# System Perspective and Lean Coordination

**The case of Open Source Software**

## Joao Castro

Committee:    Warren Seering
Eric Rebentisch
Chris Magee

# Coordination and lean

What they mean to me:

- – Lean as being effective, efficient, minimal waste

- – Coordination as being effective in managing interactions in complex environments

Can we get efficient coordination?

- In high complexity environments?

Product development is the study context

# Coordination

- Coordination is consistently cited as one of the most important factors of competitive advantage

> "The primary task of management is to get people to **work together in a systematic way**"
>
> Christensen et al, (2006), *The Tools of Cooperation and Change.* Harvard Business Review

- Literature describes many methods
  - Centralize people
  - Centralize information
  - Facilitate Communication
  - Structure Communication
  - Structure Processes

# Field exploration

- Visited and collected data, interviews from different PD companies:
  - Industrial machinery (2)          - Services (2)
  - Aerospace (2)                     - Food packaging
  - Medical equipment                 - PD consultancy


- Factors affecting coordination are varied:
  - Team size                         - Reputation for prob. solving
  - Schedule                          - Product complexity

# Where to focus?

"Ohno thought that assembly workers could probably do most of the functions of the specialists and **do them much better because of their direct acquaintance with conditions on the line.**"

Womack, Jones, and Roos, (1990)

*The Machine That Changed The World*

# Where to focus?

- Brook's Law implies that the **ideal size for a programming team is one** – a single developer who never has to stop to communicate with a colleague.

- This approach streamlines everything, and it also provides insurance that the project will retain **"conceptual integrity"**.

Rosenberg, (2007), *Dreaming in Code.*

# Where to focus?

"I also believed there was a **certain critical complexity above which a more centralized, a priori approach was required.**"

"Linus Torvalds's style of development came as a surprise. (...)The fact that this bazaar style **seemed to work, and work well**, came as a distinct shock."

*Eric Raymond, (1999) **The Cathedral and the Bazaar***

# Where to focus?

But, as Galileo is said to have murmured after officially recanting his statement that the earth moves around the sun:

## "And yet it moves!"

## What is going on here?"

von Krogh and von Hippel (2006)

*The Promise of Research on Open Source Software*

# Driving Hypothesis

Individuals in a collaborative environment and behaving autonomously are able to *efficiently* solve complex problems
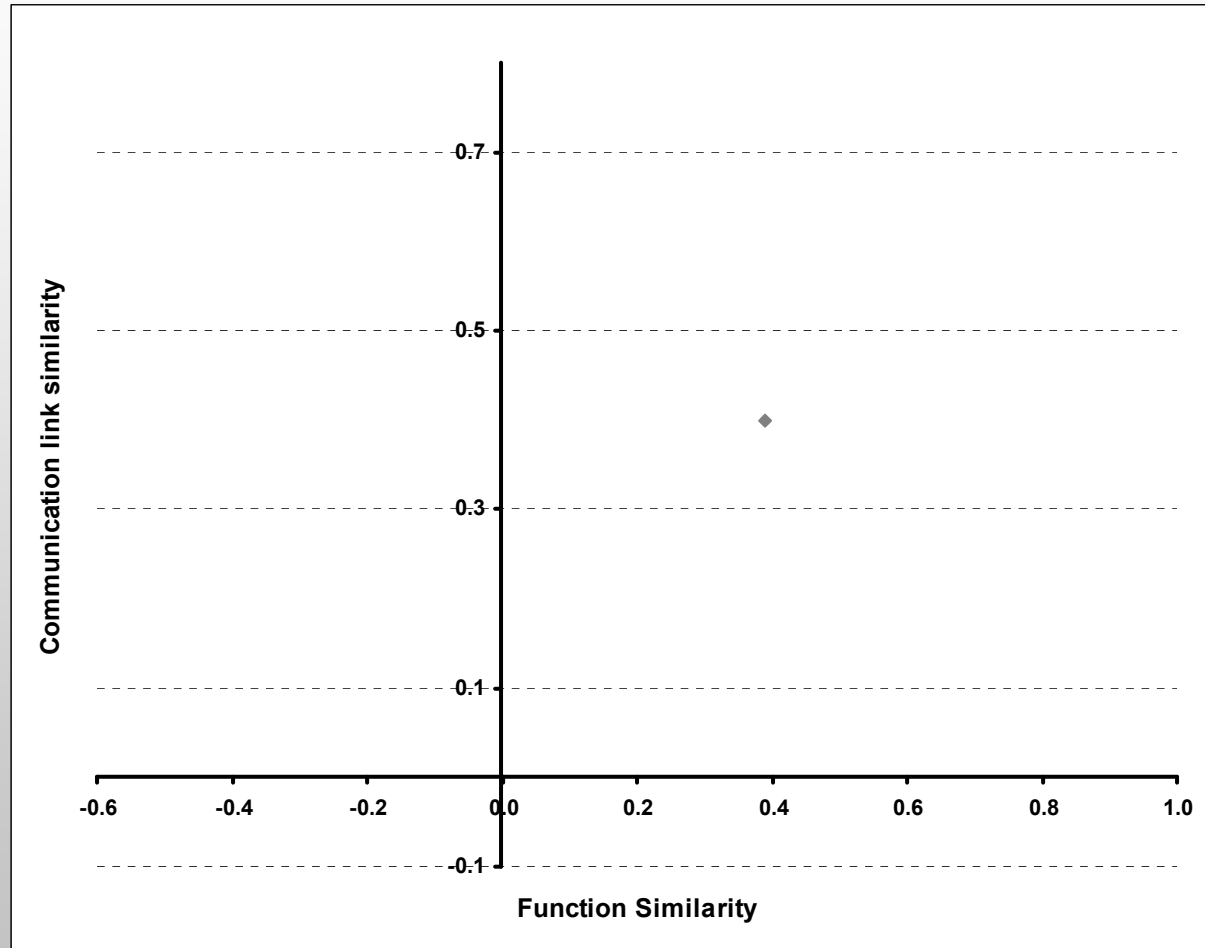
In other words:

Coordination is possible without heavy supervisory and overhead methods.

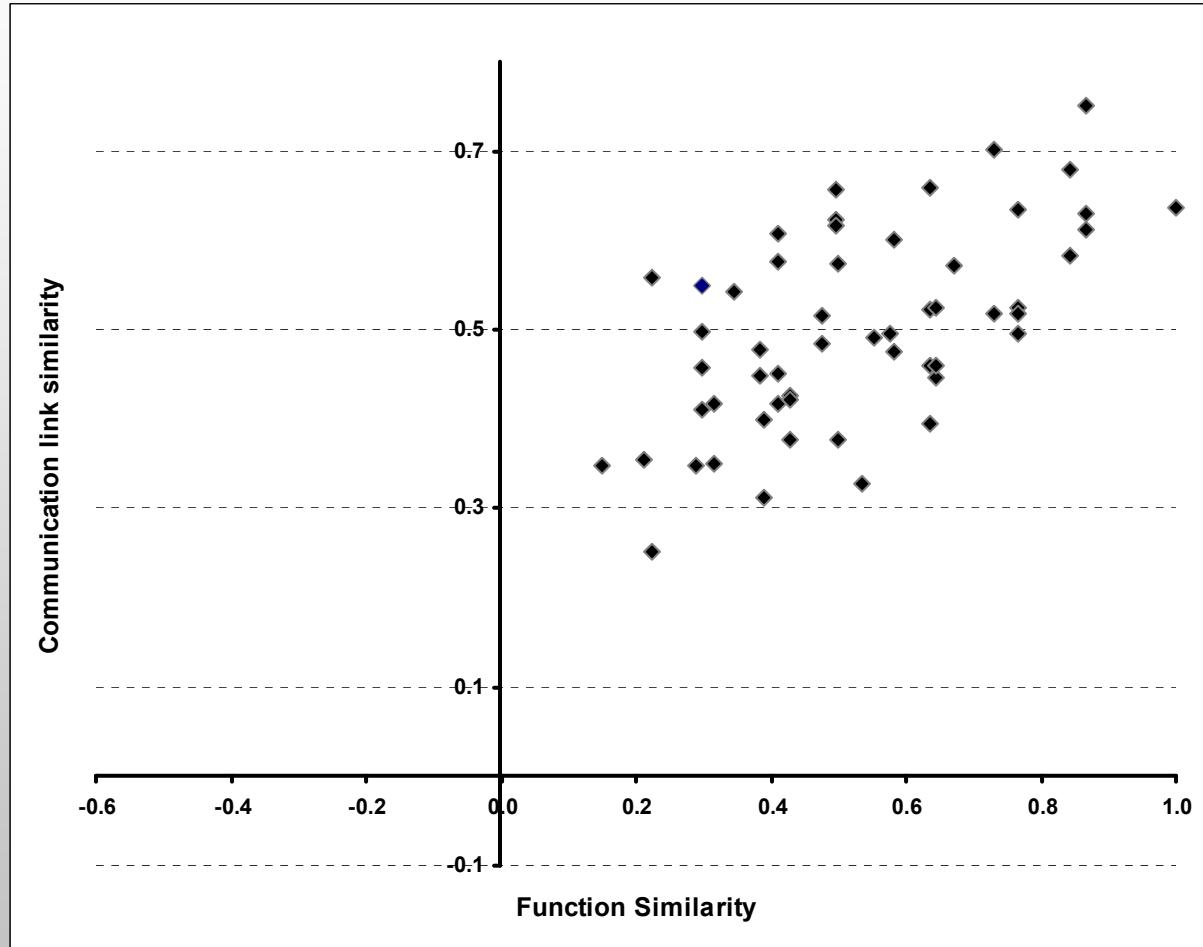# Is there a driver behind what connections are made?

- Connections in a complex project will happen

- Given two similar projects, will the connections also be similar?

- Test scenario:
  - Multiple concurrent engineering sessions
    - No barriers to communication in each session
  - Sessions have different objectives and different function areas are selected and staffed
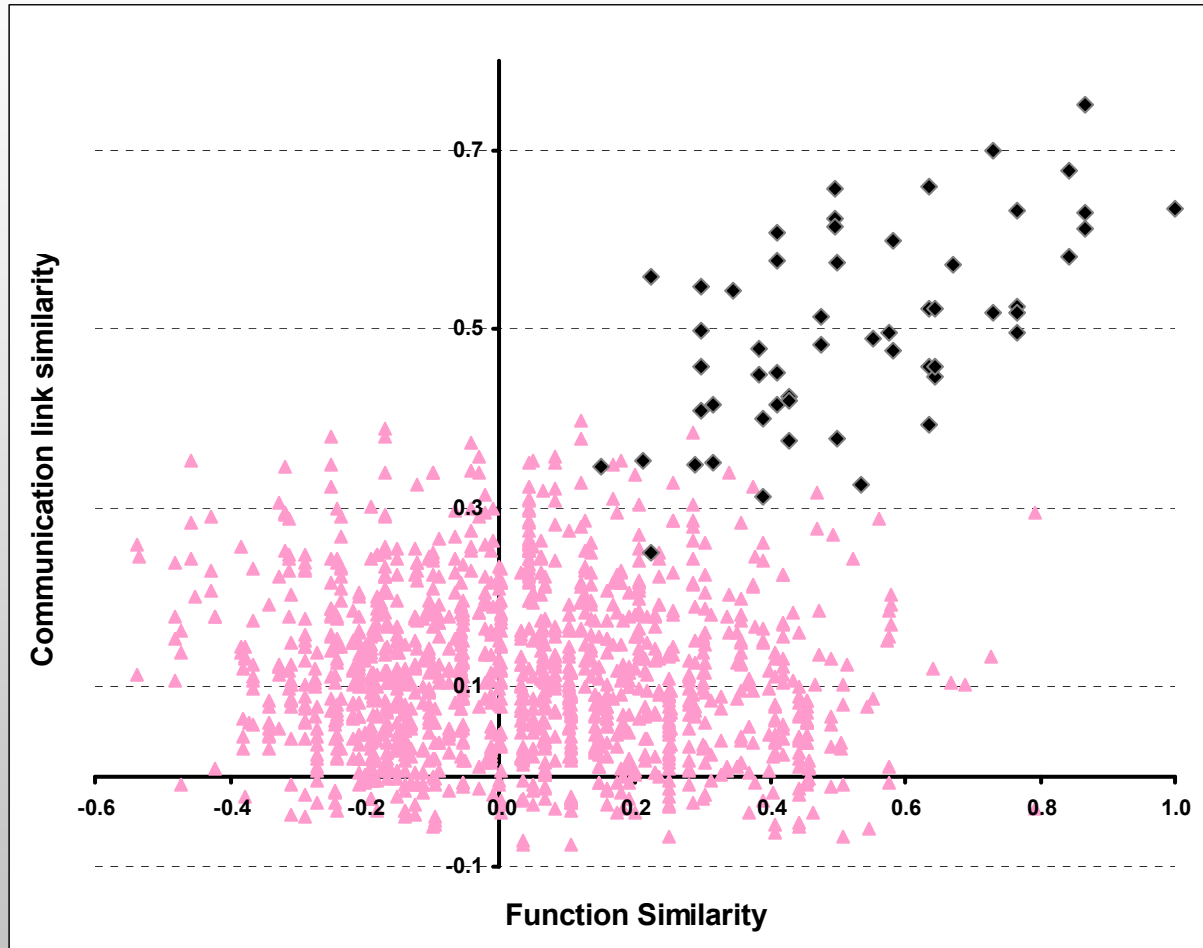
# Are connections made on purpose?



A scatter plot titled with axes "Communication link similarity" (vertical, ranging from -0.1 to 0.7) and "Function Similarity" (horizontal, ranging from -0.6 to 1.0), with a single data point plotted at approximately (0.4, 0.4).

# Are connections made on purpose?



Data for real projects provided by Mark Avnet

# Are connections made on purpose?



Data for real projects provided by Mark Avnet

# Hypotheses and questions

- Hx: Different levels of product complexity require different patterns of people communication. Number, focus and attention span vary.

- Hx: Level of oversight by a person is extremely limited (they only see what they do) when compared to the whole system.

- Hx: People's attention (as verified by their footprint) varies through time. Once something is done, they move on and do not return to it.

- Hx: Systems that operate under the freedom of participants have high redundancy communication channels.

- Hx: System critical components are verified by several people at different times.

# Finding more data

- Product complexity, component interaction and specialty interaction is context specific and varies across industries

- But, at a low-level, we can abstract to:
  - Component A <- logical interaction - > Component B

- This allows us to study the same problem in different industries and try to understand and generalize

- But to understand different behaviors, a lot of **very detailed data** from **several projects** is required
  - **Focus on one area: software**

# Why software? Fits the topic

- Software code is also made of a set of logical interactions:
    - procedures, functions, variables and objects

- Complex software is developed collaboratively by teams.
    - Each member works on a sub-part of the system that interacts
    - Members often work on code written by others
    - Teams can typically be geographically dispersed

- Software is key part in *almost* all modern complex products

# Why software? Good data for research

- Software engineering practices have excellent book-keeping methods which give us:
    - Fine grained information
    - Complete information. Long history on past projects is available
    - Uniform data over time
    - Even small projects generate large volumes of changes making it possible to detect even small effects statistically

- The data collection is nonintrusive / non-disturbing
    - doesn't require resources from project to help with the data collection

- The data collection is cheap
    - no impact on the project as this data collection is already performed

Adapted from (Mockus, Weiss et al. 2003)

- Data is ripe for processing
    - Using computer to process and analyze a lot of information

# Why software? Even better data

- The information-based nature of software products brings another benefit in that **we can track the evolution of a design over time**. (…) For a researcher, this presents an opportunity to follow the "living history" of a design, a technique that is typically not possible for physical products.

  Exploring the Structure of Complex Software Designs: An Empirical Study of Open Source and Proprietary Code (MacCormack, Rusnack, Baldwin 2005)

# Why software? Available literature

- Academic Research
  - Software development process
  - Coordination in software projects
  - Measures of software complexity
  - Visualization of software and team participation
  - Social settings, motivations and behaviors of participants
  - Social network analysis of software projects
  - System evolution
  - Case studies
- Ethnographies

TWO DOZEN PROGRAMMERS,
THREE YEARS, 4,732 BUGS,
AND ONE QUEST FOR
TRANSCENDENT
SOFTWARE

DREAMING IN CODE

SCOTT ROSENBERG
COFOUNDER OF SALON.COM

# Software Data. What it looks like

- ## Change in the code

- ## Bug correction activity

**Revision**: 4380
**Author**: morgen
**Date**: 7:11:38 PM, Tuesday, February 01, 2005
**Message**:
Added a dialog to prompt the user for new webdav account info, and fixed some typos
----
**Added** : /trunk/chandler/parcels/osaf/framework/sharing/AccountInfoPrompt.py
**Added** : /trunk/chandler/parcels/osaf/framework/sharing/AccountInfoPrompt.wdr
**Added** : /trunk/chandler/parcels/osaf/framework/sharing/AccountInfoPrompt_wdr.xrc
**Modified** : /trunk/chandler/parcels/osaf/framework/sharing/Sharing.py

| Who | When | What | Removed | Added |
|-----|------|------|---------|-------|
| pbossut@osafoundation.org | 2007-11-20 10:33:10 PST | Severity | major | norma |
| | | Component | Application | Calendar UI |
| | | Priority | P1 | P3 |
| | | Product | Chandler | Cosmo |
| | | Summary | Events belonging to several calendars are only sync'd on the calendar they were created in | Events belonging to several calendars not displayed correctly when one of the calendars is unchecked |
| | | Target Milestone | 0.7.3 | --- |
| | | Version | 0.7 | 0.8 |
| pbossut@osafoundation.org | 2007-11-20 10:34:20 PST | AssignedTo | pbossut@osafoundation.org | mde@osafoundation.org |
| aparna@osafoundation.org | 2007-11-20 11:41:10 PST | CC | | adam@osafoundation.org |
| mikea@osafoundation.org | 2007-11-30 15:02:07 PST | Target Milestone | --- | 1.0 |
| sheila@osafoundation.org | 2008-02-28 12:53:54 PST | AssignedTo | mde@osafoundation.org | travis@osafoundation.org |
| | | Target Milestone | 1.0 | Future |

# Projects of interest – Open Source

- Open Source Software is a type of software project that relies on a loose articulation between developers.

- Open source software projects are based on voluntary contributions and **involve only very light coordination** activities by a central project team

  Kogut and Metiu (2001) Open-Source Software Development and Distributed Innovation

- "What is perhaps most surprising about the process is that it **lacks many of the traditional mechanisms used to coordinate** software development, such as plans, system-level design, schedules, and defined processes."

  Mockus, Fielding et al. (2002) Two Case Studies of OSS Development: Apache and Mozilla

- "everyone, under this type of project management, is self-determining"

  Mockus and Herbsleb (2002) Why Not Improve Coordination in Distributed Software Development by Stealing Good Ideas from Open Source?

# Open Source Software

- Because of their open policies, project data is already public
    - Time to gain access to projects is cut to almost zero
    - Data is available online
    - No need to travel
    - No need to navigate NDAs

# Rewriting our Hypotheses and questions

- Hx: Different levels of product complexity require different patterns of people behavior. Number, focus and attention span vary.

- Hx: Level of oversight in code by a person is limited (the files they edit mostly reference themselves)

- Hx: Does people's attention (as verified by their footprint) varies through time. Do they come back to their older files while editing new ones?

- Hx: How well does the ensemble of perspectives cover the whole code?

- Hx: A developer engages in coordination only with those who are part of his system view

- Hx: Most time is spent on the boundary components than on independent components.

# Projects collected (so far)

| Project name | Description | Number of developers | Got code | Source code files | Lines of code | Number of functions | Number of edits in DB | Got Bugs | Number of bugs | Got ML | Total # of Msgs | Got IRC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Chandler | PIM | 44 | Yes | 5374 | 1195429 | 2560 | 14835 | Yes | 12427 | Yes | 10405 | Can get |
| Audacity | Audio editing | | Can get | | | | | No | | | | |
| Apache | Web Server | 95 | Yes | 8645 | 1382475 | 694 | 18020 | | | | | |
| Wireshark | | 40 | Yes | 12560 | 5800442 | 7374 | 26421 | Can get | | | | |
| MythTV | Media center | 41 | Yes | 20792 | 3941132 | 18091 | 17221 | Can get | | | | |
| Rsync | Backup | | Yes | 382 | 73392 | 218 | | Can get | | Can get | | |
| Git | Repository mgmt | | Yes | 1796 | 322768 | 1042 | | | | Can get | | |
| Android | Mobile OS | | Very big | | | | | | | | | |
| GNUmeric | Spreadsheet | 217 | Yes | 6583 | 4818253 | 2366 | 15827 | | | | | |
| Gimp | Image processing | 260 | Yes | 14455 | 6648419 | 7610 | 24772 | | | | | |
| Juice | Podcast | | No code | | | | | | | | | |
| Songbird | | 30 | Yes | 14433 | 1924184 | 1626 | 7954 | | | | | |

825 859 records in database so far

# Example of data from an open source project

## Chandler

- Tracked from Aug'02 to Nov'08
- 43 developers
- Source Code:
  - 14 835 changes (commits)
  - 5 347 files
  - 2 560 functions, 1 195 429 loc
- Bugs:
  - 12 427
- Mailing list:
  - 10 405 emails exchanged
- Internet chat:
  - ?

# Data analysis

- Demographics of projects

- Analysis of **individuals**
  - Footprint – what parts of the product are focused
  - Change over time

- Analysis of **team**
  - Communication network
  - Visibility overlap
  - System hand-offs

- Analysis of **product**
  - Function call graph
  - Bug duration, origin, severity

- Analysis of **product, team**
  - Overlap in communication and objects

# Some results

- Scale of review and rework
  – How many times a file is edited

- System visibility
  – How much of the product does each one see

- System overlap
  – Who worked on whose files

- Evolution of personal footprint

# Scale of review and rework

**Histogram of number of edits on a file**



Average 9 edits/file

# System visibility

## How much do they work on?

Top ten developers and average on each project

| Chandler | Apache | Gimp | Gnumeric | MythTV | Wireshark | Songbird |
|----------|--------|------|----------|--------|-----------|----------|
| 36% | 67% | 61% | 71% | 35% | 44% | 39% |
| 32% | 35% | 57% | 29% | 15% | 36% | 25% |
| 29% | 22% | 31% | 14% | 15% | 34% | 23% |
| 24% | 21% | 18% | 13% | 14% | 32% | 22% |
| 14% | 21% | 8% | 12% | 13% | 25% | 20% |
| 12% | 20% | 8% | 10% | 12% | 24% | 13% |
| 5% | 14% | 7% | 9% | 10% | 23% | 12% |
| 3% | 13% | 7% | 8% | 10% | 13% | 8% |
| 3% | 12% | 7% | 7% | 9% | 12% | 6% |
| 3% | 12% | 5% | 7% | 8% | 12% | 4% |
| 4.0% | 4.6% | 1.0% | 1.2% | 4.7% | 8.1% | 6.5% |

% of files in the project that have at least one edit by the member

# System overlap

## Who worked on whose files (1st order)



**Chandler: 43 developers**

# System overlap

## Who worked on whose files (1st order)



**Apache: 94 developers**

# System overlap

## Who worked on whose files (entire project)
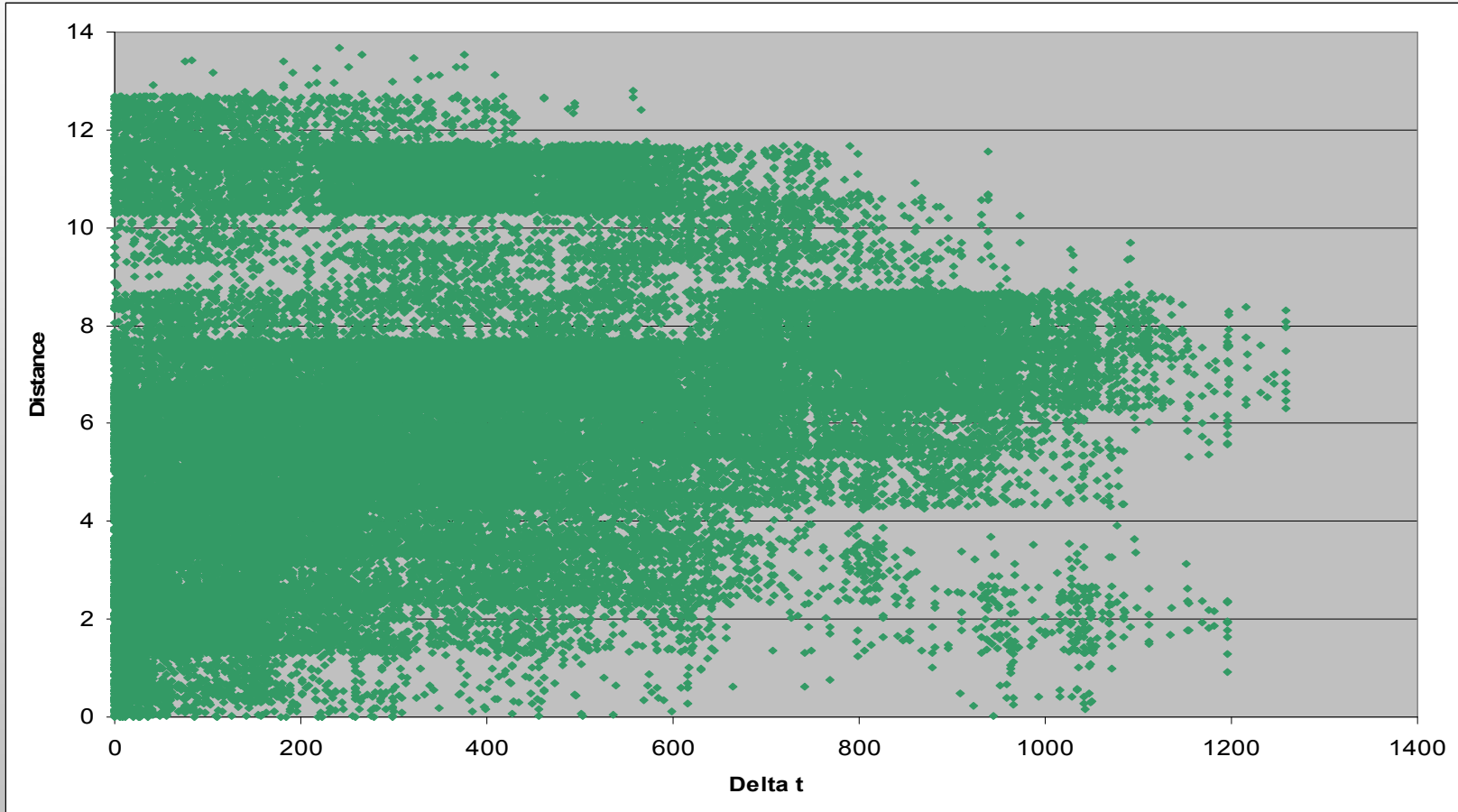
# Evolution of personal footprint

## Δt v Distance
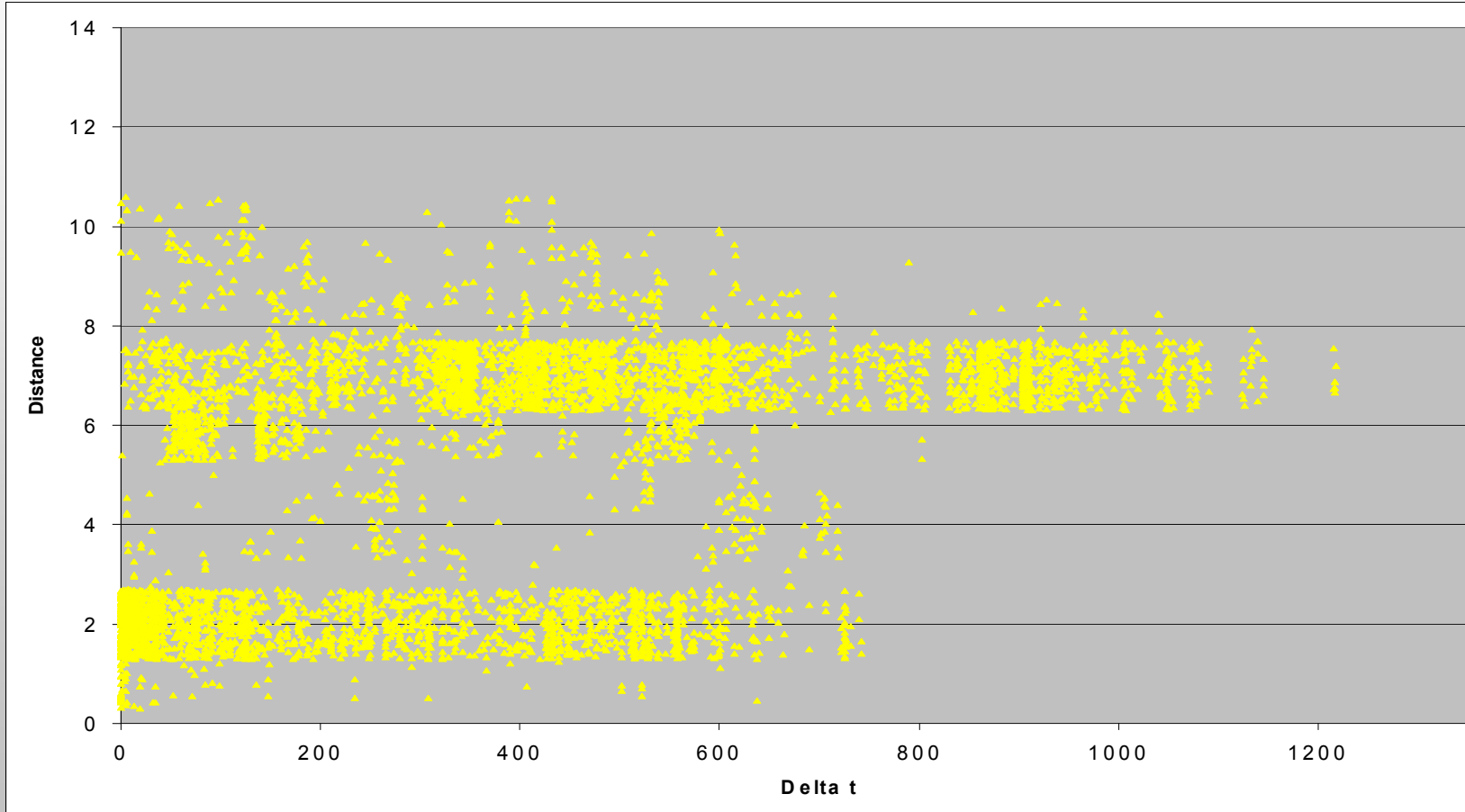
# Evolution of personal footprint

## Δt v Distance

# Evolution of personal footprint

## Δt v Distance

# Next steps

- Integrate in analysis the *bug* and *mailing list* data

- Analyze relationship between
    - Product and communication
    - Analyze problem solving over time and product structure

- Analyze using FCG instead of folder hierarchy for product structure

# Thank you