefore 36 degrees of freedom ( $f$ ). The consistency index following regression was 45.7, less than  $\chi^2(0.9,36)=47.2$ , so it is unlikely that the remaining set contains gross inconsistencies. Thus the linear model, accounting for measurement error, does have sufficient support to be considered statistically reasonable.

#### 4.3.2. Error analysis and model validation

The calculation of an Erk composition using the ELISA model, through solution of the optimization problem in Equation 4.20, is an example of inverse regression. In other words, an initial set of known composition  $\mathbf{x}$  was used to generate model estimates  $\hat{\mathbf{m}}$ , and so the model involves regression of  $\mathbf{m}$  onto  $\mathbf{x}$ . In Equation 4.20 the opposite is occurring, and  $\mathbf{x}$  is being regressed onto a known measurement  $\mathbf{m}$ . For a simple one-to-one linear model between  $m$  and  $x$  ( $m = ax + b$ ), the confidence intervals for inverse regression are relatively straightforward to determine [41]. Essentially, the model is solved to estimate some  $\hat{x}_0$  at a known  $m_0$ . The confidence intervals for  $\hat{x}_0$  are determined by solution at the bounds of the confidence interval for  $m_0$ , which is obtained from estimates in the error in the parameters  $a$  and  $b$ .

For the current arrangement, however, error analysis is complicated by the constraints provided in the inverse regression (Equation 4.20) as well as the inherent nonlinearity in the original regression, caused by cross-multiplication of parameters (Equation 4.19). For such a problem, no simple closed-form relationship exists between the standard error in the parameters,  $se(\mathbf{p})$ , and the confidence band for  $\hat{\mathbf{x}}$ . Therefore one way to estimate a reasonable range for  $\hat{\mathbf{x}}$  is to obtain a confidence interval for the parameters  $\mathbf{p}$ , then solve Equation 4.20 at the bounds of the parameter values. What remains is an estimation of  $se(\mathbf{p})$  for the nonlinear regression problem, by combining a sensitivity analysis of the model with an estimate for the data variance  $\sigma^2$ .

The optimization function in Equation 4.19 is the weighted sum-of-squares error  $S(\mathbf{p})$ , which can also be written as a matrix product:

$$S(\mathbf{p}) = \sum_i \sum_j \frac{(\hat{\mathbf{m}}_{i,j} - \mathbf{m}_{i,j})^2}{\sigma_{i,j}^2} = (\hat{\mathbf{m}} - \mathbf{m})^T \mathbf{V}^{-1} (\hat{\mathbf{m}} - \mathbf{m}) \quad (4.21)$$

where

$\hat{\mathbf{m}}$  is a (128x1) matrix of estimates calculated using Equations 4.3-4.18 at each  $x_T$ ,

$\mathbf{m}$  is a (128x1) matrix of corresponding observed measurements, and

$\mathbf{V}^{-1}$  is the inverse of the variance-covariance between the measurements. Assuming independence in the measurements then  $\mathbf{V}^{-1}$  is a (128x128) matrix with  $1/\sigma_{i,j}^2$  along the diagonal, and zeros elsewhere.

To investigate the variation of the parameters, embedded within the model  $\hat{\mathbf{m}}$ , a Taylor expansion is performed around a particular solution  $\mathbf{p}^0$  [42, 43]:

$$\hat{\mathbf{m}}_i(\mathbf{p}) \approx \hat{\mathbf{m}}_i(\mathbf{p}^0) + \sum_j \left. \frac{\partial \hat{\mathbf{m}}_i(\mathbf{p})}{\partial p_j} \right|_{\mathbf{p}^0} (p_j - p_j^0) \quad (4.22)$$

Or equivalently,

$$\hat{\mathbf{m}}(\mathbf{p}) \approx \hat{\mathbf{m}}(\mathbf{p}^0) + \mathbf{G}(\mathbf{p} - \mathbf{p}^0) \quad (4.23)$$

The sensitivity matrix  $\mathbf{G}$  in Equation 4.23 reflects how each parameter  $p_j$  may influence the estimate for a particular measurement  $m_i$ . The contents  $G_{ij}$  can be readily determined from partial differentiation of Equations 4.3-4.18 as appropriate. The results of differentiation are shown in Appendix 5. Note that the total amount of  $\text{Erk } x_T$  appears in Equations 4.3-4.18 (and thus the entries  $G_{ij}$ ), and therefore will appear as part of  $\mathbf{G}$ . Thus the true matrix  $\mathbf{G}$  will contain 16 (expressions) \* 8 (concentrations  $x_T$ ) or 128 rows, and 14 (parameter) columns.

Combining Equations 4.21 and 4.23 leads to the following approximation:

$$S(\mathbf{p}) \approx \left[ \mathbf{m} - \hat{\mathbf{m}}(\mathbf{p}^0) - \mathbf{G}(\mathbf{p} - \mathbf{p}^0) \right]^T \mathbf{V}^{-1} \left[ \mathbf{m} - \hat{\mathbf{m}}(\mathbf{p}^0) - \mathbf{G}(\mathbf{p} - \mathbf{p}^0) \right] \quad (4.24)$$

Equation 4.24 is of the same form as the classic linear regression problem, and by analogy is minimized when

$$(\mathbf{p} - \mathbf{p}^0) = [\mathbf{G}^T \mathbf{V}^{-1} \mathbf{G}]^{-1} \mathbf{G}^T \mathbf{V}^{-1} [\mathbf{m} - \hat{\mathbf{m}}(\mathbf{p}^0)] \quad (4.25)$$

By analogy with the linear problem, it can be shown that around the solution  $\mathbf{p}^0$  the parameters approximately follow a normal distribution, with variance  $[\mathbf{G}^T \mathbf{V}^{-1} \mathbf{G}]^{-1}$ , i.e. that  $(\mathbf{p} - \mathbf{p}^0)$ , and therefore  $\mathbf{p} \sim N(\mathbf{p}^0, \sigma^2 [\mathbf{G}^T \mathbf{V}^{-1} \mathbf{G}]^{-1})$  [43]. By the properties of the normal distribution, multiplication of the  $\mathbf{p}$  by a factor  $\mathbf{a}$  yields  $\mathbf{a}^T \mathbf{p} \sim N(\mathbf{a}^T \mathbf{p}^0, \sigma^2 \mathbf{a}^T [\mathbf{G}^T \mathbf{V}^{-1} \mathbf{G}]^{-1} \mathbf{a})$ , and that:

$$\frac{\mathbf{a}^T \mathbf{p} - \mathbf{a}^T \mathbf{p}^0}{s \sqrt{\mathbf{a}^T [\mathbf{G}^T \mathbf{V}^{-1} \mathbf{G}]^{-1} \mathbf{a}}} \sim t_f \quad (4.26)$$

Recall that  $f$  is the number of degrees of freedom of the system, while  $s$  is the estimate of the standard deviation of the system, calculated from  $s^2 = S(\mathbf{p})/f$ . Rearranging Equation 4.26, an approximate  $100(1-\alpha)\%$  confidence interval for the parameters around the solution can be written as

$$\mathbf{a}^T \mathbf{p}^0 \pm t_{f, \alpha/2} s \sqrt{\mathbf{a}^T [\mathbf{G}^T \mathbf{V}^{-1} \mathbf{G}]^{-1} \mathbf{a}} \quad (4.27)$$

The confidence interval for a particular parameter  $p_i$  is found by setting  $\mathbf{a}(i) = 1$  and all other elements 0, i.e.  $\mathbf{a} = [0 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]^T$ . In that case the vector product in Equation 4.27 will leave only the  $i^{\text{th}}$  diagonal element of the matrix  $[\mathbf{G}^T \mathbf{V}^{-1} \mathbf{G}]^{-1}$ , denoted  $\{[\mathbf{G}^T \mathbf{V}^{-1} \mathbf{G}]^{-1}\}_{ii}$ . Therefore the confidence interval for  $p_i$  is:

$$p_i^0 \pm t_{f, \alpha/2} s \sqrt{\{[\mathbf{G}^T \mathbf{V}^{-1} \mathbf{G}]^{-1}\}_{ii}} \quad (4.28)$$

An example set of parameter values and confidence intervals, for the data shown in Figures 4.9-4.11, is shown in Table 4.3. The aDP antibody appears to show the highest affinity, based upon the high value of  $a_{11}$ . This should not be particularly surprising given that recognition of DP-Erk, containing two phosphorylated residues in close proximity, would necessarily involve several strong charge-based interactions with an antibody. Other Erk forms, lacking such charged epitopes, would interact with an antibody using more hydrophobic or hydrogen-bonding interactions, which are relatively

weak by comparison. The relative values of  $a_{22}$  and  $a_{24}$  indicate an approximately 3-fold higher affinity of aTP for TP-Erk over NP-Erk, somewhat lower than the 10-fold difference expected from previous in vitro characterizations [12]. However, these previous studies were based on competitive ELISAs using short peptide sequences, not parallel ELISAs using entire proteins. It may be that conformational changes in protein structure due to phosphorylation state may change the availability of some residues to interact with the antibodies, and thus alter the affinities for different forms of Erk. On the other hand, aNP was expected to cross-react weakly with TP- and YP-Erk, but the results from this study indicate no statistically significant affinity. Here, the charge on tyrosine or threonine due to phosphorylation could impair interaction with aNP, which may recognize the phosphorylation site through hydrogen bonding to the unmodified residues.

**Table 4.3.** Parameter values from ELISA standards regression (data shown in Figures 4.9-4.11)

<b>Parameter</b>	<b>Value (95% CI)</b>	<b>Parameter</b>	<b>Value</b>
$a_{11}$	$16700 \pm 4800$	$b_1$	$-1500 \pm 3100$
$a_{22}$	$4400 \pm 1200$	$b_2$	$-4600 \pm 2100$
$a_{24}$	$1760 \pm 160$	$b_3$	$-6300 \pm 3700$
$a_{33}$	$4390 \pm 930$	$b_4$	$-8000 \pm 3100$
$a_{42}$	$0 \pm 1000$	$d$	$0.164 \pm 0.042$
$a_{43}$	$420 \pm 500$	$t$	$0.191 \pm 0.117$
$a_{44}$	$2000 \pm 200$	$y$	$0.208 \pm 0.049$

In many cases the uncertainty in a parameter can be as high as 30% of the optimized value, which will induce a high uncertainty later when estimating Erk composition in unknown samples. Unfortunately, this error appears to be a characteristic of the system, likely a combination between variability in the measurements themselves and nonlinearity in the response. This can be observed as scatter between the model predictions (lines) and observations (open symbols) in Figures 4.10-4.11. For example, aTP exhibits significant scatter, particularly at 10 ng total Erk for all samples, and original ErkPP (D, circles) shows curvature for aYP. Improvements in measurement accuracy, including better estimates of how the variance for each measurement varies with  $x_T$ , would certainly help. However, a much more powerful benefit would be realized by preparation of well-defined standards, to avoid assumptions necessary about the composition of samples D, T, Y, and N.

It is also interesting to note the estimates for the ErkPP composition,  $d$ ,  $t$ , and  $y$ . The regression results predict that at most 20% of the “active” sample is in fact DP-Erk, and another 40-50% may be monophosphorylated on threonine or tyrosine. The monophosphorylated species may arise during the proprietary procedure that Upstate Biotechnology uses to activate Erk, or may be an artifact of sample handling. Separate experiments using MALDI and LC-MS of ErkPP after digestion with trypsin yielded an estimate of approximately 50% phosphorylated Erk (DP-, TP-, and YP-Erk) with monophosphorylated species certainly present (data not shown). Unfortunately, it was not possible to separately resolve between TP- and YP-Erk using MS techniques, or to quantify the relative amounts of DP- and monophosphorylated Erk as a separate method to validate the regression estimates.

Recall that ultimately the ELISA model is used to estimate  $\hat{\mathbf{x}}$  from a set of measurements  $\mathbf{m}$  for an unknown sample. The error in  $\hat{\mathbf{x}}$  is approximated by solving the optimization problem (Equation 4.20) with parameters  $\mathbf{p}$  at the bounds of the confidence interval for each  $p_i$ , calculated from Equation 4.28. Therefore, an alternative approach to investigating the accuracy of the model parameters was to perform a validation experiment using mixtures of the standards samples. The predicted Erk composition could be determined from the compositions of the standards themselves, calculated from the parameters  $d$ ,  $t$ , and  $y$  in accordance with Equation 4.2. On the other hand, ELISA measurements of these same mixtures could be analyzed to give a model-predicted composition, by solving the regression problem in Equation 4.20.

The results of the validation experiment are shown in Table 4.4. Two standards samples (D, T, Y, or N) were combined together in equal parts (10 ng total Erk each) and applied to ELISA wells. In general, very good agreement is seen between the predictions from the original sample composition (“Pred”) and estimates based upon ELISA measurements (“Meas”) for DP-, YP-, and NP-Erk forms, although significant ranges can be seen for the estimates. However, the measured estimates for TP-Erk appear consistently lower than the predicted values based on the original ErkPP composition, and the confidence interval for TP-Erk does not contain the predicted value. Most likely, this is the result of an overestimation for  $a_{24}$ , perhaps because the N sample (assumed to

be 100% NP-Erk) actually contains some remaining TP-Erk. Analysis of the measurement system may attribute too much of the aTP signal to NP- and not TP-Erk, lowering the overall estimation for TP-Erk. Recall that previous experiments appear to indicate that in general, very little TP-Erk is formed during Erk phosphorylation in cycles (e.g. Figure 4.4 and [13, 14]). Therefore, an error in the estimation for TP-Erk may not significantly alter the analysis of a sample from Erk cycles, which would be dominated by DP-, YP-, and NP-Erk forms. Without a more characterized set of standards (ideally, pure in each Erk form), it is difficult to separate the potential sources of error associated with inaccuracies in sample composition as opposed to those related to describing the response of the antibodies.

**Table 4.4.** Comparison of predicted (from standards composition) and observed (estimated from measurements) Erk amounts using mixtures of standards (D, T, Y, N) samples. 10 ng of each standard (20 ng total) was added to wells, and ELISA measurements used with Equation 4.20 to calculate  $x$ . Predictions were calculated according to Equation 4.2. All entries are shown in ng Erk of each form.

	<b>DP-Erk Estimates</b>				<b>TP-Erk Estimates</b>			
	Pred.	Meas	Low	High	Pred.	Meas	Low	High
<b>D+N</b>	4.02	4.56	3.74	5.63	2.17	1.18	1.13	1.15
<b>D+T</b>	4.02	4.74	3.90	5.84	8.36	6.36	5.72	7.12
<b>T+Y</b>	0	0.10	0	0.49	6.19	5.22	4.82	5.63
<b>Y+N</b>	0	0.10	0	0.48	0	0	0	0.17
<b>N+T</b>	0	0	0	0.34	6.19	3.39	3.72	3.09
<b>D+Y</b>	4.02	4.29	3.51	5.32	2.17	1.42	1.28	1.43

	<b>YP-Erk Estimates</b>				<b>NP-Erk Estimates</b>			
	Pred.	Meas	Low	High	Pred.	Meas	Low	High
<b>D+N</b>	1.20	1.88	1.19	2.78	12.61	11.85	8.95	15.70
<b>D+T</b>	1.20	1.17	0.57	1.96	6.42	3.59	0.83	7.54
<b>T+Y</b>	5.22	5.72	4.57	7.21	8.59	5.51	2.19	10.23
<b>Y+N</b>	5.22	5.85	4.69	7.38	14.78	11.68	8.57	15.16
<b>N+T</b>	0	0.45	0	1.13	13.81	13.19	9.98	17.48
<b>D+Y</b>	6.42	5.67	4.53	7.16	7.39	2.65	0.18	6.12

It has already been noted that the current model fails to account for signal observed using aDP for the samples T, Y, and N, and aYP for samples T and N, which may arise from nonspecific antibody cross-reaction or from incomplete dephosphorylation by phosphatases. Additional modifications to the model, to account for these observations, and possibly to include the nonlinear regions as well, could perhaps provide a more successful representation in validation experiments. However, it



was felt that this introduction would complicate the sample analysis (Equation 4.20) while introducing additional parameters that might not be determined with any higher accuracy than in the current model. Thus the model, however imperfect, was maintained in its current form for analysis of immobilized and liquid-phase cycles.

#### **4.4. Activation Ratios in Erk Phosphorylation Cycles**

As described in Chapter 2, the activation ratios for a system should be dependent upon kinetic parameters as well as concentrations of each component, in particular when enzyme-substrate complexes cannot be separated from free species. Several sets of experiments were performed to investigate whether Erk activation profiles and activation ratios were consistent with predicted patterns. Liquid-phase and immobilized enzyme reactions were compared to examine the linear and nonlinear behavior of activation ratios under different conditions. Variation of the total Erk and phosphatase concentrations allowed for investigation on how changes in these parameters would be reflected in the activation ratios. Unfortunately, there was no method available for making modifications to the catalytic parameters for the individual enzymes.

Results for these experiments are discussed in the sections below and shown in Figures 4.13-4.16, where both activation profiles (normalized by total Erk) and all possible activation ratios are included. In most cases, the primary monophosphorylated intermediate was YP-Erk, with estimates for TP-Erk within the error of the measurements. Therefore although activation ratios based on TP-Erk are calculated and shown in Figures 4.13-4.16, caution must be taken in interpreting the results. Confidence intervals for Erk estimates and activation ratios are omitted for the sake of clarity, as in most cases excessive overlap between different series would cause confounding. Hence there was little statistically significant differences observed between series, but observable variations in patterns did appear and are noteworthy.

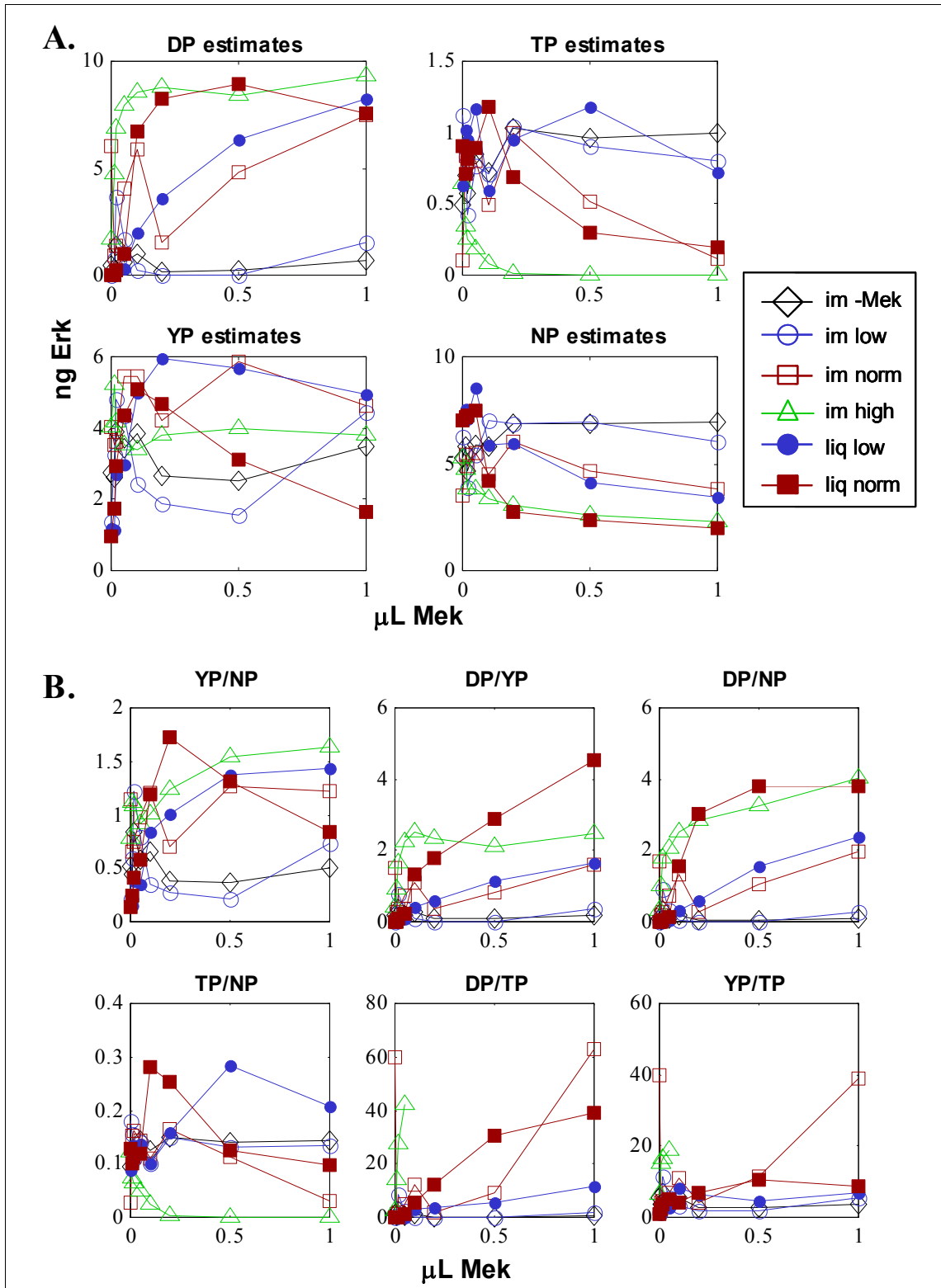
##### **4.4.1. Comparison of liquid-phase and immobilized reactions**

A critical design feature of the immobilized reaction system was to allow resolution of freely soluble substrate from that bound by enzyme. As described in Section 2.4, the activation ratio for a system may appear linear when calculated using free substrate concentrations, but appear hyperbolic for a system if enzyme-substrate

complexes are included in calculations. A test both for the immobilized reaction system and activation ratio analysis itself was to experimentally observe this prediction. Activation ratios calculated for samples taken from immobilized reactions should appear linear, since enzyme-bound substrate would remain within the wells of the reaction plate. Samples from liquid-phase reactions would contain enzyme and substrate, and would therefore be expected to yield hyperbolic activation ratios, provided that sufficient enzyme is added to the system. At low enzyme concentrations, activation ratios should appear linear regardless of which system is utilized.

Comparison of results between liquid-phase and immobilized enzyme systems are shown in Figure 4.13, where the amount of N4-Mek applied to the wells ranges from 0-0.1  $\mu\text{L}$  (“low”, circles), 0-1  $\mu\text{L}$  (“normal”, squares), or 0-10  $\mu\text{L}$  (“high”, triangles). Liquid-phase reactions are indicated by closed symbols, and immobilized enzyme reactions by open symbols. The experiment was conducted in this manner to account for differences in enzymatic activity between the two systems. Recall from Section 4.2.3 that immobilized antibodies were observed to retain approximately 5-10% of the enzyme during the binding step. Thus direct comparisons should be made between liquid-phase reactions and immobilized reactions with 10-fold higher N4-Mek concentrations (i.e. red squares vs. green circles, or blue squares vs. red circles). In general, the results for “low” concentrations of N4-Mek ( $< 0.1 \mu\text{L}$ ) were indistinguishable from those containing no enzyme at all.

However, for “normal” and “high” amounts of N4-Mek added, liquid reactions yield very similar results to corresponding immobilized reactions, both in terms of activation profiles and activation ratios. At higher N4-Mek concentration ranges, formation of DP-Erk saturates quickly, whereas a peak is observed for YP-Erk that disappears more quickly for the liquid-phase reaction. Lower ranges of N4-Mek show wider peaks for YP-Erk formation, as well as more shallow curves for appearance of DP-Erk. Together, this suggests that the enzymatic system continues to operate essentially unchanged in the immobilized phase, except for an overall lower activity.



**Figure 4.13.** A) Activation profiles and B) Activation Ratios for Erk cycles vs. normalized volume of N4-Mek added, in liquid phase (closed symbols) or immobilized using capture antibodies (open symbols). “Low” (circles) signifies a maximum of 0.1  $\mu\text{L}$  N4-Mek added, “Norm” (squares): 1  $\mu\text{L}$ , “High” (triangles): 10  $\mu\text{L}$  Mek.

In general, it should be slightly surprising that the Erk distribution from immobilized enzymes appears similar to freely soluble enzymes. Both systems use the same enzymes in similar proportions, and all other parameters are kept constant. However, the amount of substrate bound to different enzymes changes with each Erk form, and thus the amount remaining in free solution should also change. Furthermore, a small decrease is expected in recovery of all Erk forms, but as shown in Figure 4.13A samples from immobilized reactions contain the same amount of each Erk species as those arising from liquid-phase reactions. Separate ELISA data using an antibody that reacts with all Erk forms also indicated that Erk recovery from the two systems was equal (data not shown). Therefore, the activation profile data appears to suggest that enzyme-substrate complexes are not being separated from free substrate in immobilized reactions.

This conclusion is made more apparent by the hyperbolic shape for activation ratios for high concentrations of immobilized N4-Mek. As described in Chapter 2, activation ratios appear hyperbolic if total concentrations are utilized and the system is operating under saturating conditions. If the enzymes are in fact far from saturation, then activation ratios should appear linear regardless of whether total or free substrate concentrations are measured. Since activation ratios for liquid-phase reactions are hyperbolic, the system must indeed be operating with saturated enzymes. That immobilized reactions follow the same behavior indicates that the immobilized system is not operating as intended in separating enzyme-bound substrate.

The most obvious explanation for this observation is that the enzyme-substrate complexes are in fact not stable under current conditions. This conclusion would be consistent with prior difficulties in separating protein-protein complexes by filtration or electrophoresis. Interactions between Mek and Erk have been observed experimentally by co-immunoprecipitation [44]. However, these studies may have been performed under experimental conditions more conducive to protein-protein complex formation than for enzymatic reaction, and perhaps more importantly, within a cellular context where a scaffolding protein like MP-1 could assist in bringing the enzyme and substrate together [45, 46]. Whereas immunoprecipitation is performed using high salt concentrations that help bring proteins together, the reaction buffer for this system had no salt. The presence

of the reducing agent DTT in the reaction buffer, necessary for maintaining PTP1B activity, may have had adverse affects upon the Mek-Erk interaction as well.

Although less likely, it may also be that the enzyme-substrate complexes account for only a small fraction of the total substrate pool. The fraction of substrate bound by enzyme is dependent upon the ratio between the concentrations of enzyme the Michaelis constant defined for their interaction ( $[E \cdot A]/A_T \sim E_T/K_m$ ). Thus, for complexes to be negligibly small it would require that the Michaelis constants for each enzymatic step be significantly greater than the enzyme concentrations. The Michaelis constants for dephosphorylation of DP-, YP- and TP-Erk by PTP1B and PP2A could not be determined during *in vitro* reactions, but were estimated to be much greater than 1  $\mu$ M [25]. On the other hand, the interaction between Mek and Erk is more specific, with the  $K_m$  estimated at approximately 300-460 nM for NP-Erk phosphorylation and 30-50 nM for the second phosphorylation step [13, 47]. Therefore, some fraction of Erk (in particular, of YP-Erk, the likely monophosphorylated intermediate) should be abstracted by interaction with Mek if the complexes are stable. Since no change in total Erk recovery was observed in immobilized reactions, it appears more reasonable to conclude that enzyme-substrate complexes were unstable using the immobilized reaction system.

If the reaction is assumed to proceed only through a YP-Erk intermediate (thus, only one path) then the system is identical to that described in Section 2.3.4. Unfortunately, the failure to separate free species prevents the calculation of individual activation factors for each step or even verification of the linear responses expected for each intermediate step (YP/NP and DP/YP) and quadratic response expected for the overall dual activation ratio (DP/NP). Instead, the data shown in Figure 4.13B represents total activation ratios, as described in Section 2.4. After accounting for enzyme-substrate complexes in balances for total enzyme and substrate concentrations, the free activation ratios in Equations 2.35-2.36 can be converted into expressions for the total activation ratios as follows:

$$\begin{aligned} \text{TAR}_{A1} &\equiv \frac{A_T^*}{A_T} = \frac{A^* (1 + E_2/K_{m2} + E_1/K_{m3})}{A (1 + E_1/K_{m1})} \\ &= \frac{k_1 E_{1T}}{k_2 E_{2T}} \frac{K_{m2} + A^* + E_{2T} + E_{1T} \left( \frac{K_{m2} + A^*}{K_{m3} + \frac{K_{m3}}{K_{m1}} A + A^*} \right)}{K_{m1} + A + \frac{K_{m1}}{K_{m3}} A^* + E_{1T}} \end{aligned} \quad (4.29)$$

$$\begin{aligned} \text{TAR}_{A2} &\equiv \frac{A_T^{**}}{A_T^*} = \frac{A^{**} (1 + E_4/K_{m4})}{A^* (1 + E_2/K_{m2} + E_1/K_{m3})} \\ &= \frac{k_3 E_{1T}}{k_4 E_{4T}} \frac{K_{m4} + A^{**} + E_{4T}}{\left( 1 + \frac{E_{2T}}{K_{m2} + A^*} \right) \left( K_{m3} + \frac{K_{m3}}{K_{m1}} A + A^* \right) + E_{1T}} \end{aligned} \quad (4.30)$$

For the Erk experimental system,  $E_1$ ,  $E_2$  and  $E_4$  correspond to N4-Mek, PTP1B, and PP2A respectively, while  $A$ ,  $A^*$ , and  $A^{**}$  would correspond to the free amounts of NP-, YP-, and DP-Erk (a closed-form expression for TAR using only total substrate concentrations such as  $A_T^*$  is not feasible). Assuming that the phosphatases operate far from saturation, as described above, then  $K_{m2}$  and  $K_{m4}$  will dominate adjacent terms in Equations 4.29-4.30, which thus can be simplified to yield:

$$\text{TAR}_{A1} = \frac{k_1 E_{1T}}{k_2 E_{2T}} \frac{K_{m2} \left( K_{m3} + \frac{K_{m3}}{K_{m1}} A + A^* + E_{1T} \right)}{\left( K_{m1} + A + \frac{K_{m1}}{K_{m3}} A^* + E_{1T} \right) \left( K_{m3} + \frac{K_{m3}}{K_{m1}} A + A^* \right)} \quad (4.31)$$

$$\text{TAR}_{A2} = \frac{k_3 E_{1T}}{k_4 E_{4T}} \frac{K_{m4}}{K_{m3} + \frac{K_{m3}}{K_{m1}} A + A^* + E_{1T}} \quad (4.32)$$

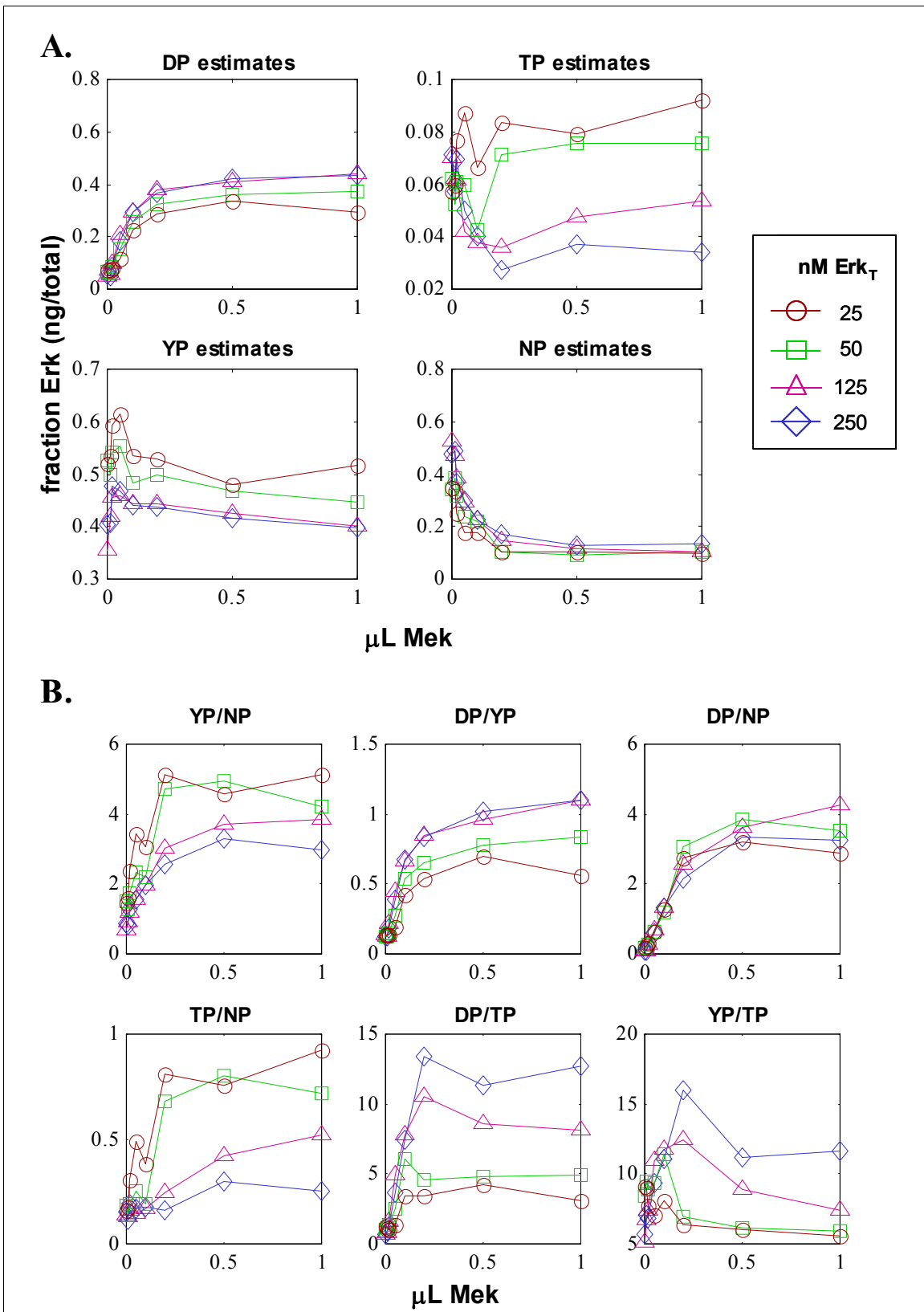
Since enzyme and substrate concentrations appear in Equations 4.29-4.32, the magnitude of total activation ratios does not reflect the kinetic parameters alone. Nevertheless, the fact that DP/YP activation ratios appear greater than YP/NP ratios seems to suggest that the second phosphorylation step is overall faster than the formation of YP-Erk. This may be due to differences in the catalytic constants and phosphatase concentrations ( $k_3/k_4 E_{4T} > k_1/k_2 E_{2T}$ ), but perhaps also because of differences in the Michaelis constants. With  $K_{m1} > K_{m3}$  and all other parameters equal it is likely that  $\text{TAR}_{A2} > \text{TAR}_{A1}$ , or in terms of Erk activation,  $\text{DP/YP} > \text{YP/NP}$ .

Without better estimates for concentrations of the enzymes, it is difficult to resolve between these possibilities. Both cases appear to require that the second phosphorylation step proceed more quickly than the first (in terms of higher net  $k_{\text{cat}}$  or lower  $K_{\text{m}}$ ). Should the first step be faster, an accumulation of YP-Erk would be observed that would lead to an overall increase in YP/NP activation ratios, and a decrease in the magnitude of DP/YP across all Mek concentrations. With the first step limiting overall reaction, then the concentration of intermediate YP-Erk would be expected to decrease relative to the other forms, and DP/YP activation ratios would increase beyond YP/NP.

#### 4.4.2. Variation of total Erk concentration

This issue was investigated further through modulation of the amount of Erk substrate. The experiments described in Section 4.4.1 were performed using a total Erk concentration of 125 nM, in between the approximate Michaelis constants for Mek,  $K_{\text{m}1}$  and  $K_{\text{m}3}$ , but less than those for the phosphatases. By varying the amount of Erk involved in liquid-phase reactions from 25-250 nM, the extent of saturation for each phosphorylation step should also change. According to Equations 4.29-4.32, the amount of Erk in each form influences the total activation ratios, and of particular interest is the distribution between YP-Erk and NP-Erk ( $A^*$  and  $A$ , respectively). Results for this case are shown in Figure 4.14. Note that the activation profiles in Figure 4.14A are normalized by the total Erk estimated for each sample, thus demonstrating how the fractional distribution of Erk forms varies with total Erk concentration.

Interestingly, the fractions of NP and DP increase with total Erk, while intermediate YP decreases relative to the others. Accordingly, the activation ratios for DP/YP increase with Erk but YP/NP decreases (DP/NP increases very slightly). These results might appear nonintuitive, since one would expect that increasing the total Erk concentration would drive both reactions forward and thus all activation ratios should increase. However, if we consider that Mek has a 10-fold higher affinity for YP, then as more total Erk is present, more of the available Mek will be devoted to conversion of YP to DP, with less available to form YP from NP, and the balance of the system tips towards the DP/YP cycle. Since the phosphatases are far from saturation, free YP will be readily dephosphorylated to NP and thus the YP/NP ratio decreases. Use of activation ratios make these results more apparent than examining the activation profiles alone.



**Figure 4.14.** A) Activation profiles (fraction Erk in each form) and B) Activation Ratios vs. volume Mek added, for Erk cycles in liquid phase, varying total concentration of Erk.

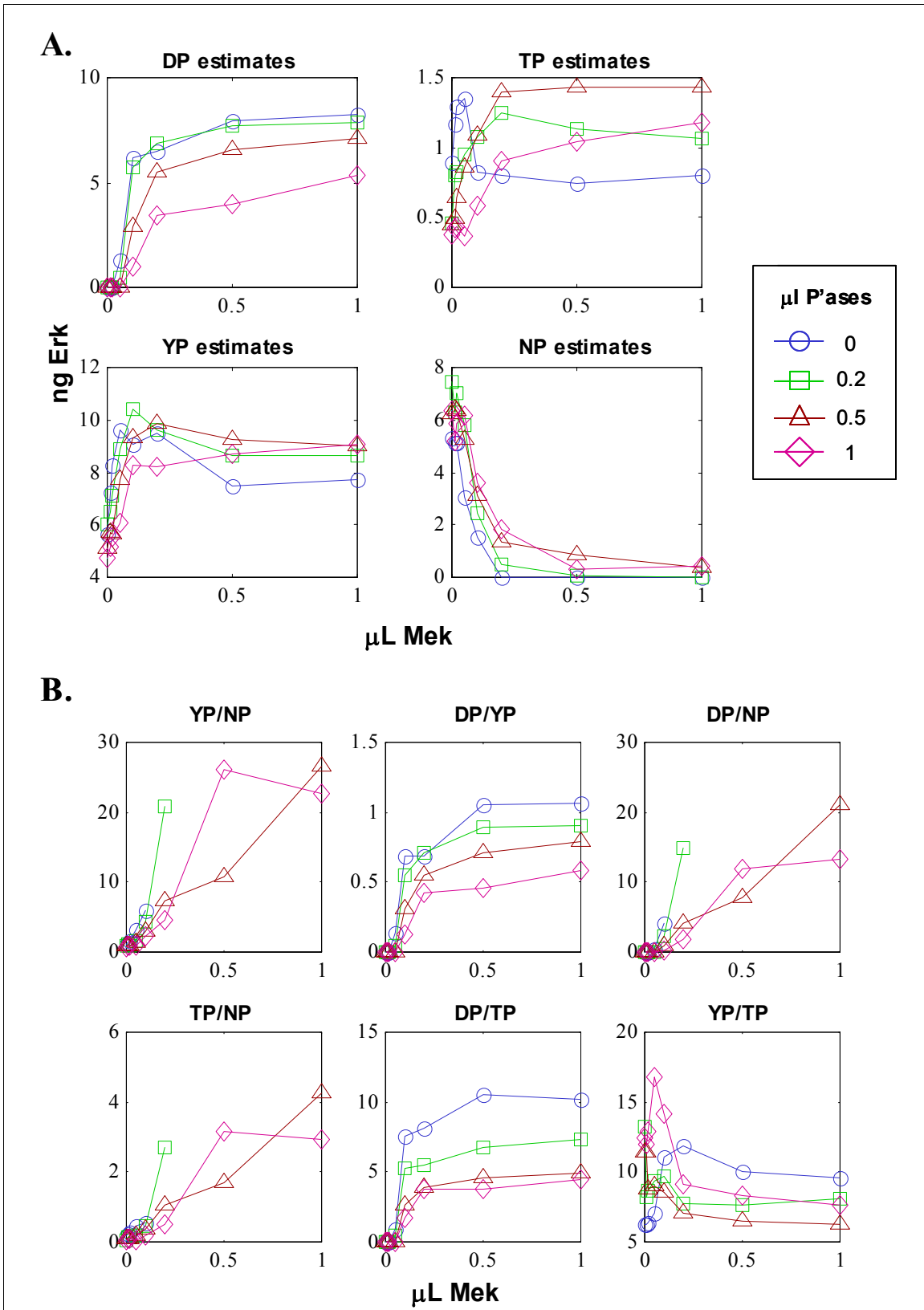


#### 4.4.3. Modulation of phosphatases

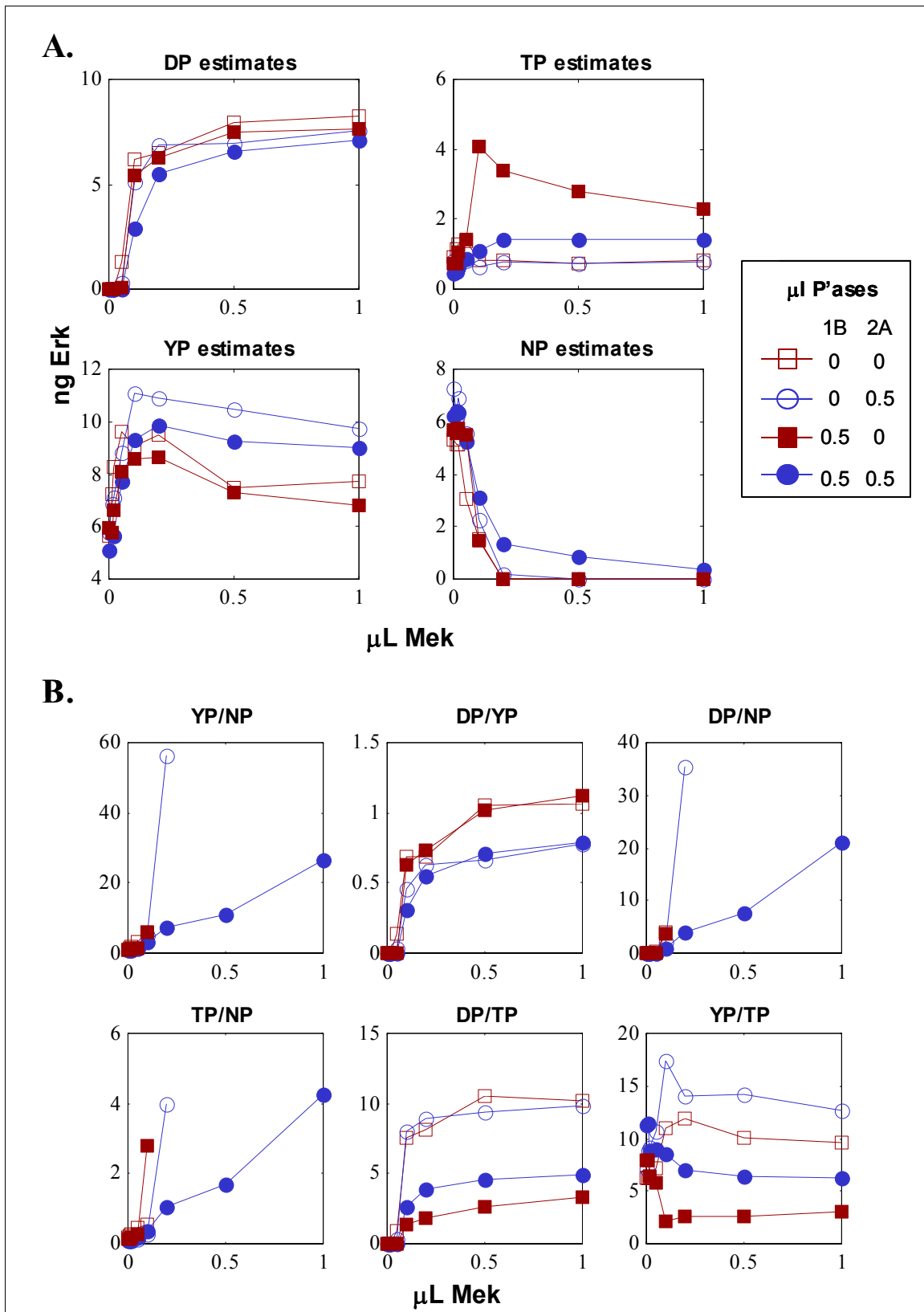
Although examination of the Erk phosphorylation system can be complicated when interactions between Erk and Mek are considered, the analysis is much more straightforward when considering the effects of each phosphatase. This is because each phosphatase will primarily encounter only one Erk species, whereas Mek binds both NP- and YP-Erk. From physical grounds it might be expected that increasing both phosphatase concentrations will tend to decrease all phosphorylated forms. The same conclusion arises from examination of Equations 4.31-4.32. Moreover, activation ratios suggest that if only one phosphatase is modulated, the activation ratio for the step regulated by the other phosphatase should not change significantly.

To verify these predictions, experiments involving liquid-phase reactions using 125 nM Erk were performed, but in this case varying from the base value of 0.2  $\mu$ L phosphatase per 50  $\mu$ L reaction. As shown in Figure 4.15, increasing both phosphatases together generally resulted in decreases in DP-Erk and increases in YP- and NP-Erk, but furthermore, decreases in the YP/NP activation ratio as well as DP/YP. This again suggests that the second phosphorylation step tends to be faster, since a fast initial phosphorylation would help to maintain levels of YP even in the presence of additional phosphatase, and so the YP/NP ratio would be expected to increase.

Each phosphatase was also added independently, to examine the effects of either enzyme. With only one phosphatase, the system is expected to center primarily around the second phosphorylation step, i.e. formation of DP-Erk either from YP-Erk or (if only PTP1B is present) TP-Erk. The activation profiles, shown in Figure 4.16A, indicate as expected that high PTP1B concentrations tend to increase the amount of TP-Erk and decrease YP-Erk, while PP2A promotes formation of YP-Erk. Curiously, YP-Erk is still formed when only PTP1B is present, suggesting perhaps that most of the PTP1B is dedicated to dephosphorylation of DP-Erk. The impact of each phosphatase is strikingly apparent when attention is placed on the activation ratios DP/YP and DP/TP. The ratio DP/YP appears to depend only upon the amount of PP2A (compare squares vs. circles), whereas DP/TP is almost entirely dependent upon PTP1B (open vs. closed symbols). Therefore even with total activation ratios, the influence of each phosphatase on Erk phosphorylation can be separated.



**Figure 4.15.** A) Activation profiles and B) Activation Ratios vs. volume Mek added, for Erk cycles in liquid phase, varying total amounts of both phosphatases.



**Figure 4.16.** A) Activation profiles and B) Activation Ratios vs. volume Mek added, for Erk cycles in liquid phase, varying total amounts of each phosphatase independently.

## **4.5. Conclusions**

The cycles produced by phosphorylation of Erk2 by the protein kinase Mek and dephosphorylation by phosphatases PP2A and PTP1B were selected as a model system for the validation of the theoretical framework described in Chapter 2. A regression scheme was coupled with measurements of Erk standards (originating from an unknown mixture obtained from Upstate Biotechnology) to prepare a model for the response of ELISAs to each form of Erk. The resulting parameters were used to estimate Erk composition in test samples. A sensitivity analysis on the ELISA model was used to estimate uncertainty in the model parameters, which could then be used to approximate the uncertainty in measured Erk composition.

Enzymes were immobilized using capture antibodies so as to facilitating the separation of free substrate from enzyme-substrate complexes following reactions. This procedure enabled immobilization of all three enzymes while maintaining activity and inhibiting nonspecific adsorption of substrate. However, immobilized enzymes exhibited only 5-10% activity when compared against their liquid form, apparently due to poor capture efficiency by the antibodies. After accounting for the differences in activity, the immobilized system behaves identically to liquid-phase reactions.

Activation ratios for Erk phosphorylation cycles were hyperbolic for both reaction systems, suggesting that enzyme-bound substrate was not separated from freely soluble species. It appears likely that the enzyme-substrate complexes were not stable under current reaction conditions, perhaps because of a lack of appropriate buffer composition or scaffold protein. This prevented examination of some of the features of activation ratio analysis, such as the predicted linear and quadratic responses and calculation of the activation factors, as described in Section 2.3.4.

Nevertheless, several qualitative characteristics including the effects of Erk and phosphatase concentrations were still realized, by analyzing the data using total activation ratios (Equations 4.29-4.32). In particular, the relative affinities of Mek for nonphosphorylated and monophosphorylated Erk forms ( $K_m$  for each phosphorylation step) could be readily established by the behavior of activation ratios in response to variation in total Erk concentration. As expected, increasing both phosphatase

concentrations decreased all activation ratios, while independent modulation of each enzyme primarily regulated only one of the two steps. The overlap of activation ratios for constant phosphatase concentrations (Figure 4.16) is a particularly striking example of how this method can isolate activation steps even when there may be additional interactions with other components.

The difficulties in performing these experiments highlight key issues regarding investigation into signaling systems, notably a lack of effective measurement techniques, standards for quantitation of samples, and well-defined methods for handling of enzymes. Certainly, these issues will need to be addressed for further detailed studies into kinetics of enzymes, in particular when calculating activation ratios where more demands are placed upon the measurements. However, they should not detract from understanding the utility of activation ratios themselves, as a different way to investigate individual signaling interactions and to interpret the data. Furthermore, activation ratio analysis should be considered in the context of analyzing larger signaling networks, where detection of interactions in the whole system is more important than precise examination of the individual steps. In that case verification of the qualitative characteristics of activation ratio analysis is more pertinent, and the ability to perform network reconstruction with “total” species measurements should be recognized.

#### **4.6. References**

1. Cobb, M.H. and E.J. Goldsmith, *How MAP kinases are regulated*. J Biol Chem, 1995. **270**(25): p. 14843-6.
2. Lewis, T.S., P.S. Shapiro, and N.G. Ahn, *Signal transduction through MAP kinase cascades*. Adv Cancer Res, 1998. **74**: p. 49-139.
3. Robinson, M.J. and M.H. Cobb, *Mitogen-activated protein kinase pathways*. Curr Opin Cell Biol, 1997. **9**(2): p. 180-6.
4. Zhou, B. and Z.Y. Zhang, *The activity of the extracellular signal-regulated kinase 2 is regulated by differential phosphorylation in the activation loop*. J Biol Chem, 2002. **277**(16): p. 13889-99.
5. Prowse, C.N. and J. Lew, *Mechanism of activation of ERK2 by dual phosphorylation*. J Biol Chem, 2001. **276**(1): p. 99-103.
6. Burack, W.R. and T.W. Sturgill, *The activating dual phosphorylation of MAPK by MEK is nonprocessive*. Biochemistry, 1997. **36**(20): p. 5929-33.
7. Khokhlatchev, A., et al., *Reconstitution of mitogen-activated protein kinase phosphorylation cascades in bacteria. Efficient synthesis of active protein kinases*. J Biol Chem, 1997. **272**(17): p. 11057-62.

8. Huang, W., et al., *Raf-1 forms a stable complex with Mek1 and activates Mek1 by serine phosphorylation*. Proc Natl Acad Sci U S A, 1993. **90**(23): p. 10947-51.
9. McLaughlin, M.M., et al., *Identification of mitogen-activated protein (MAP) kinase-activated protein kinase-3, a novel substrate of CSBP p38 MAP kinase*. J Biol Chem, 1996. **271**(14): p. 8488-92.
10. Oda, Y., et al., *Accurate quantitation of protein expression and site-specific phosphorylation*. Proc Natl Acad Sci U S A, 1999. **96**(12): p. 6591-6.
11. Gygi, S.P., et al., *Quantitative analysis of complex protein mixtures using isotope-coded affinity tags*. Nat Biotechnol, 1999. **17**(10): p. 994-9.
12. Yao, Z., et al., *Detection of partially phosphorylated forms of ERK by monoclonal antibodies reveals spatial regulation of ERK activity by phosphatases*. FEBS Lett, 2000. **468**(1): p. 37-42.
13. Haystead, T.A., et al., *Ordered phosphorylation of p42mapk by MAP kinase kinase*. FEBS Lett, 1992. **306**(1): p. 17-22.
14. Ferrell, J.E., Jr. and R.R. Bhatt, *Mechanistic studies of the dual phosphorylation of mitogen-activated protein kinase*. J Biol Chem, 1997. **272**(30): p. 19008-16.
15. Mansour, S.J., et al., *Transformation of mammalian cells by constitutively active MAP kinase kinase*. Science, 1994. **265**(5174): p. 966-70.
16. Mansour, S.J., et al., *Interdependent domains controlling the enzymatic activity of mitogen-activated protein kinase kinase 1*. Biochemistry, 1996. **35**(48): p. 15529-36.
17. Huang, W., D.S. Kessler, and R.L. Erikson, *Biochemical and biological analysis of Mek1 phosphorylation site mutants*. Mol Biol Cell, 1995. **6**(3): p. 237-45.
18. Zheng, C.F. and K.L. Guan, *Activation of MEK family kinases requires phosphorylation of two conserved Ser/Thr residues*. Embo J, 1994. **13**(5): p. 1123-31.
19. Dent, P., et al., *Expression, purification and characterization of recombinant mitogen-activated protein kinase kinases*. Biochem J, 1994. **303**(Pt 1): p. 105-12.
20. Saxena, M. and T. Mustelin, *Extracellular signals and scores of phosphatases: all roads lead to MAP kinase*. Semin Immunol, 2000. **12**(4): p. 387-96.
21. Alessi, D.R., et al., *Inactivation of p42 MAP kinase by protein phosphatase 2A and a protein tyrosine phosphatase, but not CL100, in various cell lines*. Curr Biol, 1995. **5**(3): p. 283-95.
22. Silverstein, A.M., et al., *Actions of PP2A on the MAP kinase pathway and apoptosis are mediated by distinct regulatory subunits*. Proc Natl Acad Sci U S A, 2002. **99**(7): p. 4221-6.
23. Wang, Z.X., et al., *A kinetic approach for the study of protein phosphatase-catalyzed regulation of protein kinase activity*. Biochemistry, 2002. **41**(24): p. 7849-57.
24. Zhao, Y. and Z.Y. Zhang, *The mechanism of dephosphorylation of extracellular signal-regulated kinase 2 by mitogen-activated protein kinase phosphatase 3*. J Biol Chem, 2001. **276**(34): p. 32382-91.
25. Zhou, B., et al., *The specificity of extracellular signal-regulated kinase 2 dephosphorylation by protein phosphatases*. J Biol Chem, 2002. **277**(35): p. 31818-25.

26. Mansour, S.J., et al., *Mitogen-Activated Protein (Map) Kinase Phosphorylation of Map Kinase Kinase - Determination of Phosphorylation Sites By Mass-Spectrometry and Site-Directed Mutagenesis*. Journal of Biochemistry, 1994. **116**(2): p. 304-314.
27. McCain, D.F. and Z.Y. Zhang, *Assays for protein-tyrosine phosphatases*. Methods Enzymol, 2002. **345**: p. 507-18.
28. Muller, J. and D.K. Morrison, *Assay of Raf-1 activity*. Methods Enzymol, 2002. **345**: p. 490-8.
29. Seger, R., et al., *Human T-cell mitogen-activated protein kinase kinases are related to yeast signal transduction kinases*. J Biol Chem, 1992. **267**(36): p. 25628-31.
30. Zhang, Z.Y., et al., *Determinants of substrate recognition in the protein-tyrosine phosphatase, PTP1*. J Biol Chem, 1996. **271**(10): p. 5386-92.
31. LaPorte, D.C. and D.E. Koshland, Jr., *Phosphorylation of isocitrate dehydrogenase as a demonstration of enhanced sensitivity in covalent regulation*. Nature, 1983. **305**(5932): p. 286-90.
32. Shacter, E., P.B. Chock, and E.R. Stadtman, *Regulation through phosphorylation/dephosphorylation cascade systems*. J Biol Chem, 1984. **259**(19): p. 12252-9.
33. Meinke, M.H., J.S. Bishop, and R.D. Edstrom, *Zero-order ultrasensitivity in the regulation of glycogen phosphorylase*. Proc Natl Acad Sci U S A, 1986. **83**(9): p. 2865-8.
34. Meinke, M.H. and R.D. Edstrom, *Muscle glycogenolysis. Regulation of the cyclic interconversion of phosphorylase a and phosphorylase b*. J Biol Chem, 1991. **266**(4): p. 2259-66.
35. Horiuchi, K.Y., et al., *Competitive inhibition of MAP kinase activation by a peptide representing the alpha C helix of ERK*. Biochemistry, 1998. **37**(25): p. 8879-85.
36. Seger, R., et al., *Overexpression of mitogen-activated protein kinase kinase (MAPKK) and its mutants in NIH 3T3 cells. Evidence that MAPKK involvement in cellular proliferation is regulated by phosphorylation of serine residues in its kinase subdomains VII and VIII*. J Biol Chem, 1994. **269**(41): p. 25699-709.
37. Myles, T., et al., *Active-site mutations impairing the catalytic function of the catalytic subunit of human protein phosphatase 2A permit baculovirus-mediated overexpression in insect cells*. Biochem J, 2001. **357**(Pt 1): p. 225-32.
38. Guan, K.L. and J.E. Dixon, *Evidence for protein-tyrosine-phosphatase catalysis proceeding via a cysteine-phosphate intermediate*. J Biol Chem, 1991. **266**(26): p. 17026-30.
39. Klapa, M.I., J.C. Aon, and G. Stephanopoulos, *Systematic quantification of complex metabolic flux networks using stable isotopes and mass spectrometry*. Eur J Biochem, 2003. **270**(17): p. 3525-42.
40. Romagnoli, J.A. and M.C. Sanchez, *Data processing and reconciliation for chemical process operations*. Process systems engineering ; v. 2. 2000: Academic, San Diego, Calif. ; London.
41. Draper, N.R. and H. Smith, *Applied regression analysis*. 2d ed. Wiley series in probability and mathematical statistics. 1981: Wiley, New York.

42. Seber, G.A.F. and C.J. Wild, *Nonlinear regression*. Wiley series in probability and mathematical statistics. Probability and mathematical statistics. 1989: Wiley, New York.
43. Bates, D.M. and D.G. Watts, *Nonlinear regression analysis and its applications*. 1988: Wiley, New York.
44. Bardwell, A.J., et al., *A conserved docking site in MEKs mediates high-affinity binding to MAP kinases and cooperates with a scaffold protein to enhance signal transmission*. J Biol Chem, 2001. **276**(13): p. 10374-86.
45. Burack, W.R. and A.S. Shaw, *Signal transduction: hanging on a scaffold*. Curr Opin Cell Biol, 2000. **12**(2): p. 211-6.
46. Levchenko, A., J. Bruck, and P.W. Sternberg, *Scaffold proteins may biphasically affect the levels of mitogen-activated protein kinase signaling and reduce its threshold properties*. Proc Natl Acad Sci U S A, 2000. **97**(11): p. 5818-23.
47. Ferrell, J.E., Jr., *Tripping the switch fantastic: how a protein kinase cascade can convert graded inputs into switch-like outputs [see comments]*. Trends Biochem Sci, 1996. **21**(12): p. 460-6.



# 5 CONCLUSIONS – FUTURE WORK

The field of biology has changed dramatically within the last fifty years, since the discovery of the structure of DNA and the birth of molecular biology. The capabilities to sequence genomes and computational tools to identify gene products and splice variants have given researchers a veritable “parts list” for the cell. This is only a first step (albeit a necessary one) in describing cellular behavior. Simply knowing the components of a system is not enough. Understanding the way these molecules react together, forming a responsive network, by determining the reactions that occur is also important. However, a critical final step in the study of cellular operation is determination of the *in vivo* engagement of each of these reactions, as this provides information not just on what the cell is capable of but what is actually occurring under a particular set of conditions.

This thesis represents an approach to perform both the second and third steps, namely, network structure identification and *in vivo* quantitative characterization, for the analysis of signal transduction processes. The former has traditionally been achieved using a variety of experimental techniques, normally focused on individual steps: reactions involving completely isolated species *in vitro*, heterologous expression of proteins in novel hosts, genetic and functional knockouts, etc. The latter has been attempted using a mathematical sensitivity analysis approach originally developed for study in metabolic reactions, but never applied experimentally.

Use of activation ratios as described here enables both steps to be performed simultaneously *in vivo*, without resorting to detailed mechanistic models or requiring extensive molecular perturbations even to analyze simple systems. The nature of the connection between two signaling species is reflected in the shape of activation ratio plots, thus enabling the inversion of data for activation ratios to determine the structure of their interaction. Activation factors, found from simple linear regression of data for direct interactions, yield simple quantitative measures of signal transfer.

The procedure is simple enough to perform visually, although has also been partially automated via a computational regression and evaluation protocol, which can be used to assign confidence in assignment in the face of experimental errors or insufficient data. Although generally robust to errors in estimates for the activation ratios, the pattern assignment algorithm performs more poorly when fitting sharp inverse hyperbolae with errors in the predictor (activating enzyme) concentrations.

One difficult but important extension of this work would therefore be to modify the regression procedure to allow for errors in the x-data. Essentially, this amounts to solving the standard regression problem (for a linear, hyperbolic, inverse hyperbolic model) but with additional parameters being the “true” values of the x data. In general solution of this problem requires replicate measurements at each value of x, and is facilitated if separate information about the potential variance in each x ( $\sigma_x$ ) is available [1]. Although this adds additional experimental requirements and analytical complexity to the procedure, it should provide more stability to the process of pattern assignment.

Somewhat more straightforward should be the development of an automated network reconstruction algorithm, used to connect the results of pattern assignment (a matrix **R** as described in Section 3.3.3) to a graphical structure. This would simplify assignment for simple signaling systems, and be essential for the study of larger networks. In effect, this would be a computational implementation of the algorithm described in Section 3.1, and the consistency rules in Section 3.1.3 would play an important role in error-checking, as would the Akaike weights as measures in confidence. One option would be to simply enumerate all possible networks for a particular number of components, and eliminate those that would not result in the observed **R**. However, considering the combinatorial permutations that arise as the number of species increases, a much more tractable alternative appears to be to assemble a structure starting from a known input and working down.

The process involved in development of the ELISA model, as well as the difficulties encountered in its application, illustrate a key limitation preventing quantitative analysis of signaling systems. Namely, the lack of highly accurate measurement techniques for signaling intermediates, and of well-defined standards to

help translate measured values into estimates of species concentrations. The sensitivity and selectivity of antibodies varies widely, and is rarely characterized or optimized for quantitative purposes. Although individual peptides also vary in their recovery during MS, use of an internal standard still enables quantitation of proteins with less than 10% variance [2, 3]. Use of mass spectrometry would certainly both improve the accuracy and expand the variety of proteins that could be investigated, since it would not depend upon the availability of antibodies. Nevertheless, pure standards would be required for absolute quantitation of proteins in different states. (Only relative quantitation is possible if a consistent, but unknown, sample is used as an internal control)

In the end, the most important, and most useful, continuation of this work will be continued application in a variety of experimental settings. Small, well-controlled *in vitro* studies, similar to the examination of Erk phosphorylation discussed in Chapter 4, could be used to illustrate further the capacity of activation ratios to reflect the kinetics of signaling systems. The issue of separating enzyme-bound from free substrate still remains, and may be addressed by further experiments with immobilized enzymes. Considering the problems encountered using antibody capture, it appears that a different strategy for enzyme immobilization is required. A promising option is to utilize a RNA-protein fusion, which simplifies immobilization through simple DNA-RNA base pairing, while maintaining catalytic function [4].

However, cascades and other signaling network arrangements should also be examined, to demonstrate the capacity of activation ratios to examine signaling structure. Controlled *in vivo* experiments could first be performed using mammalian signaling intermediates expressed in bacteria or insect cells to construct cascades in a novel setting [5, 6]. As methods for separation and analysis of intermediates improve, more complicated systems including yeast and higher-order eukaryotic cells could be studied, to characterize known signaling networks as well as previously undetermined systems, including a variety of multifaceted diseases such as diabetes.

## **5.1. References**

1. Seber, G.A.F. and C.J. Wild, *Nonlinear regression*. Wiley series in probability and mathematical statistics. Probability and mathematical statistics. 1989: Wiley, New York.

2. Oda, Y., et al., *Accurate quantitation of protein expression and site-specific phosphorylation*. Proc Natl Acad Sci U S A, 1999. **96**(12): p. 6591-6.
3. Gygi, S.P., et al., *Quantitative analysis of complex protein mixtures using isotope-coded affinity tags*. Nat Biotechnol, 1999. **17**(10): p. 994-9.
4. Jung, G.Y. and G. Stephanopoulos, *A functional protein chip for combinatorial pathway optimization and in vitro metabolic engineering*. (in preparation), 2003.
5. Macdonald, S.G., et al., *Reconstitution of the Raf-1-MEK-ERK signal transduction pathway in vitro*. Mol Cell Biol, 1993. **13**(11): p. 6615-20.
6. Khokhlatchev, A., et al., *Reconstitution of mitogen-activated protein kinase phosphorylation cascades in bacteria. Efficient synthesis of active protein kinases*. J Biol Chem, 1997. **272**(17): p. 11057-62.

# 6 APPENDICES

## ***Appendix 1. Single cycle model***

### A1.1. Simulation details

All models of signaling systems were constructed as a set of coupled first-order differential equations for the time-dependent change in concentration of each species, including enzyme-substrate complexes. For the single cycle described in Section 2.2, this constitutes Equations 2.8-2.13. As not all equations are independent, only the expressions for the signaling intermediate itself (2.8-2.11) were explicitly coded. Conservation relationships (Equations 2.14-2.16) were used to find the concentrations of free enzymes ( $E_1$ ,  $E_2$ ) and free unmodified intermediate ( $A$ ). At the initial conditions, all species = 0 except for  $A = A_T$ ,  $E_1 = E_{1T}$ ,  $E_2 = E_{2T}$ .

The stiff ODE solver function “ode15s” was then used to integrate the set of differential equations from  $t_0 = 0$  to  $t_{\text{final}} = 5000$  to ensure steady state. The solver internally sets the integration step size; through the parameter “tspan” results at specific time points can be returned. Default solver options were used; this includes relative tolerance ( $10^{-6}$ ) and absolute tolerance ( $10^{-3}$ ) requirements for each species at each integration step time.

The models were constructed as two files. A script file is used to create variables and assign values, create matrices to store output data, and iteratively call the ODE solver ode15s to simulate the cycle activation at each value of input stimulus (or pair of input stimuli). The equations for reaction rates are included in a separate file as a function, that returns  $\mathbf{f}(t, \mathbf{x}, \mathbf{p}) = d\mathbf{x}/dt$ . Kinetic parameters and total concentrations of species are passed as matrices  $\mathbf{p}$  to the function, as well as to the ODE solver.

## A1.2. MATLAB files

Shown below are the two files used for examination of a single cycling intermediate, as described above and in Section 2.2. Extension to other signaling arrangements was achieved by changing the script file (to add kinetic parameters, total concentrations, and handle additional species) and the function of rate equations (to add expressions for additional species or change reactions steps) as appropriate. Copies of all files, including those containing output data, are available from the author by request.

```
% Script file to prepare data for single cycle, single activation step
% Calculates steady-state distribution of forms (free active, free
% inactive, bound to each enzyme) at a range of activating enzyme
% concentrations and also varies a1 so as to modify Km1
%
% F. Javier Femenia, MIT Chemical Engineering Department
% 8/29/00

% Assignment of kinetic parameters
a = 20*ones(1,2);
d = 10*ones(1,2);
k = 10*ones(1,2);

Param = [a;d;k]';

% Total concentrations for each species
Atot = 10;
E1tot = 1;
E2tot = 1;

Totals = [Atot; E1tot; E2tot];

% Varying the concentration of input stimulus (E1tot), and the
% association constant for first reaction (a1)

npts = 27;
Elvals = logspace(-3,1,npts);
Avals = [1 2 5 10 20 50 100];
npts2=length(Avals);
[X,Y] =meshgrid(Elvals,Avals);

% Setting up necessary variables for the ode solver
% See 'help odeset' for information regarding options.
% tspan sets the integration from 0 to tfinal, such that the output
% matrix 'Res' will only have 3 rows--easier for later.

tfinal = 5000;           % Length of integration--until system at ss.
options = [];           % Default options, can make modifications later
tspan = [0 10 tfinal]; % Output of odesolver will give rows at
                        % t=0, 10, 5000

Results = zeros(npts2,npts,4); % Matrix to store output
```

```

% At each value of Eltot and a1, integrates to steady state and
% assigns the output to the matrix Results.
% The ode solver calls the file "singlefile.m" which contains the
% differential equations governing each species (form of A)

for i = 1:npts          % Outer (Eltot) loop
    for j = 1:npts2    % Inner (a1) loop
        Eltot = X(j,i)
        a(1) = Y(j,i);
        vars0 = [Atot;0;0;0];
        Param = [a;d;k]';
        Totals = [Atot; Eltot; E2tot];
        [t,Res] = ode15s('singlefile',tspan,vars0,options,Param,Totals);
        Results(j,i,:) = Res(3,:);
    end
end

% Calculate fractions of enzymes in active form (including those in
% complexes with other proteins).

A=Results(:, :, 1);
AP=Results(:, :, 2);
AE=Results(:, :, 3);
APE=Results(:, :, 4);

% If desired, create plots of activation profiles and activation ratios

figure(1)
plot(Elvals,AP./Atot)
title('Fraction A Activated for Different Values of a_1')
xlabel('E1_t_o_t')
legend('1','2','5','10','20','50','100')

figure(2)
plot(Elvals,AP./A)
title('Activation Ratios for Different Values of a_1')
xlabel('E1_t_o_t')
legend('1','2','5','10','20','50','100')

```

```

function varargout = singlefile(t,vars,flag,Param,Totals)

% Function used to calculate distribution of forms in signaling cycles.
% Based on using deterministic model for each chemical reaction
% between species as a set of first-order differential equations
% that describes the time-dependent behavior of each species.
%
% The flag is needed for the odesolver. Thus there are subfunctions to
% allow evaluation at different flags.
% Param is a matrix containing all the kinetic rate parameters, and
% Totals is a vector containing the total concentrations of all
% species.
% vars represents the variables y for which we are integrating; in this
% case is A, AP, E1_A and E2_AP

switch flag
case ''
    % Return dy/dt = f(t,y).
    varargout{1} = f(t,vars,Param,Totals);
case 'init'
    % Return initial conditions y(0).
    [varargout{1:3}] = initval(Param,Totals);

%case 'jacobian'
    % Return Jacobian matrix df/dy.
% varargout{1} = jacobian(t,y,p1,p2);
%case 'jpattern'
    % Return sparsity pattern matrix S.
% varargout{1} = jpattern(t,y,p1,p2);
%case 'mass'
    % Return mass matrix M(t) or M.
% varargout{1} = mass(t,y,p1,p2);
%case 'events'
    % Return [value,isterminal,direction].
% [varargout{1:3}] = events(t,y,p1,p2);

otherwise
    error(['Unknown flag '' flag ''.']);
end

% -----
% Normal case: Evaluates df/dt at given t and parameter values

function eqneval = f(t,vars,Param,Totals)

% Relation of local variables to input values

a = Param(:,1); % Column 1 in Param is ai
d = Param(:,2); % Column 2 is di
k = Param(:,3); % Column 3 is ki

Atot = Totals(1);
E1tot = Totals(2);
E2tot = Totals(3);

AP = vars(2);
A_E1 = vars(3);
AP_E2 = vars(4);

```



```

% Application of conservation relationships

A = Atot - (AP + A_E1 + AP_E2);
E1 = E1tot - A_E1;
E2 = E2tot - AP_E2;

% Actual system of equations (evaluation of differential equations)

eqns(1) = -a(1)*A*E1 + d(1)*A_E1 + k(2)*AP_E2;      % dA/dt
eqns(2) = -a(2)*AP*E2 + d(2)*AP_E2 + k(1)*A_E1;     % dAP/dt
eqns(3) = a(1)*A*E1 - (d(1)+k(1))*A_E1;             % dA_E1/dt
eqns(4) = a(2)*AP*E2 - (d(2)+k(2))*AP_E2;           % dAP_E2/dt

eqneval=eqns';

% -----

function [tspan,y0,options] = initval(Param,Totals)
tspan = [0 10 5000];
y0 = [Totals(1);0;0;0];
options = [];

disp('Go Bears!!');

%% -----
%
%function dfdy = jacobian(t,y,p1,p2)
%dfdy = < Insert Jacobian matrix here. >;
%
%% -----
%
%function S = jpattern(t,y,p1,p2)
%S = < Insert Jacobian matrix sparsity pattern here. >;
%
%% -----
%
%function M = mass(t,y,p1,p2)
%M = < Insert mass matrix here. >;
%
%% -----
%
%function [value,isterminal,direction] = events(t,y,p1,p2)
%value = < Insert event function vector here. >
%isterminal = < Insert logical ISTERMINAL vector here.>;
%direction = < Insert DIRECTION vector here.>;

```

## Appendix 2. Parameter values for simple signaling models

### A2.1. Linear Cascade

Total Concentrations ( $\mu\text{M}$ )

A	10
B	10
C	10
$E_1$	$10^{-6} - 100$
$E_{IA}$	10
$E_{IB}$	10
$E_{IC}$	10

Kinetic Rate Constants ( $\text{s}^{-1}$  or  $\mu\text{M}^{-1}\text{s}^{-1}$ )

Reaction	a	d	k	$K_m$
$E_1 \rightarrow A$	2	1	1	1
$E_{IA} \rightarrow A^*$	0.5	1	1	4
$A^* \rightarrow B$	4	1	1	0.5
$E_{IB} \rightarrow B^*$	0.67	1	1	3
$B^* \rightarrow C$	2.5	1	1.5	1
$E_{IC} \rightarrow C^*$	0.85	1	0.7	2

### A2.2. Converging Cycles

Total Concentrations ( $\mu\text{M}$ )

A	0.1
$E_1$	0 - 0.1
$E_2$	0 - 0.1
$E_{IA}$	0.01

Kinetic Rate Constants ( $\text{s}^{-1}$  or  $\mu\text{M}^{-1}\text{s}^{-1}$ )

Reaction	a	d	k	$K_m$
$E_1 \rightarrow A$	200	100	100	1
$E_2 \rightarrow A$	200	150	50	1
$E_{IA} \rightarrow A^*$	200	100	100	1

### A2.3. Diverging Pathways

Total Concentrations ( $\mu\text{M}$ )

A	10
B	10
$E_1$	0 - 20
$E_{IA}$	1
$E_{IB}$	1

Kinetic Rate Constants ( $\text{s}^{-1}$  or  $\mu\text{M}^{-1}\text{s}^{-1}$ )

Reaction	a	d	k	$K_m$
$E_1 \rightarrow A$	20	10	10	1
$E_{IA} \rightarrow A^*$	20	10	10	1
$E_1 \rightarrow B$	20	15	5	1
$E_{IB} \rightarrow B^*$	20	10	10	1

### A2.4. Cascades with Feedback

(<sup>†</sup>Note:  $k_{FB}$  was varied from  $0-10^4$ ,  $a_{FB}$  adjusted to maintain  $K_{FB}=100$ )

Total Concentrations ( $\mu\text{M}$ )

A	10
B	10
C	10
D	10
E	10
$E_1$	$10^{-5} - 10$
$E_{IA}$	0.1
$E_{IB}$	0.1
$E_{IC}$	0.1
$E_{ID}$	0.1
$E_{IE}$	0.1

Kinetic Rate Constants ( $\text{s}^{-1}$  or  $\mu\text{M}^{-1}\text{s}^{-1}$ )

Reaction	a	d	k	$K_m$
$E_1 \rightarrow A$	101	1	100	1
$E_{IA} \rightarrow A^*$	25.25	1	100	4
$A^* \rightarrow B$	4	1	1	0.5
$E_{IB} \rightarrow B^*$	33.67	1	100	3
$B^* \rightarrow C$	2.5	1	1.5	1
$E_{IC} \rightarrow C^*$	35.5	1	0.7	2
$C^* \rightarrow D$	3	1	2	1
$E_{ID} \rightarrow D^*$	25.25	1	100	4
$D^* \rightarrow E$	5	1	1	0.4
$E_{IE} \rightarrow E^*$	51	1	50	1
$E^* \rightarrow A$ or $A^*$	$0.11^{\dagger}$	1	$10^{\dagger}$	100

### Appendix 3. Model signaling network

Total Concentrations ( $\mu\text{M}$ )

A	10
B	10
C	10
D	10
E	10
F	10
G	10
H	10
E <sub>1</sub>	10 <sup>-6</sup> – 10
E <sub>2</sub>	10 <sup>-6</sup> – 10
E <sub>IA</sub>	10
E <sub>IB</sub>	10
E <sub>IC</sub>	10
E <sub>ID</sub>	10
E <sub>IE</sub>	10
E <sub>IF</sub>	10
E <sub>IG</sub>	10
E <sub>IH</sub>	10

Kinetic Rate Constants ( $\text{s}^{-1}$  or  $\mu\text{M}^{-1}\text{s}^{-1}$ )

Reaction	a	d	k	K <sub>m</sub>
E <sub>1</sub> → A	12	10	2	1
E <sub>IA</sub> → A*	2.75	10	1	4
E <sub>2</sub> → B	12	10	2	1
E <sub>IB</sub> → B*	11	10	1	1
A* → C	22	10	1	0.5
B* → C	11	10	1	1
E <sub>IC</sub> → C*	3.67	10	1	3
B* → D	12	10	2	1
E <sub>ID</sub> → D*	11	10	1	1
C* → E	11.5	10	1.5	1
D* → E	5.75	10	1.5	2
E <sub>IE</sub> → E*	5.35	10	0.7	2
D* → F	5.75	10	1.5	2
E <sub>IF</sub> → F*	5.5	10	1	2
A* → G	12	10	2	1
E <sub>IG</sub> → G*	11	10	1	1
A* → H	11.5	10	1.5	1
E <sub>IH</sub> → H*	11	10	1	1

Calculations for predicted  $\alpha_M^N$  based on Equation 2.51:

$$\frac{N^*}{N} = \frac{k_N M^*}{k_{IN} E_{INT}} \frac{K_{mIN} + N^*}{K_{mN}} \alpha_M^N M^* \text{ so } \alpha_M^N = \frac{k_N}{k_{IN} E_{INT}} \frac{K_{mIN} + N^*}{K_{mN}}$$

Reaction #	Reaction	maximum N*	$\alpha_M^N$ predicted	$\alpha_M^N$ observed
1	E <sub>1</sub> → A	0.6905	0.9381	0.937 ± 0.001
2	E <sub>2</sub> → B	0.5818	0.3164	0.315 ± 0.001
3	A* → C	1.0234	0.8047	0.799 ± 0.006
4	B* → C	0.7779	0.3780	0.373 ± 0.003
5	B* → D	0.4777	0.2956	0.290 ± 0.004
6	C* → E	1.7317	0.7997	0.770 ± 0.002
7	D* → E	1.7317	0.3998	0.399 ± 0.005
8	D* → F	0.5388	0.1904	0.1878 ± 0.0004
9	A* → G	0.7165	0.3433	0.336 ± 0.006
10	A* → H	0.5593	0.2339	0.230 ± 0.003
6 (E <sub>1</sub> only)	C* → E	1.5637	0.7637	0.76 ± 0.01
6 (E <sub>2</sub> only)	C* → E	1.5604	0.7629	0.98 ± 0.01
7 (E <sub>2</sub> only)	D* → E	1.5604	0.3815	1.58 ± 0.03

#### **Appendix 4. Extended cascade with noise**

Parameter values and model were identical to the extended cascade with feedback (See A2.4 for values). However, feedback was eliminated by setting the rate constants for the feedback step (a, d, and k) equal to zero. The resulting model was integrated to yield “true” values of active and inactive species (e.g. A, A<sup>\*</sup>, etc).

The MATLAB function “randn” was used to help simulate experimental error in measurements. This function produces as an output a randomly selected number x from a Normal distribution centered around zero with a variance of one ( $x \sim N(0,1)$ ). Multiplying this number by any factor a will make the standard deviation equal to a. Therefore, multiplying the true value of a measurement (say, A<sup>\*</sup>) by 0.05\*randn will produce a number that is normally distributed about the true A<sup>\*</sup>, with a standard deviation equal to 5% of the true value. This value (which may be positive or negative) can be added to the “true” value as the error, yielding an error-adjusted estimate for the measured value. This process was performed either on the activation ratios AR<sub>I</sub> or on all individual values I and I<sup>\*</sup>, from which activation ratios were then calculated. The automated regression and pattern assignment procedure was applied to the adjusted data, producing an output matrix as shown in Figure 3.10B. This approach was repeated ten times (to produce ten sets of adjusted data) and the average results are shown in Table 3.5.

## Appendix 5. Sensitivity matrix for ELISA measurement model

Results of differentiation of ELISA model (Equations 4.3-4.18) as a matrix  $G_{ij} = \partial m_i / \partial p_j$ . Columns represent parameters  $p_j$ , and rows correspond to measurements  $m_i$ .

$G_{ij}$	$\partial a_{11}$	$\partial a_{22}$	$\partial a_{24}$	$\partial a_{33}$	$\partial a_{42}$	$\partial a_{43}$	$\partial a_{44}$	$\partial b_1$	$\partial b_2$	$\partial b_3$	$\partial b_4$	$\partial d$	$\partial t$	$\partial y$
$\partial a_{DP_D}$	$dx_T$							1				$a_{11}x_T$		
$\partial a_{TP_D}$		$tx_T$	$(1-d-t-y)x_T$						1			$-a_{24}x_T$	$(a_{22}-a_{24})x_T$	$-a_{24}x_T$
$\partial a_{YP_D}$				$yx_T$						1				$a_{33}x_T$
$\partial a_{NP_D}$					$tx_T$	$yx_T$	$(1-d-t-y)x_T$				1	$-a_{44}x_T$	$(a_{42}-a_{44})x_T$	$(a_{43}-a_{44})x_T$
$\partial a_{DP_T}$								1						
$\partial a_{TP_T}$		$(d+t)x_T$	$(1-d-t)x_T$						1			$(a_{22}-a_{24})x_T$	$(a_{22}-a_{24})x_T$	
$\partial a_{YP_T}$										1				
$\partial a_{NP_T}$					$(d+t)x_T$		$(1-d-t)x_T$				1	$(a_{42}-a_{44})x_T$	$(a_{42}-a_{44})x_T$	
$\partial a_{DP_Y}$								1						
$\partial a_{TP_Y}$			$(1-d-y)x_T$						1			$-a_{24}x_T$		$-a_{24}x_T$
$\partial a_{YP_Y}$				$(d+y)x_T$						1		$a_{33}x_T$		$a_{33}x_T$
$\partial a_{NP_Y}$						$(d+y)x_T$	$(1-d-y)x_T$				1	$(a_{43}-a_{44})x_T$		$(a_{43}-a_{44})x_T$
$\partial a_{DP_N}$								1						
$\partial a_{TP_N}$			$x_T$						1					
$\partial a_{YP_N}$										1				
$\partial a_{NP_N}$							$x_T$							