# Three Essays in Macroeconomics

by

## Raphael Anton Maximilian Peter Gabriel Auer

M.A., University of Maastricht (2004)

Submitted to the Department of Economics
in partial fulfillment of the requirements for the degree of
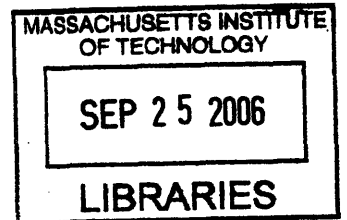
Doctor of Philosophy

at the

## MASSACHUSETTS INSTITUTE OF TECHNOLOGY

September 2006

Signature of Author..................................................................
Department of Economics
15 August 2006

Certified by .......................................................................
Xavier Gabaix
Dornbusch Career Development Associate Professor
Thesis Supervisor

Certified by .......................................................................
Daron Acemoglu
Charles P. Kindleberger Professor of Applied Economics
Thesis Supervisor

Accepted by.......................................................................
Peter Temin
Elisha Gray II Professor of Economics
Chairman, Departmental Committee on Graduate Studies

# Three Essays in Macroeconomics

by

Raphael Anton Maximilian Peter Gabriel Auer

Submitted to the Department of Economics
on 15 August 2006, in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy

## Abstract

This thesis is a collection of three essays on international trade and economic growth.

Chapter 1 analyzes the dynamic gains from trade in a Hecksher-Ohlin economy with endogenous factor accumulation. In a framework where heterogeneous workers make educational decisions in the presence of complete markets, I first show how convergence of factor rewards induces divergence of factor abundance and levels of income. When heterogeneous workers invest in schooling, higher type agents earn a surplus from their investment. By affecting educational decisions, trade influences the international distribution of this surplus. The latter effect tends to benefit richer countries disproportionately, leading to divergence of welfare when markets are opened to trade. The shift of investments to initially rich countries also leads to a global increase of the average skill premium despite a decrease of the price of skill intensive goods. I next examine whether the factor content of trade indeed does affect domestic education decisions. To establish a causal relation, I instrument for the factors embodied in actual imports by the geographic component of trade. The constructed measures of geographical proximity to skilled and unskilled labor have significant effects on domestic educational decisions. Countries that tend to be close to international supply of skilled labor have lower levels of advanced education, while the reverse is true for countries that are close to labor abundant nations. A one standard deviation difference in geographic proximity to skilled labor is associated with a difference of about 2/3 of a year of average higher education.

Chapter 2 examines why movements of relative costs brought about by exchange rate fluctuations are passed on to customers only slowly, and never to a full extent. We first develop a perfectly competitive economy featuring heterogeneity of both good qualities and of consumer valuations. In equilibrium, high valuation consumers and high quality firms are matched. The relative scarcity of different qualities leads to pricing-to-market and markups that are determined by the local toughness of competition. Our production setup features trade in intermediate goods, local assembly that is subject to decreasing returns and fixed costs of market entry. In every export market, firm entry and size decisions are determined by how local prices compare to the cost of production at home. We next analyze how changes in the real exchange rate are transmitted internationally. In the short run, the set of firms active in the export sector is fixed, but each firm accommodates changes in the exchange rate by adjusting the quantity of its exports. Due to this response of export volume to the relative cost of production, market toughness counteracts exchange rate movements, leading to partial pass-through in the short run. Due to the presence of fixed costs of market access, in the long run also the set of firms

that are actively exporting reacts to movements of the real exchange rate, with two associated consequences. Firstly, pass-through is larger than in the short run because long run export volume responds to relative costs due to changes in both the average firm size and in the number of firms. Secondly, the response of the market entry decision to changes in the relative cost of production affects only low quality firms, which fetch a relatively low price for their output. Exchange rate movements thus change the composition of actively exporting firms, with the consequence that aggregate price indexes overstate the actual extent of pass-through in the long run.

Chapter 3 further examines the seminal work of Acemoglu et al. (2001) on the effects of settler mortality on colonization policies during early imperialism. The authors build a strong case for the importance of institutions as the primary force of economic development. However, because their empirical analysis is limited to former colonies, they cannot directly distinguish their theory from the rivaling view that a country's disease environment has direct effects on economic prosperity and institutions. In this paper, using either additional historical sources or a model of the geographic determinants of disease, I first construct two measures of mortality rates including up to 36 countries that have not been colonized. I then show that mortality did affect institutional development in former colonies but not in the rest of the sample. This can only be rationalized in the context of the colonial origins theory of Acemoglu et al. Turning to disentanlge the relation between institutions and income, I sometimes find that disease environment influences income also directly and correspondingly, that institutions are somewhat less important for prosperity in my specifications than when working with a sample composed of only former colonies. Incorporating these findings, I estimate that institutions are the major determinant of long run prosperity and can explain about 50% of the observed variation of current income levels, while the direct effects of disease environment can account for about 15%.

Thesis Supervisor: Xavier Gabaix
Title: Dornbusch Career Development Associate Professor

Thesis Supervisor: Daron Acemoglu
Title: Charles P. Kindleberger Professor of Applied Economics

# Acknowledgements

# Contents

# Chapter 1

# Human Capital and the Dynamic Effects of Trade

**Summary 1** *This Chapter analyzes the dynamic gains from trade in a Hecksher-Ohlin economy with endogenous factor accumulation. In a framework where heterogeneous workers make educational decisions in the presence of complete markets, I first show how convergence of factor rewards induces divergence of factor abundance and levels of income. When heterogeneous workers invest in schooling, higher type agents earn a surplus from their investment. By affecting educational decisions, trade influences the international distribution of this surplus. The latter effect tends to benefit richer countries disproportionately, leading to divergence of welfare when markets are opened to trade. The shift of investments to initially rich countries also leads to a global increase of the average skill premium despite a decrease of the price of skill intensive goods. I next examine whether the factor content of trade indeed does affect domestic education decisions. To establish a causal relation, I instrument for the factors embodied in actual imports by the geographic component of trade. The constructed measures of geographical proximity to skilled and unskilled labor have significant effects on domestic educational decisions. Countries that tend to be close to international supply of skilled labor have lower levels of advanced education, while the reverse is true for countries that are close to labor abundant nations. A one standard deviation difference in geographic proximity to skilled labor is associated with a difference of about 2/3 of a year of average higher education.*

## 1.1 Introduction

While the static effects of international trade are well understood, there is much more limited analysis of the dynamic effects of trade, especially in the presence of accumulated factors. In essence, a large part of the existing literature on the dynamic gains from trade focuses on market failures that may become exacerbated with trade. In this paper, I show that even in a world with complete markets, trade, while benefiting all nations, may also create divergence in income per capita and in relative welfare. I then present evidence that the proposed mechanism is of statistical and economic significance.

A great deal of literature debates the gains from trade for poor nations. In a seminal article, Young (1991) shows how trade can cause countries to specialize in industries with differential learning-by-doing potential and hence be on different dynamic learning paths. Countries with low initial experience in industrial production specialize in sectors with low learning potential and may thus loose from trade. Krugman and Venables (1995) show how, in the presence of increasing returns, initial patterns of specialization tend to reinforce themselves because new firms locate close to existing industry. Other contributions, not limited to, but including Matsuyama (1991) and the new economic geography literature originating from Krugman (1991), focus on similar mechanisms of increasing returns. The model developed here differs substantially from the existing literature on the dynamic gains from trade because it does not focus on the evolution of location and productivity of different industries but rather on the endogenous formation of factor supplies. One of the most fundamental insights of the theory of international trade – Samuelson's (1948) factor price equalization theorem – establishes conditions under which foreign trade equates international returns to factors. Stated in its simplest and weakest form, trade increases the reward of domestically abundant factors. When some factors, such as physical and human capital, are in variable supply, trade thereby re-enforces initial patters of specialization[1]. In his early contribution to the theory of international trade, Ohlin (1933) discusses this dynamic feedback from trade to induced changes in skill accumulation.

"The adaptation of labour to the requirements of industries where it is employed

---

[1]Stiglitz (1970) has used this insight to argue that complete specialization is indeed extremely likely in a dynamic context.

goes so far as to become a cause of *extended* trade. Acquired no less that inherited qualities which involve differences between the productive resources of various nations lead to specialization along different lines, i.e. to international trade. Trade thus engenders more trade [...] and it tends to increase the unevenness of the international distribution of factors of production." Original emphasis, Ohlin (1933), pp. 125/126

What are the welfare consequences of trade-induced accumulation of factors? The framework of this paper draws on the insights of Findlay and Kierzkowski (1983), who propose a general equilibrium model of human capital accumulation in the presence of international trade. I depart from their model with two key assumptions. In my framework, countries are characterized by exogenously given differences of the efficiency of human capital. In autarky equilibrium, countries with a high level of human capital efficiency are characterized by a high demand for skills and hence a high level of human capital. Trade leads nations with a comparatively high level of human capital efficiency to specialize in skill intensive goods and leads to other nations providing unskilled labor services. In a dynamic context, the basic asymmetry of the model is that trade induces productive nations to specialize in a factor that can be accumulated, which increases the growth potential of the economy. Less productive nations specialize in 'raw' labor, a factor in fixed supply. Opening markets to trade therefore results in divergence of the world distribution of income. A second departure from Findlay and Kierzkowski (1983) is that I assume that, while workers are homogenous in how well they can provide unskilled labor, they differ in how well they can supply skilled labor if they chose to get an education. In the resulting equilibrium of the economy, low type workers do not accumulate skills. Higher type workers do, and while there may exist a cut-off type that is indifferent between getting an education or not, all other skilled workers earn a surplus from education. Trade induced changes in relative wages affect this surplus in a way that favors already rich and developed nations, leading to divergence of welfare[2].

The first part of the paper analyses the effects of globalization – opening markets to trade –

---

[2]This characteristic of the model is what makes human capital differenct from physical capital. Baldwin (1992) discusses the dynamic gains from trade when physical capital is accumulated endogenously. He concludes that in the absence of externalities, trade induced accumulation has no welfare consequences.

in a partial equilibrium setting, i.e. taking as a given world goods prices. In autarky equilibrium, a country with a high efficiency of human capital is skill abundant, has a low domestic price of the skill intensive good and is characterized by a high level of income and consumption. Countries with a high level of human capital efficiency are hence rich, or "developed" while skill scarce countries are "poor" in equilibrium. By the Stolper-Samuleson effect, exposing the economy to international prices leads to an increase (decrease) of the relative skilled wage if the country is more (less) skill abundant than the rest of the world. Hence, open markets increase entry into the skilled labor force in skill abundant countries and decrease them elsewhere.

I compare the path of income after opening markets to trade for skill abundant or developed countries to that of skill scarce or poor countries. This path is characterized by three distinct phases. At the moment of opening to trade, all countries benefit and the relative size of these gains is determined by global conditions of demand and supply. This static comparison is of importance since it establishes the benchmark of what would happen to income and consumption differentials after opening markets in a standard model of Hecksher-Ohlin trade with factors of production in fixed supply. Since educational investments take the form of forgone earnings, the dynamic path of income is first characterized by a phase of convergence. Poor countries send a smaller fraction of young workers to the educational sector, but their older cohorts are still relatively skilled. Thus, they experience an increase in the supply of unskilled labor, while the supply of skilled labor reacts only with a lag when new cohorts finish their schooling. Richer countries start sending a larger fraction of young workers into schooling, while their supply of skilled labor is stable for a while. The resulting medium term response of income displays not only relative but also absolute convergence of income levels. The GDP of poor countries increases, while that of richer countries decreases. This pattern prevails until the first cohort of workers, who started schooling at the moment of opening markets to trade, leaves the educational sector and enters the labor force. From then on, earlier changes in educational investment start to pay off and the GDP of rich countries increases, while the opposite is true in poor countries. The resulting long term dispersion of income is larger than at the moment of opening to trade, and also likely to be larger than it was in autarky. This result of dynamic divergence is related to Ventura (1997), who argues that open economies can avoid running into decreasing marginal product of capital by shifting the structure of their exports

11

into successively more capital and skill intensive sectors. Precisely the same mechanism that allows growth miracles in an open economy leads to divergence when countires open to trade: in autarky, factor abundance differentials are smaller than in an open economy regime. Trade results in dynamic divergence of the world distribution of income because it shifts investment to relatively developed countries.

I then turn to evolution of welfare. International capital markets enable countries to smooth consumption. Welfare changes are hence equivalent to changes in the net present value of the future flow of income. At the moment of opening markets – i.e. taken as given factor supplies – rich and poor nations benefit from trade and the size of relative gains is determined by the global scarcity of factors. In addition, all countries gain when the skill supply adjusts to the changed demand conditions under free trade. This dynamic response of the economy introduces the asymmetry between nations in the model: the dynamic gains from trade are likely to favor already rich nations. This result stems from the two margins in which the relative wage influences the surplus from education. A higher relative wage increases the income for all workers who already would have chosen schooling at lower wages. In addition, an increase in the relative wage induces more entry into the skilled labor force. In total, the net income from education – taking into consideration the opportunity cost of forgone unskilled labor – responds more than proportionally to changes in the relative wage. Skill scarce nations, in contrast, have their comparative advantage in labor, a factor that is in fixed supply and cannot be accumulated. These countries gain linearly in the increase of the unskilled wage. I develop conditions under which the total gains from trade lead to overall divergence of welfare. This tends to be the case if the global skill intensive sector is large compared to the labor intensive sector, the period of time needed to get an education is short, and the heterogeneity of workers is not too large. Summarizing, the key insight of the mechanism at work is that, while all countries gain from trade, already developed nations gain proportionally the most from trade. Trade hence results in divergence of welfare.

The next part of the paper evaluates the general equilibrium response of simultaneously opening many countries to trade. The results of this section are related to a growing literature on the skill bias of global trade. The increased exposure to international trade seems to have resulted in both a pervasive increase in the skill premium while resulting in a decrease in the price of skill intensive goods. One group of papers includes Dinopolous and Segerstrom (1999) and Gancia and Epifani (2005) and argues that the skill intensive sector is more sensitive to scale. Trade increases the market size for an average firm and hence leads to a relative expansion of the skill intensive sector. A second class of models builds on the directed technical change literature, with contributions by Acemoglu and Zilibotti (2001), Acemoglu (2003), and Gancia (2004). Here, it is a combination of unequal protection of intellectual property rights and differential factor endowments that creates technical change biased towards skilled workers. By increasing the market size for skill complementary technologies in those countries that have good intellectual property rights protection, trade increases the skill bias of global technology.

The current paper presents a new channel for why trade is skill biased and leads to a global expansion of the skill intensive sector. The mechanism does not rely on how trade influences technology, but on how trade influences the international location of human capital. Trade equates goods prices across the world, and the dynamic response of education decisions tends to concentrate human capital in countries that can use skills efficiently. With the average skilled worker working in a country with a higher level of human capital augmenting technology, the output of skill intensive goods increases. This results in a decrease of the price of skill intensive goods. The expansion of the skill intensive sector takes place slowly as new cohorts enter the labor force. Some countries start the process of globalization as exporters of the skill intensive good, but successively become importers of the latter. Despite the decrease in the price of the skill intensive good, I show that an open economy is skill biased. This is a consequence of two related mechanisms. At the moment of opening to trade, the skill premium increases in human capital abundant countries, while it decreases in skill scarce countries. The arithmetic average of the skill premium - weighted by relative supply - hence increases with trade. Dynamically, there exists another effect leading to further skill bias. The supply of human capital decreases in countries that are skill scarce and increases elsewhere, resulting in a further increase in the arithmetic average of the skill premium. The results of the model in general equilibrium hence

explain why a globalizing world is characterized by both a decreasing price of the skill intensive good while at the same time resulting in a pervasive increase in the skill premium.

Finally, I present empirical evidence that the skilled labor embodied in current trade flows is indeed a significant factor for domestic education decisions. The empirical strategy first constructs measures of geographic proximity to international supply of skilled and unskilled labor. I then show that, conditional on the level of a country's development, these measures are a significant determinant of investment in human capital and average years of education in the workforce.

A sharp and testable prediction of the Hecksher-Ohlin trade theory made by Vaneck (1968) is that trade can be reduced to the net factor content it embodies. Extensive research efforts have been aimed at establishing the empirical validity of this Hecksher-Ohlin-Vanek (HOV) prediction, with mixed success. While earlier studies (Bowen, Leamer Sveikauskas (1984) and Trefler (1995)) have struggled to show that trade embodies a sizeable net factor content, there has been substantial progress in estimating the factor content of trade when adjusting for productivity differences (David and Weinstein (2001), Antweiler and Trefler (2002)). Davis and Weinstein show that even among the homogenous group of ten wealthy OECD countries, the net factor content of trade is typically equivalent to 10 percent of national endowments. Trefler (2002), building on an observation of Conway (2002) uses data that also covers a significant number of poor countries to show that the dimension in which the flow of embodied factors follows mostly closely the direction of the HOV theory is skilled labor.

Despite these apparent success of these studies, the underlying very restrictive assumptions of the Hecksher Ohlin model of trade (no transportation costs, no productivity differences) makes the precise estimation of the HOV prediction very difficult. A weaker test of the theory of comparative advantage uses bilateral trade flows and their embodied factor content: comparing all bilateral trade flows, on average, skill abundant nations should be net exporters of skilled labor. This arguably weaker prediction receives strong support in several studies (see Debaere (2003) and Choi and Krishna (2004)). Romalis (2004) shows how countries capture a larger global market share in sectors that intensively use their abundant factors. Romalis also establishes that the predictions of Hecksher Ohlin theory hold qualitatively in the context of Krugman's (1980) model of monopolistic competition and transport costs. The finding that

14

endowments shape the bilateral factor content of trade and Romalis' results are the starting point for the empirical section of this paper. I use the bilateral trade data from the World Trade Database and the US productivity matrix to construct the factor content of bilateral trade. I then use national factor endowment data from Barro and Lee (1994) as well as geographic data to develop a gravity model relating the size of the factor content of bilateral trade to distance, to population size and – most importantly – to the abundance of the respective factor in the exporting nation. In line with the results of the current literature, I find that one can predict the factor content of bilateral trade rather well using geographical data and factor supply differentials.

A potential problem with estimating the relation between the observed factor content of trade and domestic education decisions is the establishing causality: even in a static version of the model with fixed supply of skilled and unskilled labor would a measure of the factor content of trade be correlated with domestic education levels. To deal with this endogeneity, I only use the information of a country's geographic proximity to international supply of skilled and unskilled labor to instrument for the observed factor content of trade. I do not use any domestic information except population size. In this way I isolate the component of international trade that is not stemming from domestic supply and demand, but exclusively from the factor supply of other nations. This empirical strategy is related to Frankel and Romer (1999), who isolate the geographic component of trade to establish a causal relation between trade and growth. Similarly, my constructed measures reflect how much skilled and unskilled labor other countries are likely to export to a given nation, and I subsequently test whether this measure of geographic proximity to skilled and unskilled labor has significant effects on domestic education decisions and the stock of human capital.

I first evaluate effects my measures have on education levels. The unconditional correlation between instrumented factor imports and education levels is significant, but only marginally. I therefore condition on levels economic development, and show that education is strongly affected by geographical proximity to skilled and unskilled labor given the level of a nation's development. In accordance with my theory, levels of higher education (average years in the population) are negatively affected by proximity to international skilled labor, while proximity to unskilled labor has a positive effect on education. Interestingly, the same channel is not

present for levels of primary education (again in average years in the population). Rather, when I find a relation between proximity to international supply of human capital and levels of primary education, I find the reverse relationship: countries that trade a lot with skill abundant nations tend to accumulate more primary education. The next question of interest is whether the relation between trade and domestic education has become more pronounced with the increasing importance of international trade over the last 40 years. I test this formally by evaluating changes of the level of education from 1960 to 1990. Again, I find that proximity to skilled labor had a negative effect on changes of average advanced education, while this channel is not present for primary education. In total, I conclude that the theory proposed in this paper is supported by the data, and using back-of-the-envelope calculations I also show that these effects are quite sizeable: between two otherwise identical nations, a one standard deviation difference in geographic proximity to skilled labor is associated with a difference of 0.6 - 0.7 in the average number years of higher education per worker, which corresponds to a difference in GDP in the order of magnitude of 5%.

The structure of the paper is as follows. Section 1.2 develops a general equilibrium model where heterogeneous and finitely lived workers invest in their human capital. Section 1.3 characterizes the resulting autarky equilibrium. Section 1.4 establishes the path of income as well as the welfare effects of opening a small economy to trade. Section 1.5 endogenizes world prices and establishes the skill bias of world trade. Section 1.6 describes the path of world development after globalizing markets in a world composed of many nations. Section 1.7 presents the empirical results and Section 1.8 concludes.

## 1.2 Preferences, Production Relations and Demography

This section describes the economic environment. The model is formulated in continuous time, which is indexed by $t$ ($t \geq 0$). The world economy consists of many small countries that are indexed by $i$. Each country $i$ has mass 1 of identically and infinitely lived households. Each household is composed of a mass of heterogenous and finitely lived workers. I describe the formation of skills below. Households make the education decisions for workers, which is described below. Households have stable preferences over consumption that are additive, time

16

separable and exhibit a constant rate of time preference.

$$V(t,i) = \int\limits_{t}^{\infty} U(C_{\tau,i}) e^{-\delta(\tau-t)} d\tau \qquad (1.1)$$

I assume that $U$ is strictly increasing, strictly concave and twice continuously differentiable, with $U'(0) = \infty$. Infinite marginal utility at $C_{\tau,i} = 0$ is assumed for convenience so that the economy is never on a path where investment is 0 for all times. A standard budget constraint applies, which restricts the net present cost of the path of consumption to being at most as big as the net present value of future income. Let $Y_{i,t}$ denote a country's production. The budget constraint of each household is given by

$$\int\limits_{t}^{\infty} C_{\tau,i} e^{-\int_{t}^{\tau} r_{\nu} d\nu} d\tau \leq \int\limits_{t}^{\infty} Y_{i,\tau} e^{-\int_{t}^{\tau} r_{\nu} d\nu} d\tau + B_{i,t} \qquad (1.2)$$

The interest rate $r_t$ is not country specific, i.e. well developed global capital markets exist. $B_{t,i}$ denotes the net asset position of country $i$.[3]

Final output $Y$ is defined over a constant elasticity of substitution (CES) aggregate of a skill intensive and a labor intensive good. Denoting the amount of the labor intensive intermediate good used in production by $X_{l,i}$ and the amount of the human capital intensive good by $X_{h,i}$, final output in country $i$ is given by[4]

$$Y_i = \left( X_{l,i}^{\beta} + X_{h,i}^{\beta} \right)^{\frac{1}{\beta}} \qquad (1.3)$$

The final good is produced competitively. The elasticity of substitution between the two intermediate goods is constant and equal to $(1-\beta)^{-1}$. Throughout the analysis, I assume that the intermediate goods are gross substitutes.

**Assumption 1.** $0 < \beta < 1$

Assumption 1 implies that price effects are not too strong. In equilibrium, a human capital

---

[3] This implies that final output can always be traded so that countries can borrow, lend and repay to each other.

[4] For simplicity, (1.3) omitts the distribution parameters normally present in the CES production function.

abundant economy is characterized by a low price of skill intensive goods but still larger total expenditures on skill intensive goods than a labor abundant economy. Autor et al. (1998) have estimated the elasticity between skilled and unskilled labor directly. They conclude that it is unlikely to fall outside the interval $[1, 2]$, which in this model corresponds to $0 < \beta < 0.5$. I denote the prices of the two intermediate goods in country $i$ by $p_{l,i}$ and $p_{h,i}$. Normalizing the price of the final good to unity implies

$$p_{h,i}^{-\frac{\beta}{1-\beta}} + p_{l,i}^{-\frac{\beta}{1-\beta}} = 1 \tag{1.4}$$

The relative price differs across countries when there is no international trade. The two intermediate goods are produced from two factors, human capital and"raw" unskilled labor. Human capital $H_i$ can be used to produce the skill intensive good using a linear transformation technology. Labor $L_i$ can be used to produce the labor intensive good using a linear transformation technology. I sometimes refer to these two goods as the skill intensive sector and the labor intensive sector respectively. While raw labor can be used equally efficiently in all countries, I assume that the effectiveness of human capital depends on some exogenously given, country-specific parameter $A_i$ that is stable over time.[5]

I denote the output of the skill intensive good in country $i$ by $Y_{h,i}$ and the output of the labor intensive good by $Y_{l,i}$.

$$Y_{l,i} = L_i \quad \text{and} \quad Y_{h,i} = A_i H_i \tag{1.5}$$

The two intermediate goods are produced competitively. There are no factors of production other than human capital and labor. Equation (1.5) incorporates the simplification that production in each sector requires either only unskilled labor or only human capital. Ventura (1997) also evaluates the case of goods that require both factors at different intensities. He concludes that the basic results are unchanged if countries have similar enough factor supplies such that factor price equalization holds when countries trade.

I now turn to the supply of skilled and unskilled labor. Each household consists of a mass of

---

[5]These cross country differences in $A_i$ can be seen as stemming from differences in the institutional setup of a country, see Caselli and Coleman (2005). Appendix B endogenizes the level of technolgy.

heterogenous and finitely lived workers. Per household and unit of time, a mass of $\delta$ workers is born. Young workers are of type $\theta$ and can spend time educating themselves. If they choose to get an education, they enter the labor force after a fixed period of time $T$ and start supplying $0 \le \alpha \le 1$ units of unskilled labor and $\theta$ units of skilled labor. Workers that do not get an education supply one unit of skilled labor from their moment of birth. For each type $\theta$ and at each moment of time, households decide whether the worker does get an education or not. Let $h(t, i, \theta)$ denote the education decision for a worker of type $\theta$ in country $i$ at time $t$. $h(t, i, \theta)$ equals 1 if the worker gets an education and 0 otherwise. There is no cost of education other than time spent in school. Also, there is no utility from getting an education or working. After entering the Labor force, all agents face a constant and age-independent rate of death $\delta$. This convenient structure of the life cycle ensures that the size of a country's working population and the demographic composition are constant along any stationary equilibrium.

Types are distributed equally in all households and countries with a Pareto density function with shape parameter $\frac{1}{1-\eta}$ and scale parameter $\eta c$.

$$F(\theta) = 1 - \left(\frac{\eta c}{\theta}\right)^{\frac{1}{1-\eta}} \tag{1.6}$$

The parameter restrictions $0 < \eta < 1$ and $0 < c$ as well as the lower bound of $\eta c \le \theta$ apply. A lower $\eta$ is associated with more heterogenous workers. The scale parameter in (1.6) is chosen such that $\eta$ does not affect the average type and it is always true that $E(\theta) = c$. With this formulation, a decrease of $\eta$ is a mean preserving spread of the distribution of types. In equilibrium, $\eta$ affects how different countries are and also determines the surplus from education.

Human capital and unskilled Labor of different workers are perfectly substitutable. Due to this and the fact that all workers die with equal probability, the supply of human capital and labor is completely described by the size of the current Labor force. The total supply of human capital is given by the sum over past education decisions adjusted for types, the probability of survival and whether a worker is currently schooling or working.

$$H_{i,t} = \delta \int_{-\infty}^{t} e^{-(t-(T+\tau))\delta} \int_{\theta} f(\theta)\, \Upsilon_{t,i,\tau} h(t, i, \theta)\, \theta d\theta d\tau \tag{1.7}$$

19

Where $\Upsilon_{\tau,i}$ denotes the indicator function that equals 1 if a worker has left school and 0 otherwise. Since education is restricted to take place at the beginning of an individuals' life, $\Upsilon_{\tau,i}$ takes the value 1 whenever $\tau \leq t - T$. Similarly, the supply of labor takes into consideration that some agents are currently at school.

$$
\begin{aligned}
L_{t,i} = {} & \delta \int_{-\infty}^{t} e^{-\delta(t-\tau)} \int_{\theta} f(\theta)(1 - h(t,i,\theta))\, d\theta d\tau \\
& + \delta \int_{-\infty}^{t} \alpha e^{-(t-(T+\tau))\delta} \int_{\theta} f(\theta)\, \Upsilon_{t,i,\tau} h(t,i,\theta)\, d\theta d\tau
\end{aligned}
\tag{1.8}
$$

Supply of services from labor $L$ comes from two groups: unskilled workers and skilled workers who have finished their education.

## 1.3 Autarky Wage Patterns

This section establishes the equilibrium in a closed economy. Before solving for the stationary equilibrium path of the economy in autarky, I establish the instantaneous competitive equilibrium. Thereafter, I establish a stationary equilibrium and explain the origin of income and consumption differences in autarky.

**Definition 1** *A feasible autarky allocation in country i given the supply of labor (1.8) and the supply of human capital (1.7), consists of functions $[h(t,i,\theta), Y_{i,t}, C_{t,i}]$ that satisfy (1.5) and (1.2) such the integral over (1.1) is finite and well defined. A resource constraint restricting input use in (1.3) to $X_{l,i} \leq Y_{l,i}$ and $X_{h,i} \leq Y_{h,i}$ applies.*

At each point in time $t$, there are perfectly competitive spot markets for the two intermediates and the final good. Non-satiation of the instantaneous utility together with the strictly positive marginal product of inputs in (1.3) ensures that all inequalities hold. I first establish the instantaneous equilibrium given factor supplies. For simplicity, I drop time subscripts $t$ unless there is danger of confusion. I denote the wage of raw labor by $w_{l,i}$, the factor return of one unit of human capital by $w_{h,i}$ and the relative wage by $w_i \equiv \frac{w_{h,i}}{w_{l,i}}$. Profit maximization

by competitive final goods producers (1.3) relates the relative price of intermediate goods to relative input use. Also, I denote the relative prices of the skill intensive good in country $i$ by $p_i$.

$$p_i \equiv \frac{p_{h,i}}{p_{l,i}} = \left(\frac{Y_{h,i}}{Y_{l,i}}\right)^{-(1-\beta)} \tag{1.9}$$

Intermediate goods are produced using a linear transformation technology and (1.9) also determines the relative wage.

$$w_i = A_i^\beta \left(\frac{H_i}{L_i}\right)^{-(1-\beta)} \tag{1.10}$$

The relative wage is increasing in the efficiency of technology but decreasing in the relative abundance if human capital. Since the price of the final good is normalized to 1, the relative price $p_i$ alone pins down $p_{l,i}$ and $p_{h,i}$ and consequently also wages.

Each household chooses a strategy taking the strategy of other households in the economy as given. A strategy for a household is a subset of each cohort of workers that are sent to the educational sector and the intertemporal consumption decision. I evaluate first the education decision $h(t, i, \theta)$ of each household. Since there exist perfect capital markets, each household maximizes the net present flow of labor income from each worker. Denote by $N(t, i, \theta, h)$ the net present value of the lifetime income that a worker of type $\theta$ born at $t$ in country $i$ receives when the education decision is $h(t, i, \theta)$. Income is discounted to the point of birth $t$ of the respective worker and equal to

$$N_{t,i}(\theta, h) = \begin{cases} \int\limits_t^\infty w_{l,\tau,i} e^{-\int_t^\tau \delta + r(\nu)d\nu} d\tau & \text{if } h(t, i, \theta) = 0 \\ \int\limits_{t+T}^\infty (\theta w_{h,\tau,i} + \alpha w_{l,\tau,i}) e^{+\delta T - \int_t^\tau \delta + r(\nu)d\nu} d\tau & \text{if } h(t, i, \theta) = 1 \end{cases} \tag{1.11}$$

The effective cost of education is giving up the unskilled wage from time $t$ to $t + T$ and a share $(1 - \alpha)$ of unskilled labor income thereafter. The benefit is the additional income equal to $\theta$ times the skilled wage from time $t + T$ on. Along any path of the economy, (1.11) leads to a single crossing property of the type and the education decision of a household. If it is optimal for a household to choose $h(t, i, \theta) = 1$, then the same is true for any other type $\theta' > \theta$. Therefore, there exists a cutoff level $\bar{\theta}_{i,t}$ such that all types $\theta \geq \bar{\theta}_{i,t}$ get an education and all other types do not. The main sections of the paper are concerned with across-countries comparison of the

21

aggregate gains from trade. I therefore define the aggregate net present income from the current cohort of workers $I_{t,i}$. Total income is equal to the integration of the maximal income (1.11) over types. This defines the discounted flow of income from the current generation of workers, which is of mass $\delta$.

$$I_{t,i} \equiv \delta \int_{\theta} f(\theta) \max_{h(t,i,\theta)} N_{t,i}(\theta, h)\, d\theta \tag{1.12}$$

There is no aggregate uncertainty in this economy. Given (1.12) for past, present and future generations, the household has a separate consumption decision. Optimization of intertemporal utility (1.1) subject to (1.2) yields a familiar result for the slope of the consumption process.

**Definition 2** *A competitive static equilibrium, given by the initial stock of human capital (1.7), labor (1.8) and $A_i$ consists of a feasible allocation of functions for $[c(\tau, i), T(t, i), r(t), p(x_i)]$ such that (1.9) and (1.10) hold, $h(t, i, \theta)$ maximizes lifetime income for all cohorts (1.12) and the path of consumption maximizes (1.1) subject to (1.2).*

I next consider the existence and uniqueness of a stationary equilibrium (SE) in autarky. Let an "$A$" superscript denote expressions along such a stationary equilibrium, in which the relative price is constant and equal to $p_i^A$, the relative wage is a function of $A_i$ and $p_i$ and the interest rate is stable. Households choose a cutoff level $\overline{\theta}_i^A$ and, since there is no technological progress, output and consumption are constant. Convergence to a stationary equilibrium is established easily because investment and intertemporal consumption decisions are independent. First, evaluate the cutoff condition (2.28) along any path of development. A single household has no influence on the relative wages or interest rates. Even if is optimal to school all types of workers, there is still a well defined and finite supply of unskilled and skilled labor for any path of wages and interest rates that leads to a finite net discounted value of income. Arbitrage considerations ensure a non-negative rate of interest at all moments of time. A nonzero interest rate combined with a positive rate of death $\delta$ implies that the discounted value of income is finite for any worker. Hence intertemporal income of a household is always defined. By standard arguments, time separable and concave preferences combined with a constant rate of time preference lead to a constant interest rate of $r = \rho$ along any path where income is stable. If $\rho > 0$, a unique and stable stationary equilibrium exists in which the choice of the cutoff point is a constant function of the interest rate and the autarky wages $w_{h,i}^A$ and $w_{l,i}^A$. Evaluating the entry condition

(1.11) at the worker of type $\theta = \overline{\theta}_i^A$ who is indifferent between going to school or not, this cutoff level solves

$$w_{l,i}^A = e^{-\rho T} \left( \overline{\theta}_i^A w_{h,i}^A + \alpha w_{l,i}^A \right) \tag{1.13}$$

Given the optimal choice of $\overline{\theta}_i^A$, one can solve for the maximal net present value of income from the present cohort of workers, which if given by (1.12) in autarky. Along any path of the economy with constant wages and cutoff level $\overline{\theta}_i$, I denote the net present value of income from the current cohort of workers by $I\left(\overline{\theta}_i, w_{l,i}, w_{h,i}\right)$. Without assuming any specific distribution of types, it is always possible to express the net present income of a cohort of workers depending exclusively on the two wages. Evaluated at $\overline{\theta}_i^A$, the total income discounted to the point of birth of a generation of workers is equal to

$$I\left(\overline{\theta}_i^A, w_{l,i}, w_{h,i}\right) = \frac{\delta}{\rho + \delta} \left( 1 + e^{-\rho T} \frac{w_{h,i}}{w_{l,i}} \int_{\overline{\theta}_i^A}^{\infty} f(\theta) \left( \theta - \overline{\theta}_i^A \right) d\theta \right) w_{l,i} \tag{1.14}$$

For any relative wage $w_i = \frac{w_{h,i}}{w_{l,i}}$, income is at least equal to $\frac{\delta}{\rho+\delta} w_{l,i}$. There are $\delta$ young workers who could start working right away and earn the unskilled wage forever, where the future is discounted at rate $\rho + \delta$ to account for the probability of death. Secondly, for any $w_{h,i} > 0$, there may exist high type agents that find it worthy to get an education. The marginal worker of type $\theta = \overline{\theta}_i^A$ just breaks even on his educational investment, but for all workers of higher type $\theta$, the possibility to get educated increases their lifetime income. It is important to note that the aggregate surplus from having access to an education, which is represented by the second term in (1.14), is more than proportionally increasing in the relative wage $w_i$: if the relative wage increases, there are two margins in which net income from education is affected. The increased relative wage benefits all worker proportionally that would have chosen to get educated at lower wages. In addition, if the relative wage increases, the optimal cutoff level $\overline{\theta}_i^A$ decreases, hence benefiting the additional entrants (weakly). An increase of the relative wage - given the unskilled wage - hence results in a more than proportional increase in the net income from education. In the case of no heterogeneity of workers (this corresponds to $\eta \to 1$ in the specific case of the Pareto distribution (1.6)) there is no surplus from education. In this case, the model becomes very similar to that of Findlay and Kierzkowski (1983) and all

workers earn the unskilled wage (1.16). Intrinsic cross country differences in $A_i$ hence do no longer matter because different workers earn different wages, but exclusively through general equilibrium effects that influence the unskilled wage.

I now solve for general equilibrium prices, wages and level of income (1.14) in the case of the Pareto distribution of types (1.6) of an economy in autarky. For the rest of the paper, I also assume that $\alpha = 1$. This assumption allows a closed form solution of the supply of unskilled labor $L_i^A$ and the resulting relative wage.

In the autarky stationary equilibrium the only source of cross-country variation is $A_i$. Solving the supply of labor (1.8) and human capital (1.7) for the constant cutoff $\overline{\theta}_i^A$, factor supply is given by

$$L_i^A = 1 \quad \text{and} \quad H_i^A = \lambda^{\frac{1}{\beta}} A_i^{\frac{\eta\beta}{1-\eta\beta}} \tag{1.15}$$

where $\lambda \equiv \eta^{\frac{\beta\eta}{1-\eta\beta}} \left(e^{\rho T} - 1\right)^{-\frac{\beta\eta}{1-\eta\beta}} c^{\frac{\beta}{1-\eta\beta}}$. In equilibrium, the higher a country's relative efficiency of human capital $A_i$, the more skill abundant a country is. With the supply of factors given, prices (1.9) and consequently wages (1.10) are determined uniquely. In autarky, skill abundant countries have a lower relative price of the skill intensive good, but still a higher relative wage.

The relative abundance of factors, technology and the normalization of the final good (1.4) relate the equilibrium unskilled wage $w_{l,i}^A$ to the level of domestic skill complementary technology $A_i$.

$$w_{l,i}^A = \left(1 + \lambda A_i^{\frac{\beta}{1-\eta\beta}}\right)^{\frac{1-\beta}{\beta}} \tag{1.16}$$

A country that is characterized by a high $A_i$ has a low autarky price of the skill intensive good. Because the normalization of the final good relates relative and absolute prices one to one, the price of the labor intensive good is high in these countries. Since each unit of raw labor can produce one unit of the unskilled good it thus receives a high wage. I denote stationary output by $Y\left(\overline{\theta}_i, w_{l,i}, w_{h,i}\right)$, which in autarky is equal to

$$Y\left(\overline{\theta}_i^A, w_{l,i}^A, w_{h,i}^A\right) = \left(1 + \lambda A_i^{\frac{\beta}{1-\eta\beta}}\right) w_{l,i}^A \tag{1.17}$$

In equilibrium, a country that is characterized by a high efficiency of human capital has a high level of net income (1.18), i.e. it is "rich". The stationary net present income (1.12) of

young cohort of workers is equal to the total income from skilled labor plus the net income from human capital.

$$I\left(\overline{\theta}_i^A, w_{l,i}^A, w_{h,i}^A\right) = \frac{\delta}{\rho + \delta}\left(1 + e^{-\rho T}\left(1 - \eta\right)\lambda A_i^{\frac{\beta}{1 - \eta\beta}}\right)w_{l,i}^A \qquad (1.18)$$

High efficiency countries have a high level of net income and are rich. Because of the convenient Pareto distribution of types, the *net* income from human capital is equal to a fraction $e^{-\rho T}\left(1 - \eta\right)$ of the *total* income from skilled labor services.

How does the heterogeneity of workers influence the lifetime income of a cohort of workers? Consider first the case of homogenous types $(\eta \to 1)$, in which all workers earn $w_{l,i}^A$. The model then becomes very similar to that of Findlay and Kierzkowski (1983). All workers earn the unskilled wage (1.16) and technology differences matter only through relative supply and price effects: a country with high $A_i$ is characterized by a high supply of human capital and hence a lower price of the skill intensive good. A low price of the skill intensive good implies a high price of the labor intensive good and consequently a high unskilled wage. Consider now the case of a decrease in $\eta$, i.e. a mean preserving spread of the distribution of types. In autarky equilibrium a low $\eta$ is associated with a large share of surplus as a fraction of total revenue of the skill intensive sector.[6]

More important than the effects $\eta$ has on absolute levels of income and output is the impact it has on relative cross-country differences. Nations intrinsically only differ with respect to their level of human capital augmenting technology $A_i$. The heterogeneity of workers guides how differences in technology translate into differences of income and factor abundance. If types are similar, small differences in human capital efficiency translate into large differences of relative factor abundance and income. If workers are very heterogenous, differences in $A_i$ translate into only moderate differences of factor endowment: the more spread the distribution of types is, the lower is the density of workers at any point along the distribution $f(\theta)$. For a given intrinsic difference in $A_i$ and therefore in the relative demand for factors and in the

---

[6]For given wages and therefore cutoff level $\overline{\theta}_i$, the supply of skilled workers (1.7) is lower if types are more heterogenous. Although the expected value of the distribution of types is unaffected by $\eta$, the truncated expected value (that is the expected value given that the type is higher than $\overline{\theta}_i$) actually increases with $\eta$. This effect is captured in the value of $\lambda$.

cutoff point $\bar{\theta}_i$, the resulting international dispersion of relative factor supply is large if the distribution of workers is homogenous.

Cross country differences are influenced by the elasticity of substitution between the skill and labor intensive intermediate goods. Consider first the case of $\beta$ bigger than, but close to 0. In this case, price effects in (1.3) are offsetting differences in technology and countries have nearly identical factor supplies. Countries thus only differ in their level of technology and hence output.[7] A higher beta is associated with weaker prices effects and thus increasingly pronounced cross country differences in autarky factor supply. In the case of $\beta = 1$ the production of the final good (1.3) is linear in inputs used, relative input prices are fixed and therefore international factor abundance levels are very different. The level of $\beta$ also determines the size of gains from trade, which are derived in the next section.

## 1.4 Trade and the Evolution of Income and Welfare

The notion that exchange - if it happens - must benefit all involved parties is an axiomatic insight of economic theory, and the same should be true for exchange between countries, international trade. But how are these gains from trade split up between nations at different stages of their economic development? This section establishes the gains from opening to trade. This is done in a partial equilibrium setting taking as given world prices. Global prices are derived in the next section. I focus on relative effects that occur to 'poor' and 'rich' countries. The structure of the present section is the following. First, as a benchmark model, I establish the gains from trade that would prevail in a world where education decisions are fixed at autarky levels. This is equivalent to welfare effects in a standard Heckscher Ohlin model of trade with factors of production in fixed supply. In this static setting, a country gains from trade because it is different from the rest of the world.

I show that the initial gains from trade are likely to lead to neutral gains from trade that favour neither developed nor developing countries and leave the relative dispersion of income unchanged. Dynamically, one has to distinguish between income divergence and divergence of

---

[7]In the case of $\beta = 0$, (1.3) takes the Cobb Douglas shape with expenditure share of $1/2$ for each sector. As is well known, the factor augmenting productivity $A_i$ is in this case equivalent to a Hick's neutral productivity level of $\sqrt{A_i}$.

welfare. I first describe the evolution of income. After opening markets, there is a phase of convergence of income that reflects the increased investment activity in richer countries and the decrease in education in other countries. After a period of time $T$ the increased investment in human capital translates into again diverging income. I establish that the steady state of an open world is characterized by larger differences in human capital abundance and also larger output differences than in a world of closed economies. I then turn to establish the evolution welfare. There are always additional efficiency gains that occur to countries because the education decision can adjust to international prices. However, because of the way in which trade affects the surplus from education, there is always dynamic divergence of welfare compared to the moment just after opening to trade. Finally, I develop conditions for when trade leads to absolute divergence of welfare compared to autarky and argue that these are likely to hold in reality.

Assume a small country $i$ has a level of human capital efficiency of $A_i$ and is in its autarky stationary equilibrium. At point in time $\tau^*$, markets are unanticipatedly opened to trade with a large world that is characterized by $A_w$ and a resulting relative price of the skill intensive good $p_w = \lambda^{-\frac{1-\beta}{\beta}} A_w^{-\frac{(1-\beta)}{1-\eta\beta}}$. $A_w$ will be endogenized in the next section. Instantaneously after opening to trade, output of country $i$ is given by autarky factor supplies (1.7) and (1.8), but valued at international prices.

$$Y\left(\overline{\theta}_i^A, w_{l,w}, A_i p_w\right) = \left(1 + \lambda \left(\frac{A_i}{A_w}\right)^{\frac{1}{1-\eta\beta}} A_w^{\frac{\beta}{1-\eta\beta}}\right) \left(1 + \lambda A_w^{\frac{\beta}{1-\eta\beta}}\right)^{\frac{1-\beta}{\beta}} \tag{1.19}$$

Opening to trade has two effects on income: it influences both relative wage $w_i$ and the unskilled wage $w_{l,i}$. These two effects always work in opposite direction. If a country is more skill abundant than the rest of the world ($\frac{A_i}{A_w} > 1$), it benefits from trade because the relative wage $w_i$ increases, but at the same time looses from trade because the unskilled wage decreases. The opposite is true for a country $j$ that is less skill abundant than the rest of the world.

It is important to point out that net effect always results in an increase of output. This can formally be shown by evaluating the first order condition of the ratio of (1.19) divided by (1.17) with respect to $A_i$. The minimum level of this ratio is equal to 1 and occurs at $A_i = A_w$. A country that happens to have autarky prices that are equal to the rest of the world is not

affected by trade. Statically, all other countries strictly gain from trade. Evaluating the second order condition of the above ratio establishes that countries that are more different from the rest of the world gain relatively more from trade. The intuition for this result follows from standard trade theory. Each country faces a concave frontier of how much it can supply of the two factors and because there are no market failures, the current supply is on and not inside this frontier. Statically, factor supply is fixed, but trade can change the relative price. At any relative price, the input constraint of final goods producers under trade passes through the current factor supply (1.5), is tangent to the concave factor supply frontier and hence encompasses the latter. Trade enables producers to a strictly larger set of input bundles, and since production isoquants are convex, output increases.

A second point of interest is whether at the moment of trade, it is poor or rich nations that benefit relatively more. A statement on convergence or divergence involves comparing income differences before and after opening to trade, i.e. four different levels of income. To establish the direction of relative gains from trade, I evaluate income differences for two small economies, a country form the north ($n$) and a country from the south ($s$). I assume that the North is skill abundant compared to the rest of the world, so that that $A_n = (1 + \gamma) A_w$, where $\gamma > 0$. South is skill scarce and I assume that $A_s = (1 + \gamma)^{-1} A_w$. $N$ and $S$ are hence symmetrically different from the rest of the world. If for every pair of countries defined in this way there is divergence of output, I speak about uniform relative divergence.

**Definition 3 (Uniform relative Di- and Convergence)** *Let $n$ and $s$ be two small countries with $A_n = (1 + \gamma) A_w = (1 + \gamma)^2 A_s$. There is uniform relative divergence (convergence) of output if trade results in an increase (decrease) of relative income differentials for every $\gamma > 0$ and for every $A_w$.*

The definition of uniform divergence at hand allows to make statements about how the world income distribution would evolve for a given level of world relative output and prices. At the moment, statements of convergence or divergence will be made for each pair of countries. If, for all of these hypothetical pairs, opening to trade differences in output and net present income are increased, one can make statements of the world distribution of income.

The appealing feature of the definition at hand is that it helps to establish for which range

of world prices there will be divergence when opening to trade. For now, this comparison is of course a hypothetical one since world prices have to be determined endogenously. This problem is tackled in the next section. The following lemma establishes instantaneous effects from trade.

**Lemma 1 (Static Output Effects of Trade)** *Consider the moment of opening to trade $\tau^*$. There is uniform relative convergence (divergence) of output if the global size of the labor intensive sector is smaller (bigger) than the human capital intensive one.*

**Proof.** Appendix A establishes that

$$\frac{Y\left(\overline{\theta}_n^A, w_{l,w}, A_n p_w\right) \Big/ Y\left(\overline{\theta}_s^A, w_{l,w}, A_s p_w\right)}{Y\left(\overline{\theta}_n^A, w_{l,n}^A, w_{h,n}^A\right) \Big/ Y\left(\overline{\theta}_s^A, w_{l,s}^A, w_{h,s}^A\right)} \begin{cases} \geq 1 \text{ if } \lambda A_w^{\frac{\beta}{1-\eta\beta}} \leq 1 \\ < 1 \text{ if } \lambda A_w^{\frac{\beta}{1-\eta\beta}} > 1 \end{cases}$$

It is also true that if $\lambda A_w^{\frac{\beta}{1-\eta\beta}} > 1$, the skill intensive sector is larger in terms of output and revenue than the labor intensive sector. ■

If the skill intensive sector is large there is divergence, i.e. poor nations gain more from trade if their sector of specialization is relatively unimportant. This result seems striking at first sight, but thinking in terms of wages offers a good intuition. If $\lambda A_w^{\frac{\beta}{1-\eta\beta}} > 1$ the gains for unskilled labor are relatively large because labor is a globally scarce factor. Poor countries that export labor hence benefit more from trade than do rich countries. Mankiw et al. (1992) estimate that the global expenditure share of the human capital is about as big as the one on pure labor services. A similar comparison can be made from the calculations of Hall and Jones (1999): estimates suggest that the two sectors are of about the same size. Hence, trade is in a static sense neither likely to favour poor nor rich nations and results in uniform gains from trade.

How is the path of income affected after opening to trade? Throughout the following analysis, I denote open economy expressions with an "$O$" superscript. The new optimal cutoff level for the education decision is hence denoted by $\overline{\theta}_i^O$ and solves

$$w_{l,w} = e^{-\rho T}\left(\overline{\theta}_i^O A_i p_{h,w} + w_{l,w}\right) \tag{1.20}$$

In the period from $\tau^*$ to $\tau^* + T$, any country $n$ with $\frac{A_n}{A_w} > 1$ has a level of education investment

that is larger than in autarky. During this period of time, the increased rate of schooling decreases the level of output in these countries: the supply of human capital is still fixed at autarky levels because skilled workers born before $\tau^*$ enter the Labor market.

At the same time, an increased number of young workers chooses to get an education, leading to a temporary decrease of $L_i$ and consequently output. At point in time $\tau^*$, the rate of change of income is equal to equal to

$$\frac{\partial Y_i}{\partial t} |_{t=\tau^*} = \delta w_{l,w} \left( F\left(\overline{\theta}_i^O\right) - F\left(\overline{\theta}_i^A\right) \right) \tag{1.21}$$

Because the supply of skilled labor lags school enrollment rates by a period of time $T$, (1.39) is positive for a poor country with $\frac{A_i}{A_w} < 1$: in these economies, investment decreases instantaneously at $\tau^*$, leading to a temporary 'overshooting' of output[8]. Only after point in time $\tau^* + T$ does the increased investment in human capital start to pay of as workers that started their schooling at $\tau^*$ enter the skilled labor force. After this point in time, there is divergence of output.

The resulting long term level of output is given by

$$Y\left(\overline{\theta}_i^O, w_{l,w}, A_i p_w\right) = \left(1 + \lambda \left(\frac{A_i}{A_w}\right)^{\frac{1}{1-\eta}} A_w^{\frac{\beta}{1-\eta\beta}}\right) \left(1 + \lambda A_w^{\frac{\beta}{1-\eta\beta}}\right)^{\frac{1-\beta}{\beta}} \tag{1.22}$$

For any country $n$ with $\frac{A_n}{A_w} > 1$ the long term level of output is necessarily bigger than the one prevailing at the moment of opening to trade. This reflects the increased investment activity compared to autarky. Similarly, the long term level of output under trade for any country $s$ with $\frac{A_s}{A_w} < 1$ is necessarily smaller than the one prevailing just after autarky. The following proposition summarizes trade-induced changes of output after opening to trade.

**Proposition 1 (Trade and the Dynamics of Income)** *Let $n$ and $s$ be two small countries with $A_n = (1 + \gamma) A_w = (1 + \gamma)^2 A_s$. There is uniform relative divergence of output comparing the output just after opening to trade (1.19) to the one in the stationary equilibrium under*

---

[8]Depending on the rate of death $\delta$ and the time required for education $T$ this effect can be very pronounced, and even lead to temporary reversals of income levels of rich and poor nations.

*free trade (1.22). There is also uniform relative divergence of output comparing the output in autarky stationary equilibrium (1.17) to the stationary equilibrium under free trade (1.22).*

**Proof.** see Appendix A ■

The results of diverging output after opening to trade are straightforward. Trade increases investment rates in rich countries while it decreases them in poor countries. Naturally, an open world is characterized by a more stratified distribution of incomes.



Figure 1: Increases in educational expenditures and thus forgone income in richer countries result in a period of intial convergence (from $\tau^*$ to $\tau^* + T$). Because educational investments start to pay of after $\tau^* + T$, therafter, income diverges. The long run dispersion of income levels is larger than just after opening markets to trade.

Figure 1 displays the path of output for two countries $n$ and $s$ as they have been defined previously. The path of outcome displays first convergence by more pronounced divergence. What does the preceding analysis imply for welfare considerations? Because international capital markets exist, at each moment of time, each household simply consumes a fraction $\rho$ of its complete net present value of future flows of income. Therefore, changes in welfare are equivalent to changes in the net present value of income from all cohorts of income. The

comparison is simple for workers that have made their education decision before $\tau^*$. Since their education decision is sunk, the increase of output due to trade is equivalent to the increase of net present income for this group of workers.

For young workers, there are two questions of interest. The first is whether they gain from trade and the second is whether they gain more than they would have if the education choice had not adjusted. First evaluate the net present value of income for cohorts of workers born at or after $\tau^*$ if the cutoff point had not changed from its autarky level $\overline{\theta}_i^A$.

$$I\left(\overline{\theta}_i^A, w_{l,w}, A_i p_w\right) = \frac{\delta}{\rho+\delta}\left(1 + \left(\left(\frac{A_i}{A_w}\right)^{\frac{1}{1-\eta\beta}} - \eta\left(\frac{A_i}{A_w}\right)^{\frac{\beta}{1-\eta\beta}}\right)\lambda e^{-\rho T}A_w^{\frac{\beta}{1-\eta\beta}}\right)w_{l,w} \quad (1.23)$$

Compare this to the level of net present income that the same cohort of workers get from adjusting to the new optimal cutoff level $\overline{\theta}_i^T$.

$$I\left(\overline{\theta}_i^O, w_{l,w}, A_i p_w\right) = \frac{\delta}{\rho+\delta}\left(1 + (1-\eta)\left(\frac{A_i}{A_w}\right)^{\frac{1}{1-\eta}}\lambda e^{-\rho T}A_w^{\frac{\beta}{1-\eta\beta}}\right)w_{l,w} \quad (1.24)$$

(1.23) is the also net present value that a worker born just before $\tau^*$ gets.

**Lemma 2 (Gains From Trade)** *For all $A_i$ and any $A_w$, there are gains from trade also when the cutoff remains at $\overline{\theta}_i^A$. There are additional gains from trade when $\overline{\theta}_i$ adjusts optimally.*

**Proof.** Compare (1.18) , (1.23) and (1.24). It is easily established that

$$I\left(\overline{\theta}_i^O, w_{l,w}, A_i p_w\right) \geq I\left(\overline{\theta}_i^A, w_{l,w}, A_i p_w\right) \geq I\left(\overline{\theta}_i^A, w_{l,i}^A, w_{h,i}^A\right)$$

With equality if $A_i = A_w$  ∎

What happens to relative levels? The following proposition establishes whether there is divergence of net present income.

**Proposition 2 (Post Opening Divergence)** *Let $n$ and $s$ be two small countries with $A_n = (1+\gamma)A_w = (1+\gamma)^2 A_s$. It is always the case that comparing $I\left(\overline{\theta}_i^O, w_{l,w}, A_i p_w\right)$ to $I\left(\overline{\theta}_i^A, w_{l,w}, A_i p_w\right)$, there is uniform relative divergence. There is uniform relative divergence of $I\left(\overline{\theta}_i^O, w_{l,w}, A_i p_w\right)$ and $I\left(\overline{\theta}_i^A, w_{l,i}^A, w_{h,i}^A\right)$ iff*

$$e^{-\rho T} - 1 + \eta e^{-\rho T}\lambda A_w^{\frac{\beta}{1-\eta\beta}} > 0 \quad (1.25)$$

**Proof.** Evaluate the ratio of (1.24) to (1.23) for two countries $N$ and $S$.

$$\frac{I\left(\overline{\theta}_n^O, w_{l,w}, A_n p_w\right) \Big/ I\left(\overline{\theta}_s^O, w_{l,w}, A_s p_w\right)}{I\left(\overline{\theta}_n^A, w_{l,w}, A_n p_w\right) \Big/ I\left(\overline{\theta}_s^A, w_{l,w}, A_s p_w\right)}$$

If $\gamma = 0$, this ratio equal 1. For any $\gamma > 0$, this ratio can be shown to be bigger 1. The second claim involves a similar comparison of (1.24) to (1.18). The equivalent ratio can shown to be bigger 1 for any $\gamma > 0$ if (1.25) holds. ∎

The preceding proposition establishes whether the net present income of young workers diverges when opening to trade. The household receives additional income from old cohorts of workers that were born before $\tau^*$. To establish whether the total net present income of the economy diverges, one has to evaluate the total relative increase in consumption, which is a combination of contributions from generations born before $\tau^*$ and from younger cohorts born thereafter. The total net present value of all future income of country $i$ is given by two flows of income. First, there is a flow of $Y\left(\overline{\theta}_i^A, w_{l,w}, A_i p_w\right)$ from old cohorts or workers $\tau^* + T$, the size of old cohorts stays constant at 1 but it decreases at rate $\delta$ thereafter. In addition, starting from $\tau^*$, each moment of time a new cohort of workers of mass $\delta$ is born, receiving a net income of $I\left(\overline{\theta}_i^O, w_{l,w}, A_i p_w\right)$. Consumption smoothing implies that the household consumes a fraction $\rho$ of its net wealth.

The new level of consumption after opening markets to trade is hence given by

$$C_i^T = \left(\left(1 - e^{-\rho T}\right) + \rho \left(\rho + \delta\right)^{-1} e^{-\rho T}\right) Y\left(\overline{\theta}_i^O, w_{l,w}, A_i p_w\right) + I\left(\overline{\theta}_i^O, w_{l,w}, A_i p_w\right) \quad (1.26)$$

When is consumption, and therefore also welfare, likely to diverge?

**Proposition 3 (Trade and Divergence of Welfare)** *Let $n$ and $s$ be two small countries with $A_n = (1 + \gamma) A_w = (1 + \gamma)^2 A_s$. Opening to trade results in uniform divergence of welfare iff*

$$(1 - \beta) e^{-\rho T} - 1 + (1 - \beta) \eta e^{-\rho T} \lambda A_w^{\frac{\beta}{1 - \eta \beta}} > 0 \quad (1.27)$$

**Proof.** In autarky, households would consume (1.17). Again evaluating whether the gains

from trade are bigger for a skill abundant country $N$ than for a skill scarce country $S$, this is true for any $\gamma > 0$ if (1.27) holds. ∎

How likely is trade leading to divergence under realistic parameter values? Consider first the conditions for post opening divergence of net present income (1.25). If the duration of education is sufficiently short or $\rho$ approaches 0, there is always divergence. This result is straightforward: as $e^{-\rho T}$ goes to 1, workers do not have to invest much in order to become skilled. Any human capital accumulation that is induced by trade hence leads to large net gains for human capital abundant countries. If $e^{-\rho T}$ is is substantially below one, there is a significant cost of education. In this case, rich countries are likely to gain more from trade than poor nations if the global skill intensive sector is large compared to the labor intensive sector and if the heterogeneity of workers is small. Why does heterogeneity play this role? Again, the same mechanism that controlled how different countries are in autarky influences how sensitive the supply of skilled labor is to changes in the relative wage induced by trade. Consider a developed country $(A_i > A_w)$. If workers are heterogeneous, for a given change in the wage only a moderate number of additional workers enters the skilled sector. The increase in net income (i.e. in surplus) is only moderate, as well. In contrast, if workers are homogenous, a small increase in the wage induces a sizable entry in the skilled labor supply and consequently a larger increase in the surplus from education.[9]

The condition for total divergence of welfare is similar to the one for post opening divergence. Different countries are more likely to diverge if the time of schooling is short, the human capital intensive sector is relatively important and if workers are more homogenous. In addition, the elasticity of substitution now guides relative divergence. Since the skill supply reacts only slowly to changed demand conditions, there is less likely to be divergence. It is again noteworthy that empirical estimates of beta are small (see Autor et al. (1998)) so that conditions (1.25) and (1.27) are similar.

---

[9]An interesting benchmark is when all workers are identical. In this case trade induces complete specialization and the gains from trade are the following. Workers in poor nations receive the global unskilled wage, while workers in rich nations receive $A_i / A_w$ times the global unskilled wage. Because identical workers earn the unskilled wage, the gains from trade are fundamentally different and depend again only on wage and price effects of trade. Due to this, Lemma 1 (Static Gains From Trade) also describes conditions under which there is con- and divergence in the case of homogenous workers.

## 1.5 General Equilibrium

The analysis so far has shown that given world prices, the dynamic response of education investments leads to divergence of income. It has also established that given world prices, opening the economy to trade might result in overall divergence of welfare. But how are the world prices determined when the global supply of human capital adjusts to factor prices? The general equilibrium effects are indeed different than a Heckscher-Ohlin model with factors of production in fixed supply would predict. Trade reduces education investments in some countries and raises it in other nations. The key feature of interest is that trade induces accumulation of human capital in countries that can use skills relatively efficiently and disaccumulation elsewhere. The average skilled worker hence works in a more skilled country in an open than in a closed economy world. The global output of skill intensive goods increases with trade and therefore decreases the price of these goods. Additionally, both static and dynamic effects of opening markets to trade lead to a skill bias. This part of the paper hence establishes that when human capital accumulation is endogenous, globalization increases the skill premium despite falling prices of human capital intensive goods.

For this section of the paper, I order all countries $i$ by their relative human capital effectiveness $A_i$. I assume that this measure is distributed with probability density function $g(A_i)$. This distribution might also include large countries, and hence the following analysis also encompasses the typical two-country North and South case. I assume that countries are not too different, so that there exists no country that would only have skilled workers in equilibrium.[10] A global competitive equilibrium follows the definition of equilibrium in a closed economy. The global resource constraint restricts total input use to be at most as big as global output of the two intermediate goods.

$$\int_i X_{l,i} di \leq \int_i Y_{l,i} di \quad \text{and} \quad \int_i X_{h,i} di \leq \int_i Y_{h,i} di \tag{1.28}$$

At the moment of opening to trade, global relative supply of factors is given by steady state

---

[10]In the long run equilibrium, this restriction is equivalent to $A_{MAX} < \eta c \left( e^{\rho T} - 1 \right) p_w^{-1}$.

autarky levels.

$$\frac{Y_{h,w}}{Y_{l,w}}\Big|_{t=\tau^*} = \lambda^{\frac{1}{\beta}} \int_{A_i} g\left(A_i\right)\left(A_i\right)^{\frac{1}{1-\eta\beta}} dA_i \tag{1.29}$$

I denote the average world level of human capital efficiency at the moment of opening to trade by $A_w^A$.

$$A_w^A \equiv \left(\int_{A_i} g\left(A_i\right)\left(A_i\right)^{\frac{1}{1-\eta\beta}} dA_i\right)^{1-\eta\beta} \tag{1.30}$$

Instantaneously after $\tau^*$, countries with $A_i > A_w^A$ accumulate further human capital, while other nations disaccumulate. In a stationary equilibrium, each country chooses a level of human capital dependent on its level of $A_i$ and on global prices $H_i^O = \lambda^{\frac{2-\beta}{\beta}} A_i^{\frac{\eta}{1-\eta}} p_w^O$, resulting in a total level of global output of

$$\frac{Y_{h,w}^O}{Y_{l,w}^O} = \lambda^{\frac{1}{\beta}} \left(\int_{A_i} g\left(A_i\right)\left(A_i\right)^{\frac{1}{1-\eta}} dA_i\right)^{\frac{1-\eta}{1-\eta\beta}} \tag{1.31}$$

Similarly to the definition of the average world level of human capital efficiency at the moment of opening to trade $A_w^A$, I denote the average long run global level of human capital efficiency in an open world by $A_w^O$.

$$A_w^O \equiv \left(\int_{A_i} g\left(A_i\right)\left(A_i\right)^{\frac{1}{1-\eta}} dA_i\right)^{1-\eta} \tag{1.32}$$

What is the difference between (1.30) and (1.32)? The key mechanism is that in autarky, the general equilibrium response of prices dampens differences in the supply of human capital: a nation that is characterized by a low $A_i$ has a high price of the skill intensive good, thereby increasing demand. In an open economy, all countries face the same price and cross country differences in the supply of human capital are thus more pronounced. The next proposition establishes the net effect of this concentration of skills.

**Proposition 4 (Expansion of the Skill Intensive Sector)** *The stationary equilibrium under trade is characterized by a larger world production of skill intensive goods than in autarky.*

**Proof.** The dynamic relative supply of skill intensive goods (1.31) is larger than the static

37

one (1.29) if the following inequality holds.

$$\left( \int_{A_i} g\left(A_i\right) \left(A_i\right)^{\frac{1}{1-\eta}} dA_i \right)^{1-\eta} > \left( \int_{A_i} g\left(A_i\right) \left(A_i\right)^{\frac{1}{1-\eta\beta}} dA_i \right)^{1-\eta\beta}$$

By Assumption 1 and the general means inequality, this is always true. ∎

Trade, to a first order, "shifts" skilled workers from low $A$ to high $A$ countries. While trade also reduces skill abundance and the supply of the skill intensive good in some nations, it raises the supply in exactly those countries that can use them very efficiently. This concentration results in an expansion of the skill intensive sector.

A higher relative output implies a lower relative price of the skill intensive good, and one might suspect that trade therefore lowers the average skill premium. Interestingly, the opposite is the case. In the context of the present model, skill bias is not easily established, since some countries may see their skill premium increase with trade, but in other countries there may be a decrease of the relative skilled wage. I therefore define skill bias in an average sense.

**Definition 4 (Pervasive Skill Bias)** *Trade is pervasively skill biased if the arithmetic average of the relative wage of human capital increases with trade.*

There are two questions of interest. First, is there pervasive skill bias at the moment of opening to trade? Second, is there additional skill bias along the dynamic path of the global economy? The following proposition answers both of these questions.

**Proposition 5 (The Skill Bias of Global Trade)** *There is pervasive skill bias at $\tau^*$. The dynamic response of educational investment results in further skill bias.*

**Proof.** Compare the arithmetic average of the skill premium before, at the moment of and long after opening to trade. It is both true that

$$\int_{A_i} g\left(A_i\right) H_i^A A_i p_w \left( \frac{Y_{h,w}}{Y_{l,w}} |_{t=\tau^*} \right) dA_i > \int_{A_i} g\left(A_i\right) H_i^A w_i^A dA_i$$

and

$$\int_{A_i} g\left(A_i\right) H_i^O A_i p_w^O dA_i > \int_{A_i} g\left(A_i\right) H_i^A A_i p_w \left( \frac{Y_{h,w}}{Y_{l,w}} |_{t=\tau^*} \right) dA_i$$

These two inequalities are satisfied by the general means inequality. ∎

The same mechanism that is responsible for the output increase of the skill intensive sector is responsible for the skill bias of trade. Consider first the moment of opening to trade. When goods prices are equalized, the wage increases in skill abundant countries, while it decreases in skill scarce countries. Since all countries have an equal endowment of unskilled labor, this channel is not present for unskilled sector. Additionally, the dynamic response of human capital amplifies the initial skill bias. Trade induces skill accumulation in high wage countries and de accumulation elsewhere. The arithmetic average of the wage hence increases further. This result is related to Dinopolous and Segerstrom (1999) and especially to Gancia and Epifani (2005): while these authors argue that the skill intensive sector is more sensitive to scale, here it is the fact that the factor used intensively in one industry can be accumulated and this, on the aggregate, leads to the industry to having a higher growth potential.

## 1.6  Trade and the Path of Global Development

The two preceding sections have established what the consequences of opening a small economy to trade are given world prices, and how these world prices are determined. However, what is the path of globalization when many nations simultaneously open their markets to trade? This section establishes the path of global development. The main insights are very similar to that of Section 5 evaluated at $A_w^O$ determined in (1.32). Also, there may exist a group of intermediately developed countries that starts the process of globalization as exporters of the skill intensive good and successively become importers of the latter. This results stems from the expansion of the skill intensive sector and the resulting decline of skill intensive goods prices caused by trade.

At the moment of opening to trade, world prices are given by the static world supply (1.29). Lemma 1 (Static Output Effects of Trade) describes the evolution of the world income distribution at the moment of trade, which converges if at the moment of opening to trade, the size of the world skill intensive sector is larger than that of the unskilled sector. The long run distribution of income is given by Proposition 2 (Trade and the Dynamics of Income) and the long term level of world average technology $(A_w^O)$, which is pinned down by (1.32). That is, there is long term divergence of the world distribution of income around $A_w^O$.

What happens to investment and the path of income after $\tau^*$? Each household in each country considers the same evolution of goods prices and the education decision leads to an optimal cutoff level that solves (for any $t \geq \tau^*$)

$$\int_{t+T}^{\infty} \left( \overline{\theta}_{i,t} A_i p_{h,\tau,w} + p_{l,\tau,w} \right) e^{+\delta T - \int_t^\tau \delta + r(\nu) d\nu} d\tau = \int_t^\infty p_{l,\tau,w} e^{-\int_t^\tau \delta + r(\nu) d\nu} d\tau \qquad (1.33)$$

Because all countries evaluate the same evolution of prices, it is easy to make *relative* statements about the cutoff level $\overline{\theta}_{i,t}$. Define the level of technology for which the dynamic path of wages (1.33) evaluated at $t \geq \tau^*$ leads to exactly the same cutoff level as in autarky by $A_t^D$. Comparing (1.33) to (1.30) and (1.31), for any $t \geq \tau^*$, $A_{t,w}^D$ satisfies

$$A_w^A < A_{\tau^*,w}^D \leq A_w^O$$

After opening markets to trade, all countries with $A_i < A_{t,w}^D$ decrease their educational investment while others increase their level of investment. There is thus instantaneous convergence of income around $A_{t,w}^D$. How does the distribution of world income evolve thereafter? First evaluate the evolution of relative output for countries that are more technologically advanced than the rest of the world in long term equilibrium $A_w^O$. At any point in time $t > \tau^*$ these countries invest more in skills than in autarky. Thus, again the GDP of these countries initially decreases until period $\tau^* + T$ and then increases. The opposite is true for all countries with $A_i < A_{\tau^*,w}^D$. Along any point in time $t > \tau^*$ these countries invest less in skills than in autarky. The dynamics for the countries with $A_{\tau^*,w}^D < A_i < A_w^O$ are more interesting. This group starts the process of globalization being an exporter of the skill intensive good but successively becomes an importer of the latter. That is, as $A_{t,w}^D$ converges to $A_w^O$, a larger and larger fraction has lower investment in education than in autarky.

What are the welfare consequences of globalizing markets? Consider first countries that do not change their pattern of specialization i.e. all $A_{\tau^*,w}^D < A_i < A_w^O$. Proposition 5 (Trade and the Divergence of Welfare) establishes conditions under which welfare diverges for a given *constant* level of world prices. In fact, the world price of skill intensive goods is higher than in the long run equilibrium with $A_w^O$. This creates additional temporary gains for exporters

of skill intensive goods: if (1.27) holds for $A_w = A_w^O$, there is divergence around $A_w^O$. In fact, divergence is more likely than if prices would instantaneously jump to their long term levels. Consider now this group of countries that has $A_{r^*,w}^D < A_i < A_w^O$. Some of these countries might experience smaller dynamic than static gains from trade. For example consider – if it exists – a country with $A_i = A_w^O$. This country starts as an exporter of the skill intensive sector but does not trade in the long run. By Lemma 1 (Static Gains From Trade), this nations gains statically from trade, but it does not gain from trade in the long run. Putting things together, globalization results in long term divergence of income around $A_w^O$, the long term average level of technology.

## 1.7 Empirical Evidence

This section presents empirical evidence that the factor content of trade flows is indeed affects education decisions. I first construct measures of geographic proximity to international supply of skilled and unskilled labor. I then show that, conditional on the level of a country's development, these measures are a significant determinant of investment in human capital and average years of education in the workforce.

While researchers have struggled to establish the empirical validity of the theory of comparative advantage (Bowen, Leamer Sveikauskas (1984) and Trefler (1995), but see also David and Weinstein (2001)), research has been more successful in showing that bilateral trade flows follow the direction implied by Hecksher Ohlin trade theory quite closely. This arguably weaker prediction receives strong support in several studies (see Debaere (2003) and Choi and Krishna (2004)). Romalis (2004) shows how countries capture a larger global market share in sectors that intensively use their abundant factors. Romalis also establishes that the predictions of HOV theory hold qualitatively in the context of Krugman's (1980) model of monopolistic competition and transport costs. The tests of the bilateral factor content of trade and the findings of Romalis are the starting point for the empirical section of this paper. I use the bilateral trade data from the World Trade Database and the US productivity matrix to construct the factor content of bilateral trade. I then use national factor endowment data from Barro and Lee (1994) as well as geographic data to develop a gravity model relating the size of the factor

content of trade to distance between pairs of countries, to country size and most importantly to the abundance of the respective factor in the exporting nations. By running gravity equations that relate the factor content of trade flows to geograpical distance and factor supplies, I estimate the effect factors embodied in imports have on domestic education decisions. In line with the findings of the current literature, I find that one can predict the factor content of trade rather well using geographical data and factor supply differentials. To establish a causal relationship, I only use the information of a country's geographic proximity to international supply of skilled and unskilled labor to instrument for the observed factor content of trade. In this way I isolate the component of international trade that is not stemming from domestic supply and demand, but exclusively from the factor supply of other nations. This empirical strategy is related to Frankel and Romer (1999), who isolate the geographic component of trade to establish a causal relation between trade and growth. Similarly, my measure reflects how much skilled and unskilled labor other countries are likely to export to a given nation, and I subsequently test whether this measure of geographic proximity to skilled and unskilled labor has significant effects on domestic education decisions and the stock of human capital.

My findings are that the unconditional correlation between my the instrumented factor content are correlated to education levels, but not strongly. This finding can be explained by the fact that rich and poor nations are often geographically grouped together. I thus condition on the level of economic development directly or indirectly by using geographic variables. Conditional on these measures, I show that education is strongly affected by my geographical proximity to skilled and unskilled labor: levels of higher education (average years in the population) are negatively affected by proximity to skilled labor, and positively to proximity to unskilled labor. The same channel is not present for levels of primary education. Rather, when I find a significant relation, it is of the reverse nature: countries that trade a lot with skill abundant nations tend to accumulate more primary education. Another question of interest is whether this relation between trade and domestic education has become more pronounced with the increasing importance of international trade over the last 40 years. I test this formally by evaluating changes in education from 1960 to 1990, with similar results: proximity to skilled labor had a negative effect on changes of advanced education, while this channel is not present for primary education.

### 1.7.1 Constructing the Instrument

In a first step, I use bilateral trade flow data from Feenstra et al. (1997) and the US productivity matrix (obtained from Trefler and Antweiler (2002)) to translate 3 digit industry level trade flows into the net factor content of trade. I then use factor endowment data to predict the actual factor content of trade by using geographic data and national endowment levels. Let $i \in I$ index countries and let $o \in I-1$ index countries of origin, let $f \in F$ index factors and let $j \in J$ index industries. Let $V$ be the $(I, F)$ matrix summarizing international factor endowments, so that $V_{i,f}$ denotes the endowment of factor $f$ in country $i$. Let $A^*$ be the $(F, J)$ unit requirement matrix that is common to all countries once factor augmenting productivity differences are taken into account.[11] Finally, let $M_{i,o}$ be the $(J)$ vector of good imports from country $o$ to country $i$. A country's imports can be converted into the factor content by applying the unit requirement matrix. I denote the factor content of bilateral imports from country $o$ to country $i$ by $FCT_{i,o}$, which is given by

$$FCT_{i,o} = A^* \times M_{i,o} \tag{1.34}$$

The strong HOV hypothesis would boil down to testing whether the net vector of the factor content of trade (i.e. imports to country $i$ minus exports from country $i$) equals the national endowment minus a constant share $s_i$ of world endowment for each factor.

$$\sum_{o \in I-1} FCT_{i,o} - \sum_{i \in I-1} FCT_{o,i} = V_i - s_i \sum_{i \in I} V_i \tag{1.35}$$

Rather than testing (1.35), I relate the factor content of bilateral trade (1.34) to the relative abundance of factors at home and abroad. The weakest test is whether the relative supply of skills at home and abroad affects the relative factor content of imports. My first regression, represented in Table 1, tests the following relation.

---

[11] Antweiler and Trefler (2002) also allow for country-specific factor augmenting productivity differences, obtained from factor prices. However, lack of precise wage data leads them to uniformly adjust for all types of labor and for not skilled and unskilled worker productivity separately. The fit of my model that focuses on *relative* levels of factor supply is therefore not affected by their adjustments.

$$\frac{FCT_{i,o,HK2}}{FCT_{i,o,HK2} + FCT_{i,o,HK1}} = \alpha + \beta_1 \frac{V_{i,HK2}}{V_{i,HK2} + V_{i,HK1}} + \beta_2 \frac{V_{i,HK2}}{V_{i,HK2} + V_{i,HK1}} + \varepsilon \qquad (1.36)$$

$HK1$ and $HK2$ are my measures of skill and labor abundance corresponding to the number of people in the workforce with at least a high school equivalent degree ($HK2$) and those that did not finish high school ($HK1$). The dependent variable is the relative factor content of imports (to home). In regression (1) of Table 1, I first estimate equation (1.36) with $\beta_2$ restricted to 0, i.e. only considering the relative supply of skilled labor in the exporting country. Countries that are more skill abundant, on average, tend to export more skill intensive goods and hence have a larger relative skilled factor content compared to the unskilled factor content of trade. The coefficient equals 12 %, i.e. if a country, ceteris paribus, has 10 percentage point more skilled workers, the export composition is 1.2 percentage points more skill intensive. Regression (2) of Table 1 adds domestic skill supply, which is negatively related to the skill intensity of imports. Countries that are more skill abundant, on average, tend to export relatively less skill intensive goods and hence have a lower relative factor content of skilled labor. The coefficient equals about -6%: domestic skill abundance is a less important determinant of the factor content of imports than the skill abundance of the exporting country.

Regressions (3) and (4) of Table 1 replicate this finding for a different measure of skilled labor that only counts workers with finished secondary or higher education as skilled and all other workers as unskilled. The results are robust to this different definition a "skilled" worker. Equations (1) to (4) demonstrate that foreign, and to a lesser extent domestic, factor abundance is a major determinant of the factor content of trade. Even unconditional on any other information, relative endowment differences have significant effects on the relative factor content of imports. However, many other variables, such as the size of the export and the import markets or geographical barriers determine the absolute amount of a certain factor embodied in actual trade flows. I incorporate these in the next step.

|  | (1) Relative FCT I HK2/(HK1+HK2) | (2) Relative FCT I HK2/(HK1+HK2) | (3) Relative FCT II (TTA+STA)/TOTAL | (4) Relative FCT II (TTA+STA)/TOTAL |
|---|---|---|---|---|
| Foreign Skill Supply I | 0.116 (0.009)** | 0.109 (0.009)** |  |  |
| Home Skill Supply I |  | -0.058 (0.013)** |  |  |
| Foreign Skill Supply II |  |  | 0.051 (0.005)** | 0.047 (0.005)** |
| Home Skill Supply II |  |  |  | -0.031 (0.005)** |
| Observations | 3055 | 3055 | 3055 | 3055 |
| R-squared | 0.08 | 0.1 | 0.07 | 0.1 |

Robust standard errors in parentheses, observations clustered within importing nation
* significant at 5% level; ** significant at 1% level

Table 1: Relative Factor Supply and Relative Factor Content

To get a finer measure of the factor content of trade, I now run two separate gravity equations that relate (always in logarithms) geographical distance, country size measured in population and factor abundance to the magnitude of the factor content of imports between each pair of countries.[12] Because I want to focus on the causal effect that does not depend on domestic factors (other than country size), I do not include domestic skill abundance in the regressions. I next thus estimate a gravity style equaiton of the factor content of trade.

$$Log(FCT_{i,o,HK1}) = \alpha + \gamma_1 Log\,(dist_{i,o}) + \gamma_2 Log\,(POP_i) + \gamma_3 Log\,(V_{o,HK1}) + \overline{\varepsilon} \qquad (1.37)$$

$$Log(FCT_{i,o,HK2}) = \alpha + \tilde{\gamma}_1 Log\,(dist_{i,o}) + \tilde{\gamma}_2 Log\,(POP_i) + \tilde{\gamma}_3 Log\,(V_{o,HK2}) + \tilde{\varepsilon} \qquad (1.38)$$

---

[12] As is often noted (see for example Helpman, Melitz and Rubinstein (2004)) gravity equations are often biased because zero-trade observations are ignored when estimation a log-normalized model. In contrast to industry specific trade observations, where there exist many country pairs that do not trade in a specific industry, there are very few country pairs that do not trade at all. Therefore the 0-trade bias is not severe in my regressions.

| | (1) Log FCT HK1 (1992) | (2) Log FCT HK1 (1972) | (3) Log FCT HK2 (1992) | (4) Log FCT HK2 (1972) | (5) Log FCT HK2 (1992) |
|---|---|---|---|---|---|
| Log Distance | -1.312 (0.075)** | -1.261 (0.066)** | -1.333 (0.087)** | -1.174 (0.076)** | -1.334 (0.079)** |
| Log Pop Home | 0.76 (0.113)** | 0.677 (0.141)** | 0.691 (0.099)** | 0.613 (0.123)** | 0.704 (0.070)** |
| Log HK1 Foreign | 0.887 (0.035)** | 0.827 (0.035)** | | | |
| Log HK2 Foreign | | | 0.715 (0.025)** | 0.714 (0.030)** | 0.743 (0.024)** |
| Home Skill Supply I | | | | | 3.185 (0.543)** |
| Observations | 3055 | 2942 | 3055 | 2942 | 3055 |
| R-squared | 0.44 | 0.36 | 0.32 | 0.26 | 0.37 |

Robust standard errors in parentheses, observations clustered within importing nation

* significant at 5% level; ** significant at 1% level

Table 2: A Gravity Model of the Factor Content of Trade

Table 2 describes the results from estimating the gravity-style equations (1.37) and (1.38). Column (1) estimates model (1.37) using trade and endowment data from 1992. As expected, distance between two nations decreases trade flows and therefore also the embodied factor content. Secondly, the size of the importing market, as captured by the population in the importing nation increases the trade flow, but with an elasticity of less than 1, because bigger nations tend to trade less in percentage terms of their GDP. The last regressor I add is the supply of unskilled labor in the exporting country.[13] On average, a country whose unskilled labor force is 1 percentage point higher tends to export 0.89 percentage point more unskilled labor embodied in the goods it exports.

Comparing the $R^2$ in Table 1 and Table 2, the log estimation improves the total fit of the regression substantially. In Equation (3) of Table 2, I estimate the skilled labor content of trade in 1992. In accordance with the results of David and Weinstein (2001), factor endowments correlate strongly and in the right order of magnitude (elasticity of about 3/4) with the observed trade data in a gravity specification. Equations (2) and (4) replicate equation (1) and (3)

---

[13] From here on, I report measures that count all workers with at least a highschool degree as skilled workers. The results with other measures of skilled and unskilled labor are comparable.

respectively, but using trade and endowment data from 1972. The relationship seems to have been stable over time. Finally, in Equation (5), I add skill supply of the importing nation to check how good one can predict the factor content of trade when also including domestic information. This model is not used in the later analysis, but shows an interesting finding: the elasticity of the factor content of imports with respect to relative domestic skill supply is estimated at 3.19 - a value that is very large and positive. If a country is *more skill abundant*, it tends to *import more skill intensive goods*. However, the previous models evaluating relative factor content (Table 1) have established that a *relatively more skill abundant* nation's imports tend to be *relatively less skill intensive*. The explanation for this discrepancy is the positive correlation between human capital abundance and the volume of a nation's trade. Countries with a more educated workforce also tend to be richer and thus trade more.

I now predict equations (1) and (3) of Table 2. After taking exponents, for each country I sum over all its trading partners and then divide by the country's population size to end up with 2 measures for each country. From equation (1) I predict the measure $PredictedFCT\_HK1$, and from equation (3) I predict the measure $PredictedFCT\_HK$. $PredictedFCT\_HK1$ corresponds to a country's geographical proximity to unskilled labor. For example, a value of 0.001 means that given the geographic location of this country, on average the unskilled worker content of trade is equal to 0.001 per head of the population. $PredictedFCT\_HK2$ corresponds to the geographical proximity of a country to skilled labor, and a value of 0.001 corresponds to an average skilled worker content of trade equal to 0.001 per head of the population. I end up with estimates for 64 countries that are listed in Table 6 of Appendix C. Graph 2 in Appendix C also shows that the variables $PredictedFCT\_HK1$ and $PredictedFCT\_HK2$ are well behaved and that there are no outliers. I now check whether these measures of geographical proximity to unskilled labor influence domestic levels of human capital.

### 1.7.2 Skill Content and Domestic Education Decisions

The preceding gravity model correlates well with actual factor content of trade, and using the constructed measures $PredictedFCT\_HK1$ and $PredictedFCT\_HK2$ allows to make causal statements relating trade and education decisions. The main results of the empirical section are reported in Table 3. As dependent variables, I use education data in 1990 from the Barro and Lee dataset. I focus on average years of education in the population aged 15 an older. As independent variables, I use my two instruments that exploit the variation in the geographical proximity to skilled and unskilled labor.

Model (1) in Table 3 displays the relation between average years of higher education (secondary and tertiary) in the workforce and my measures of proximity to skilled labor. Proximity to unskilled labor is a significant positive determinant of education, while proximity to skilled labor is a (not significant) negative determinant of education. The third last row of Table 3 reports the P-value for a test of equality for the coefficients on $PredictedFCT\_HK1$ and $PredictedFCT\_HK2$ which is rejected for model (1). While this is encouraging, the results are not very strong, and the coefficient on proximity to skilled labor is not significant. Considering the fact that skill abundant nations tend to be clustered together (e.g. Western Europe), as is the case with skill scarce nations (e.g. Africa) this is not surprising: this clustering biases the coefficient on $PredictedFCT\_HK2$ upwards and the coefficient on $PredictedFCT\_HK1$ downwards. Therefore, I now check whether my measures have an effect on the formation of human capital conditional on the level of development of a nation. I first condition on the Log of GDP per head in the population in 1990. The results are presented in column (2) of Table 3. Closeness to unskilled labor has a positive impact on higher education, while closeness to skilled labor has negative effect on higher education. Both effects are highly significant and also the test on the equality of the coefficients $PredictedFCT\_HK1$ and $PredictedFCT\_HK2$ is strongly rejected. Conditional on a nations levels of development, there is a very strong relation between the factor content of trade and domestic education. A concern with this regression is the endogeneity of GDP per head and the level of education. I therefore use Distance from Equator, or Latitude as an instrument for GDP. This measure is strongly related to the level of a countries development (see for example Hall and Jones (1999)), but at the same time it is

48

clearly exogenous. The results are reported in model (3) of Table 3,[14] and show that indeed, the previous relationship is robust to instrumenting the actual level of geography with geographic variables. I provide further evidence that the relationship between the factor content of imports and the formation of human capital is robust by conditioning on other geographic measures or on direct measures of development such as the quality of institutions. Table 5 in Appendix D lists additional models with different control variables.

A bigger potential concern is that my two measures of proximity to international factor supply simply capture the direct effect openness has on economic prosperity and therefore education. I hence also include the measure of geographic openness constructed by Frankel and Romer (1999). They construct a nation's potential for trade by estimating gravity models relating the volume of a nation's trade to its proximity to international population. I add their variable in model (4). While other coefficients remain significant, geographic openness itself has no impact on levels of education. Although the volume of trade is an important determinant for economic growth, it is the factor content and not the possibility for trade itself that is influencing the formation of human capital. Model (1) of Table 6 in Appendix D shows that this finding is robust to different control variables. In these models, I condition on levels of economic development by including $KGPTemp$, the proportion of the population living in temperate climate and $ME$, a measure of the geographic potential for malaria (Kiszewski et al. (2004)). The relation between levels of higher education and proximity to skilled labor is robust and negative, while the opposite is true for proximity to unskilled labor. Model (5) and (6), however show that this is not true for levels of primary education. Neither unconditionally (5), nor conditionally on the level of a nations development (Latitude; (6)) is there a significant relation. Again, I demonstrate that these findings are robust to different sets of control variables in Table 6 of Appendix D. This finding is in line with my proposed theory, since the unskilled labor content of trade also embodies a sizeable amount of primary education: workers who nearly finished highschool have high levels of primary education, but are counted as unskilled workers when I construct my instrument. Therefore, proximity to unskilled embodies primary education and should have a (weak) negative effect on education.

---

[14] The reduced form is shown here. A full model instrumenting $Log\,(GDP/Pop)$ by $Latitude$ displays equally significant results.

| | (1) Higher Education (avg. years) | (2) Higher Education (avg. years) | (3) Higher Education (avg. years) | (4) Higher Education (avg. years) | (5) Primary Education (avg. years) | (6) Primary Education (avg. years) |
|---|---|---|---|---|---|---|
| Predicted FCT EK1 | 238.8 (77.7)* | 162.4 (47.2)** | 154.7 (46.1)** | 157.5 (46.8)** | 86.2 (76.2) | 34.7 (62.8) |
| Predicted FCT EK2 | -86.3 (63.5) | -129.8 (33.6)** | -127.4 (33.4)** | -117.2 (33.6)** | 27.6 (56.7) | -36 (62.5) |
| Log GDP90/Pop | | 0.768 (0.0079)** | | | | |
| Latitude | | | 4.94 (1.01)** | 4.91 (0.99)** | | 3.98 (1.03)** |
| Geog. Openness | | | | -0.0052 (0.015) | | 0.0033 (0.025) |
| F(H2 ≠ EK1) | (5.77)* | (12.41)** | (15.46)** | (11.96)** | (0.20) | (0.29) |
| Observations | 63 | 63 | 60 | 60 | 63 | 60 |
| R-squared | 0.2 | 0.68 | 0.483 | 0.484 | 0.134 | 0.292 |

Robust standard errors in parentheses
Unskilled and Skilled Factor content per Population are measured in millions; * significant at 5% level; ** significant at 1% level

Table 3: The Factor Content of Trade and Domestic Education Levels

Table 4 evaluates the hypothesis that the increasing magnitude of trade flows in the period

50

since World War II has intensified the exposure to international markets and hence effects of trade are more pronounced today than they were before the period of globalization we have witnessed over the past 40 years. I repeat earlier regressions with changes of the level of average schooling from 1960 to 1990 as dependent variable. The general picture is very similar to the preceding results: higher education has been affected negatively in countries close to skill abundant nations, and has been affected positively when other nations tended to provide unskilled labor. The opposite is true for effects on primary education. In models (1) and (2) of Table 4, I condition on the logarithm of GDP per head in 1990, while in models (3) to (6) I condition on purely geographic information. Interestingly, I now find a significant relation between changes in primary education and my measures of proximity to factors. The coefficients are of the opposite sings as they are for higher education. Being close to exporters of labor intensive goods reduces primary education, while being close to exporters of skill intensive goods increases primary education.

In total, the empirical section demonstrates that both levels of education today as well as changes of education from 1960 to 1990 have significantly been affected by geographical proximity to skilled labor, and they have been affected in a way that is very consistent with the proposed theory. A simple calculation shows that the describe effects are sizeable: the standard deviation of $PredictedFCT\_HK2$ is equal to 0.0054 while the coefficient of $PredictedFCT\_HK2$ on levels of education lies between -117 and -130 (Conditional on Development, see Table 4). Comparing two otherwise identical nations, a one standard deviation difference in geographic proximity to skilled labor is hence associated with a difference between 0.63 and 0.70 in the average number of years of higher education per worker. This is large, especially considering that in my sample, the average level of higher education is only 2.3 years. While economists disagree on the precise estimate increased schooling has on income, most researchers would agree that an additional year of education increases a worker's lifetime income by more than 6%. A lower bound for the effects of the proposed theory is hence that a one standard deviation increase of proximity to skilled labor implies at least a 4-5% decrease of GDP. Again, this is quite sizeable when compared to standard estimates of the effect of trade to GDP.

| | Higher Education (change 1960-90) | Primary Education (change 1960-90) | Higher Education (change 1960-90) | Primary Education (change 1960-90) | Higher Education (change 1960-90) | Primary Education (change 1960-90) |
|---|---|---|---|---|---|---|
| Produced FCT HK1 | 71.6 (34.6)* | -75.2 (23.3)** | 53 (29.6) | -71.1 (23.4)** | 79.8 (33.7)* | -80.4 (24.5)** |
| Produced FCT HK3 | -48.6 (28.5) | 34.7 (21) | -33.6 (11.8) | -7.7 (23.9) | -46.9 (36.1) | 31.1 (21.1) |
| Log GDP90/cap | 0.301 (0.067)** | -0.12 (0.06) | | | | |
| Latitude | | | 2.17 (0.75)** | -0.61 (0.46) | | |
| ME | | | | | 0.77 (0.24)** | -0.46 (0.21)* |
| KGPT cap | | | | | -0.008 (0.025) | -0.003 (0.02) |
| Group. Openness | | | 0.001 (0.007) | 0.019 (0.01) | | |
| F(HK2 ≠ HK1) | (4.36)* | (5.17)* | (2.53) | (1.94) | (3.58) | (6.71)** |
| Observations | 62 | 62 | 59 | 59 | 56 | 56 |
| R-squared | 0.383 | 0.27 | 0.324 | 0.318 | 0.33 | 0.28 |

Robust standard errors in parentheses
Unskilled and Skilled Factor content per Population are measured in millions; * significant at 5% level; ** significant at 1% level

Table 4: The Factor Content of Trade and Changes in the Level of Education

52

# 1.8 Conclusions

Recent contributions to the literature of economic growth have argued that factor accumulation is a key ingredient for long term economic success. Countries that have sustained high rates of growth did so because of their high levels of savings and investment in human capital, see for example Young (1995) and Mankiw et al. (1992). Other nations stagnated precisely because their institutional setups hindered private savings and investment, see for example Hall and Jones (1999) and Acemoglu et al. (2005). To evaluate whether trade has sizeable and first order effects on economic performance, it is therefore essential to show how exposure to international prices influences factor accumulation. This paper starts from one of the most fundamental insights of the theory of international trade – Samuelson's (1948) factor price equalization theorem. Focusing on the role of human capital, I show that when countries trade, the response of factor supply to factor price equalization has important dynamic effects that lead to divergence of the world distribution of income and also tend to cause divergence of welfare, even in the absence of market failures. Furthermore, the present framework can explain why trade leads to a pervasive increase of the skill premium while at the same time resulting in a decline of the price of skill intensives goods. The basic asymmetry of the model is not specialization in different sectors that have varying growth potential, but the specialization in factors, which causes differential effects of trade. When some factors, such as physical and human capital, are in variable supply, trade thereby re-enforces initial patters of specialization. Poor countries gain from trade because they experience an increase of the unskilled wage, yet loose from trade because fewer workers decide to enter the skilled labor force and therefore, the surplus from education decreases. The opposite is true in rich countries. However, being a net supplier of human capital that can be accumulated increases the growth potential of the economy, while specialization in labor intensive sectors means specialization in a factor in fixed supply.

In addition to showing how trade can result in divergence of income, I evaluate how trade affects relative welfare differentials. The dynamic response of the economy that introduces the asymmetry between nations in the model. This result stems from the two margins in which the relative wage influences the surplus from education: a higher relative wage increases the income for all workers that already would have chosen schooling at lower wages. In addition, an

53

increase in the relative wage induces more entry into the skilled labor force. In total, the surplus from education responds more than proportionally to changes in the relative wage. Skill scarce nations, in contrast, have their comparative advantage in a factor that is in fixed supply and cannot be accumulated. The key insight of the mechanism at work is that while all countries gain from trade, it is already developed nations that gain proportionally the most from trade. Trade hence results in dynamic divergence of welfare.

I next present empirical evidence that the skilled labor embodied in current trade flows is indeed a significant factor for domestic education decisions. The empirical strategy first constructs measures of the skill content of trade. I then show that, conditional on the level of a country's development, these measures are a significant determinant of investment in human capital and average years of education in the workforce. A potential problem with estimating the relation between the observed factor content of trade and domestic education decisions is the problem of causality: even in a static version of the model with fixed supply of skilled and unskilled labor would a measure of the factor content of trade be correlated with domestic levels of education. To deal with this endogeneity, when constructing my measures of skill content, I only use the information of a country's geographic proximity to international supply of skilled and unskilled labor to instrument for the observed factor content of trade, and do not use the domestic relative factor supply itself. In this way I isolate the component of international trade that is not stemming from domestic supply and demand, but exclusively from the factor supply of other nations. This empirical strategy is related to Frankel and Romer (1999), who isolate the geographic component of trade to establish a causal relation between trade and growth. Similarly, my measures reflect how much skilled and unskilled labor other countries are likely to export to a given nation, and I subsequently test whether this measure of geographic proximity to skilled and unskilled labor has significant effects on domestic education decisions and the stock of human capital. My findings are that conditional on measures of economic developments, education is strongly affected by my measures of geographical proximity to skilled and unskilled labor. In accordance with the proposed theory, levels of higher education are negatively affected by proximity to international to skilled labor, while proximity to unskilled labor has a positive effect on education. I show that the same channel is not present for levels of primary education. Rather, when I find a significant relation between proximity to international supply of human

capital and levels of primary education, I find the reverse relationship: countries that trade a lot with skill abundant nations tend to accumulate more primary education. Another question of interest is whether the relation between trade and domestic education has become more pronounced with the increasing importance of international trade over the last 40 years. I test this formally by evaluating changes of the level of education from 1960 to 1990, with similar results: proximity to skilled labor had a negative effect on changes of average advanced education, while this channel is not present for primary education. A simple calculation shows that the describe effects are sizeable. Comparing two otherwise identical nations, a one standard deviation difference in geographic proximity to skilled labor is hence associated with a difference between 0.6 to 0.7 in the average number of years of higher education per worker. This is large, especially considering that in my sample, the average level of higher education is only 2.3 years.

## 1.9 Appendix A Proofs

**Proof.** Of Lemma 1 (Static Output Effects of Trade)

Evaluate autarky output (1.17) to output at $\tau^*$ (1.19) for two countries $A_n = (1 + \gamma) A_w = (1 + \gamma)^2 A_s$. The pre opening ratio of output is equal to

$$Y\left(\overline{\theta}_n^A, w_{l,w}, A_n p_w\right) \Big/ Y\left(\overline{\theta}_s^A, w_{l,w}, A_s p_w\right) = \frac{\left(1 + \lambda \left(\frac{A_i}{A_w}\right)^{\frac{1}{1-\eta\beta}} A_w^{\frac{\beta}{1-\eta\beta}}\right)}{\left(1 + \lambda \left(\frac{A_i}{A_w}\right)^{-\frac{1}{1-\eta\beta}} A_w^{\frac{\beta}{1-\eta\beta}}\right)}$$

The post opening ratio of output is equal to

$$Y\left(\overline{\theta}_n^A, w_{l,w}, A_n p_w\right) \Big/ Y\left(\overline{\theta}_s^A, w_{l,w}, A_s p_w\right) = \frac{\left(1 + \lambda A_n^{\frac{\beta}{1-\eta\beta}}\right)^{\frac{1}{\beta}}}{\left(1 + \lambda A_s^{\frac{\beta}{1-\eta\beta}}\right)^{\frac{1}{\beta}}}$$

When is there absolute divergence? The relative level of output has a form of has a form of $\left(\frac{z^k x + 1}{z^{-k} x + 1}\right)^{\frac{1}{k}}$, where $z = 1 + \gamma$. The autarky level of $k$ is higher than under autarky. If

$$\frac{\partial}{\partial k} \left(\frac{z^k x + 1}{z^{-k} x + 1}\right)^{\frac{1}{k}} > 0$$

holds, there is divergence. First note that when $x = 1$, $\left(\frac{z^k x + 1}{z^{-k} x + 1}\right)^{\frac{1}{k}} = z$. If $x < 1$, $\left(\frac{z^k x + 1}{z^{-k} x + 1}\right)^{\frac{1}{k}} < z$. Rewriting this expression

$$\frac{\partial}{\partial k} \exp\left[\frac{1}{k} \log\left(z^k x + 1\right) - \frac{1}{k} \log\left(z^{-k} x + 1\right)\right]$$

Where the inner exponents need to be rewritten in log-exponential form, too. This can be shown to be

$$\left(\frac{z^k x + 1}{z^{-k} x + 1}\right)^{\frac{1}{k}} \left(\frac{\partial}{\partial k}\left(\frac{1}{k} \log\left(z^k x + 1\right)\right) - \frac{\partial}{\partial k}\left(\frac{1}{k} \log\left(z^{-k} x + 1\right)\right)\right)$$

Omitting terms that are positive it remains to be shown that

$$\frac{\partial}{\partial k}\left(\frac{1}{k}\log\left(z^k x + 1\right)\right) - \frac{\partial}{\partial k}\left(\frac{1}{k}\log\left(z^{-k} x + 1\right)\right) > 0$$

Putting things together, multiplying by $k^2$ this is true, if

$$\log\left(z\right)\left[\frac{z^k x}{(z^k x + 1)} + \frac{z^{-k} x}{(z^{-k} x + 1)}\right] - \log\left(\frac{z^k x + 1}{z^{-k} x + 1}\right)^{\frac{1}{k}} > 0 \qquad (1.39)$$

Recalling that when $x = 1$, $\log\left(\frac{z^k x + 1}{z^{-k} x + 1}\right)^{\frac{1}{k}} = \log(z)$ so (1.39) is equal to 0. In addition, (1.39) takes the value 0 if $x \to \infty$ and if $x = 0$. Evaluating the slope of (1.39) with respect to x, it can be shown that at levels of $x$ just below 1,(1.39) is decreasing, implying that (1.39) is bigger 0 for any $1 > x > 0$. In addition, it can be shown there exist at most 3 levels of x where (1.39) equals 0. Hence, (??) is increasing in $k$ whenever $1 > x > 0$. ∎

**Proof.** of Proposition 2 (Trade and The Dynamics of Income)

To establish the two claims of the proposition, compare the relative ratio of output for two countries N and S in autarky (1.17), just after opening to trade (1.19) and in the stationary equilibrium with trade (1.22). It is both true that for any $\gamma > 0$ , the following inequalities hold.

$$\frac{Y\left(\overline{\theta}_n^O, w_{l,w}, A_n p_w\right)}{Y\left(\overline{\theta}_s^O, w_{l,w}, A_s p_w\right)} \bigg/ \frac{Y\left(\overline{\theta}_n^A, w_{l,w}, A_n p_w\right)}{Y\left(\overline{\theta}_s^A, w_{l,w}, A_s p_w\right)} \qquad (1.40)$$

First note that if $\gamma = 0$, this ratio is equal to 1. Now evaluate the first and second derivative of (1.40) with respect to $\gamma$. The first derivative is positive at $\gamma = 0$, while the second derivative is positive for any $\gamma$. Hence, for any $\gamma > 0$ (1.40) takes a value larger than 1. A proof along the same lines establishes that for any $\gamma > 0$ comparing (1.17) to (1.22) result in dynamic divergence, i.e.

$$\frac{Y\left(\overline{\theta}_n^O, w_{l,w}, A_n p_w\right)}{Y\left(\overline{\theta}_s^O, w_{l,w}, A_s p_w\right)} > \frac{Y\left(\overline{\theta}_n^A, w_{l,n}^A, w_{h,n}^A\right)}{Y\left(\overline{\theta}_s^A, w_{l,s}^A, w_{h,s}^A\right)}$$

∎

# 1.10 Appendix B Endogenous Technology

The previous results have been derived under the assumption of exogenous technology differences. This section endogenizes technology in a two sector model of endogenous growth. It is established that all results are amplified.

Assume that the production of intermediate goods is modified alla Romer (1990) in the two sector version of Acemoglu (1998). Technology is local and in each country, the production function combines factor specific differentiated input goods and the respective factor. Each of these input goods is produced using a linear transformation of the respective intermediate good[15]. I denote the amount of each input good used in the labor intensive sector by $i_{L,i}$ and the one used for production of the skill intensive good by $i_{H,i}$. The net output of each intermediated (denoted by $\widetilde{Y}_{j,i}$ for $j \epsilon [L,H]$) given by

$$\widetilde{Y}_{L,i} = \left( \int_0^{N_{L,i}} i_{L,i}^\gamma di \right) L_i^{1-\gamma} - \int_0^{N_{L,i}} i_{L,i} di - R_{L,i} \tag{A1}$$

$$\widetilde{Y}_{H,i} = \left( \int_0^{N_H} i_{H,i}^\gamma di \right) H_i^{1-\gamma} - \int_0^{N_{H,i}} i_{H,i} di - R_{H,i} \tag{A2}$$

$(1 < \gamma < 1)$ Where $R_{L,i}$ and $R_{H,i}$ are the flows of R&D expenditures that are used to invent new blueprints in each sector.[16] I assume that innovation in sector j uses only the respective intermediate good as input to produce new innovaitons. Furthermore, as in Jones(1995) , innovation becomes more difficult as the current level of innovation is higher. Denoting the flow-cost of innovation in terms of the respective intermediate goods in sector $j$ and country $i$ by $\vartheta (N_{j,i})$, such that the absolute cost is $P_{j,i} \vartheta (N_{j,i})$

$$\vartheta (N_{j,i}) = N_{j,i}^\mu$$

with $\mu > 0$, innovation in each country essentially runs into decreasing returns. I now characterize countries not by their intrinsic difference in technology, but by their difference in their

---

[15]To be sure to make no confusions: there are now two sorts of input goods used two produce two intermediate goods used to produce the final good.

[16]Epifani and Gancia, in one oftheir formulations show how trade can lead to skill bias when the elasticity f substituion $(1 - \gamma)^{-1}$ between varieties is larger in the skill intensive sector that nin the labor intensive one.

educational sector. That is, some countries are essentially better at educating their workforce. The demographic structure is essentially unchanged, except that in country $i$, a skilled worker of type $\theta$ now supplies $s_i\theta$ units of skilled labor if she chooses to get educated. $s_i$ is the country specific efficiency of the educational system that is given exogenously. Models of endogenous investments in technology that are financed with profits from monopolistic competition features two related market failures: because each input good monopolist cannot prices discriminate, it charges a constant markup prices hence producing a suboptimal amount. For given levels of technology, the production of a country is thus suboptimal. More importantly, the same lack of ability to price discriminate also leads to the monopolist not capturing the full social surplus from her invention. There is thus also suboptimal entry into the input producing sector, with important dynamic consequences for technology, output and welfare. Because innovators face a constant demand elasticity, they charge a price of $1/\gamma$ times their marginal costs. Each innovator in the L sector hence sells $i_{L,i} = \gamma^2 L_i$ units while a firm in the H sector sells $i_{H,i} = \gamma^2 H_i$. Free entry into the market for input goods hence implies that

$$N_{L,i} = \left((1-\gamma)\gamma^2 L_i\right)^{\frac{1}{\mu}} \quad \text{and} \quad N_{H,i} = \left((1-\gamma)\gamma^2 H_i\right)^{\frac{1}{\mu}}$$

The net output in each sector is hence given by

$$\widetilde{Y}_{L,i} = \left((1-\gamma)\gamma^2\right)^{\frac{1}{\mu}} \gamma^{2\gamma} L_i^{\frac{1+\mu}{\mu}} - \int_0^{N_{L,i}} i_{L,i}di - R_{L,i} \tag{A1}$$

$$\widetilde{Y}_{H,i} = \left(\int_0^{N_H} i_{H,i}^\gamma di\right) H_i^{1-\gamma} - \int_0^{N_{H,i}} i_{H,i}di - R_{H,i} \tag{A2}$$

and the relative wage is given by

$$w_i = \left(\frac{N_{H,i}}{N_{L,i}}\right)^\beta \left(\frac{H_i}{L_i}\right)^{-(1-\beta)} = \left(\frac{H_i}{L_i}\right)^{\frac{\beta}{\mu}-(1-\beta)}$$

Because relative technology is increasing in factor abundance, the relative wage may now be increasing in the supply of skilled labor. The steady state education supply in each country is a function of the wage and the domestic schooling technology $s_i$.

$$\frac{H_i}{L_i} = (\eta c)^{\frac{1}{1-\eta}} \eta \left(e^{\rho T} - 1\right)^{-\frac{1}{1-\eta}} w_i^{\frac{\eta}{1-\eta}} s_i^{\frac{1}{1-\eta}}$$

$s_i$ matters more than proportional, because it influences both the cutoff and the average level of education per skilled worker.

$$\left(\frac{H_i}{L_i}\right)^{1-\beta\eta\frac{1+\mu}{\mu}} = (\eta c)^{\eta} \eta^{1-\eta} \left(e^{\rho T} - 1\right)^{-1} s_i$$

An non explosive equilibrium closed economy requires that $\beta\eta\frac{1+\mu}{\mu} < 1$, and the condition $\eta\frac{1+\mu}{\mu} < 1$ is required for a non-explosive open equilibrium. Otherwise, the result of the model with endogenous education and an intrinsic difference in the efficiency in the educational system is equivalent to the model with fixed technology, except that because also technology adjusts, countries tend to be more dissimilar.

## 1.11 Appendix C List of Data





Figure 2: Relationship Between Geographical Estimates and the Actual Factor Content of Trade

| COUNTRY | CODE | Actl FCT HK1/POP | Pred. FCT HK1/POP | Actual FCT HK2/POP | Pred. FCT HK2/POP |
|---|---|---|---|---|---|
| Argentina | ARG | 0.0033684 | 0.0013843 | 0.0076796 | 0.0019838 |
| Australia | AUS | 0.0156595 | 0.0012524 | 0.0353459 | 0.0015204 |
| Austria | AUT | 0.050929 | 0.0079536 | 0.1049795 | 0.0102248 |
| Belgium | BEL | 0.0995534 | 0.0242058 | 0.2002868 | 0.0319474 |
| Bangladesh | BGD | 0.0001472 | 0.0006098 | 0.0002498 | 0.0005991 |
| Bolivia | BOL | 0.0013662 | 0.0021585 | 0.0029329 | 0.0025738 |
| Brazil | BRA | 0.000902 | 0.0008657 | 0.0020244 | 0.0008593 |
| Canada | CAN | 0.0340874 | 0.0126008 | 0.0768641 | 0.0038954 |
| Chile | CHL | 0.0048128 | 0.0014096 | 0.0107315 | 0.0014921 |
| Cameroon | CMR | 0.0005447 | 0.0013528 | 0.0011327 | 0.0013688 |
| Colombia | COL | 0.001603 | 0.0017553 | 0.0035777 | 0.0013072 |
| Costa Rica | CRI | 0.0062644 | 0.0044054 | 0.0112671 | 0.003127 |
| Germany | DEU | 0.0371908 | 0.0078331 | 0.0737739 | 0.0086087 |
| Dominican Republi | DOM | 0.0045665 | 0.003405 | 0.0070124 | 0.0020678 |
| Ecuador | ECU | 0.0018717 | 0.0026027 | 0.0043515 | 0.0027665 |
| Egypt | EGY | 0.0011397 | 0.0018414 | 0.0022832 | 0.0026529 |
| Spain | ESP | 0.0189396 | 0.0034027 | 0.0393566 | 0.0041596 |
| Ethiopia | ETH | 0.0001348 | 0.0007075 | 0.0002709 | 0.0007102 |
| Finland | FIN | 0.0280207 | 0.0080966 | 0.0606272 | 0.010228 |
| France | FRA | 0.0338046 | 0.009455 | 0.0704005 | 0.0108103 |
| UK | GBR | 0.0299144 | 0.008869 | 0.0627145 | 0.010267 |
| Ghana | GHA | 0.0004863 | 0.0013363 | 0.0010198 | 0.0012835 |
| Greece | GRC | 0.0156624 | 0.0036018 | 0.0312232 | 0.004544 |
| Guatemala | GTM | 0.0018572 | 0.0038483 | 0.0035795 | 0.0024711 |
| Honduras | HND | 0.0018829 | 0.0052114 | 0.0032362 | 0.005305 |
| India | IND | 0.0001068 | 0.0005961 | 0.0002446 | 0.0008205 |
| Ireland | IRL | 0.051388 | 0.0105353 | 0.1065766 | 0.0140363 |
| Iceland | ISL | 0.0439459 | 0.0118418 | 0.0901199 | 0.0150537 |
| Israel | ISR | 0.0212666 | 0.0031558 | 0.047085 | 0.0046584 |
| Italy | ITA | 0.0230816 | 0.0040044 | 0.0474788 | 0.0052702 |
| Jamaica | JAM | 0.005729 | 0.0080753 | 0.0099525 | 0.0042925 |
| Japan | JPN | 0.0109529 | 0.0020878 | 0.0206181 | 0.0019955 |
| Korea, Rep. | KOR | 0.0106769 | 0.0026641 | 0.0241156 | 0.0020616 |
| Sri Lanka | LKA | 0.0010173 | 0.0009491 | 0.001901 | 0.0009943 |
| Morocco | MAR | 0.0019483 | 0.0027749 | 0.0038536 | 0.0037612 |
| Madagascar | MDG | 0.0001902 | 0.0007833 | 0.00041 | 0.000801 |
| Mexico | MEX | 0.0053249 | 0.0031339 | 0.0115914 | 0.0013737 |
| Malawi | MWI | 0.0001881 | 0.0009287 | 0.000401 | 0.0009002 |
| Malaysia | MYS | 0.0108268 | 0.0012462 | 0.0263758 | 0.0014133 |
| Nigeria | NGA | 0.0004729 | 0.0008495 | 0.0010941 | 0.000735 |
| Netherlands | NLD | 0.0735976 | 0.0164973 | 0.147273 | 0.0210376 |
| Norway | NOR | 0.0476122 | 0.0098605 | 0.1006144 | 0.0124308 |
| New Zealand | NZL | 0.0170064 | 0.001469 | 0.038523 | 0.0013423 |
| Pakistan | PAK | 0.0003959 | 0.003445 | 0.0008485 | 0.0051765 |
| Panama | PAN | 0.0385122 | 0.0055632 | 0.0857173 | 0.0059807 |
| Peru | PER | 0.0009234 | 0.0013595 | 0.0018967 | 0.001142 |
| Philippines | PHL | 0.0012117 | 0.0014731 | 0.0026486 | 0.0017009 |
| Papua New Guine | PNG | 0.0017484 | 0.001458 | 0.0037988 | 0.0010718 |
| Portugal | PRT | 0.0258084 | 0.0062916 | 0.0508878 | 0.0090776 |
| El Salvador | SLV | 0.0019679 | 0.0051129 | 0.0039094 | 0.0037858 |
| Suriname | SUR | 0.0059145 | 0.0047717 | 0.0126213 | 0.0047342 |
| Sweden | SWE | 0.0411911 | 0.0076858 | 0.0865051 | 0.0099012 |
| Syrian | SYR | 0.0010638 | 0.0013792 | 0.0022904 | 0.0019241 |
| Thailand | THA | 0.0045238 | 0.0026697 | 0.0105152 | 0.0037086 |
| Tunisia | TUN | 0.0059709 | 0.004736 | 0.0111689 | 0.006899 |
| Turkey | TUR | 0.002042 | 0.0019111 | 0.004691 | 0.002062 |
| Tanzania | TZA | 0.0002024 | 0.000708 | 0.0004511 | 0.0006767 |
| Uruguay | URY | 0.0051781 | 0.0057908 | 0.0110467 | 0.0093274 |
| US | USA | 0.0148143 | 0.0016287 | 0.0313321 | 0.0016086 |
| Venezuela, | VEN | 0.0044967 | 0.0022322 | 0.0101182 | 0.0018017 |
| South Africa | ZAF | 0.0026293 | 0.0005744 | 0.005827 | 0.0005227 |
| Zambia | ZMB | 0.0004482 | 0.0007637 | 0.0008866 | 0.0007648 |
| Zimbabwe | ZWE | 0.0006618 | 0.0010386 | 0.0013359 | 0.0010024 |

Table 6: List of Constructed Factor Content of Trade

# Bibliography

[1] Acemoglu, D. (2003) "Patterns of Skill Premia", The Review of Economic Studies, 70(2): 199-230.

[2] Acemoglu, D., Johnson, S. and Robinson J.A. (2005) "Institutions as the Fundamental Cause of Long-Run Growth", Handbook of Economic Growth, forthcoming

[3] Acemoglu, D. and Ventura, J. (2002) "The World Income Distribution", The Quarterly Journal of Economics, 117(1): 659-694.

[4] Acemoglu, D. and Zilibotti, F. "Productivity Differences", The Quarterly Journal of Economics, 116(1): 563-606.

[5] Antweiler, W and Trefler, D (2002) "Increasing Returns and All That: A View from Trade", The American Economic Review, 92(1)), pp. 93-119.

[6] Autor, D.H., L.F. Katz, and A.B. Krueger (1998) "Computing Inequality: Have Computers Changed the Labor Market?", The Quarterly Journal of Economics, 113(4): 1169-1213.

[7] Barro, R. and Lee, J. W.(1996) "International Measures of Schooling Years and Schooling Quality", The American Economic Review P&P, 86(2), 218-223.

[8] Baldwin, R.E. (1992) "Measurable Dynamic Gains from Trade", The Journal of Political Economy, 100(1): 162-174

[9] Caselli, F. and Coleman J. (2005) "The World Technology Frontier", The American Economic Review, forthcoming.

[10] Choi, Yong-Seok and Krishna, Pravin (2004) "The Factor Content of Bilateral Trade: An Empirical Test", The Journal of Political Economy, 112(4): 887-914.

[11] Davis, D. R. and Weinstein, D.E. (2001) "An Account of Global Factor Trade", The American Economic Review, 91(5): 1423-1453.

[12] Debaere, Peter (2003) "Relative Factor Abundance and Trade", The Journal of Political Economy, 111(3): 589-610

[13] Dinopoulos, E and P. Segerstrom (1999) "A Schumpeterian Model of Protection and Relative Wages", The American Economic Review, 89(3): 450-472.

[14] Epifani, P. and Gancia, G. (2005) "Increasing Returns, Imperfect Competition and Factor Prices", mimeo, CREI.

[15] Feenstra, Robert C., Robert E. Lipsey and Harry P. Bowen (1997) "World Trade Flows, 1970-1992, with Production and Tariff Data," NBER Working Paper 5910.

[16] Findlay, R. and Kierzkowski, H (1983) "International Trade and Human Capital: A Simple General Equilibrium Model", The Journal of Political Economy, 91(6): 957-978.

[17] Gancia, G. (2005) "Globalization, Divergence and Stagnation", mimeo, CREI.

[18] Anthony Kiszewski, Andrew Mellinger, Andrew Spielman, Pia Malaney, Sachs Jeffrey, and Sonia Ehrlich Sachs (2004)."A Global Index of the Stability of Malaria Transmission" American Journal of Tropical Medicine and Hygiene, 70(5), pp. 486-498

[19] Krugman, P. (1980) "Scale Economies, Product Differentiation, and the Pattern of Trade", The American Economic Review, 70(5): 950-959.

[20] Krugman, P. (1991) "Increasing Returns and Economic Geography", The Journal of Political Economy, 99(3): 483-499.

[21] Krugman, P. and Venables, A.J (1995) "Globalization and the Inequality of Nations", The Quarterly Journal of Economics, 110(4): 857-880.

[22] Mankiw, N.G, Romer, D and Weil, D "A Contribution to the Empirics of Economic Growth", The Quarterly Journal of Economics, 107(2): 407-437.

[23] Matsuyama, K. (1991) "Increasing Returns, Industrialization, and Indeterminacy of Equilibrium", The Quarterly Journal of Economics, 106(2):617-650.

[24] Ohlin, B. (1933) "Interregional and International Trade", Harvard University Press, Cambridge.

[25] Romalis, J. (2004) "Factor Proportions and the Structure of Commodity Trade", The American Economic Review, 94(1): (67-97).

[26] Stiglitz, J.E. (1970) "Factor Price Equalization in a Dynamic Economy", The Journal of Political Economy, 78(3): 456-488.

[27] Trefler, D (1995) "The Case of the Missing Trade and Other Mysteries", The American Economic Review, 85(5), pp. 1029-1046.

[28] Trefler, D (2002) "The Case of the Missing Trade and Other Mysteries: Reply", The American Economic Review, 92(1), pp. 405-410.

[29] Vanek, Jaroslav (1968), "The Factor-Proportions Theory: The N-Factor Case", Kyklos, 21(4): 749-756.

[30] Ventura, J. (1997) "Growth and Interdependence", The Quarterly Journal of Economics, 112(1): 57-84.

[31] Young, A. (1991) "Learning by Doing and the Dynamic Effects of International Trade", The Quarterly Journal of Economics, 106(2): 369-405.

[32] Young, A. (1995) "The Tyranny of Numbers: Confronting the Statistical Realities of the East Asian Growth Experience", The Quarterly Journal of Economics, 110(3): 641-680.

# Chapter 2

# Quality, Pricing to Market and Entry - The Short and Long Run of Exchange Rate Pass-Through (with Thomas Chaney)

**Summary 2** *Why are movements of relative costs brought about by exchange rate fluctuations passed on to customers only slowly, and never to a full extent? In this paper, we first develop a perfectly competitive economy featuring heterogeneity of both good qualities and of consumer valuations. In equilibrium, high valuation consumers and high quality firms are matched. The relative scarcity of different qualities leads to pricing-to-market and markups that are determined by the local toughness of competition. Our production setup features trade in intermediate goods, local assembly that is subject to decreasing returns and fixed costs of market entry. In every export market, firm entry and size decisions are determined by how local prices compare to the cost of production at home. We next analyze how changes in the real exchange rate are transmitted internationally. In the short run, the set of firms active in the export sector is fixed, but each firm accommodates changes in the exchange rate by adjusting the quantity of its exports. Due to this response of export volume to the relative cost of production, market toughness counteracts exchange rate movements, leading to partial pass-through in the short*

*run. Due to the presence of fixed costs of market access, in the long run also the set of firms that are actively exporting reacts to movements of the real exchange rate, with two associated consequences. Firstly, pass-through is larger than in the short run because long run export volume responds to relative costs due to changes in both the average firm size and in the number of firms. Secondly, the response of the market entry decision to changes in the relative cost of production affects only low quality firms, which fetch a relatively low price for their output. Exchange rate movements thus change the composition of actively exporting firms, with the consequence that aggregate price indexes overstate the actual extent of pass-through in the long run.*

## 2.1 Introduction

Why are movements of relative costs brought about by exchange rate movements passed on to customers only slowly, and never to a full extent? While research has greatly advanced our understanding of the direction and volume of trade flows and also of the location and productivity of industry, the often used assumption of constant markups has left us with little understanding of how firms set prices in the international economy.

In this paper, we first develop a competitive economy where the relative scarcity of goods of different qualities leads to pricing-to-market and markups that are determined by the local tightness of competition. We then analyze how changes in the relative cost of production affect prices and markups in the short and in the long run. While changes in the intensive margin lead to limited pass-through in the short run, the long run change of the set of firms actively exporting induces a further adjustment of market conditions to changes in the relative cost of production and also a shift in the composition of exporting good towards goods of different quality and prices, leading to a large measured degree of aggregate pass-through.

In our set up, which is derived from Landier and Gabaix (2006),[1] we think of industries as rather narrowly determined sectors such as German cars, Italian designer fashion or French red wines in which the country as a whole has some monopoly power, but there are potentially many domestic firms that compete in the sector. Although all firms in the industry are producing

---

[1]Landier and Gabaix (2006) explain how the matching of manager talent to firms of different size can explain the extremely large salaries for the CEO's of publicly listed companies.

the same good as their local competitors, the firms differ in the intrinsic quality of the good they produce, i.e. we think that certain goods such as a Porsche, an Armani shirt or a bottle of Chateau Dauzac are valued higher than others, such as a Volkswagen, a Benetton shirt or a lesser wine. Similarly, consumers are heterogenous with respect to how much they value quality or having a good from a certain industry at all. While some customers may put only a very small weight on driving a good car, others do find having a prestigious car very important and are hence willing to pay high prices for such goods. In the resulting equilibrium, we show that high quality consumers and high quality producers are matched together. The novel implication of our model of preferences is that the change in the prices from one quality to another depends on the valuation that the respective firms are matched with. The latter matching itself depends on relative supply and demand: in markets characterized by large supply, a firm that produces a good of a given quality is matched with a lower valuation consumer than in a market charaterized by low supply. Crowded markets are hence characterized by lower average prices and also by smaller differences in the price of low and high quality goods.

Summarizing, the key feature of our model of preferences is that markets are competitive, but within any industry, there is heterogeneity in both the quality of goods and in the valuation of goods by consumers. In equilibrium, high valuation customers are matched with high quality goods, and the relative density of consumer valuations and good's qualities determines the mapping of qualities into prices. The two basic results of this formulation are that firstly, markups are positive for all but the worst quality goods even though markets are competitive and that secondly, markups respond to relative supply and demand. In a market characterized by entry of a large number of firms the average firm is matched with a lower valuation customer and thus receives a lower average price.

Our model of preferences, it's underlying economic intuitions and also the resulting consequences for pass-through differ radically from existing frameworks that rely on Dixit and Stiglitz (1977) preferences of a CES demand, such as the ones used by Krugman (1980) or Melitz (2003). Our approach also differs from the other models of monopolistic competition reviewed by Krugman (1987) or the ones used by Dornbusch (1987), Melitz and Ottaviano (2005) or Atkeson and Burstein (2005).

The first difference between the existing literature and our framework is that we are able to explain pricing to market in a competitive framework, where markups are positive, determined by market conditions and not dependent on any form of market power. In any model of monopolistic competition, the markup is a function of the elasticity of substitution and the behavior of this markup hence depends exclusively on changes of the curvature of demand. Constant Elasticity of Substitution preferences – the workhorse of current international economic – imply constant markups and hence full exchange rate pass through. Setups using models of Cournot competitions, such as the one used by Dornbusch or by Atkeson and Burstein, imply more limited short run pass through, but the results hinge on few firms being active in a given market and also imply nearly full exchange rate pass through in realistic calibrations. Our model disentangles prices from the elasticity of substitution and thus allows for more sizeable price movements in a realistic setting.

The second key difference between our setup and the existing literature is the way in which we introduce firm heterogeneity. While we rely on Melitz's (2003) seminal insight that firms are heterogenous and market access is subject to fixed access costs, we assume that firms are equally productive, but differ in their inherent "quality." This subtle change results in findings for pricing to market and exchange rate-pass through that are in stark contrast to the composition effects present in the current literature. For simplicity, analyze the long run consequences of an increase in the relative wage in home in Melitz (2003).[2] In addition to an immediate one-to-one response of prices, in the long run less productive firms leave the export sector. Since less productive firms charge higher prices, the observed long run pass-through is smaller than the short run effect, which is in stark contrast to the empirical findings. In contrast, following the same increase in the relative wage at home, our model implies that quality firms that leave the export sector. Since low quality firms also receive low prices for their goods, in our model the average observed long run exchange rate pass-through is larger than the short run response, hence better matching both anecdotal and empirical evidence.

We nest our model of preferences in an international economy setting with many countries and many markets. Each country has a monopoly in one industry and it is endowed with a number of firms that are competing in the respective sector at home and abroad. Trade

---

[2]But see Baldwin (1988) and Baldwin and Krugman (1989)

allows these companies to expand internationally and each firm can do so by paying a fixed access fee for each market it wants to be active in. We assume that production takes place in two steps. First intermediate goods are produced at home and may then be shipped to international markets. Intermediate goods then have to be assembled and distributed locally, which requires additional inputs and can only be done in the final destination market. Setting up our production economy in this way is motivated by the often noted empirical regularity that the manufacturing terms of trade measured at the border are far less volatile than the real exchange rate, while local producer prices fluctuate with the exchange rate (see for example Atkeson and Burstein (2005)). Goldberg and Campa (2006) decompose the sources of this counter-intuitive result into the price changes of imported inputs and the sensitivity of local distribution margins to changes in the exchange rate. They find that indeed local distribution margins are a sizeable factor when trying to understand why exchange rate pass through is so limited and we hence allow for this channel of local adjustment. In our setup, if a firm is exporting to a specific market it has access to a local assembly plant and distribution network. Assembly and distribution at the local level is characterized by increasing marginal costs due to capacity constraints, hence leading to a well defined quantity decision of each firm.[3]

In general equilibrium, we show that countries that are "close" to the rest of the world (i.e. low average transport cost) are characterized by massive entry of firms, larger average firm size and also by entry of firms with lower average good quality than countries that are more distant from the rest of the world. Because markets are more crowded and the average quality of goods is lower, average prices are low in countries with large entry of foreign firms. We next analyze the short and the long-run consequences of an unanticipated one time appreciaiton of the real exchange rate, or an increase in domestic relative wages. We point out three distinct channels why measured exchange rate pass through is present in the data and point out short and long run effects.

In the short run, the set of firms active in the export sector is fixed, but each firm accommodates changes in the relative cost brought about by a change in the exchange rate by

---

[3]In a way our specificaiton of inputs being assembled or distributed locally matched the empirical facts presented by Atkeson and Kehoe very well. In our formulation, the actually traded goods (input goods) fluctuate with wage or exchange rate movements one to one, while the increasing cost of assembly feature ensures that aggregate price of good that are actually sold to final users reacts far less to these movements.

adjusting the quantity of its exports. Since the latter decrease when the home currency appreciates, export markets get relatively less crowded and thus prices measured in the foreign currency increase, leading to partial exchange rate pass-through in the short run. In the long run, the range of firms that are actively exporting changes, because in the presence of fixed costs of market access some low quality firms may no longer find it profitable to export at all. While the short run change in the intensive margin (volume of exports per firm) affects all firms equally, this change in the extensive margin affects only low quality firms, with two associated consequences. Firstly, after an appreciation less firms are active in the export sector and these markets are thus characterized by relatively lower demand and consequently upward pressure on prices. Due to this further adjustment of relative demand to changes in the cost of production, long run pass through – although smaller than 100% – is larger than in the short run.

In addition, the long run change in the set of firms being actively exporting also changes the average composition of exported goods, leading to even larger measures pass-through when evaluating aggregate data: an appreciation of the home currency drives out low quality firms that receive a low price for their good. The observed aggregate price is hence averaged over a set of higher priced firms, leading to an overestimation of long term pass through when using aggregate data. Incorporating this finding, we show that a researcher estimating pass through in the long run might actually arrive at the conclusion that long run pass through is equa to 100%.

In total, the present study shows what the determinants of short and long run pass-through are, and also the precise mechanisms at work. We furthermore pont out that in addition to two real effects present in our model, the exit and entry of low cost firms leads to an overestimation of long run exchange rate pass through because the composition of firms is always shifted to accomodate exchange rate fluctuations.

The structure of the paper is the following. In the next section, we introduce our model of preferences and production. We analyze the static equilibrium in Section 2 and the stationary equilibrium in Section 3. Finally, we present our main results of long and short run pass-through in Section 4 and Section 5 concludes.

71

## 2.2 The Economy

The model presented in this paper displays heterogeneity in both preferences and in good's quality. The resulting associative matching pairs high utility customers with high quality firms. We assume that each country has a monopoly in a certain industry, but many firms compete in that industry. In every market, prices are thus determined competitively but market tightness and thus markups depends on the relative mass of firms entering and the size of each firm. Although each single firm takes prices as a given, exchange rate movements are passed on to customers partly because short as well as long run supply respond to prices thereby mitigating initial cost movements.

### 2.2.1 Preferences

Let $j$ index an industry and assume that each country is endowed with a monopoly in exactly one industry. Since each industry is unambiguously associated with one country, $j$ also indexes the country of origin of a certain good (i.e. the exporting country or home in one industry). In our set up, we think of these industries as rather narrowly determined sectors such as German cars, French Wine or Italian designer fashion in which the country has some monopoly power, but there are potentially many firms that compete. Although all firms in the industry are producing the same good as their local competitors the firms differ in the intrinsic quality of the good they produce, i.e. we think that certain goods such as a Porsche are valued higher than others, such as a Volkswagen. Similarly, we think of consumers being heterogenous with respect to how much they value having a good from a certain industry at all. While some customers may put only a very small weight on driving a good car, others do find having a prestigious car very important and are hence willing to pay high prices. A second assumption we impose on preferences is that goods are indivisible and each consumer can at most use one good in each sector, i.e. you either have a car or you do not. Let $i$ index the market under consideration, i.e. the importing country. Thus each narrowly defined market is described by the index $i, j$, for example French wines sold in the US.

We analyze the utility of consumers in country $i$ who have access to goods from all countries $j \varepsilon J$[4] and an outside good $A$ that is produced domestically. We denote the valuation of a single customer in country $i$ for a good in sector $j$ by $v_{i,j}$ and index the different qualities of a good from country $j$ by $q_j$. Let $\Pi_{i,j}(v,q)$ denote the integer strategy of consumers mapping the decision of a costumer in country $i$ with valuation $v$ to buy a certain good of quality $q$ in sector $j$. The utility of the consumer is given by

$$U = \sum_{t=0}^{\infty} \frac{C_{i,t}^{\gamma-1}}{\gamma-1} \delta^t \tag{2.1}$$

Dropping the time subscript $t$, the instantaneous utility $C_i$ is given by

$$C_i = \sum_j \left( \max_{q_j \epsilon Q_j} [v_i q_j \min[1, \Pi_{i,j}(v,q)]] \right) + a A_i \tag{2.2}$$

Where $A$ is the quantity of the alternative good bought by the consumer. If a consumer decides to buy a good in a certain industry, she gets a utility that is proportional to both her valuation $v$ in that industry and the quality of that good $q$. The preferences formulated here hence display a complementarity between high valuation customers and high quality goods that in equilibrium leads to a well defined matching of firms and customers. The inner minimization in (2.2) reflects the fact that if a consumer buys more than one unit of from a certain firm, she cannot use the amount in excess of 1. Similarly, the maximization over firms in (2.2) reflects the fact that if the consumer buys from two firms in a given industry $j$, she can only use one. The consumer cannot save[5] and thus is subject to a budget constraint.

$$I_i \leq \sum_j \Pi_{i,j}(v,q) P_{i,j}(q_j) + P_i^A A \tag{2.3}$$

In what follows below, we assume that $P_i^A = 1$, implying that prizes are normalized by the outside good which also translates linearly into instantaneous utility. We also assume that

---

[4] $J$ includes the good that country $i$ has a monopoly in.

[5] Since we assume that income is large enough such that any consumer in equilibrium always buys a positive amount of the outside good $A$ and the outside good enters the instantenous utility linearly, allowing for savings would allow the consumer to smooth $C_i$ by altering the amount of the outside good consumed, but would not change any of the pricing results obtained below.

$I_i$ is large enough such that in equilibrium, all consumers buy a strictly positive amount of the alternative good. We also assume that at each point in time and in any market $i, j$, the distribution of taste shocks is given by a Pareto distribution.

$$F_v(v) \sim 1 - \left(\frac{v}{\overline{v}}\right)^{-\lambda_v} \tag{2.4}$$

$\overline{v}$ is the scale parameter of the distribution of valuation shocks. A higher $\overline{v}$ is associated with a higher average valuation and thus a higher average price of each good. The shape parameter $\lambda_v$ measures the fatness of the tails and is a measure of relative demand for high quality goods. In equilibrium, a smaller $\lambda_v$ is associated with a large demand for high quality goods and thus a higher price for these goods. We assume that $\lambda_v > 1$.

## 2.2.2  Production, Shipping and Distribution

We adopt an economic formulation of production that features trade in intermediate goods and local assembly of final goods. While the production of intermediary goods takes places in the home market, these inputs are subsequently shipped to each market and assembled and distributed locally.

Each economy is endowed with $N_j$ firms that are all active in industry $j$. Each firm is endowed with certain characteristics of the good it produces that can be summarized in the quality $q$ of its output. We assume that $q$ is time invariant and we also assume that the quality shocks are distributed equally in all countries with a Pareto distribution function.

$$F_q(q) \sim 1 - \left(\frac{q}{\overline{q}}\right)^{-\lambda_q} \tag{2.5}$$

$\overline{q}$ is the scale parameter of the distribution of quality shocks and a higher $\overline{q}$ is hence associated with higher average quality and equilibrium prices. The shape parameter $\lambda_q$ is a measure of how many high quality goods there are. In equilibrium, a larger $\lambda_q$ is associated with relatively low supply of high quality goods and thus a relatively high prices for these goods. We assume that $\lambda_q > 1$. Because in equilibrium, every firm's strategy is a function of the quality of its good, we do not index firms.

Firms can potentially sell their good in any market, but may choose not to do so because of

fixed costs of market access. We also assume that keeping the firm itself alive is costless, so that even a firm that does not even choose to access its home market is still alive. Following Melitz (2003), we assume that in order to be active in a certain market, a firm has to pay an up-front cost $Z_{i,j}$. We assume that this entry cost has to be paid stochastically: in the end of every period, every firm is drawn randomly for every market $i$ with constant and time independent probability $\delta < 1$ and may decide to enter market $i$. If the firm decides to do so, it has access to a local assemble plant or distribution infrastructure $\overline{D}_{i,j}$ until it is randomly selected again. If a firms opts not to pay the fixed access fee, it is inactive in the respective market but may decide to become active when it is randomly picked again in the future. Payment of the fixed access fee hence enables the firm to be active in one market for an average duration of $\frac{1}{1-\delta}$ periods. We assume that the cost of access is industry and country specific and we also assume that this access cost has to be paid in units of the outside good A.

The second key feature of our model of production is that we require a two step production function of goods, where intermediate goods are produced at home and may then be shipped to international markets. Intermediate goods then have to be assembled and distributed locally, which requires additional inputs and can only be done in the final destination market. Setting up our production economy in this way is motivated by the often noted empirical regularity that the manufacturing terms of trade measured at the border are far less volatile than the real exchange rate while local producer prices fluctuate with the exchange rate (see for example Atkeson and Burstein (2005)). Goldberg and Campa (2006) decompose the sources of this counter-intuitive result into the price changes of imported inputs and the sensitivity of local distribution margins to changes in the exchange rate. They find that indeed local distribution margins are a sizeable factor when trying to understand why exchange rate pass through is so limited and we hence allow for this channel of local adjustment.

In our setup, a firm uses $1/A_j$ units of domestic labor (i.e. labor in country $j$) to produce one unit of the intermediate good. It can then ship these intermediate goods to any market, where they are assembled locally and can subsequently be sold. We assume that assembly uses two inputs, a fixed amount infrastructure $\overline{D}_{i,j}$ and the intermediate good that is shipped from home. Denote the quantity of intermediate goods shipped by a firm of quality $q$ from country $j$ to $i$ by $I_{i,j}(q)$ and the final amount assembled by the same firm in that market by the same

firm by $S_{i,j}(q)$.

If a firm wants to be able sell $S_{i,j}(q)$ units in market $i,j$, it has to ship a sufficient amount of intermediate goods for both direct use in assembly as well as to build instantaneous assembly capacity. Denote the amount of intermediate goods used in local assembly by $I_{i,j}^D(q)$. The instantaneous production function of the assembly and distribution capacity is of Cobb Douglas type

$$D_{i,j}(q) = (1 + 1/\eta) I_{i,j}^D(q)^{\eta/(\eta+1)} \overline{D}_{i,j}^{1/(\eta+1)} \qquad (2.6)$$

Where $\overline{D}_{i,j}$ is the initially installed fixed amount of infrastructure that is assumed to be a constant and $\eta > 0$ is measure of the factor shares of inputs and fixed infrastrucutre. In equilibrium, as will be established below the elasticity of supply with respect to the wage is equal to $\eta$. In order to be able to assemble (and distribute) $D_{i,j}(q)$ units of the good, the firm needs both the fixed infrastructure of $\overline{D}_{i,j}$ as well as $I_{i,j}^D(q)$ intermediate goods that depreciate every period since they are used in the process of assembly.

Secondly, to produce one unit of the final good, the firm needs sufficient assembly capacity as well as again one unit of the intermediate good. Denote the amount of intermediate goods used as inputs for the production of the final good by $I_{i,j}^A(q)$. Output of final goods $S_{i,j}(q)$ is given by the minimum of the assembly capacity or the actual amount of inputs used for assembly.

$$S_{i,j}(q) = \min\left[I_{i,j}^A(q), D_{i,j}(q)\right] \qquad (2.7)$$

Where feasibility requires that the total amounts of inputs that arrive in market $i,j$ must be less that the total amount used directly in production or to build assembly infrastructure. Shipping is subject to iceberg transportation cost and thus

$$I_{i,j}(q)/\tau_{i,j} < I_{i,j}^A(q) + I_{i,j}^D(q) \qquad (2.8)$$

In what follows below we assume that any firm that in active in a certain market can instantaneously adjust its supply decision $S_{i,j}(q)$ in a given market it is active in. The key feature of local assembly (2.6) is that $\overline{D}_{i,j}$ is fixed, implying an increasing effective marginal cost schedule. We interpret this formulation as the short run capacity constraints that each firm faces on a

local level. Since all firms produce less when their costs of production rises, they reduce their production following an appreciation of their home currency, leading to limited pass through in the short run. We next turn to the static analysis of the economy, but for reasons of clarity, we first we summarize the timing of our economy.

- Each country is endowed with a monopoly in 1 industry and $N_j$ firms. Each firm is endowed with a quality draw $q$

- Past decisions have determined the entry decisions for all firms in all countries $j$ into all countries $i$.

- Firms simultaneously choose production of intermediate goods, the size of shipments to each market and the quantity of local assembly for each market they are active in.

- Given the valuations of customers and the supply decisions, prices are determined.

- Firms realize their profits, some firms get randomly drawn and decide whether to (re-)enter certain markets they were drawn for.

## 2.3 Analysis

We next analyze the static equilibrium of our economy. It is important to point out that our timing is crucial to establish the existence of an equilibrium: because production and shipping decisions have been made before prices are determined, the matching of consumers to goods is ex post easy to determine. To solve for most of our pricing results, we do not need to rely on the two closed form distributions of taste and quality shocks (2.5) and (2.4).

**Definition 5** *(Static Competitive Equilibrium) A static competitive equilibrium in market $i, j$ is defined by a firm strategy mapping prices into a supply decision for all active firms in market $j$ and an entry strategy in country $i$ for all firms $j \varepsilon J$ and a mapping of consumer valuations and qualities into the buying strategy $\Pi_{i,j}(v, q)$. There is an associated matching of qualities and valuations $v : v_{i,j}(q)$ and a unique mapping of the quality of each good into prices $p : p_{i,j}(q)$.*

The model of preferences is sufficiently rich such that the analysis is complicated at first sight, and we thus first establish in the following Lemma that we can analyze preferences in each market on their own and independent of all other markets.

**Lemma 3** *(Indirect utility function). Instead of evaluating (2.1), (2.2) and (2.3), one can evaluate the indirect utility function in each market $i, j$*

$$\widetilde{U}_{i,j}(v) = \max\left[0; \max_q \left[vq - ap_{i,j}(q)\right]\right] \tag{2.9}$$

*Proof. First note that since we assume that $A \geq 0$ for all consumers each unit of income has a marginal utility of $a$. A consumer of valuation $v$ hence only buys a good in industry $j$ if there exists at least one good in that market $i$ such that $vq - ap_{i,j}(q) \geq 0$. Secondly, since she can only use 1 good in that industry, she chooses the good that maximizes $vq - ap_{i,j}(q)$.* ∎

Evaluating the indirect utility function (2.9) greatly simplifies the analysis because buying decisions in one industry are not affected by choices in other industries. Before we determine prices, we need to know what the actual matching of customers and goods is. Here it noteworthy that in order for our equilibrium to be determinate, we need that production and assembly decisions are already made once the matching is made. We assume that producers correctly anticipate the matching, rely on this structure of the timing of the economy to determine the matching of the economy.

**Lemma 4** *(Assortative Matching) For all parameters of the model that allow any transaction to take place, a stable matching exists. In each industry and country a buyer with valuation $v_{i,j}(q)$ is matched with seller with quality $q$, and all buyers with valuation above $v_{i,j}(q)$ are matched with sellers with quality above $q$. The matching function is a mapping of $v : v_{i,j}(q)$ such that*

$$N_j \int_q^\infty S_{i,j}(q) f_q(x)\, dx = L_i \int_{v_{i,j}(q)}^\infty f_v(v)\, dx \tag{2.10}$$

**Proof.** Imagine that there exists a transaction schedule where (2.10) is violated and two customers 1 and 2 of valuations $v_1 < v_2$ buy goods of qualities $q_1$ and $q_2$ respectively, where $q_1 > q_2$. The two consumers have no influence on prices, and there thus exists a Pareto

improvement in which the two trade their goods and customer 2 compensates 1 for the loss in quality, while still making a net gain of $(v_2 - v_1)(q_1 - q_2)$. ∎

The assortative matching is a very strong result that greatly eases the analysis below. There may exist a valuation that is low enough such that the consumer decides not to buy when matched. Denote the lowest valuation customer that is actually buying by $v_{i,j}^{\min}$. The matching has to hold for every $q$ and $v$ and thus market clearing determines $v_{i,j}^{\min}$.

**Corollary 6** *(Market Clearing) Denote the lowest quality firm of sector $j$ that is present in market $i$ and actively selling by $q_{i,j}^{\min}$. Denote the lowest valuation customer that is buying in the same sector by $v_{i,j}^{\min}$. Market Clearing implies*

$$N_j \int_{q_{i,j}^{\min}}^{\infty} S_{i,j}(q) f_q(x) \, dx = L_i \int_{v_{i,j}^{\min}}^{\infty} f_v(v) \, dv \qquad (2.11)$$

Because matching is well defined associating each customer to one specific good, prices can be solved for. The crucial insight is that in a competitive market the local scarcity goods determines the price schedule. Consider a given customer who is matched with a certain good of quality $q$, but can potentially also buy goods of other quality. Each producer takes as a given the prices of his competitors, and just sets his price such that the consumer marginally prefers his good. In equilibrium, the relative supply and demand of goods at every quality $q$ determines the shape of the pricing schedule. We establish next the two main results of our model of pricing.

**Proposition 7** *(Pricing Schedule) Denote the equilibrium function in market $i,j$ that maps a good's quality $q$ into it's price by $p_{i,j}(q)$. The indirect utility function (2.9) and the assortative matching lead to a unique relationship between the prices of goods of different qualities such that*

$$p_{i,j}(q) = a^{-1} \int_{q_{i,j}^{\min}}^{q} v_{i,j}(x) \, dx + C_{i,j} \qquad (2.12)$$

*Where the constant of integration is determined below. Also is always true that*

$$\frac{\partial^2 p_{i,j}(q)}{\partial^2 q} = \frac{\partial v_{i,j}(q)}{\partial q} \geq 0$$

79

**Proof.** Consider a customer of given valuation $v_{i,j}(q)$ who is matched with good $q$ but considers buying a variety of a slightly higher quality $q\prime$. The price differential that would make her indifferent between the two equals $v_{i,j}(q)(q\prime - q)$. That is for both firms to have nonnegative demand, the price differential has to satisfy

$$p_{i,j}(q') - p_{i,j}(q) = v_{i,j}(q)(q\prime - q)$$

Or in the limit when firms are comparatively similar such that $q\prime - q \to 0$

$$\frac{\partial p_{i,j}(q)}{\partial q} = v_{i,j}(q)$$

Integration of the above expression yields the price schedule (2.12).

Finally, implicit derivation of matching (2.10) implies that $\frac{\partial v_{i,j}(q)}{\partial q} \geq 0$, thus $\frac{\partial^2 p_{i,j}(q)}{\partial^2 q} \geq 0$. The latter inequality holds strictly if valuations are distributed such that $f_v(v) > 0$ for the whole support of $v$. $\blacksquare$

We determine the constant of integration in (2.12) below. The price schedule (2.12) is a central equation that is responsible for some of our main results for the exact nature of exchange rate pass through. The key importance is that the slope of the pricing schedule is increasing in the valuation that the respective consumer is actually matched with. When there are less customers, the average firm is matched with a lower valuation customer, and the whole pricing schedule is lower. As we show below, the relative density of valuations and quality draws determine the steepness of the pricing schedule. If high valuations are relatively scarce ($\frac{\lambda_v}{\lambda_q}$ is small) a given high quality good is matched with a very high valuation, and hence fetches a high price. Consider, for example a market with homogenous consumers ($\lambda_v \to \infty$) and assume that the constant of integration is 0 in that case (this would indeed be an equilibrium consequence). This case results in the familiar case that a 1% higher quality good is associated with a 1% higher price. In the case of no heterogeneity of consumers, prices are a linear function of quality, but they respond stronger when higher quality goods are matched with higher valuation goods.

We next solve for the constant of integration, which depends on actual supply. Although matching takes place after the supply decision, firms anticipate the resulting prices and produce accordingly. Evaluating the first order condition we can express the input requirement function

80

in terms of the effective marginal cost schedule, which we denote by $c_{i,j}\left(S_{i,j}\left(q\right)\right)$. Also, we introduce the effective real wage $\omega_{i,j} \equiv \frac{e_{i,j}\omega_{i,j}}{w_i}$ between country $i$ and $j$. Producing one unit of intermediate good has a cost of $\frac{\omega_{i,j}}{A_j}$ measured in terms of $i$ currency. The production decision over (2.6), (2.7) and (2.8) hence solves for the effective marginal cost schedule

$$c_{i,j}\left(S_{i,j}\left(q\right)\right) = \theta^{-1}\left(\frac{S_{i,j}\left(q\right)}{\overline{D}_{i,j}}\right)^{1/\eta}\frac{\tau_{i,j}\omega_{i,j}}{A_j} + \frac{\tau_{i,j}\omega_{i,j}}{A_j} \tag{2.13}$$

Reflecting the two parts of the production setup introduced above, the cost of production has two parts. He first term in (2.13) reflects the cost of local assembly and distribution, which is increasing due to the fixed infrastructure $\overline{D}_{i,j}$. The second part reflects the fact that production requires one unit of intermediate good. Statically, the entry decision of a firm is sunk and a strategy of a firm is thus the function of the price it gets in a certain market and the cost of its input into an output decision. The optimally condition is that shipments of inputs are such that

$$c_{i,j}\left(S_{i,j}\left(q\right)\right) = p_{i,j}\left(q\right) \tag{2.14}$$

Since marginal costs are increasing in the short run, a higher equilibrium price is associated with a larger output. We solve for the output decision in closed form in the next section, but first characterize the constant of integration in (2.12).

**Proposition 8** *(Lowest Observed Price) Denote the lowest quality firm that has paid the access cost in market $i,j$ by $\widetilde{q}_{i,j}$, and the lowest quality good that is actually sold by $q_{i,j}^{\min}$, where $q_{i,j}^{\min} \geq \widetilde{q}_{i,j}$. The price of the lowest quality good actually sold $p_{i,j}\left(q_{i,j}^{\min}\right)$ is given by one of the following three cases.*

*A (Weak) Excess Potential Demand: $v_{i,j}^{\min} \geq \overline{v}$, $q_{i,j}^{\min} = \widetilde{q}_{i,j}$.*

$$v_{i,j}^{\min}q_{i,j}^{\min} - ap_{i,j}\left(q_{i,j}^{\min}\right) = 0 \tag{2.15}$$

*B Excess Potential Supply: $v_{i,j}^{\min} = \overline{v}$, $q_{i,j}^{\min} > \widetilde{q}_{i,j}$*

$$p_{i,j}\left(q_{i,j}^{\min}\right) = \min_{S_{i,j}(q)}\left[\frac{\int_0^{S_{i,j}(q)} c_{i,j}\left(x\right)dx}{S_{i,j}\left(q\right)}\right] \tag{2.16}$$

$$v_{i,j}^{\min} q_{i,j}^{\min} - a \min_{S_{i,j}(q)} \left[ \frac{\int_0^{S_{i,j}(q)} c_{i,j}(x)\,dx}{S_{i,j}(q)} \right] = 0 \qquad (2.17)$$

**Proof.** In all three cases, market clearing (2.11) determines the relation between $v_{i,j}^{\min}$ and $q_{i,j}^{\min}$. Consider first case A. For given production amounts for any firm $S_{i,j}(q)$, market clearing (2.11) can imply that $N_j \displaystyle\int_{q_{i,j}^{\min}=\widetilde{q}_{i,j}}^{\infty} S_{i,j}(q) f_q(x)\,dx < L_i$ so supply does not suffice to satisfy all potential demand. The lowest observed price is then determined by the valuation of the lowest actual buyer who competes with a set of marginally lower valuation customers for the good of quality $\widetilde{q}_{i,j}$ and is thus breaking even, implying (2.15).

Consider next Case B, where all potential customers by one good but there is still potential excess supply. Firms anticipate that not all companies can ex post sell their good and hence some firms do not produce. The marginal firm that actually produces faces a set of marginally lower quality competitors and thus produces at the quantity that minimizes average variable costs and also fetches a price equal to $p_{i,j}\left(q_{i,j}^{\min}\right) = \min_{S_{i,j}(q)} \left[ \frac{\int_0^{S_{i,j}(q)} c_{i,j}(x)dx}{S_{i,j}(q)} \right]$.

Finally consider case C, where the minimum average cost of production is high enough such that some potential customers actually do not buy. Although there is a mass of potential customers, some valuations are just too low such that firms cannot produce cheap enough. The marginal customer just breaks even, while the marginal firm produces at the minimum of average variable costs. The latter case can be an equilibrium whenever

$$\overline{v}\widetilde{q}_{i,j} - a \min_{S_{i,j}(q)} \left[ \frac{\int_0^{S_{i,j}(q)} c_{i,j}(x)\,dx}{S_{i,j}(q)} \right] < 0$$

■

We argue below that in any steady state with a positive access cost only case A is the relevant one, since only in this case the lowest valuation firm actually sells a positive quantity and is induced to enter the market in the first place.[6]

---

[6]A fourth case where all consumers buy and all firms sell a positive amount exists only when market clear *exactly* in that equilibrium, which is gernerically not true. The definition of Case A hence includes the case where $v_{i,j}^{\min} = \overline{v}$, $q_{i,j}^{\min} = \widetilde{q}_{i,j}$ and is thus termed Weak Excess Potential Demand.

For given prices, the firm's entry decision can be analyzed. The entry decision of each firm depends on the expected profits it fetches in market $i, j$ and the respective access cost. We assume that there is an industry and country specific entry cost $Z_{i,j}$. In each period that the firm is active it makes an operating profit $P_{i,j,t}(q) S_{i,j,t}(q) - \int_0^{S_{i,j,t}(q)} c_{i,j,t}(s) \, ds$. The marginal firm $\tilde{q}_{i,j}$ that still chooses to access the market from the next period on breaks even on it investment and $\tilde{q}_{i,j}$ solves

$$Z_{i,j} = E_t \left[ \sum_{t=0}^{\infty} \beta^t \delta^t \left( P_{i,j,t}(\tilde{q}_{i,j}) S_{i,j,t}(\tilde{q}_{i,j}) - \int_0^{S_{i,j,t}(\tilde{q}_{i,j})} c_{i,j,t}(s) \, ds \right) \right] \qquad (2.18)$$

We next turn to the stationary equilibrium of the economy.

## 2.4 Stationary Equilibrium

In this section, we characterize the steady state of the economy. In the next section, we shock an economy in steady state with an unanticipated exchange rate shock and analyze short and long run exchange rate pass through. In our setup, a Stationary Equilibrium is a Competitive Equilibrium characterized by a constant firm Entry Strategy such that all firms always enter market $i, j$ if $q \geq q_{i,j}^{\min}$ and never otherwise. Along any stationary equilibrium, aggregate consumption in each industry is constant and the interest is thus equal to the discount rate. We first establish the following Lemma to know how the minimum observed price is determined.

**Lemma 5** *(Market Conditions in Steady State) Assume that $Z_{i,j} > 0$. Then in any stationary equilibrium $v_{i,j}^{\min} > \bar{v}$ and $q_{i,j}^{\min} = \tilde{q}_{i,j}$.*

**Proof.** Assume that either B or C are the relevant case in a steady state. Then, the lowest active firm in market $i, j$ with quality $\tilde{q}_{i,j}$ actually does not sell. In a steady state the price for a good of given quality is constant and a firm that does not sell at present will never sell. Whenever such a firm comes up for repaying the entry cost it would choose not to and leave the market. Hence, while the path of history may have induced an initial mass of low quality firms to have a distribution network in place in market $i, j$ but not be active, this mass converges to 0 with rate of $\delta$ per period and is not present in a stationary equilibrium. ∎

In any steady state, all firms that have accessed the market are also selling (i.e. $q_{i,j}^{\min} = \tilde{q}_{i,j}$) and the competition of the lowest customer that actually buys a good with marginally lower valuation customers determines the minimum observed price in market $i, j$.

$$p_{i,j}\left(q\right) = a^{-1} \int_{q_{i,j}^{\min}}^{q} v_{i,j}\left(x\right) dx + a^{-1} v_{i,j}^{\min} q_{i,j}^{\min} \tag{2.19}$$

Recall our explicit distribution of preference shocks (2.4) and quality draws (2.5). To find an equilibrium, we guess and verify a price schedule. Prices are a function of the cost of production, the relative number of customers and firms and the quality of each good. We are only able to solve for our model in closed form when we assume a specific relation between the cost of market entry and the size of the local assembly and distribution infrastructure $\overline{D}_{i,j}$. We discuss a slight modification of our model of preferences that enables a closed form solution in a more general case in Appendix A.

**Proposition 9** *(The Steady State Pricing Schedule ) Assume that the fixed access cost of entering a market satisfies $\frac{Z_{i,j}}{\overline{D}_{i,j}} = (1 - \beta\delta) \left(\frac{\lambda_q - \eta}{\lambda_q + \lambda_v}\right)^{-(\eta+1)} / (\eta + 1)$*
*Then, in a stationary equilibrium, the price schedule mapping a good's quality into a unique price is given by*

$$p_{i,j}\left(q\right) = \Lambda \left(\frac{N_j}{L_i}\right)^{-\frac{1}{\lambda_v + \eta}} \left(\frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)^{\frac{\eta}{\lambda_v + \eta}} q^{\frac{\lambda_q + \lambda_v}{\lambda_v + \eta}} + \frac{\tau_{i,j}\omega_{i,j}}{A_j} \tag{2.20}$$

*where $\Lambda^{1 - \frac{1}{\lambda_v}\eta} \equiv a^{-1}\overline{vq}^{-\frac{\lambda_q}{\lambda_v}} \left(\overline{D}_{i,j}\theta^{\eta}\frac{\lambda_v + \eta}{\lambda_q - \eta}\frac{\lambda_q}{\lambda_v}\right)^{-1/\lambda_v}\frac{\lambda_v + \eta}{\lambda_q + \lambda_v}$ is a constant that depends on the shape and mean of the distribution of preferences and quality draws and on the parameters of the marginal cost function (2.13).*

**Proof.** To establish the price schedule, we guess and verify that (2.20) is an equilibrium. At each moment, each of the firms from country $j$ that are active in market $i$ are maximizing profits and hence producing up to the point where marginal costs are equal to the price of their good. The supply decision of every firm (2.14) is a function of the price of the respective good

$$S_{i,j}\left(q\right) = \overline{D}_{i,j} \left(p_{i,j}\left(q\right) \left(\frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)^{-1} - 1\right)^{\eta} \tag{2.21}$$

Since higher quality firms fetch a higher price for their good but have the same cost function, they are large in equilibrium. Incorporating the output of each firm conditional on its quality, the general matching function (2.10) solves for the equilibrium matching

$$v_{i,j}(q) = a \frac{\lambda_v + \eta}{\lambda_q + \lambda_v} \Lambda \left( \frac{N_j}{L_{i,j}} \right)^{-\frac{1}{\lambda_v}} \left( \frac{\tau_{i,j}\omega_{i,j}}{A_j} \right)^{\frac{\eta}{\lambda_v + \eta}} q^{\frac{\lambda_q - \eta}{\lambda_v + \eta}} \tag{2.22}$$

Applying the previously obtained results for the slope of prices (2.12) from Proposition 1, we can confirm our initial guess up to the constant of integration.

$$p_{i,j}(q) = \Lambda \left( \frac{N_j}{L_{i,j}} \right)^{-\frac{1}{\lambda_v}} \left( \frac{\tau_{i,j}\omega_{i,j}}{A_j} \right)^{\frac{\eta}{\lambda_v + \eta}} q^{\frac{\lambda_q + \lambda_v}{\lambda_v + \eta}} + C_{i,j} \tag{2.23}$$

To solve for the constant of integration, we need to know the entry decision of firms. In a stationary equilibrium, the firms follow a cutoff rule such that all firms of $q > \widetilde{q}_{i,j}$ enter the respective market, where the cutoff firm makes 0 profits from entering a given market $i, j$. In steady state, the entry decision (2.18) simplifies to

$$\frac{Z_{i,j}}{1 - \beta\delta} \frac{\tau_{i,j}\omega_{i,j}}{A_j} = P_{i,j}(\widetilde{q}_{i,j}) S_{i,j}(\widetilde{q}_{i,j}) - \int_0^{S_{i,j}(\widetilde{q}_{i,j})} c_{i,j}(s) \, ds \tag{2.24}$$

Again noting that in steady state every firm that has chosen to access a market also sells a positive amount, market clearing (2.11) implies that in steady state

$$\left( q_{i,j}^{\min} \right)^{\frac{\lambda_q + \lambda_v}{\lambda_v + \eta}} = \Lambda^{-1} \left( \frac{1 + \eta}{1 - \beta\delta} \frac{Z_{i,j}}{\overline{D}_{i,j}} \right)^{1/(1+\eta)} \left( \frac{\tau_{i,j}\omega_{i,j}}{A_j} \right)^{-\frac{\eta}{\lambda_v + \eta}} \left( \frac{N_j}{L_{i,j}} \right)^{\frac{1}{\lambda_v}} \tag{2.25}$$

The steady state level $q_{i,j}^{\min}$, the implied consumer valuation (2.22) and the associated minimum price (2.19) determine the minimum price and thus the constant of integration $C_{i,j}$, which equals $\frac{\tau_{i,j}\omega_{i,j}}{A_j}$. ∎

How do market conditions determine prices? We have to take into consideration two parts, lowest observed price $p_{i,j}\left( q_{i,j}^{\min} \right)$ as well as the pricing schedule for higher quality goods. Consider first how the slope of the pricing schedule is affected by the tightness of competition. When there are relatively more potentially active firms $N_j$ for a given number of potential consumers $L_i$ markets are more crowded and the average valuation is lower leading to lower

price in equilibrium. The same is the case when the cost of production $\left(\frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)$ is low such that each firm produces a larger amount for a given price of its good. The toughness of competition affects both the slope and the average of the pricing schedule. The pricing slope $\frac{\partial p_{i,j}(q)}{\partial q}$ is steeper in more crowded markets. Because of our that valuations and quality are distributed Pareto, all measures of market crowing (number of firms, average production) shift the price schedule with constant elasticity.

Secondly, consider how the constant of integration is determined. This effect works through the entry decision in the market. The latter condition and the specific the effective cost function (2.13) and the assumption on the form of the entry cost $Z_{i,j}$ leads to the constant of integration $C_{i,j}$ always being equal to 0. In equilibrium, a firm makes per period operating profits equal to a constant fraction of sales minus the direct cost of inputs $I_{i,j}^A(q)$ such that profits are a function of $S_{i,j}(q)\left(p_{i,j}(q) - \frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)$. In equilibrium, the lowest active firm just breaks even, and since $S_{i,j}(q)$ is a pure function of prices relative to costs, the following must hold for some constant $\Theta$.[7]

$$\overline{D}_{i,j}\left(p_{i,j}\left(q_{i,j}^{\min}\right)\left(\frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)^{-1} - 1\right)^{\eta}\left(p_{i,j}\left(q_{i,j}^{\min}\right) - \left(\frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)\right) = \left(\frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)\Theta \quad (2.26)$$

Since the right hand side is linear in the wage, the left hand side is too, such that when $\frac{\tau_{i,j}\omega_{i,j}}{A_j}$ changes, the equality still holds. This can only be the case when $p_{i,j}\left(q_{i,j}^{\min}\right)A_j/\tau_{i,j}\omega_{i,j}$ is constant for any $\omega_{i,j}$. Since discounted profits have to equal the access cost, thus tying down the minimum price as a linear function of the costs of inputs $\frac{\tau_{i,j}\omega_{i,j}}{A_j}$ one-to-one.

## 2.5 Exchange Rate Pass-Through

We are now able to evaluate the short and long run impacts of an unanticipated one time change in relative costs of production brought about the a fluctuating exchange rate. We first need to establish our precise concept of exchange rate pass-through.

**Definition 6** *(Pass Through Elasticity) Define the Price Pass Through Elasticity as the per-*

---

[7]The presence of the term $\left(\frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)$ on the right hand side of (2.26) reflects the assumption that the fixed acces cost has to be paid in intermediate goods.

*centage change of the equilibrium price for a constant quality good q conditional on the firm still being active in that market.*

$$\sigma_{p_{i,j}(q),\omega_{i,j}} \equiv \frac{\partial p_{i,j}(q)}{\partial \omega_{i,j}} \frac{\omega_{i,j}}{p(q)}$$

*Define the producer price index (PPI) pass through elasticity as*

$$\sigma_{PPI_{i,j},\omega_{i,j}} \equiv \frac{\partial PPI_{i,j}}{\partial \omega_{i,j}} \frac{\omega_{i,j}}{PPI_{i,j}}$$

*Where the producer price index is the quantity weighed average price of final goods.*

$$PPI_{i,j} \equiv \frac{\int_{q_{i,j}^{\min}}^{\infty} f_q(s) p_{i,j}(s) S_{i,j}(s) ds}{\int_{q_{i,j}^{\min}}^{\infty} f_q(s) S_{i,j}(s) ds}$$

We introduce two different concepts of exchange rate pass-through. Both are measured in the standard elasticity sense, i.e. the percentage response following a given change of the exchange rate or wages at home or abroad. As a benchmark one can assume either the standard competitive economy or a standard model of monopolistic competition alla Krugman (1980). In both cases the pass through elasticity constant over time, 100% for both $\sigma_{p_{i,j}(q),\omega_{i,j}}$ and $\sigma_{PPI_{i,j},\omega_{i,j}}$. In our setup, in contrast exchange rate is generically different from 100%, different in the short and long run and also not equal for $\sigma_{p_{i,j}(q),\omega_{i,j}}$ and $\sigma_{PPI_{i,j},\omega_{i,j}}$.

## 2.5.1 Short Run Pass Through

**Lemma 6** *(Small Change in the Wage Rate) For a small movement of the exchange rate, the instantaneous exchange rate pass through of the measured PPI $\sigma_{PPI_{i,j},\omega_{i,j}}$ and of each single price $\sigma_{p_{i,j}(q),\omega_{i,j}}$ is constant for all qualities q and equal to*

$$\sigma_{PPI_{i,j},\omega_{i,j}} = \sigma_{p_{i,j}(q),\omega_{i,j}} = \frac{\eta}{\lambda_v + \eta}$$

**Proof.** Assume that steady state wage is $\omega_{i,j}$ leading to steady state cutoff $q_{i,j}^{\min*}$ given by (2.24). At impact of the change in the wage, $q_{i,j}^{\min*}$ is fixed, so if firms didn't adjust their quantity, prices would be unchanged, leading to 0 pass through. However, the quantity each

firm produces responds with elasticity $\eta$. This immediately counteracts the initial change of the wage rate because with less, the average firm is matched with a higher customer and hence achieves a higher price measured in foreign currency. In total, the change of the whole price schedule in response to a small change of the exchange rate is given by

$$\frac{\omega_{i,j}}{p\left(q_{i,j}^{\min}\right)} \frac{\partial p\left(q_{i,j}^{\min}\right)}{\partial \omega_{i,j}} \Big|_{q_{i,j}^{\min} \approx q_{i,j}^{\min *}, \omega_{i,j} \approx \omega_{i,j}} = \eta/\left(\lambda_v + \eta\right)$$

$$\frac{\omega_{i,j}}{p\left(q\right) - p\left(q_{i,j}^{\min}\right)} \frac{\partial \left(p\left(q\right) - p\left(q_{i,j}^{\min}\right)\right)}{\partial \omega_{i,j}} \Big|_{q_{i,j}^{\min} \approx q_{i,j}^{\min *}, \omega_{i,j} \approx \omega_{i,j}} = \eta/\left(\lambda_v + \eta\right)$$

Consequently,$\sigma_{PPI_{i,j},\omega_{i,j}} = \sigma_{p_{i,j}(q),\omega_{i,j}} = \eta/\left(\lambda_v + \eta\right)$ ∎

Why is the short run pass-through elasticity to a first order exactly equal to $\eta/\left(\lambda_v + \eta\right)$ for all firms equally? Imagine that while other firms face constant costs, a single firm of quality $q$ faces an increase in cost. The single firm is not price taker, but it can accommodate the change in its cost by adjusting the quantity of its output. Along the steady state price schedule $p_{i,j}\left(q\right)^*$, the firms output responds to the cost its faces with an elasticity of $\eta$:

$$\frac{\omega_{i,j}}{S_{i,j}\left(q\right)} \frac{\partial S_{i,j}\left(q\right)}{\partial \omega_{i,j}} \Big|_{p_{i,j}(q) = p_{i,j}(q)^*} = -\eta$$

When the wage change does not affect only one firm, but all together, this change in the total amount produced affects the equilibrium price schedule. If a small increase in the cost of production implies a large adjustment of the typical firms output, a lot of the pricing changes will be passed on to customers, because low valuation customers are driven out of the market since the quantity supplies contracts. While the first round effect of a given percentage increase in the wage is that all firms adjust their quantity with an elasticity of $\eta$. Then since total output expands and prices the iterated equilibrium effect implies and elasticity of $\frac{\eta}{\lambda_v + \eta}$ which is smaller than the first round effect $\eta$. We can also make statements about a larger change in the wage rate. We can no longer solve for the prices in closed form, but we can establish that the effects must be smaller than to a first order.

**Corollary 10** *(Exchange Rate Pass Through in the Short Run) For large movements of the exchange rate the short run pass-through elasticity for the measured PPI $\sigma_{PPI_{i,j},\omega_{i,j}}$ and for*

*each single price* $\sigma_{p_{i,j}(q),\omega_{i,j}}$ *is smaller than* $\eta/(\lambda_v + \eta)$ *for larger movements of the exchange rate.*

**Proof.** For the first part of the proposition, see Lemma 4. Consider next a sizeable increase in the cost of production. The magnitude of the percentage change in production for a larger change in the exchange rate is smaller than to a first order

$$\frac{\omega_{i,j}}{S_{i,j}(q)} \frac{\partial S_{i,j}(q)}{\partial \omega_{i,j}} \bigg|_{q_{i,j}^{\min *}, \omega_{i,j} > \omega_{i,j}^*} > -\eta$$

Since the first round change in quantity is larger than what a first order would imply we can establish that

$$0 < \frac{\omega_{i,j}}{p(q)} \frac{\partial p(q)}{\partial \omega_{i,j}} \bigg|_{q_{i,j}^{\min *}, \omega_{i,j} > \omega_{i,j}^*} < \frac{\eta}{\lambda_v + \eta}$$

■

For extreme values of the exchange rate fluctuation, some firms start being crowded out of the market. The observed minimum price adjusts leading to $q_{i,j}^{\min} < \widetilde{q_{i,j}}$ (I.e. case B or C from Proposition 2) and the minimum price is of

$$p_{i,j}\left(q_{i,j}^{\min}\right) = \min_{S_{i,j}(q)} \left[ \frac{\int_0^{S_{i,j}(q)} c_{i,j}(x)\,dx}{S_{i,j}(q)} \right] = 0$$

In the short run, we can hence establish that the pass through elasticity is at most $\eta/(\lambda_v + \eta)$, and somewhat smaller when the exchange moves more. We next establish the long run pass-through.

### 2.5.2 Long Run Pass-Through

With our convenient assumption on the form of the entry cost $Z_{i,j}$, long run exchange rate pass through can be pinned down precisely. We establish in this section that the pricing schedule itself adjusts to exchange rate movements due to two channels. Firstly the adjustment of the quantity that every firm produces adjusts, and therefore prices move against exchange rates. Secondly, because a mass of firms enters and leaves a market when the exchange rate changes, exchange rate movements are transmitted even more in the long run than in the short run.

A further consequence of the long run change in the set of firms being active exporters is

that the composition of active exporters changes, leading to even larger measured pass through when evaluating aggregate data. An appreciation of the home currency drives out low quality firms that receive a low price for their goods. The observed aggregate price is hence averaged over a set of higher priced goods, leading to an overestimation of long term pass through when using aggregate data. Incorporating this finding, we show that a researcher estimating pass through in the long run might actually arrive at the conclusion that long run pass through is equal to 100%. In most of our results presented here, we evaluate the steady state pricing schedule (2.20).

The following Proposition Summarizes our results for the long run pricing schedule.

**Proposition 11** *(Exchange Rate Pass-Through in the Long Run) Assume that $\frac{Z_{i,j}}{D_{i,j}} = (\eta + 1)^{(\eta+1)}(1 - \beta\delta)$. Then, in a stationary equilibrium a change in the relative cost of production $\omega_{i,j}$ is passed on to the price of a each single variety $q$ with a long run elasticity of*

$$\frac{\eta}{\lambda_v + \eta} \leq \sigma_{p_{i,j}(q),\omega_{i,j}} < 1$$

*Where $\sigma_{p_{i,j}(q),\omega_{i,j}}$ is decreasing in $q$ and $\underset{q\to\infty}{Lim}\left(\sigma_{p_{i,j}(q),\omega_{i,j}}\right) = \frac{\eta}{\lambda_v + \eta}$.*
*In contrast, the measured long run elasticity of the aggregate PPI index is given by*

$$\sigma_{PPI_{i,j},\omega_{i,j}} = 1$$

**Proof.** Consider first the elasticity of the constant and the pricing slope in (2.20)

$$\frac{\omega_{i,j}}{p(q)-\frac{\tau_{i,j}\omega_{i,j}}{A_j}}\frac{\partial\left(p(q)-\frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)}{\partial\omega_{i,j}} = \frac{\eta}{\lambda_v + \eta} \quad \text{and} \quad \frac{\omega_{i,j}}{\frac{\tau_{i,j}\omega_{i,j}}{A_j}}\frac{\partial\frac{\tau_{i,j}\omega_{i,j}}{A_j}}{\partial\omega_{i,j}} = 1$$

That is, if it were not for the constant of integration, up to a first order the price exchange rate pass-through would be the same in the long and short run. The reason this is the case is that the exit of low quality firms affects the location but not the shape of the pricing schedule (2.12). For each variety, the pass through elasticity is equal to

$$\sigma_{p_{i,j}(q),\omega_{i,j}} = \frac{\eta}{\lambda_v + \eta} + \frac{\lambda_v}{\lambda_v + \eta}\frac{\frac{\tau_{i,j}\omega_{i,j}}{A_j}}{p_{i,j}(q)}$$

Since the price is always larger than $\frac{\tau_{i,j}\omega_{i,j}}{A_j}$ and increasing in $q$, $\frac{\eta}{\lambda_v+\eta} \leq \sigma_{p_{i,j}(q),\omega_{i,j}} < 1$ for all $q$.

Also $\sigma_{p_{i,j}(q),\omega_{i,j}}$ converges to $\eta/(\lambda_v+\eta)$ for very large $q$.

Next consider the measured PPI and its elasticity to the exchange rate

$$PPI_{i,j} = \widetilde{\Lambda}\left(\frac{N_j}{L_i}\right)^{-\frac{1}{\lambda_v+\eta}}\left(\frac{\tau_{i,j}\omega_{i,j}}{A_j}\right)^{\frac{\eta}{\lambda_v+\eta}}\left(q_{i,j}^{\min}\right)^{\frac{\lambda_q+\lambda_v}{\lambda_v+\eta}\eta} + \frac{\tau_{i,j}\omega_{i,j}}{A_j}$$

[8] The $PPI_{i,j}$ is a correct measure of price changes in the short run, when the entry is fixed, but incorporating the endogenous entry decision, we find that in steady state,

$$PPI_{i,j} = \left(\frac{\widetilde{\Lambda}}{\Lambda}\left(\frac{N_j}{L_i}\right)^{-\frac{1}{\lambda_v+\eta}} + 1\right)\frac{\tau_{i,j}\omega_{i,j}}{A_j}$$

The elasticity of the $PPI_{i,j}$ with respect to changes in the relative cost of production is equal to 1. is aggregated over a different set of goods, to overstate exchange rate pass through. ∎

There are two channels through which changes in the relative wage rate are transmitted Consider first the multiplicative term $(\tau_{i,j}\omega_j)^{\frac{\eta}{\lambda_v+\eta}}$. When a given firm that is active in the certain market is faced with an increase in the price, it raises its production with an elasticity of $\eta$. Because all firms do the same, the relative increase in supply shifts the matching between consumers and goods towards a lower price. This effect uniformly shift the prices upwards, and it does so with an elasticity of $\frac{\eta}{\lambda_v+\eta} < 1$. A second effect changes the range of firms actually being active in that market, which also affects the general price schedule. This effect is subtle and works through the outside option consumers have at hand. The crucial difference between being the lowest firm actually producing and being a firm that produces is that while the lowest firm producing can cream of all of the consumer surplus, a firm that sits in a continuum of competition can only extract that part of consumer surplus that it is better than the competition. In our setup, the latter relation implies a price slope of $p'_{i,j}(q) = a^{-1}v_{i,j}(q)$. In contrast, if $q_{i,j}^{\min}$, the lowest firm that accesses the market increases from $q_{i,j}^{\min}$ to $\widetilde{q_{i,j}^{\min}}$, not only will the firm command a higher price because it has a lower valuation, but it will also be matched with higher valuation customer, implying $\frac{\partial p_{i,j}\left(q_{i,j}^{\min}\right)}{\partial q_{i,j}^{\min}} = a^{-1}v_{i,j}\left(q_{i,j}^{\min}\right) + a^{-1}\frac{\partial v_{i,j}\left(q_{i,j}^{\min}\right)}{\partial q_{i,j}^{\min}}q_{i,j}^{\min}$ which

---

[8] $\widetilde{\Lambda} = \frac{\frac{\lambda_q+\lambda_v}{\lambda_v+\frac{\theta}{1-\theta}}\frac{\theta}{1-\theta}-\lambda_v}{\frac{\lambda_q+\lambda_v}{\lambda_v+\frac{\theta}{1-\theta}}\frac{1}{1-\theta}-\lambda_v}\Lambda$

is larger than $p'_{i,j}(q)$. This effect shifts up the whole pricing slope, i.e. is reflected in the constant of integration in (2.20).

## 2.6 Conclusion

The contribution of this paper is to explain how – in the presence of complete markets and perfect competition – cost changes brought about by movements of the real exchange rate are transmitted internationally in a competitive market setup and what factors explain the timing of pass-through. In addition, we point out that exchange rate pass-through measured from aggregate data may often be overstated, because the aggregation is over a different set of firms.

We first develop a perfectly competitive economy featuring heterogeneity of both good qualities and of consumer valuations. In equilibrium, high valuation customers and high quality firms are matched, and the relative scarcity of goods of different qualities leads to pricing-to-market and markups that are determined by the local tightness of competition. We analyze how of changes in the relative cost of production affect prices and markups in the short and in the long run. In the short run, the set of firms active in the export sector is fixed, but each firm accommodates changes in the relative cost brought about by a change in the exchange rate by adjusting the quantity of its exports. Since the latter decrease when the home currency appreciates, export markets get relatively less crowded and thus prices measured in the foreign currency increase, leading to partial exchange rate pass-through in the short run. In the long run, the range of firms that are actively exporting changes, because in the presence of fixed costs of market access some low quality firms may no longer find it profitable to export at all. While the short run change in the intensive margin (volume of exports per firm) affects all firms equally, this change in the extensive margin affects only low quality firms, with two associated consequences. Firstly, after an appreciation fewer firms are active in the export sector and these markets are thus characterized by relatively lower demand and consequently upward pressure on prices. Due to this further adjustment of relative demand to changes in the cost of production, long run pass through – although smaller than 100% – is larger than in the short run.

A further consequence of the long run change in the set of firms being active exporters is that the average composition of firms changes, leading to even larger measures pass through

when evaluating aggregate data. An appreciation of the home currency drives out low quality firms that receive a low price for their goods. The observed aggregate price is hence averaged over a set of higher priced firms, leading to an overestimation of long term pass through when using aggregate data. Incorporating this finding, we show that a researcher estimating pass-through in the long run might actually arrive at the conclusion that long run pass through is equal to 100%. These results differ drastically from the existing literature since we firm heterogeneity into our economy through good quality and not through productivity. In our model, firms producing higher priced goods are more profitable, because only high quality goods can command high prices. Our setup hence predicts that following an increase in costs it is low quality and thus low priced firms that leave a given markets. In the long run, the composition effect hence tends to counteract the initial exchange rate movement, which is the opposite what models with heterogeneity in productivity would predict.

## 2.7 Appendix A Endogenizing the outside option of customers

In this appendix we propose a slight modification of our model of preferences such that the model always has a closed form solution, also outside the steady state. While we do not think that the assumption that $\frac{Z_{i,j}}{D_{i,j}} = (1 - \beta\delta)\theta^\eta (1 - \theta) \left(\frac{\lambda_q + \lambda_v}{\eta - \lambda_q}\right)^{\frac{1}{1-\theta}}$ captures any generality, it is also not a very restrictive assumption since it does not depend on any market conditions but merely restricts the constant of integration in (2.20) to be equal to 1. However we here offer an alternative specification no preferences that is not subject to a poential critique that the results hinge on this assumption made in the main text.

The slight modification we propose to our setup is that instead of having a fixed valuation draw in a certain industry, customers can pay a fixed cost $T_{i,j}$ that will give them a new valuation draw. We interpret this set up as valuations not being a fixed characteristic of a consumers, but as the fluctuating occurrence of consumer needs at various moments during their lifetime. The minimum valuation draw that actually buys in our setup is indifferent between buying now and paying the fixed cost $T_{i,j}$ for a new draw. The minimum price is hence not determined by competition of low valuation customers for the lowest quality good actually sold, but by the average consumer surplus that a consumer expects when she decides to get a new draw.

Denote the expected net value from having a valuation draw in market $i,j$ by $V_{i,j}(v)$. All customers above $v \geq v_{i,j}^{\min}$ choose to buy the good, but lower valuation customers choose to obtain an new valuation draw, i.e. the equilibrium cutoff is determined on the one side by

$$V_{i,j}\left(v_{i,j}^{\min}\right) = v_{i,j}^{\min} q_{i,j}\left(v_{i,j}^{\min}\right) - ap_{i,j}\left(q_{i,j}\left(v_{i,j}^{\min}\right)\right) \tag{2.27}$$

But the cutoff customer is indifferent between getting a new draw or not and thus buying not, or paying the cost $T_{i,j}$ to obtain a new valuation and thus it is also true that

$$V_{i,j}\left(v_{i,j}^{\min}\right) = V_{i,j}\left(v_{i,j}^{\min}\right) P\left(v < v_{i,j}^{\min}\right) + E\left[\int_{v_{i,j}^{\min}} f_v(v)\left(v_{i,j}q_{i,j}(v) - ap_{i,j}\left(q_{i,j}(v)\right)\right) dv\right] - T_{i,j}$$
$$\tag{2.28}$$

Since $P\left(v < v_{i,j}^{\min}\right) = 1 - \left(\frac{v_{i,j}^{\min}}{\overline{v}}\right)^{-\lambda_v}$

$$\left(\frac{v_{i,j}^{\min}}{\overline{v}}\right)^{-\lambda_v} V_{i,j}\left(v_{i,j}^{\min}\right) = E\left[\int_{v_{i,j}^{\min}} f_v\left(v\right)\left(v_{i,j}q_{i,j}\left(v\right) - ap_{i,j}\left(q_{i,j}\left(v\right)\right)\right)dv\right] - T_{i,j}$$

The consumer knows that they will be allocated to the following rule

$$q_{i,j}\left(v, v_{i,j}^{\min}, q_{i,j}^{\min}\right) = \begin{cases} q_{i,j}^{\min} \ if \ v < v_{i,j}^{\min} \\ q_{i,j}\left(v\right) = \widetilde{\Lambda}\left(\frac{N_j}{L_{i,j}}\right)^{\frac{1}{\lambda_v}\frac{\lambda_v+\eta}{\lambda_q-\eta}}\left(\tau_{i,j}\omega_j\right)^{-\frac{\eta}{\lambda_q-\eta}}v^{\frac{\lambda_v+\eta}{\lambda_q-\eta}} \ if \ v \geq v_{i,j}^{\min} \end{cases}$$

[9] We will now solve for entry. Assume that again prices are given and so are valuations (that is the consumer who enters knows what to expect and the entry decision of other consumers. A given consumer takes as given the entry decision of other consumers and hence $v_{i,j}^{\min}$ and $q_{i,j}^{\min}$. Let the decision variable of the consumer be $\widetilde{v_{i,j,l}}$ . he consumer will always be matched according to the matching function (2.22) if he follows on a cutoff rule such that $\widetilde{v_{i,j,l}} \geq v_{i,j}^{\min}$. Potentially, each consumer $l$ could choose a cutoff variable $\widetilde{v_{i,j,l}} < v_{i,j}^{\min}$ and in some instances, he will then be matched with the lowest quality firm. The net payoff conditional on a certain cutoff level is hence kinked at $\widetilde{v_{i,j,l}} = v_{i,j}^{\min}$ because form there on, the consumer will also be matched with higher firms. In equilibrium, there is thus a strong economic force that leads to $\widetilde{v_{i,j,l}} = v_{i,j}^{\min}$ for all consumers.

Solving for the Left hand side in equilibrium

$$V_{i,j}\left(v_{i,j}^{\min}\right) = \overline{v}^{\lambda_v}\widetilde{\Lambda}^{-\lambda_v}A\left(\widetilde{\Lambda}\left(\frac{N_j}{L_{i,j}}\right)^{-\frac{1}{\lambda_v}}\left(\tau_{i,j}\omega_j\right)^{\frac{\eta}{\lambda_v+\eta}}\left(q_{i,j}^{\min}\right)^{\frac{\lambda_q-\eta}{\lambda_v+\eta}}\right)^{1-\lambda_v}$$

Solving for the equilibrium value determines (2.27)

$$\int_{v_{i,j}^{\min}} f_v\left(v\right)v^{\frac{\lambda_q+\lambda_v}{\lambda_q-\eta}}dv = \frac{\lambda_v\left(\lambda_q-\eta\right)}{\left(\lambda_q+\lambda_v\right)-\lambda_v\left(\lambda_q-\eta\right)}\overline{v}^{\lambda_v}\left(v_{i,j}^{\min}\right)^{\frac{\lambda_q+\lambda_v}{\lambda_q-\eta}-\lambda_v}$$

---

[9] $\widetilde{\Lambda} = \left(a\frac{\lambda_v+\frac{\theta}{1-\theta}}{\lambda_q+\lambda_v}\right)^{-\frac{\lambda_v+\frac{\theta}{1-\theta}}{\lambda_q-\frac{\theta}{1-\theta}}}\Lambda^{-\frac{\lambda_v+\frac{\theta}{1-\theta}}{\lambda_q-\frac{\theta}{1-\theta}}}$

Market clearing implies $v_{i,j}^{\min} = a \frac{\lambda_v + \eta}{\lambda_q + \lambda_v} \Lambda \left( \frac{N_j}{L_{i,j}} \right)^{-\frac{1}{\lambda_v}} (\tau_{i,j} \omega_j)^{\frac{\eta}{\lambda_v + \eta}} \left( q_{i,j}^{\min} \right)^{\frac{\lambda_q - \eta}{\lambda_v + \eta}}$ . Inserting this impression in the minimum price leads to the a price that is always as in (2.12), even outside the steady state. However, since all consumer buy in equilibrium total market demand is fixed and equal to $L_i$ and some terms such as $\left( \frac{N_j}{L_{i,j}} \right)$ cancel in equilibrium.

# Bibliography

[1] Atkeson, A. and Burstein, A. (2006) "Trade Costs, Pricing to Market, and International Relative Price," MIMEO, Department of Economics, UCLA, 2006.

[2] Baldwin, R. (1988) "Hysteresis in import prices: The beachhead effect" The American Economic Review Vol. 78, No. 4, Pages 773-85.

[3] Baldwin, Richard and Krugman, P. (1989) "Persistent trade effects of large exchange rate shocks" The Quarterly Journal of Economics, Vol. 104, pages 635-654, 1989.

[4] Campa, J. M.and Golberg, L. (2005) "Exchange Rate Pass Through into Import Prices" The Review of Economics and Statistics, Vol. 87, No. 4, Pages 679-690, November 2005.

[5] Campa, J. M.and Golberg, L. (2006) "Distribution Margins, Imported Inputs, and the Sensitivity of the CPI to Exchange Rates" National Bureau of Economic Research, Working Paper No. 12121, March 2006.

[6] Dixit, A.V. and Stiglitz, J.E. (1977) "Monopolistic Competition and Optimum Product Diversity" The American Economic Review, Vol. 67, No. 3. , pp. 297-308. June 1977

[7] Dornbusch, R. (1987) "Exchange Rates and Prices" The American Economic Review, Vol. 77, No. 1., pp. 93-106. March 1987.

[8] Gabaix, X. and Landier, A. (2006)"Why Has CEO Pay Increased So Much?" MIMEO, Department of Economics, MIT, 2006.

[9] **Ghironi, F. and Melitz, M. (2005)** "International Trade and Macroeconomic Dynamics with Heterogeneous Firms" The Quarterly Journal of Economics, Vol. 120, No. 6 Pages 1695 1725, August 2005

[10] **Krugman, P. (1980)** "Scale Economies, Product Differentiation, and the Pattern of Trade" The American Economic Review, Vol. 70, No. 5., pp. 950-959. December 1980.

[11] **Krugman, P. (1987)** "Pricing to market when the exchange rate changes." In: S. Arndt – J. Richardson (eds.): Real Financial Linkages Among Open Economies, Cambridge, MA: MIT Press.

[12] **Melitz, M. (2003)** "The Impact of Trade on Intra-Industry Reallocations and Aggregate Industry Productivity" Econometrica, Vol. 71, No. 6 Pages 1695 1725, November 2003.

[13] **Melitz, M. and Ottaviano, G. (2005)** "Market Size, Trade, and Productivity," MIMEO, Department of Economics, Harvard University, October 2005.

# Chapter 3

# Colonial and Geographic Origins of Comparative Development

**Summary 3** *With their seminal work on the effects of settler mortality on colonization policies during early imperialism, Acemoglu et al. (2001) build a strong case for the importance of institutions as the primary force of economic development. However, because their empirical analysis is limited to former colonies, they cannot directly distinguish their theory from the rivaling view that a country's disease environment has direct effects on economic prosperity and institutions. In this paper, using either additional historical sources or a model of the geographic determinants of disease, I first construct two measures of mortality rates including up to 36 countries that have not been colonized. I then show that mortality did affect institutional development in former colonies but not in the rest of the sample. This can only be rationalized in the context of the colonial origins theory of Acemoglu et al. (2001). Turning to disentanlge the relation between institutions and income, I sometimes find that disease environment influences income also directly and correspondingly, that institutions are somewhat less important for prosperity in my specifications than when working with a sample composed of only former colonies. Incorporating these findings, I estimate that institutions are the major determinant of long run prosperity and can explain about 50% of the observed variation of current income levels, while the direct effects of disease environment can account for about 15%.*

## 3.1 Introduction

What are the fundamental causes of long run economic growth? With their seminal article, Daron Acemoglu, Simon Johnson and James Robinson (2001) (AJR from here on) build a strong case for the importance of institutions as the primary force of economic development. The authors argue that differences in the disease environment of former colonies led to very different institutional policies pursued by the colonizing nations. In places unfavorable to European physiology, the main objective of the colonizers was to extract as much valuables as possible by corrupting local institutions. In contrast, when chances of survival where high, European settlers came in large numbers and the focus of the colonizers was to produce rather than to extract leading them to install institutions geared towards facilitating production and ensuring property rights. The authors support their argument with empirical evidence on the effect of European settler mortality rates on the quality of early and current institutions, on today's level of income directly and indirectly with institutional quality instrumented by settler mortality. The effects of early disease environment on economic prosperity are statistically significant, economically large and robust to the inclusion of geographic controls and further instruments for institutional quality.

Throughout their analysis, the authors' biggest concern is to ensure that their empirical approach is capturing exclusively institutional channels of comparative development. As they note, their basic empirical approach suffers from the potential problem that "[...] mortality rates of settlers could be correlated with the current disease environment, which may have a direct effect on economic performance" (p.1371). A large and certainly influential literature, see for example Bloom and Sachs (1998), Gallup et al. (1998) and Diamond (1997), has argued that current climate, disease environment and other geographical features have a strong direct impact on economic outcomes. Mortality in the 18th century could hence be correlated with current prosperity because it is a measure of current prevalence of disease. In accordance with the theory of Lipset (1960), the relation between institutional quality and mortality could then be an indirect consequence of good economic outcomes. Furthermore, it is not a priory clear that disease environment has no direct influence on institutions through its impact on life expectancy, education decisions and work culture.

To address these critiques, AJR take two separate strategies to test for the presence of a

direct effect of disease environment on economic outcomes. In a first, they add geographic variables such as latitude or prevalence of malaria to their regressions and show that once the influence of settler mortality is taken into account, these geographical variables have little effect on comparative development. Secondly, they include other instruments for institutions, such as institutional quality in 1900, and use overidentification tests to check whether settler mortality has a direct effect on current performance. They conclude that their instrumentation strategy is valid, and that the channel through which early disease environment affected development is indeed indirect via early institutions.

While the authors make the best use of the data at their disposal, their empirical strategy suffers from the basic problem that historical studies such as Curtin (1989) have only collected mortality rates for countries that have been colonized, but not for the rest of the world. As I argue in this study, a rigorous test of their hypothesis is not whether or not early disease environment did influence development in former colonies, but rather whether it influenced prosperity *differentially* in former colonies compared to other nations. This differential effect disease environment had on development is also the central theme of the qualitative argument of Acemoglu et al. (2003). The authors argue that while disease environment may well have also directly affected development, the latter effect is likely to be of a smaller magnitude than the institution building channel. In what follows below, I aim to quantify this assessment.

Settler mortality rates – an aggregate measure of many diseases that are potentially detrimental to growth – are a very good proxy for current geography and hence any relationship between mortality and economic outcomes is also in accordance with the geography theory of development.[1] Because mortality rates may well have large direct effects on economic development, both the institutional and the geographic view of development are observationally equivalent in the empirical setup of AJR. The validity of the institutional hypothesis hence depends crucially on the validity of AJR's additional instruments for institutions, which are measures of early institutional quality (constraints on the executive in 1900, constraints on the executive at the time of independence and democracy in 1900) and the size of European settlements in 1900. Arguably, the first set of instruments is substantially less endogenous

---

[1]For example, the mortality rate in AJR's sample of 64 countries is correlated with current life expectancy with a coefficient of -.75, which is substantially stronger than the correlation between current life expectancy and other variables that are often used as proxy for geography, for example latitude.

than current institutions and levels of income, yet it is not unreasonable to argue that early institutions are directly related to current economic prosperity.[2] Furthermore, using the size of European settlements in 1900 as an instrument can be questioned on the same grounds as instrumenting with settler mortality itself: maybe European settlers came only in big number to geographically favorable areas and these areas continue to prosper today.

Summarizing, while I think that AJR provide a substantial amount of evidence for their proposed theory, it is worth the while to examine their theory further by testing its sharpest prediction: *if the institution-building hypothesis is valid, early disease environment should have had an additional effect on development in former colonies.* If this is the case, such a difference of how geography affected comparative development can only be explained in the context of institutions set up during colonization. Furthermore – and the two views are not exclusive in the empirical setup presented below – if early disease environment has an effect that is common to all nations, this reflects the direct impact of geography on prosperity.

The key insight of this paper is that also AJR's hypothesis is ultimately an argument relating geography to economic outcomes, but through an indirect channel that only applies to a subset of countries. Among others, this has been noted by Easterly and Levine (2003).

"The geography hypothesis stresses that the disease environment directly influences productivity. Settler mortality measures the disease environment as European settlers arrived and thereby provides an exogenous indicator of "germs." The AJR (2001) institution hypothesis, instead, stresses that (i) initial endowments shape the long-lasting institutions created by European conquerors and (ii) these long-lasting institutions continue to shape economic development today."

Easterly and Levine (2003), p. 12

In the language of Easterly and Levine, geography determined the prevalence of "germs," which influenced European settler mortality, thereby affected early institutions and current economic outcomes. A test of the geography view of development is whether any geographic variable – including AJR's settler mortality – is a significant determinant development in all

---

[2]In fact, this is indeed the opinion of AJR, see their footnote 31.

countries equally. In contrast, the crucial test of the institutional view of development is whether geography had an additional effect on development of institutions in former colonies.

The mental exercise of this paper is closely related to McArthur and Sachs (2001), who argue that "if a disease environment associated with high mortality rates in the early $19^{th}$ century let European powers to develop predatory political institutions rather than developmental [...] institutions, it seems far fetched to argue that the disease burden itself played no adverse role in the process of economic development" (p.6). McArthur and Sachs (2001) and Sachs (2003) thus add more countries[3] and additional geographic information to the empirical setup of AJR to show that also variables other than mortality matter for comparative development. In contrast, the current study argues that to distinguish between these two theories of development, it is only of indicative use to show that certain variables are significant determinants of prosperity while others are not. Both views predict a close relationship between geography and economic outcomes. The key difference is how geography shaped prosperity differentially in each group of countries. This empirical approach is related to the work of Rigobon and Rodrik (2005) who – among other things – identify the relationship between institutions and prosperity by the heteroscedasticity between the group of former colonies and other nations. Rather than relying on the difference of the variability of economic outcomes between those two groups, I exploit the differential relation between mortality and institutional outcomes to establish a causal link between institutions and income.

In order to do this, I first need to construct additional mortality data for nations that have not been colonized. At first sight, enlarging the sample of AJR to also include a substantial number of such countries seems impossible because countries like Japan have never been settled by Europeans and there is thus no historical estimate of "settler mortality." It is, however, important to point out that also AJR do not measure the mortality of the average settler, but that of soldiers. This data is also available for countries that have never been colonized and I am hence able to construct a historical measure of mortality for 28 nations from a variety of military records and collections of vital statistics in the respective times. I merge this data with

---

[3] Interestingly, McArthur and Sachs (2001) add France and Britain to the setup of AJR. The addition of these two non-colonies substantially weakens the relation between mortality and institutions, wich the authors interpret as a failure of the institutional hypothesis of development. On the contrary, the present study argues that adding independent nations as is done by Mc Arthur and Sachs *should* weaken the overall significance if AJR's hypothesis is true.

AJR's original data and thus obtain a first measure of mortality for 91 countries. To enlarge the sample size further and address any potential bias in my additional historic mortality data, I also construct a second data series by exploiting the close relationship between geography and disease. In a first step, I use AJR's historical mortality data and a range of geographic information to estimate the relation between early disease environment and geographic variables such as latitude, temperature or elevation that are measured precisely and are available for a large set of countries. Using this model of the geographic determinants of mortality, I am then able to construct a measure of early disease environment for a large sample of 150 nations, including 36 countries that have never been colonized. The variable I construct does not measure the actually observed mortality rates, but rather the potential mortality of settlers had a country been colonized.[4]

For each of the two samples constructed, I then test for the importance of direct and indirect effects of early disease environment on economic development. I repeat the first stage specifications of AJR relating mortality to institutional development. However, instead of testing only for the significance of mortality by itself, I test for both the significance of a main effect of mortality (i.e. for geography) and for the significance of mortality interacted with a colony dummy (i.e. for early institution building). As I show below, only this way of estimation can causally disentangle the different views of development. In the first stage specifications, I find that while the main effect is not a significant determinant of institutions, the interaction of mortality and a colony dummy is strongly negatively correlated with institutional outcomes. Hence, the results support the institution building hypothesis of AJR with no evidence for any common effect of geography on institutions.

Turning to disentangle the relation between income and institutions, I exploit the differential effect disease environment had on institutional development to identify the relation between institutions and levels of income. Because the identification restriction – that any additional effect mortality has on prosperity works through the institutional channel – does not require the main effect of mortality to be equal to 0 in the second stage estimation, I am able to test whether geography did directly influence development. Indeed, I find (although not significant

---

[4]The use of this hypothetical measure of settler mortality is the appropriate to test of the validity of AJR's theory since they argue that knowledge of high prevalence of disease alone was enough to deter migration to certain countries.

in all specifications) a negative direct effect of early disease environment on current levels of prosperity, hence also supporting the geography view of development. Correspondingly, because I attribute some of the correlation between mortality and current levels of income to direct effects of geography rather than the institutional channel, I find a smaller overall importance of institutions. For example, in a typical specification,[5] I find that for a former colony, a one percent difference in disease environment is associated with a 1.06 percentage point difference in income per capita. In an estimation similar to AJR that uses only colonies, in order to identify the system, all of this difference is attributed to institutions. In contrast, when I estimate the relation in my larger sample where I can also allow for a direct effect of germs on current level of income levels, I estimate that the same one percent reduction of early disease environment leads to a 0.79 percentage increase in income through the associated improvement of early institutions and a 0.27 percentage point direct effect of disease environment on income levels. That is, although my coefficients are somewhat smaller than in an estimation using only former colonies, I still find that the institutional channel is the major determinant of prosperity while geography has smaller but still sizeable direct effects on development. This finding is shown to be robust to alternative ways of constructing the measure of early disease environment, changes in the sample, addition of geographic controls and over-identification tests.

What do these results imply for the overall importance of institutions and geography as the long run determinants of prosperity? Again, the findings of this paper point towards the importance of both factors, where institutions are the major but not exclusive ingredient for development. In my main specifications, I estimate that if one where to improve the current disease environment of a country by 1/100 of a standard deviation, this would ultimately result in an income increase of about 0.33%. In contrast, if one where to increase the quality of institutional outcomes by 1/100 of a standard deviation, the associated change of income would be 1.54%.

---

[5]The examples in the following two paragraphs refer to specifications from Tables 3 and 5, Columns 5 and 7.

The outline of the paper is the following. In the next section, I construct my two measures of mortality. In Section III, I set up the econometric framework and develop the main test of this paper. Section IV presents the first stage results relating institutions to settler mortality. Section V presents the instrumental variable estimation disentangling the relation between institutions and income. I check for the robustness of my findings in section VI and section VII concludes.

## 3.2    Data Description

In this section, I construct two mortality data samples. In the first, I rely on AJR's original data and – using additional historical sources – augment their dataset by a number of countries that have not been colonized. Arguably, my historical sources are less precise than the values reported in Curtin (1998), which is AJR's main source. I thus also construct an alternative sample using AJR's mortality rates and a range of geographic data to estimate a model of the geographic determinants of settler mortality during the period of imperialism. Since geographical information is available for almost all countries of the world, I am able to use the close relation between temperature, humidity and elevation on the one and disease on the other side to predict a variable of early disease environment for a sample of 114 former colonies and 36 countries that have not been colonized.

### 3.2.1    Historical Sources

In their unpublished appendices, AJR report mortality rates for a few additional nations. They have direct data for China, France and Britain from Curtin (1989) and construct a mortality estimate for Afghanistan, Thailand and Korea. AJR include Ethiopia in their sample of former colonies. In addition to the short period of occupation, Italy invaded Ethiopia in 1936, at a time where medicine had advanced leading to low mortality rates for western soldiers. Hence, even if the period of occupation did have some institution-building effect, it should have been of very limited scope and I therefore classify Ethiopia as an independent country. On the other side, while some researchers would regard Liberia as an independent nation I keep it in the group of former colonies. The reason for doing so is that the country was subject to heavy

slave trading throughout almost all early modern history with the same associated patterns of early institution building as was the case in former colonies. I make sure that my results are not driven by these changes and present alternative classifications in Appendix B.[6]

For these 7 countries, I end up with exactly the same estimate of European settler mortality AJR use, namely the annualized probability of death for a Caucasian soldier stationed in the respective country. This data is not available for other countries and I hence rely on estimates of the mortality of a male person in the age cohort of 20 to 40 years of age. For several countries, I have a direct estimate of this rate. Bengtsson et al. (2004) report mortality data by sex and age cohort for Italy, Japan, Sweden and Belgium, which I also use for Luxembourg. Steckel and Floud (1997) report such a value for Germany (Baden Wuertemberg). For other countries, I rely on data by Mitchell (1998), who reports the size of different age cohorts by sex at various points in time. From his data, I am able to construct additional mortality rates for Austria, Denmark, Czech and Slovak Republik, Finland, Greece, Hungary, the Netherlands, Norway, Poland, Portugal, Romania, Spain and Switzerland. I am careful to check this data for changes in size of the respective nations and migration and do not include the available data for Ireland because of large scale emigration.

I my historical sample, I include in total 91 countries, of which 63 are from AJR's original data.[7] Because I count Ethiopia as an independent nation, I end up with a sample of 62 former colonies and 29 independent nations, of which 5 (Japan, Thailand, Ethiopia, South Korea and China) lye outside Europe. In Appendix B, I show in more detail the construction of my historical mortality series.

### 3.2.2 Early Disease Environment and Geography

Curtin (1989) carefully recorded the causes of death for French and British soldiers listing a wide variety of diseases. The resulting mortality rate is an aggregate of a country's disease environment and hence a very good measure of what scholars of the geography school of development would use as a direct instrument for economic outcomes. In this section, I exploit

---

[6] Ethiopia and Liberia are poor and charaterized by high prevalence of disease. The proposed clasifications hence weaken the evidence for the "colonial origins" hypothesis in the case of Ethiopia, but strengthen it in the case of Liberia. However all the effects are mild, see Appendix B.

[7] Of AJR's 64 countries, Malta is missing some geographic data so that I have dropped it from the dataset.

the close relation between geography and germs to construct a measure of early disease environment for a large sample of countries, also including some that have never been settled by Europeans. It is important to note that the use of this hypothetical measure of settler mortality is in accordance with the institution building hypothesis of AJR. The authors provide evidence that widespread knowledge of disease prevalence – often posted in newspaper articles – alone was enough to deter migration to these places. As Acemoglu et al. (2005) argue, "by its nature, *potential* settler mortality is often not observed. In places where the potential settler mortality was high, large numbers of settlers did not go, and it is difficult to obtain comparable measures of this mortality" (original emphasis; p. 6). In the original article, the authors tackle this problem by using mortality estimates collected mostly from soldiers. Soldiers are a more homogenous group than the general settler population and also suffer less from the endogenous selection because they were stationed by their imperialist powers with little concern for chances of survival. In a first step, I thus use AJR's mortality rates to estimate the relation between geography and early disease environment.

The prevalence of many diseases is strongly influenced by the geographic characteristics of a country, which leads AJR to extrapolate certain mortality rates to countries with similar geographic features. Rather than arguing that certain countries are sufficiently similar so that one can extrapolate the mortality rate, I directly estimate a model of geography and disease environment for the countries in AJR's original dataset. To formalize this relationship, let $M_i$ denote the mortality of Caucasian soldiers in country $i$ and let $X_i$ denote the vector of geographic variables that affect development through any channel. Regardless of whether a country actually has been colonized or not, the following relationship between geography $X_i$ and disease environment holds for all countries $i$.

$$M_i = X_i'\mu_s + \nu_{M,i} \tag{3.1}$$

My measures of geography, such as average temperature, humidity and dummies for the incidence of deserts or mountains and are collected by Parker (1997). I present the estimation of my main model of the geographic determinants of disease (3.1) in Table 1.

|  | Dependent is Ln Mortality |
| --- | --- |
| Avg. Temperature | 0.096 (0.043)* |
| Temp. / max Humidity | -0.076 (0.023)** |
| Min. Rain | -0.01428 (0.00276)** |
| Max Rain | 0.00123 (0.0007) |
| Mediterranean y/n | -1.392 (0.292)** |
| Savanna y/n | 0.775 (0.192)** |
| Temperate forest y/n | -1.348 (0.264)** |
| Mountains y/n | -0.648 (0.298)* |
| Observations | 63 |
| R-squared | 0.70 |

Robust standard errors in parentheses
* significant at 5% level; ** significant at 1% level

Table 1: Geography and Mortality

In Table 1, I analyze the impact of temperature, humidity, seasonal variation of rainfall, natural vegetation and elevation, and refer the reader interested in how additional factors influence mortality to Appendix A. Higher average temperature is associated with high levels of disease. The dependent variable in Table 1 is the logarithm of mortality, so that a country with a 1 Degree Celsius warmer climate is characterized by a 10% increase in mortality. In contrast, the coefficient of high temperature at maximum humidity is negative, implying that cold and moist climate is associated with higher levels of mortality. I next evaluate the impact of rainfall and its seasonal variation. Areas with pronounced dry (low minimum monthly rain) or wet seasons (high maximum monthly rain) are characterized by a high mortality rate, although the coefficient of maximum rain is not significant. It is also noteworthy that the first coefficient is of a larger magnitude than the second, implying that a country with uniformly more rain throughout the year is characterized by lower prevalence of germs. Mirroring this finding of low

variation in climate being associated with healthier living conditions, Mediterranean climate and the natural incidence of temperate forests are associated with lower prevalence of disease. Also elevated areas are less prone to disease, and a mountain dummy is significantly negative. In contrast, countries with stretches of Savanna are characterized by higher mortality.



Figure 1: Settler Mortality and Early Disease Environment

All coefficients in Table 1 are significant at the 1 percent level, except maximum rain, which is (barely) not significant and average temperature, which is significant at the 5 percent level. The overall fit of the model is very good with an $R^2$ of 70%, and a F-score in excess of 19. For further examination of the constructed geographic model of disease, Figure 1 displays the in-sample relation between the fitted measure of early disease environment (Y-Axis) from the estimation of (3.1) in Table 1 and AJR's original's data (X-Axis), displaying again an overall

good fit and no major outliers.[8]

From these findings, I conclude that geography is indeed very closely connected to the prevalence of disease in the $18^{th}$ century. Thus, using the relationship between geography and settler mortality in Table 1, I predict my geographic measure of mortality for both the original countries of AJR plus 87 additional countries, of which 36 have not been colonized. In the analysis below, I shall refer to this measure as "early disease environment." Paralleling the definition of "settler mortality" in AJR, the term "early disease environment" refers the logarithm of the *potential* annualized probability of death for Caucasian males in the age cohort of 20 to 40 years of age.

**Summary Statistics**

| Variable | Mean | Std. Dev. | Min | Max | Observations |
|---|---|---|---|---|---|
| Ln Settler Mortality, by AJR | 4.676203 | 1.240061 | 2.145931 | 7.986165 | 63 |
| Ln Settler Mortality, Collected | 3.968353 | 1.54842 | 1.196948 | 7.986165 | 91 |
| Disease Environment | 4.377385 | 1.187361 | 1.115383 | 7.25028 | 150 |
| Ln GDP per Capita 2000 | 7.608087 | 1.687515 | 4.484354 | 11.15503 | 150 |
| Rule of Law 1996-2004 | 0.011962 | 1.022666 | -1.911856 | 2.137062 | 150 |
| Control of Corruption 1996-2004 | 0.0388283 | 1.026773 | -1.590134 | 2.464298 | 150 |

**Pairwise Correlation Diagram**

| | AJR Mortality | Collected Mortality | Disease Environment | Ln GDP | Rule of Law | Regulatory Quality |
|---|---|---|---|---|---|---|
| Ln Settler Mortality AJR | 1 | | | | | |
| Ln Settler Mortality Colllected | 1 | 1 | | | | |
| Disease Environment | 0.8373 | 0.8168 | 1 | | | |
| Ln GDP per Capita 2000 | -0.6683 | -0.6973 | -0.6197 | 1 | | |
| Rule of Law 1996-2004 | -0.6523 | -0.6674 | -0.5278 | 0.8429 | 1 | |
| Regulatory Quality 1996-2004 | -0.574 | -0.5963 | -0.4651 | 0.7408 | 0.9019 | 1 |

Table 2: Data Summary and Pairwise Correlation Diagram

The resulting measure of disease environment for 150 countries is listed in Table 9, where I also list my other data in detail. To make sure that the way of predicting early disease environment is robust, I present a variety of alternative specifications in Appendix *A* and show that the results obtained with my main measure of disease environment are robust to these alterations. In Appendix *A* I also address a recent critique of AJR's data by Albouy (2005) and show that my way of using geographic information to instrument for mortality is robust to using his

---

[8]The fitted value for Guinea seems to deviate somewhat from its actual mortality rate. I demonstrate in Appendix B that in - or exclusion of Guinea does not change the results of this study.

alternative data series. Interestinly, while his and AJR's data differ substantially, the mortality rates are estimated by geography in a very similar way and the resulting measures of early disease environment are thus very similar.

In Table 2, I present summary statistics and the pairwise correlation diagram for my two measures of mortality, AJR's original mortality estimate and current economic outcomes. Because AJR's primary measure of institutions, Protection Against Expropriation from the International Risk Country Guide is not available for the large set of countries used in this study, I use the 1996 to 2004 averages of Rule of Law and alternatively of Regulatory Quality from Kauffman et al. (2005) as proxies of institutional quality. Both variables are measured on a continuos scale ranging from −2.5 to +2.5, with higher values associated with better outcomes.[9] The number of observations is 91 in the specifications using historical data and 150 when using my measure of early disease environment. The two measures of mortality correlate highly and significant far beyond the 1 percent level with each other, AJR's original estimate and also with current levels of economic prosperity measured by GDP or institutional outcomes.

Using the data constructed so far, I next develop a test to disentangle AJR's hypothesis from alternative theories.

## 3.3   Determinants of Economic Prosperity - Theory

The key insight of this paper is that also AJR's hypothesis is ultimately an argument relating geography to economic outcomes, but through an indirect channel that only applies to a subset of countries. If the hypothesis of AJR is true, geography determined the prevalence of disease, which influenced European settler mortality and thereby affected early institutions and current economic outcomes. In contrast, the geographic view predicts that geography determined the disease environment of a country, which had direct consequences on culture, technology, productivity investments in human capital and thus on income and on institutions. It is important to point out that the geographic view predicts both a possible direct effect of disease environment on income *and* on institutions. Loosely speaking, the direct effect of disease environment

---

[9]When a country was missing an institutional score in a certain year, I averaged over the other years. Kauffman et al (2005) standardize all variables, but since I only use 150 of the available 204 countries and average over time, the mean in my sample differs slightly from 0 and my standard deviation is not exactly equal to 1.

on institutions is very close to the "cultural" view of development. Schematically, geography might have affected development through two channels.

$$\text{Geography} \Rightarrow \text{Disease Environment} \stackrel{all\ nations}{\Rightarrow} \begin{cases} \text{Technology} \Rightarrow \text{Income} \Rightarrow \text{Institutions} \\ \text{"Culture"} \Rightarrow \text{Institutions} \Rightarrow \text{Income} \end{cases}$$

This contrasts to the institutional view of AJR predicting a relation between geography and prosperity only in former countries.

$$\text{Geography} \Rightarrow \text{Disease Environment} \stackrel{only\ colonies}{\Rightarrow} \text{European Settlements} \Rightarrow \text{Institutions} \Rightarrow \text{Income}$$

I.e. in total, disease environment potentially has three distinct effects on development: a common effect on institutional quality, a common effect on income directly and an additional effect on institutional development in former colonies. Only with the large sample of countries constructed in the previous section, can this differential effect be tested for.

In Figures 2A and 2B, I display the relationship between mortality and institutional outcomes, measured by the Rule of Law from Kauffman et al. (2005), where a higher value corresponds to better enforcement of laws. In the lower figure (2A), I display the relationship between the geographic measure of settler mortality and institutions for former colonies, reproducing AJR's basic finding of a strong negative relation between the two variables in that group of countries. In contrast, In the upper figure (2B), I display the same relation, but for the group of independent nations. Graphical inspection suggests that there is definitely no negative relation between the two variables. I reproduce the same finding for my historic measure of mortality in Figure 3A and 3B, which is presented at the end of the paper.

I next establish this observation formally and develop the respective tests that enable to clearly distinguish between the geographic and institutional channel of development.
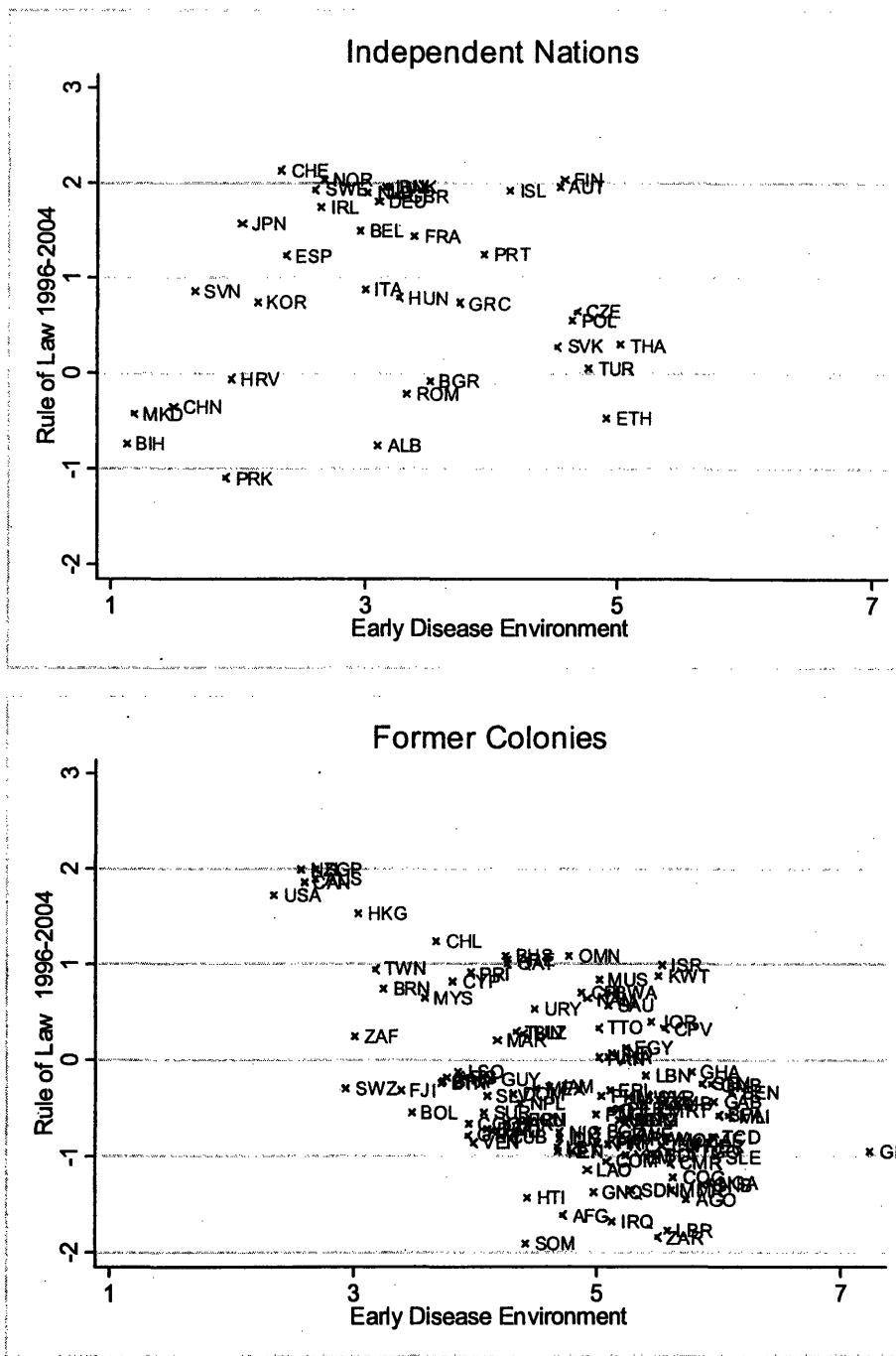
Figure 2: Disease Environment and Institutions

Throughout the analysis below, let $Y_i$ denote the logarithm of GDP per capita and denote the measure of institutional quality in country $i$ by $R_i$. In the empirical implementation, $R_i$ is measured by either Rule of Law or by Regulatory Quality. Because the course of history has shaped economic development differentially, in some instances, I set up models for either only colonies or only for independent nations. To avoid confusion, I denote any model that is valid only in former colonies by a $C$ superscript and any model that is valid only in non-colonies by a $N$ superscript. Models that are valid for all countries do not carry a superscript. Consider first the determinants of current levels of income, which is determined by the quality of a countries' institutions as well as the countries' disease environment $M_i$.[10]

$$Y_i = \widetilde{\lambda}_Y + \widetilde{\alpha}R_i + \widetilde{\vartheta}_Y M_i + \widetilde{\varepsilon}_{y,i} \tag{3.2}$$

While there is a common model of the determinants of current levels of income, the determinants of institutional outcomes are different reflecting the different historical experiences of former colonies versus other nations. In countries that have not been colonized, the quality of institutions potentially depends on disease environment through channels of "culture." Additionally, in accordance with theories put forward by Lipset (1960) and others, income $Y_i$ itself may influence institutional outcomes.

$$R_i^N = \widetilde{\lambda}_R^N + \widetilde{\vartheta}_R M_i + \widetilde{\beta}Y_i + \widetilde{\nu}_{Ri}^N \tag{3.3}$$

It is clear that the system of equations (3.3) and (3.2) is not identified and can thus not be estimated without additional information.[11] What are the determinants of institutions in former colonies? In addition to any direct effects geography might have had, there is an additional channel that works through the size of settlements of Europeans in Colony $i$, denoted by $S_i^C$. Neglecting the transition from old to current institutions, the model determining today's

---

[10]Numerous alternative variables may also directly affect growth. I neglect these variables in the theoretical analysis below, but demostrate that the results are robust to the inclusion geographic, historic and sociologic controls in Section 5.

[11]Ironically, a convinced scholar of the geographic view of development could assert that disease environment has no direct impact on institutional development (i.e. $\widetilde{\vartheta}_R = 0$) and thus identfy (3.3) and (3.2) with the opposite exclusion restriction as do AJR.

institutional quality in a former colony is hence

$$R_i^C = \widetilde{\lambda}_R^C + \widetilde{\theta}_R S_i^C + \widetilde{\vartheta}_R M_i + \beta Y_i + \widetilde{\nu}_{Ri}^C \tag{3.4}$$

Where the size of European settlements is determined by the expected mortality rate in the respective country.

$$S_i^C = \widetilde{\lambda}_S^C + \widetilde{\theta}_M^C M_i + \widetilde{\nu}_{Si}^C \tag{3.5}$$

The model in (3.4) might suffer from endogeneity, as settlers might choose to move to areas with better institutions and hence only reduced forms of (3.4) and (3.5) will be estimated, with the use of mortality directly instead of European settlements. In such a reduced form estimation in a sample of former colonies, a significant effect of settler mortality on institutional outcomes could thus either originate from the direct effect of mortality that is common to all countries ($\widetilde{\vartheta}_R$) or through the channel AJR hypothesize, i.e. through $\widetilde{\theta}_M^C \widetilde{\vartheta}_R$ being significant. More importantly, even if we have strong a priory reasons to believe that $\widetilde{\vartheta}_R$ is close to 0, the system of equations (3.2), (3.4) and (3.5) is still not identified unless also $\widetilde{\vartheta}_Y$ equals 0.

Consider, in contrast, a model with joint estimation of the determinants of institutions in former colonies (3.4) and in other nations (3.3). Stacking colonies and non-colonies and introducing a dummy $C_i$ that is equal to 1 for former colonies yields the joint model of geography, disease environment and institutions. The reduced form model of the determinants estimation of is

$$R_i = \lambda_R + \lambda_R' C_i + \vartheta_R M_i + \theta_R (M_i C_i) + X_i' \gamma_R + \nu_{Ri} \tag{3.6}$$

Where the following relations between the coefficients in (3.4), (3.5), (3.3) and in (3.6) hold: $\lambda_R = \frac{\widetilde{\lambda}_R^N}{1-\widetilde{\alpha}\widetilde{\beta}}$ , $\lambda_R' = \lambda_R^C - \lambda_R$ and $\vartheta_R = \frac{\widetilde{\vartheta}_R}{1-\widetilde{\alpha}\widetilde{\beta}}$ and $\theta_R = \frac{\widetilde{\theta}_R \widetilde{\theta}_M}{1-\widetilde{\alpha}\widetilde{\beta}}$.[12] The interpretation of the coefficients is the following. $\vartheta_R$ captures the potential direct effects of geography on institutional development, while $\theta_R$ captures the 'institution building' effect of early disease environment, which is exclusively present in former colonies. $\theta_R$ tests the joint hypothesis of AJR that mortality affected European settlement policies and that European settlement

---

[12]Furthermore $\nu_{Ri} = \frac{(1-C_i)\widetilde{\nu}_{Ri}^N + C_i(\widetilde{\nu}_{Ri} + \widetilde{\theta}_R \widetilde{\nu}_{Si})}{(1-\alpha\beta)}$. This formulation makes clear that there may well be heterscedasticity between the two groups of countries, as is assumed by Rigobon and Rodrick (2005). All results presented below are hence estimated with robust standard errors in both first and second stages of the estimation.

116

policies affected early institution building and current institutional outcomes. The model in (3.6) also allows for a different intercept for colonies and non-colonies. The coefficient of the colony dummy has to be interpreted with care, since mortality $M_i$ is measured on a logarithmic scale. $M_i$ can be scaled by taking a logarithm of mortality with a different base, which does not change any coefficient except the colony dummy itself.

More important than offerering a test of colonial origins versus direct effects of geography on institutions, is that the stacked model (3.6) offers a clear identification restriction on the relation between mortality and institutions that relies on the interaction effect of mortality and a colony dummy. With the first stage (3.6), the reduced form second stage is

$$Y_i = \lambda_Y + \lambda'_Y C_i + \alpha R_i + \vartheta_Y M + \varepsilon_{y,i} \qquad (3.7)$$

where $\lambda_Y = \frac{\tilde{\lambda}_Y + \tilde{\alpha}\lambda_R}{1 - \tilde{\alpha}\beta}$, $\lambda'_Y = \frac{\tilde{\alpha}\lambda'_R}{1 - \tilde{\alpha}\beta}$, $\alpha = \frac{\tilde{\alpha}}{1 - \tilde{\alpha}\beta}$. The coefficient $\vartheta_Y = \frac{\tilde{\vartheta}_Y}{1 - \tilde{\alpha}\beta}$ tests for the potential direct effect of disease environment on income differences, i.e. for geography. The first stage includes a main effect of both the colony dummy and of mortality as well as an interaction effect of the two. Since there are only two endogenous variables, the system can thus be identified while still allowing for a direct effect of disease environment (and a colony dummy) in the second stage estimation.

I next turn to estimating the first (3.6) and second stage (3.7) results successively. Again, it is important to point out that the exclusion restriction is different than the one in AJR. Rather than postulating that disease environment exclusively has consequences on early institution building, I require that any additional effect disease environment had on development in former colonies over and above the effect it had in independent nations works exclusively through institutions.

## 3.4 Colonization, Disease and Institutions

| Dependent<br><br>Sample | 1<br>Rule of<br>Law<br>*Colonies* | 2<br>Rule of<br>Law<br>*Other* | 3<br>Rule of<br>Law<br>*Full Sample* | 4<br>Regulatory<br>Quality<br> | 5<br>Rule of<br>Law<br>*Colonies* | 6<br>Rule of<br>Law<br>*Other* | 7<br>Rule of<br>Law<br>*Full Sample* | 8<br>Regulatory<br>Quality<br> |
|---|---|---|---|---|---|---|---|---|
| Settler Mortality | -0.486<br>(0.089)** | 0.128<br>(0.241) | 0.128<br>(0.238) | 0.086<br>(0.207) | | | | |
| Mortality* Colony y/n | | | -0.614<br>(0.254)* | -0.482<br>(0.219)* | | | | |
| Colony y/n | | | -1.544<br>(0.567)** | -1.024<br>(0.485)* | | | -1.357<br>(0.294)** | -1.235<br>(0.262)** |
| Disease Environment | | | | | -0.517<br>(0.077)** | 0.128<br>(0.162) | 0.128<br>(0.159) | 0.192<br>(0.15) |
| DE* Colony y/n | | | | | | | -0.645<br>(0.177)** | -0.68<br>(0.167)** |
| Constant | 2.057<br>(0.438)** | 0.714<br>(0.613) | 0.714<br>(0.605) | 0.671<br>(0.531) | 2.197<br>(0.388)** | 0.484<br>(0.554) | 0.484<br>(0.546) | 0.147<br>(0.52) |
| Observations | 62 | 29 | 91 | 91 | 114 | 36 | 150 | 150 |
| R-squared | 0.44 | 0.01 | 0.52 | 0.41 | 0.30 | 0.02 | 0.40 | 0.34 |

Robust standard errors in parentheses; * significant at 5%; ** significant at 1%

Table 3: First Stage Results

The basic results for estimating the determinants of institutions are presented in Table 3. Throughout Columns 1 to 4 of Table 3, I use the historical measure of mortality, while in Column 5 to 8, I present the corresponding results using the geographic measure of disease environment constructed from Table 1.

I first estimate the relationship between mortality and institutional quality separately for each group of countries, corresponding to models (3.4) and (3.3). In Column 1, I repeat the basic regression from AJR and confirm that in the group of former colonies, there is a significant and sizeable negative relation between mortality and institutional outcomes. This is not the case when the same specification is estimated for the group of non-colonies (Column 2). In fact, the correlation is even (insignificantly) positive. I next present the results for the stacked dataset estimating the combined model of institutional quality (3.6) in Column 3 of Table 3.[13] I repeat

---

[13] Throughout Table 3 and in all first stage estimations presented below, I have normalized the interaction term by the average observed level of disease environment in former colonies. This does not affect any coefficient of the presented regression other than the colony dummy. With the normalization at hand, the latter has the interpretation of the average effect of colonization.

my main specification with "Regulatory Quality" as dependent variable in Column 4. Again, while I find that the interaction of mortality with a colony dummy is a significant negative determinant of institutions the main effect is not significant and even of the opposite sign. The results of Columns 3 and 4 are the first basic tests of the colonial origins hypothesis of AJR. The interaction of disease environment and the colony dummy is significant, while the main effect is not. *Germs did indeed affect institutions differentially in former colonies compared to the rest of the world.* This differential effect can not be explained by the direct influence geography has on development, however it can be rationalized in the differential origins of institutions in former colonies and in independent nations.

In addition to a significant interaction, the main effect could also be significant, i.e. implying that geography does affect institutions in all nations through a common channel. This is not supported by the data and the main effect is even positive, though insignificant.[14] I show below that germs do have a direct effect on development, but that it does not work through institutions. I next repeat these finding using my alternative measure of settler mortality, early disease environment. Paralleling the results with the historical measure of settler mortality, this measure of disease is strongly negatively related to institutional outcomes in the group of former colonies (Column 5) while this is not the case in the rest of the sample (Column 6). Correspondingly, I find in the joint model allowing for a main and interaction effect in Column 7 that while the main effect is insignificant and even of the wrong sign, the interaction is significant and negative. I repeat this specification with my alternative measure of institutions, Regulatory Quality, with identical results in Column 8. In Appendix A, I present a variety of other geographically constructed measures of early disease environment again showing the robustness of this finding.

Before proceeding to examine the second stage effects of institutions on economic prosperity, I present some robustness checks in Table 4. An immediate critique that could be brought forward is that the results presented so far are driven by a nonlinear direct impact of geography on institutions. As can easily be made out from Figures 2A and 2B, former colonies are on

---

[14]In some sense, these two findings invalidate Rodrik (2004)'s argument that AJR have found an instrument for institutions without "providing an appropriate explanation" (p.4). Mortality is an instrument for instituions only in the group of former colonies, which can be rationalized only in the context of the channels AJR argue for.

average characterized by high levels of settler mortality. Correspondingly, a positive interaction $\theta_R$ could hence also have resulted from the fact that the relation between disease and prosperity is stronger for higher values of early disease environment. I hence add a mortality square term to the main specification in Columns 1 and 2 of Table 4.[15] In both specifications, the interaction coefficient is changed little and significant. Moreover, the coefficient of square mortality is positive, implying that the relation between disease environment and institutions *becomes weaker for high values of mortality.* Even if the data exhibits a nonlinearity, this would bias the interaction term of mortality towards being positive, not towards being negative. Indeed, comparing the magnitude of the interaction coefficient in Columns 3 and 4 of Table 3 to the magnitude of the corresponding coefficients in Columns 1 and 2 of Table 4, the differential effect becomes larger when a squared mortality regressor is added to the specification.[16] I have also added higher order terms to the estimation of the main effect of mortality, with identical results: taking the first order condition of a polynomial up to fourth degree, the relation between mortality and institutions becomes weaker rather than stronger over the range of observed values of mortality. I repeat this robustness check for my geographic measure of disease environment in Columns 5 and 6. Again, I find that inclusion of a mortality square term increases the importance of the interaction effect of mortality.

I conclude that even if there is some kind of non-linearity in the data, this would imply a lower slope for former colonies than in the rest of the sample and in the worst case, the results in my specifications underestimate the impact of colonization on institutional quality.

Are the results driven by the inclusion of certain groups of countries? A first worry might be that inclusion of African and especially Sub-Saharan countries that have very low scores of institutional development and an adverse disease environment is solely responsible for the findings presented so far. I hence exclude all countries that lye on the African tectonic plate from the sample in Columns 3 and 7.

---

[15] All values of the logarithm of settler mortality and of early disease environment are postive.

[16] In Column 1, the coefficient of disease environment of settler mortality to the square is significant, but that is not a relevant statistic since the linear and the square term are related one to one. An F test of the joint hypothesis that these two coefficients are both equal to 0 is not rejected at the 5% level.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| | | | | Dependent Variable | | | | |
| | Rule of Law | Regulatory Quality | Rule of Law | Rule of Law | Rule of Law | Regulatory Quality | Rule of Law | Rule of Law |
| | Full Sample | | w/o Africa | Europe Ind. | Full Sample | | w/o Africa | Europe Ind. |
| Settler Mortality | -0.497 | -0.297 | 0.203 | 1.371 | | | | |
| | (0.383) | (0.349) | (0.275) | (0.203)** | | | | |
| Mortality^2 | 0.102 | 0.063 | | | | | | |
| | (0.035)** | (0.034) | | | | | | |
| Mortality* Colony y/n | -0.985 | -0.709 | -1.127 | -1.857 | | | | |
| | (0.303)** | (0.242)** | (0.293)** | (0.222)** | | | | |
| Disease Environment | | | | | 0.01 | -0.003 | 0.273 | 0.24 |
| | | | | | (0.425) | (0.418) | (0.164) | -0.18 |
| DE^2 | | | | | 0.019 | 0.031 | | |
| | | | | | (0.061) | (0.06) | | |
| DE* Colony y/n | | | | | -0.696 | -0.765 | -1.042 | -0.756 |
| | | | | | (0.253)** | (0.232)** | (0.206)** | (0.196)** |
| Colony y/n | -1.864 | -1.22 | -2.057 | -4.799 | -1.4 | -1.305 | -1.849 | -1.726 |
| | (0.672)** | (0.535)* | (0.667)** | (0.571)** | (0.332)** | (0.318)** | (0.346)** | (0.319)** |
| Observations | 91 | 91 | 60 | 86 | 150 | 150 | 87 | 143 |
| R-squared | 0.55 | 0.43 | 0.51 | 0.63 | 0.4 | 0.34 | 0.43 | 0.46 |

Robust standard errors in parentheses; * significant at 5%; ** significant at 1%

Table 4: Robustness of the First Stage Relation

The exclusion of African countries does not weaken the results, instead the interaction is larger and more significant in this specification, while the main effect is negligible. Alternatively to excluding African countries I also added a dummy for African countries, which is not significant.

I next check whether the inclusion non-European countries that I classified as independent nations are driving my results. I hence exclude Japan, Thailand, Ethiopia, South Korea and China in Column 4 and in addition Afghanistan and North Korea in Column 8, with unchanged results. I present further robustness check in Appendix B where I drop further groups of nations from the estimation, and also change the classification of certain countries as colonies or independent nations. Similary to the findings presented in Table 4, my results are shown to be robust to these changes.

The findings of this section represent the first main result of this paper: early disease environment is a crucial determinant of institutional quality in former colonies, while it is not an important factor for countries that have not been colonized, hence supporting the institutional

121

theory of development of AJR. Furthermore, I find no evidence of any effect of mortality on institutions that is common to all nations, hence showing that geography has no effect on institutions other than through the indirect influence it had in the process of early institution building in former colonies. I next examine the effects of institutions on economic prosperity directly, where I use the findings so far to identify the endogenous relation between institutions and development.

## 3.5 Mortality, Institutions and Economic Performance: IV Results

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| | | | Second Stage: Dependent is Ln (GDP 2000 /Population) | | | | | |
| | Colonies | Other | Full Sample | | Colonies | Other | Full Sample | |
| **Rule of Law** | 1.771 | -0.136 | 1.373 | | 2.04 | -0.637 | 1.509 | |
| | (0.217)** | (3.00) | (0.317)** | | (0.226)** | (2.81) | (0.211)** | |
| **Regulatory Quality** | | | | 1.751 | | | | 1.429 |
| | | | | (0.435)** | | | | (0.224)** |
| Settler Mortality | | | -0.193 | -0.167 | | | | |
| | | | (0.141) | (0.16) | | | | |
| Disease Environment | | | | | | | -0.274 | -0.355 |
| | | | | | | | (0.083)** | (0.092)** |
| Colony y/n | | | 0.291 | -0.036 | | | 0.38 | 0.097 |
| | | | (0.24) | (0.217) | | | (0.229) | (0.219) |
| Constant | 7.62 | 9.224 | 8.147 | 7.952 | 7.716 | 9.541 | 8.502 | 9.022 |
| | (0.131)** | (3.053)** | (0.589)** | (0.751)** | (0.128)** | (2.510)** | (0.384)** | (0.459)** |
| Observations | 62 | 29 | 91 | 91 | 114 | 36 | 150 | 150 |

First Stage see Table 3. Both Stages Are Estimated Robust; Eicker/Huber/White/sandwich standard errors in parentheses
* significant at 5%; ** significant at 1%

Table 5: Basic IV Results

I now turn to the estimation of the second stage equation (3.7). It is important to point out the exclusion restriction in the instrumental variable specifications, which is that the interaction term of the colony dummy times the measure of mortality only affects institutions and has no further direct effects on income. Paralleling the specifications in the preceding section, I first estimate (3.7) without the full set of geographic covariates and check for their influence later. I am interested in whether the instrumented measures of institutional quality are significant

determinants of income and in whether there is an additional effect of disease environment on income directly (corresponding to $\vartheta_Y$ in (3.7)). If the latter is not the case, the measure of disease environment used by AJR has no effect other than through early institution building. If, on the other side $\vartheta_Y$ is different from 0, some of the effects AJR measure does not stem from the institution building effect, but from the direct effect of disease environment on prosperity and income.

I present the main results of the second stage estimation in Table 5, where each column is using the instrument from the corresponding column in Table 3. I first repeat the separate estimation for colonies and non-colonies in Columns 1 and 2. As is to be expected from the results presented above, the instrumented quality of institutions is highly significant in the group of former colonies, while this is not the case in the group of independent nations where the instrument has no power. For former colonies, the coefficient of Rule of Law is highly significant and estimated at 1.771. Incorporating the first stage results from Table 3, the effect of a 1% decrease of mortality for a colony is estimated to be a 0.86% increase in income per capita (1.771 times 0.486%). In Columns 5 and 6, I again estimate the model separately for the two groups, this time using the alternative mortality data to instrument for institutions. Also in this specification, the instrument is highly significant for former colonies, but not at all for other nations. The estimated coefficient implies that for a former colony, a 1% lower level of disease environment is associated with a 1.06% increase in GDP per capita. These results are comparable to AJR's basic specification, who estimate (in their Table 4, panel B and A, Column 1) a 0.61 coefficient of mortality on average protection from expropriation risk and one of 0.94 for the latter measure of institutional quality on GDP.

I next examine the stacked model in Columns 3 and 7, where the estimation exploits the interaction of mortality and a colony dummy instead of the relying on the main effect of set-tler mortality to instrument for institutions. Because the main effect of mortality or disease environment is not used as an instrument, I can test for a potential direct effect of germs on income levels in the second stage estimation. In both the model using settler mortality (Column 3) and the one using disease environment (7), the data confirms AJR's primary results that institutions are a major determinant of development. The estimated coefficient is highly significant and the associated economic importance is huge. If a country where to improve its

123

institutional score of Rule of Law by one standard deviation, it is predicted to increase income by a factor of around 4.[17]

While confirming that institutions installed during colonization are a major determinant of current income levels, I however find that the estimated coefficients of the quality of institutions are smaller when using the interaction effect rather than the main effect of mortality to instrument for institutions, and the drop is quite sizeable. Comparing Columns 1 and 3, the coefficient of Rule of Law drops from 1.771 to 1.373, while comparing Columns 4 and 7, it drops from 2.04 to 1.509. The reason for this drop is that the estimated coefficient for the direct main effect of disease is estimated to be negative. While the coefficient is not significant for the historical measure of mortality, it is significant for the geographic measure of disease environment in Column 7. Why is the importance of institutions smaller than in an estimation using only colonies? Consider first the models in Columns 1 and 3 using the historical measure of mortality. In order to identify the system, the specification in Column 1 imposes the restriction that the direct coefficient of disease environment on income per capita is equal to 0. This restriction is not present in models 3 and mortality is estimated to have a direct effect on incomes with a coefficient of $-0.193$. While it is important to again point out that the latter coefficient is not significant when using the historical measure of disease, for both samples the change of the point estimate of instituional quality is not negligible

Consider first a ceteris paribus 1% decrease of mortality in a former colony in the stacked model (3.6) including both independent nations and former colonies. In Column 3 of Table 3, a 1% decrease in mortality is associated with an increase if the score of rule of law by 0.00486 points, where the total effect is the sum of interaction ($-0.614$) and main (+0.128) effect. In Column 3 Table 5, this improvement is associated with an indirect (institution building) effect of a 0.486 times 1.373 percentage points higher level of per capita income. In addition, the change of settler mortality is associated with a direct increase of income levels by 0.193%, hence resulting in a total 0.86% increase of income for a 1% drop in mortality. Consider now the same ceteris paribus 1% decrease of mortality in a former colony in the model including only colonies (3.4). In Column 1 of Table 3, a 1% decrease in mortality is again associated

---

[17]The standard deviation of Rule of Law is equal to 1.02. If a country improves its score of Rule of Law by one standard deviation, it is predicted to have an income of $Exp(1.373 * 1.02)$ (Column 3) or $Exp(1.509 * 1.02)$ (Column 7) times its current level.

with an increase of the score of Rule of Law by 0.00486 points. However, because the direct effect of mortality in Column 1 of Table 5 is restricted to 0 in order to identify the system, the estimation attributes all of the 0.86% GDP increase to changes of institutional quality, and hence estimates a coefficient of 1.771, which satisfies $0.486 * 1.771 = 0.486 * 1.33 + 0.193$. The importance of institutions is overstated by around a fourth in the model that encompasses only former colonies.

It is important to point out that while the overall importance of settler mortality for economic prosperity remains unchanged, the interpretation is different. Consider the results from Columns 1 and 3 and consider the effects of a one standard deviation decrease of settler mortality, which in the sub-sample of former colonies is equal to 1.24. For a former colony, a one standard deviation change in early disease environment is associated with a 1.24 times 0.486 point change in the score for Rule of Law (see Table 3 Column 1 or 3). In the estimation that only uses former colonies this is associated with a 191% increase in estimated income stemming from early institution building in the process of colonization. In contrast, the estimation in Column 3 arrives at an institutional effect of 129% with the additional 62% difference being associated with direct effects of disease environment on development.

The results are of comparable magnitude for my alternative measure of mortality estimated from geographic data, khowever this time the main effect of disease on environment is significant. Again, this goes along with a correspondingly smaller effect of institutional quality than in Column 5. Incorporating the findings from the first stage in Table 3, early disease environment and the associated institution building effect is the dominant determinant of development in former colonies. The standard deviation of early disease environment in the sub-sample of former colonies is equal to 0.92. For a former colony, a one standard deviation change in early disease environment is associated with a 0.92 times 0.517 point change in the score for Rule of Law (Table 3 Column 3). This results in about a 106% increase in income due to the early institution building effect. In addition, I estimate that there is a direct geographic effect of mortality resulting in a 29% increase of GDP. I hence again find that while the effect of institutions is the main factor for economic prosperity, the results are smaller than AJR predict because disease environment has a direct effect on development. Comparing the magnitude of the two factors I find that about one fourth of the combined effect of disease environment on

125

development is associated with direct effects of development, while the other three quarters are associated with the institution building channel of development. This finding is again confirmed when using Regulatory Quality as measure of institutional quality in Columns 4 and 8.

The results of this section establish the second main result of this paper. Institutions are indeed the major determinant of long run economic prosperity, but because I find a direct effect of disease on income, the associated magnitude is somewhat smaller than the estimation in a sample composed of only colonies would suggest. My findings – that are only significant for the geographic measure of disease and not when using my historical measure of mortality – suggest that disease environment has sizeable direct effects on development. For example, if it were possible to improve Turkey's geography so that it had a disease environment equivalent to that of Spain (from 4.76 to 2.36), it is predicted to be 1.9 times as rich as it is today. In the next section, I check for the robustness of the results presented so far.

## 3.6  Robustness Checks for Disease Environment

### 3.6.1  Additional Controls

In this section, I check whether the results presented so far are robust to the inclusion of additional geographic, historic or sociologic controls. Because AJR already present a wide range of robustness test for their historical measure of settler mortality, these are not reported here.

In Table 6, I analyze the impact of additional geographic variables and of measures of internal conflict (Fractionalization along several dimensions). In all specifications, I use my main measure of disease environment from Table 1 and also use Rule of Law as dependent variable. I report the first stage results in Panel A and second stage results in Panel B of Table 6. In all the robustness checks that follow below, I find that in the first stage, the interaction of the colony dummy with early disease environment is significant, while the main effect is not. In the second stage, I find that instrumented Rule of Law is a highly significant and important determinant of growth and hence conclude that my instrumentation strategy is robust to the inclusion of the range of controls presented in this section. I also always find that disease environment has a significant direct effect on income.

In Column 1, I add "malaria ecology" from Kiszewski et al. (2004) to the estimation. Because malaria ecology uses the prevalence of certain mosquito species to determine the potential for rather than the actual prevalence of malaria, it is less endogenous than other measures of malaria. However, because also the prevalence of mosquitoes might be endogenous to institutional quality, the results have to be interpreted with care. With this caveat in mind, Column 1 shows that while malaria has no effect on institutions, it does correlate significantly negative with economic outcomes directly. A one standard deviation difference of Malaria Ecology ($\sigma_{ME} = 6.8$) is associated with a 28% difference in GDP per capita. In next add Latitude measured in degrees distance from the equator to my main specification in Column 2. While latitude has no significant impact on development directly, it has a sizeable impact on institutional quality. Through the institutional channel, a one standard deviation difference of Latitude ($\sigma_{Latitude} = 16.2$) is associated with a 67% difference in GDP per capita. This finding points towards some unexplained variation in institutions in line with the "social capital" hypothesis of Hall and Jones (1999).

Frankel and Romer (1999) have argued that openness is good for growth (see Rodrik et al. (2004), though) and I hence include a measure of geographic openness in Column 3. Dollar and Kray (2003) have argued that joint tests of the influence of institutions and trade are suffering from a weak instruments problem because of the multicollinearity of openness and institutions. Because of the latter critique, I use a landlocked dummy – which is not significantly correlated to my measure of disease environment – instead of Frankel and Romer (1999)'s measure of geographic openness. The landlocked dummy influences income directly and with a large magnitude: the impact of having access to open water is associated with an 85% increase in GDP per head, while it has no significant direct effect on institutional quality. This effect is surprisingly large, yet comparable in size to the findings of Frankel and Romer.

To further check for the influence of endowments on economic outcomes, I next include a dummy for oil rich countries that equals 1 if the oil reserves of a country are larger than $5,000$ barrels per capita in Column 4. The latter is associated with a large direct effect on income, yet not on institutions. Finally, I turn to measures of internal conflict that were introduced to the literature by Mauro (1996). Alessina et al. (2004) have constructed three measures of fractionalization along religious, linguistic and ethnic lines, which I include in the analysis in Columns 5 to 7. All three measures are significant determinants of development. However, while religious and linguistic fractionalization influence income directly, fractionalization along ethnic lines seems to work through institutional quality rather than directly. The finding of Column 5 that Ethnic fractionalization matters primarily for institutional outcomes yet does not matter for income directly mirrors the results of Mauro, and confirms that his instrumentation strategy that uses ethnic fractionalization as an instrument for institutions. In contrast, the other two measures of fractionalization from Alessina et al. matter directly for development but not for institutional outcomes.

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| **Panel B: Dependent is Ln (GDP 2000 /Population)** | | | | | | | |
| **Rule of Law** | 1.213 (0.249)** | 1.513 (0.273)** | 1.429 (0.216)** | 1.396 (0.216)** | 1.488 (0.254)** | 1.299 (0.197)** | 1.446 (0.195)** |
| Disease Environment | -0.222 (0.089)* | -0.276 (0.073)** | -0.285 (0.087)** | -0.273 (0.086)** | -0.271 (0.077)** | -0.239 (0.070)** | -0.324 (0.081)** |
| Colony y/n | 0.153 (0.21) | 0.366 (0.287) | 0.328 (0.215) | 0.275 (0.228) | 0.377 (0.221) | 0.293 (0.188) | 0.416 (0.206)* |
| Malaria Ecology | -0.041 (0.013)** | | | | | | |
| Latitude (degrees) | | -0.001 (0.014) | | | | | |
| Landlocked y/n | | | -0.642 (0.175)** | | | | |
| Religious frac | | | | | | | -0.839 (0.310)** |
| Languistic frac | | | | | | -0.993 (0.262)** | |
| Ethnic frac | | | | | -0.119 (0.44) | | |
| Oil > 5,000 b / Cap y/n | | | | 0.584 (0.148)** | | | |
| R-squared | 0.77 | 0.73 | 0.76 | 0.77 | 0.73 | 0.78 | 0.75 |
| **Panel A: First Stage: Dependent is Rule of Law** | | | | | | | |
| Disease Environment | 0.131 (0.159) | 0.16 (0.132) | 0.124 (0.161) | 0.127 (0.159) | 0.134 (0.142) | 0.129 (0.156) | 0.144 (0.164) |
| DE* Colony y/n | -0.594 (0.196)** | -0.532 (0.157)** | -0.635 (0.179)** | -0.639 (0.177)** | -0.565 (0.168)** | -0.644 (0.179)** | -0.656 (0.181)** |
| Colony y/n | -1.33 (0.300)** | -0.602 (0.275)* | -1.348 (0.299)** | -1.354 (0.294)** | -1.141 (0.281)** | -1.304 (0.299)** | -1.392 (0.305)** |
| Malaria Ecology | -0.008 (0.012) | | | | | | |
| Latitude (degrees) | | 0.029 (0.007)** | | | | | |
| Landlocked y/n | | | -0.113 (0.16) | | | | |
| Religious frac | | | | | | | 0.237 (0.263) |
| Languistic frac | | | | | | -0.176 (0.249) | |
| Ethnic frac | | | | | -0.801 (0.310)* | | |
| Oil > 5,000 b / Cap y/n | | | | 0.042 (0.14) | | | |
| Observations | 148 | 150 | 150 | 150 | 148 | 145 | 150 |
| R-squared | 0.4 | 0.49 | 0.41 | 0.40 | 0.43 | 0.42 | 0.41 |

Both Stages Are Estimated With Robust Standard Errors in Parentheses; * significant at 5%; ** significant at 1%

Table 6: Additional Geographic Controls

### 3.6.2 Other Instruments and Overidentificaiton Tests

In this section, I introduce further instruments for institutions and check whether the instrumented institutional scores are mutually consistent, i.e. I test the overidentified system. The additional instruments I use are scores of Democracy in the early $19^{th}$ century (the Polity Score & Constraints on the Executive from the Polity IV database), ethnolinguistic fractionalization and legal origin dummies.

The results are reported in the now familiar way, where Panel A of Table 7 reports the first stage results with the new instruments included. In all specifications, the interaction stays a valid instrument, while also the added instruments are significant (or at least one of the added dummies in the case of legal origins). Panel B reports the instrumented second stage, again with similar results like earlier specifications. Institutions are a main determinant of economic performance, while disease environment has a significant effect on economic development. In Panel C, I report a Hausman overidentification test with the null hypothesis that the new instrument and the interaction of a colony dummy and disease environment are mutually consistent instruments. The Hausman test rejects in no specification. Because there may well be heteroscedasticity in the sample, I also report a Hansen C-test, which rejects in one case (see below).

I first include the Polity Score as an additional instrument. This measure lies between -10 and +10 and is higher for more democratic societies. In Column 1 of Table 8, I include all countries with an available Polity score in 1900 to 1910 and take the earliest available score. Few countries that still exist today were independent at that time and there are thus only 55 observations. Still, both instruments (polity and interaction) are significant, and Polity in 1900 is a an important determinant of current institutions. I repeat this estimation in Column 2, but this time include all countries with an available Polity score before 1961, hence yielding 106 countries with similar results as in the previous specification. In Columns 3 and 4 of Table 8, I repeat this exercise, but I focus on a sub-indicator of the Polity database, "Constraint on the Executive". This variable lies between 1 and 7 and measures whether superimposed structures and rules effectively constrain the executive. I again include first all countries that have a Polity

score in 1900 to 1910 and still exist today (Column 3) and then those with a polity score before 1961 (Column 4). In both regressions, constraints on the executive as well as the interaction are significant determinants of institutions. I next include Ethnic Fractionalization from Alesina et al. (2004) in the specification. This specification is motivated by the finding from the previous section that ethnic fractionalization matters for development only through institutions. Finally, I turn to the "Legal Origins" dummies form La Porta et al.(1999). Socialist origin is associated with lower institutional outcomes (Column 6). The estimation in Column 7, which includes a British and a French legal origin dummy, yields the only model where the Hansen C-test of overidentification is rejected. However, the C-test also rejects a model where I only include the British and French legal origin dummy and not my interaction instrument. The reason that the C-test rejects in Column 7 is hence that two legal origins dummies are mutually inconsistent instruments. In Column 8, I add all three legal origins (hence only missing the German and Scandinavia origin dummies), to show that the preceding results were due to the fact that British and French legal origins countries tend to be non-soviet and again the overidentification is not rejected.

Overall, I conclude that the results presented so far are robust to both the inclusion of geographic controls and also to the inclusion of additional instruments. A difference between mine and AJR's result is however that I do find that other variables, such as a landlocked dummy or latitude do have sizeable additional explanatory power for economic outcomes.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| | | | **Panel B: Dependent is Ln (GDP 2000 /Population)** | | | | | |
| **Rule of Law** | 0.909 | 1.234 | 0.895 | 1.224 | 1.542 | 1.413 | 1.189 | 1.261 |
| | (0.318)** | (0.309)** | (0.372)* | (0.315)** | (0.168)** | (0.128)** | (0.198)** | (0.105)** |
| Disease Environment | -0.299 | -0.344 | -0.306 | -0.348 | -0.262 | -0.305 | -0.377 | -0.354 |
| | (0.274) | (0.157)* | (0.301) | (0.164)* | (0.072)** | (0.079)** | (0.098)** | (0.076)** |
| Colony y/n | 0.047 | 0.213 | 0.038 | 0.206 | 0.391 | 0.319 | 0.176 | 0.222 |
| | (0.208) | (0.242) | (0.22) | (0.24) | (0.224) | (0.184) | (0.192) | (0.178) |
| R-squared | 0.76 | 0.76 | 0.76 | 0.76 | 0.72 | 0.74 | 0.66 | 0.75 |
| | | | **Panel C: Overidentification Tests (P Values)** | | | | | |
| HaumanTest | 0.9483 | 0.6002 | 0.6215 | 0.9782 | 0.9956 | 0.9635 | 0.2543 | 0.4382 |
| C-Test | 0.08229 | 0.50338 | 0.08608 | 0.46455 | 0.79038 | 0.49875 | (0.04501)* | 0.31297 |
| | | | **Panel A: First Stage Without Interaction: Dependent is Rule of Law** | | | | | |
| Disease Environment | 0.01 | -0.093 | 0.023 | -0.072 | 0.134 | -0.001 | 0.087 | 0.031 |
| | (0.262) | (0.194) | (0.268) | (0.196) | (0.142) | (0.129) | (0.162) | (0.108) |
| DE* Colony y/n | -0.618 | -0.374 | -0.64 | -0.397 | -0.565 | -0.501 | -0.56 | -0.478 |
| | (0.264)* | (0.21) | (0.290)* | (0.212) | (0.168)** | (0.150)** | (0.180)** | (0.134)** |
| Colony y/n | -1.486 | -1.06 | -1.321 | -0.994 | -1.141 | -1.512 | -1.527 | -1.318 |
| | (0.449)** | (0.292)** | (0.449)** | (0.299)** | (0.281)** | (0.281)** | (0.326)** | (0.262)** |
| Polity 1900 | 0.061 | 0.03 | | | | | | |
| | (0.017)** | (0.012)* | | | | | | |
| Constraint on Exec. 1900 | | | 0.153 | 0.092 | | | | |
| | | | (0.053)** | (0.035)* | | | | |
| Ethnic Fractionalization | | | | | -0.801 | | | |
| | | | | | (0.310)* | | | |
| Socialist Legal Origin | | | | | | -1.141 | | -1.691 |
| | | | | | | (0.167)** | | (0.182)** |
| British Legal Origin | | | | | | | 0.546 | -0.571 |
| | | | | | | | (0.214)* | (0.209)** |
| French Legal Origin | | | | | | | 0.257 | -0.865 |
| | | | | | | | (0.194) | (0.191)** |
| Observations | 55 | 106 | 55 | 106 | 148 | 150 | 150 | 150 |
| R-squared (1st stage) | 0.54 | 0.48 | 0.53 | 0.49 | 0.43 | 0.51 | 0.43 | 0.55 |

Both Stages Are Estimated With Robust Standard Errors in Parentheses; * significant at 5%; ** significant at 1%

Note: In Column 7 the C-Test also rejects the British Legal and the French Legal Origins dummy being mutually consistent.

Table 7: Additional Instruments and Overidentification Tests

## 3.7 Conclusion

The contribution of this paper is to provide a sharp test to disentangle the colonial origins theory proposed by Daron Acemoglu, Simon Johnson and James Robinson from the potential direct impact of disease environment on prosperity. While the geography view of development predicts a relation between mortality rates and development for all countries equally, the institutional view applies only to the subset of former colonies. Using either additional historical sources or a model of the geographic determinants of disease, I first construct two measures of mortality rates including up to 36 countries that have not been colonized. I then show that mortality did affect institutional development exclusively in former colonies, which can only be rationalized in the context of the colonial origins theory of AJR.

Exploiting the differential role mortality played in the process of creating today's institutions in former colonies and in other nations, I am able to show that indeed, the theory provided by AJR is supported by the data and that institutions are the main determinant of economic success. However, I also find that disease environment influences income directly and correspondingly, that institutions are somewhat less important for prosperity in my specifications than when working with a sample composed of only former colonies. As a rough rule of thumb, I estimate that for a former colony, a one percent lower settler mortality is associated with a one percent higher income per capita. Of this 1% total difference, I attribute three quarters of a percent to the early institution-building channel, while the other fourth is a direct consequence of disease on productivity and income.

What do these results imply for the overall importance of institutions and geography as the long run determinants of prosperity? Again, the findings of this paper point towards the importance of institutions as the main ingredient for development. In a typical specification, I estimate that if one where to improve a country's current disease environment by 1/100 of a standard deviation, this would result in an increase of income between 0.3% and 0.4%. In contrast, if it were possible to increase the quality of institutional outcomes by 1/100 of a standard deviation, the associated change of income would be in the order of 1.4% to 1.6%.

Of course, these results have to be interpreted with care when thinking about the direction development efforts should be directed to. While research has struggled to develop good strategies against corruption and bad governance, programs targeted at specific diseases such

133

as the Guinea Worm, HIV or Malaria can achieve great success at a relatively low cost. If it were possible – for example by the programs reviewed in Sachs (2005) – to decrease the prevalence of disease, this could well be associated with large economic gains through the associated increases in human capital accumulation and direct effects on productivity.[18] However, the results presented in this study also suggest that policies concerned with achieving sustained catch-up of poor nations and long run convergence of international income levels have to focus on improving institutions, which are the main channel of economic development.

---

[18]For example, if it were possible to eradicate certain diseases such that the mortality of men in the cohort of 20 to 40 years of age is reduced by a third, the results of this study suggest long run gains in the order of about 15% of current income levels.

## 3.8 Appendix A: Alternative Measures of Disease Environment

The novelty of the specifications in this paper is that I am able to test for a differential effect of mortality for former colonies and other countries because I construct an estimate of mortality for nations that have not been colonized. The quality of my historical measure of mortality is limited by the availability of historical sources and the results thus have to be interpreted with care. Yet, how good is the quality of my measure of early disease environment and how sensitive are the results to the way it was constructed?

In this section, I first present alternative specifications relating geography to AJR's mortality rate and I show that that the results are robust to using alternative variables to instrument for disease. Thereafter, I address a recent critique of AJR's data by Albouy (2005), who constructs an alternative mortality rates and finds no strong relation between his series and economic outcomes. Interestingly, while AJR's and his mortality rates differ substantially, both are affected by geography in a very similar way and consequently, my predicted measure of early disease environment yields very similar results when using Albouy's instead of AJR's mortality rates as the initial mortality variable.[19]

I present the different models of the geographic determinants of disease in Table A1. I first turn to how different geographic variables relate to AJR's mortality series. In Column 1, present a basic estimation, which includes only measures of the importance of dry (low minimum monthly rain) or wet seasons (high maximum monthly rain) as well as dummies for mediteranean climate, the incidence of savanna or temperate forest and for the incidence of mountains. In I next add latitude measured in average degrees distance from the equator to the previous model in Column 2. Conditional on the other geographic variables in the estimation, latitude has no impact on mortality rates. This is to be expected since distance from equator itself should not affect disease environment, while humidity or the incidence of deserts are associated with real effects on the human physique. In Column 3, I thus drop Latitude and add measures of temperature and humidity. Throughout Table A1, I have standardized the dependent variable and hence I again present my estimation of my main model of mortality in Column 3 of Table

---

[19] A potential further critique to estimating model (3.1) in the sample of AJR is that these nations are characterized by high rates of mortality and I then predict for many nations with rather low prevalence of disease. It is important to point out that this selection by the dependent variable does not create a bias in the estimation of mortality.

135

A1, where all coefficients and standard deviations are equal to the coefficients of Table 1 divided by the standard deviation of AJR's mortality data (1.24). High temperature (absolute degrees Celsius) increases mortality levels. A measure of the maximum temperature when humidity is at its maximum level of the day is also a significant determinant of mortality. However, conditional on average temperature, higher temperature at maximum humidity decreases mortality – i.e. moist and cold weather increases morbidity. Based on an F-test, this model performs the best and is used as the main specification in the main text. I have also checked for the influence of coastal area and average distance from coast. Similarly to latitude, once the effect of humidity and temperature are taken into account, the latter variables do not affect mortality.

Rather than examining the impact of further geographic variables, I also estimate the model with a different dependent variable, the standardized logarithm of mortality as collected by Albouy (2004). This is motivated by Albouy's critique that AJR's data suffer from "a number of inconsistencies, questionable judgments, and mistakes in the mortality data [which] induce an artificial correlation of mortality with expropriation risk and GDP per capita." Acemoglu et al. (2005), in turn, raise the same concerns about his alternative data series.

| | Source | 1<br>Basic | 2<br>Latitude | 3<br>Main | 4<br>Albouy Camp. |
|---|---|---|---|---|---|
| | | **Standardized Ln (Mortality)** | | | |
| | | Acemoglu et al. (2001) | | | Albouy (2005) |
| Latitude | | | -0.0014<br>(0.098) | | |
| Avg. Temperature | | | | 0.078<br>(0.035)* | 0.063<br>(0.044) |
| max Temp/ max Humidity | | | | -0.061<br>(0.019)** | -0.056<br>(0.022)* |
| Min. Rain | | -0.011<br>(0.002)** | -0.011<br>(0.002)** | -0.012<br>(0.002)** | -0.005<br>(0.003) |
| Max Rain | | 0.00104<br>(0.00047)* | 0.00102<br>(0.00049)* | 0.00099<br>(0.00056) | 0.00152<br>(0.00064)* |
| Mediterranean y/n | | -0.738<br>(0.21863)** | -0.72392<br>(0.25339)** | -1.12219<br>(0.23559)** | -0.92302<br>(0.31706)** |
| Savanna y/n | | 0.474<br>(0.155)** | 0.463<br>(0.195)* | 0.625<br>(0.155)** | 0.463<br>(0.189)* |
| Temperate forest y/n | | -0.974<br>(0.239)** | -0.947<br>(0.304)** | -1.087<br>(0.213)** | -1.025<br>(0.396)* |
| Mountains y/n | | -0.746<br>(0.185)** | -0.744<br>(0.186)** | -0.522<br>(0.240)* | -0.406<br>(0.333) |
| Observations | | 63 | 63 | 63 | 63 |
| R-squared | | 0.64 | 0.64 | 0.7 | 0.48 |
| F-Test | | 15.77** | 14.23** | 19.23** | 8.99** |

Robust standard errors in parentheses; * significant at 5% level; ** significant at 1% level

Table A1: Different measures of Disease Environment

Rather than taking any position on the quality of the data, I now repeat the analysis with Albouy's alternative data as dependent variable. The crucial insight of doing so is that although the historical measure of mortality might be measured with error that is correlated with economic outcomes, the strategy pursued here of instrumenting mortality with its geographic component gets rid of this possible source of bias. Although my results point toward Albouy's measure being measured with more noise because they are less strongly related to geography, the main results - that early disease environment matters more for development in former colonies - is a robust finding that does not depend on which mortality estimate is used in the estimation.

I next predict the four models of Table A1. While it is not surprising that models 1 to 3 correlate highly, it is noteworthy that when I predict the relation between geography and AJR's data on the one side and between geography and Albouy's data on the other side (i.e Columns 3 and 4 in Table 4), the two resulting measures of early disease environment are very similar. For example while the temperate forest dummy carries a coefficient of −1.09 in the estimation using AJR's data as dependent, the coefficient is estimated to be −1.03 when using Albouy's data in Column 4. Correspondingly, while the coefficient of correlation of the respective historical mortality rates in the samples of AJR and Albouy is equal to 0.87, the correlation between the two predicted series is estimated at 0.93.

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Model of Mortality | Basic | Latitude | Main | Albouy Cam |
| **Panel B: Dependent is Ln (GDP 2000 /Population)** | | | | |
| **Rule of Law** | 1.78 | 1.783 | 1.509 | 1.514 |
| | (0.355)** | (0.369)** | (0.211)** | (0.202)** |
| Disease Environment | -0.235 | -0.231 | -0.34 | -0.395 |
| | (0.178) | (0.188) | (0.103)** | (0.116)** |
| Colony y/n | 0.516 | 0.519 | 0.38 | 0.355 |
| | (0.316) | (0.318) | (0.229) | (0.234) |
| R-squared | 0.67 | 0.67 | 0.73 | 0.73 |
| | | | | |
| **Panel A: First Stage: Dependent is Rule of Law** | | | | |
| Disease Environment | -0.066 | -0.084 | 0.159 | 0.226 |
| | (0.244) | (0.248) | (0.197) | (0.226) |
| DE* Colony y/n | -0.614 | -0.601 | -0.799 | -0.985 |
| | (0.263)* | (0.266)* | (0.220)** | (0.262)** |
| Colony y/n | -1.077 | -1.056 | -1.357 | -1.386 |
| | (0.299)** | (0.304)** | (0.294)** | (0.282)** |
| Observations | 150 | 150 | 150 | 150 |
| R-squared | 0.40 | 0.40 | 0.40 | 0.39 |

Both Stages Are Estimated With Robust Standard Errors in Parentheses;
* significant at 5%; ** significant at 1%

Table 2A: Results with Alternative Measures of Disease Environment

Using these four alternative models of disease environment, I repeat my basic specification that allows for a main and interaction effect of early disease environment as well as the colony dummy. The results are presented in Table A2, Panel A and B. In all four specifications,

the interaction of early disease environment and the colony dummy is highly significant for institutions, while the main effect is not significant and of alternating sign. Also in the second stage estimation in Panel B, I again find for all four measures of early disease environment that institutions as well as disease environment directly are significant determinants of institutions. I conclude that my results are not driven by the selection of a specific geographic model of the determinants of disease or by using a specific historical measure of disease.

## 3.9 Appendix B: More Alternative Samples

As discussed in the main text, I now present a couple of additional robustness checks that show that my results are not dependent on the way I classified countries as former colonies or independent nations.

Two key concerns - that the results are exclusively driven by either the inclusion of African countries or by the inclusion of non-European nations that are classified as independent nations have already been addressed in Table 4. I next turn to exclude all countries that cannot unambiguously be classified as either colonies or independent nations in Column 1 of Table A3. The estimation in total excludes 15 countries, which are Afghanistan, Albania, Cambodia, China, Ethiopia, Japan, North and South Korea, Kuwait, Liberia, Oman, Saudi Arabia, Taiwan, Thailand and the United Arab Emirates. Again, the results hold up and in the first stage the interaction effect is significant, which is not the case for the main effect. Also in the second stage the results hold up in the restricted sample of 135 countries and the associated coefficients of instrumented institutional quality and the direct effect of mortality are comparable to my basic specification.

In the main text, I define a country as former colony if it was subject to substantial foreign influence during imperialist times and it is not adjacent to that foreign power whose influence it was under. For example, this leads me to classify Ethiopia as an independent nation, although it formally was an Italian colony from 1936 to 1941. On the other side, I count Liberia as a colony because it was subject to substantial foreign influence in the forms of Slave trade and limited British occupation.[20] While the specification in Column 1 makes clear that the results are not driven by inclusion of these countries, I now show that even when classifying certain countries in the opposite way, the results hold up. I report the results of the estimation with Ethiopia counted as a former colony in Column 2 and with Liberia as an independent nation in Column 3, with unchanged results. Finally, as a last robustness check I drop Guinea from the sample in Column 4. This is motivated by the fact that Guinea was the only country that could be seen as an outlier when predicting my measure of early disease in Table 1. Again, I

---

[20] In one of the unfortunately not too rare twists of history, the freed Afro-American Slaves that arrived in Liberia after 1817 saw themselves superior to the native population and engaged in extractive activities and battles with the indigenous population in a fashion similar to the policies pursued by the Caucasian settlers in neighbouring countries.

find that this change in sample does not weaken my results.

| | 1<br>Restricted<br>Sample | 2<br>Afghanistan<br>as Colony | 3<br>Liberia<br>as Independent | 4<br>Dropping<br>Guinea |
|---|---|---|---|---|
| **Panel B: Dependent is Ln (GDP 2000 /Population)** | | | | |
| Rule of Law | 1.563 | 1.621 | 1.509 | 1.536 |
| | (0.243)** | (0.261)** | (0.211)** | (0.203)** |
| Disease Environment | -0.254 | -0.216 | -0.274 | -0.278 |
| | (0.100)* | (0.111) | (0.083)** | (0.084)** |
| Colony y/n | 0.428 | 0.359 | 0.38 | 0.428 |
| | (0.247) | (0.221) | (0.229) | (0.219) |
| R-squared | 0.73 | 0.71 | 0.73 | 0.73 |
| **Panel A: First Stage: Dependent is Rule of Law** | | | | |
| Disease Environment | 0.065 | 0.037 | 0.128 | 0.128 |
| | (0.199) | (0.172) | (0.159) | (0.157) |
| DE* Colony y/n | -0.6 | -0.576 | -0.645 | -0.684 |
| | (0.214)** | (0.188)** | (0.177)** | (0.176)** |
| Colony y/n | -1.385 | -1.087 | -1.358 | -1.322 |
| | (0.331)** | (0.327)** | (0.295)** | (0.292)** |
| Observations | 135 | 150 | 150 | 149 |
| R-squared | 0.47 | 0.37 | 0.40 | 0.40 |

Both Stages Are Estimated With Robust Standard Errors in Parentheses;
* significant at 5%; ** significant at 1%
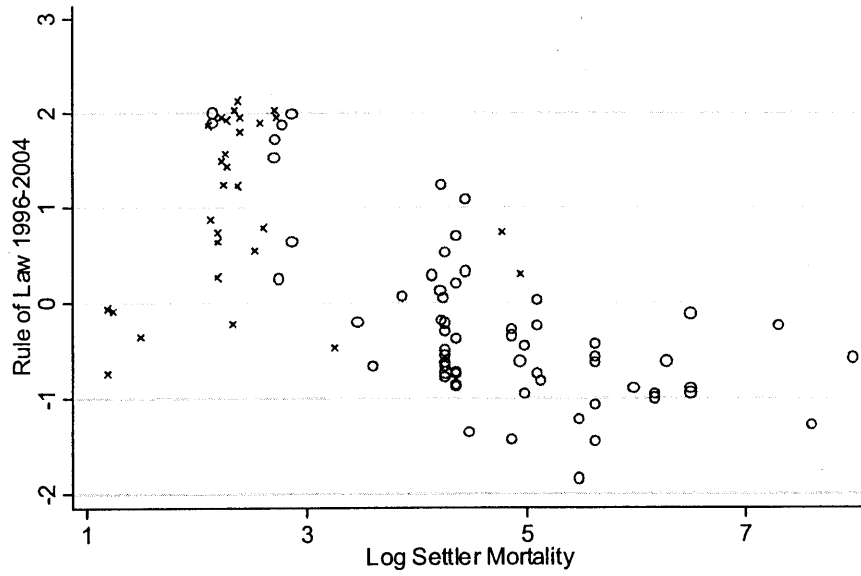
A3 Results with Different Classifications

# Bibliography

[1] Acemoglu, Daron; Johnson, Simon and Robinson, James A. "The Colonial Origins of Comparative Development: An Empirical Investigation" The American Economic Review, December 2001, 91(5), pp. 1369-1401.

[2] Acemoglu, Daron; Johnson, Simon and Robinson, James A. "Reversal of fortune: Geography and institutions in the making of the modern world income distribution" Quarterly Journal of Economics, November 2002, 117 (4), pp. 1231-1294

[3] Acemoglu, Daron; Johnson, Simon and Robinson, James A. "Disease and Development in Historical Perspective," Journal of the European Economic Association Papers and Proceedings, April 2003, 1, pp. 397-405

[4] Acemoglu, Daron; Johnson, Simon and Robinson, James A. "A response to Albouy's "A Reexamination Based on Improved Settler Mortality Data" MIT Department of Economics Working Paper, March 2005.

[5] Albouy, David. "The Colonial Origins of Comparative Development: A Reexamination Based on Improved Settler Mortality Data," December 2004, University of California, Berkeley.

[6] Alesina, Alberto; Devleeschauwer, Arnaud; Easterly, William; Kurlat, Sergio and Wacziarg, Romain. "Fractionalization," National Bureau of Economic Research, Working Paper 9411, 2003.

[7] Bloom, David E. and Sachs, Jeffrey D. "Geography, Demography, and Economic Growth in Africa." Brookings Papers on Economic Activity, 1998, (2), pp. 207-73.
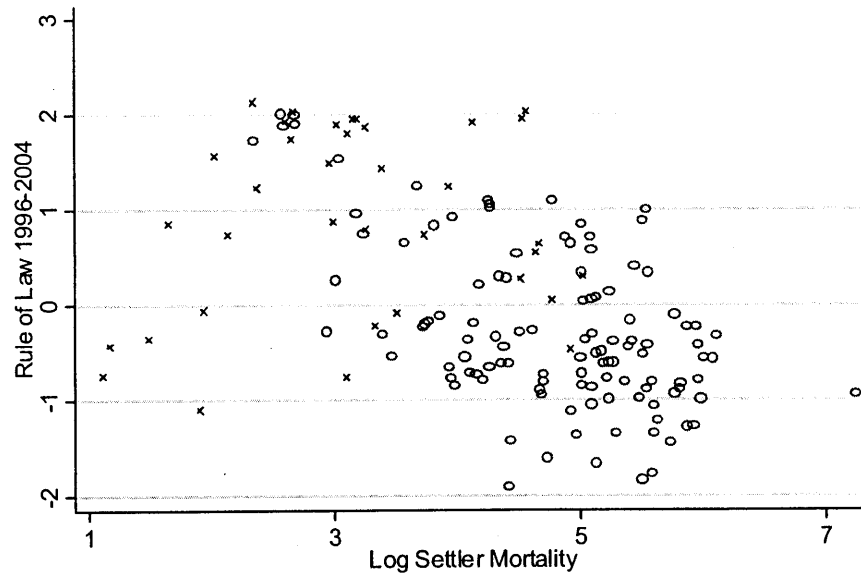
[8] **Curtin, Philipp D.** *Death by migration: Europe's encounter with the tropical world in the 19th Century.* New York: Cambridge University Press, 1989.

[9] **Dollar, David, and Aart, Kraay.** "Institutions, Trade and Growth: Revisiting the Evidence," The World Bank, Working Paper, January 2003.

[10] **Easterly, William and Levine, Ross.** "Tropics, germs, and crops: the role of endowments in economic development," January 2003, Journal of Monetary Economics, 50(1), pp. 3-39.

[11] **Gallup, John L.; Mellinger, Andrew D. and Sachs, Jeffrey D.** "Geography and Economic Development." National Bureau of Economic Research, Working Paper No. 6849, 1998.

[12] **Kaufmann, Daniel; Kraay, Art and Mastruzzi, Massimo.** "Governance Matters IV: Governance Indicators for 1996-2004" World Bank Policy Research Working Paper Series, May 2005, No. 3630

[13] **Kiszewski, Anthony; Mellinger, Andrew; Spielman, Andrew; Malaney, Pia; Ehrlich Sachs, Sonia and Sachs, Jeffrey D.** "A Global Index of the Stability of Malaria Transmission," American Journal of Tropical Medicine and Hygiene, May 2004, 70(5), pp. 486-498.

[14] **La Porta, Rafael; Lopez-de-Silanes, Florencio; Shleifer, Andrei and Vishny, Robert W.** "Law and Finance," The Journal of Political Economy, December 1998, 106(6), pp. 1113-1155.

[15] **La Porta, Rafael; Lopez-de-Silanes, Florencio; Shleifer, Andrei and Vishny, Robert W.** "The Quality of Government," Journal of Law, Economics and Organization, December 1998, 15(1),pp. 222-279.

[16] **Lipset, Seymour Martin.** *The Political Man. The Social Basis of Politics,* Garden City, N.Y: Doubleday, 1960.

[17] **Mauro, Paolo.** "Corruption and Growth," The Quarterly Journal of Economics, August 1995, 110(3), pp. 681-712.

[18] **McArthur, John W. and Sachs, Jeffrey.** "Institutions and Geography: Comment on Acemoglu, Johnson and Robinson," National Bureau of Economic Research, Working Paper No. 8114, February 2001.

[19] **Mitchel, B.R.** *International Historical Statistics Europe 1750 - 1993*, New York: Stockton Press, Fourth Edition, 1998.

[20] **Parker, Philip M.** *National cultures of the world: A statistical reference*, Cross Cultural Statistical Encyclopedia of the World, Vol. 4 Westport, CT: Greenwood Press, 1997.

[21] **Rigobon, Roberto and Rodrik, Dani.** "Rule of Law, Democracy, Openness and Income: Estimating the Interrelationships", July 2005, The Economics of Transition, 13(3), pp. 533-64.

[22] **Rodrik, Dani.** "Getting Institutions Right", April 2004, Kennedy School of Government, Harvard University.

[23] **Rodrik, Dani; Subramanian, Arvind and Trebbi, Francesco.** "Institutions rule: The primacy of institutions over geography and integration in economic development", Journal of Economic Growth, June 2004, 9 (2), pp. 131-165.

[24] **Sachs, Jeffrey D.** "Institutions Don't Rule: Direct Effects of Geography on Per Capita Income" National Bureau of Economic Research, Working Paper No. 9490, February 2003.

[25] **Steckel, Richard H. and Floud, Roderick.** *Health and Welfare during Industrialization.* National Bureau of Economic Research Project Report, Chicago: University of Chicago Press, 1997

## Independent Nations



## Former Colonies



Figure 3A and 3 B: Settler Mortality and Institutions in the Two Subsamples

## Mortaltiy and Institutions



## Diesease Environment and Institutions



Germs and Institutional Quality