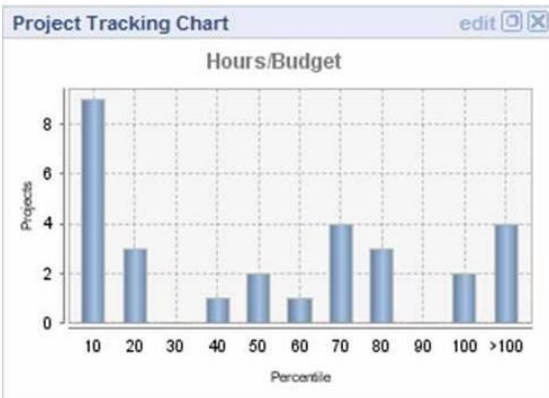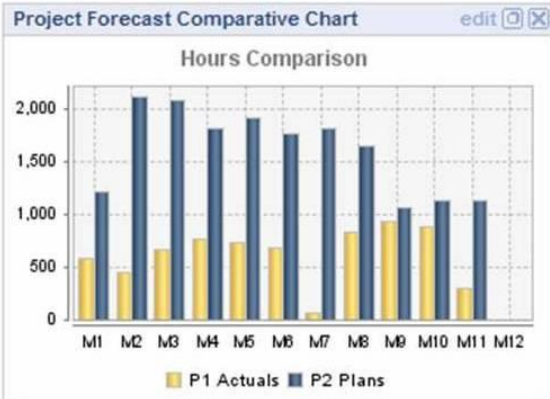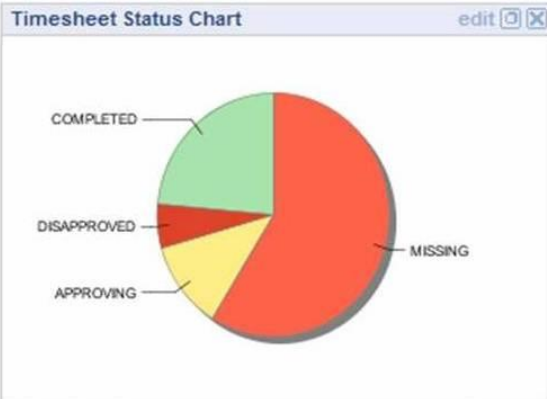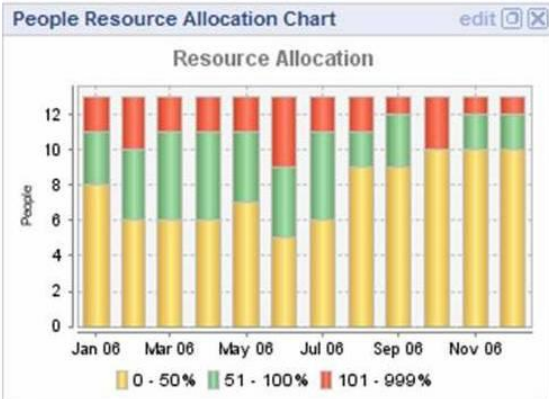# Validating the Use of Topic Models for Software Evolution
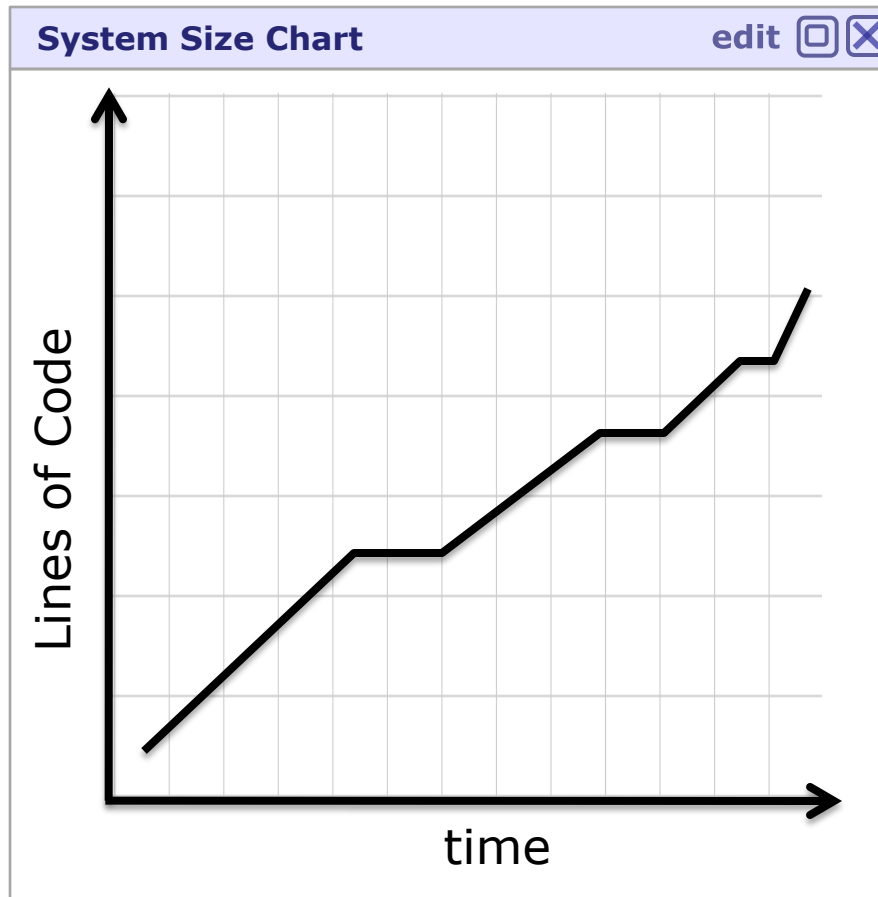
**Stephen W. Thomas**, **Bram Adams,**
**Ahmed E. Hassan, and Dorothea Blostein**

**Software Analysis and Intelligence Lab (SAIL)**
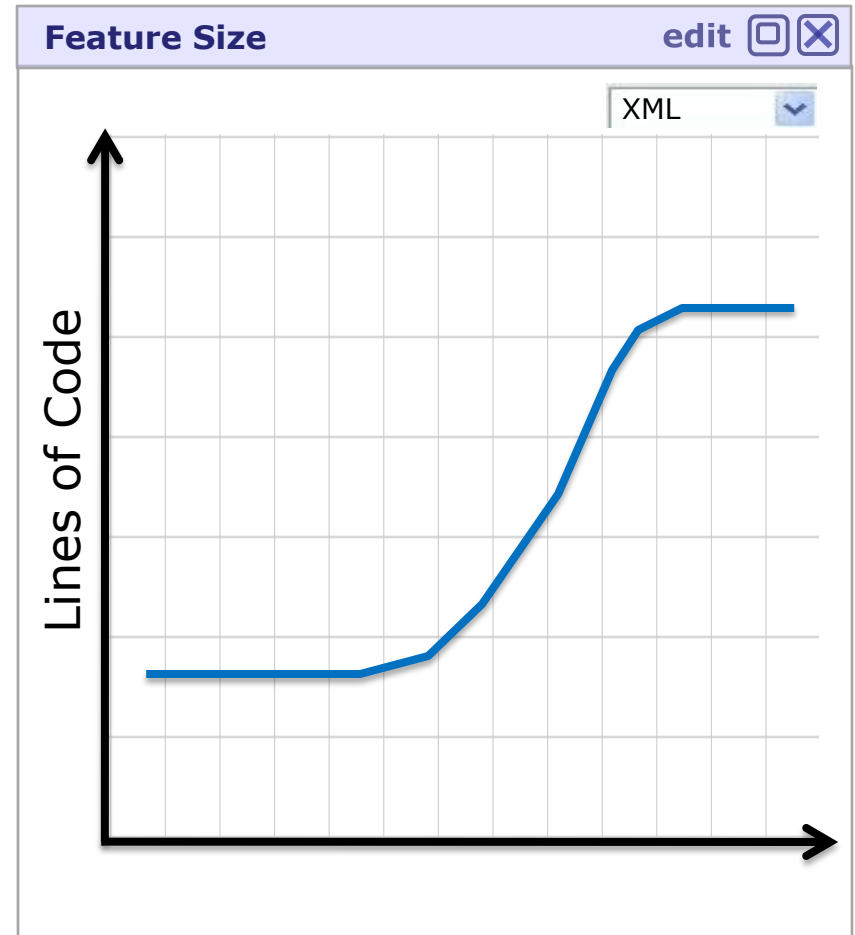**School of Computing, Queen's University, Canada**

**People Resource Allocation Chart** edit □ ☒
Resource Allocation
0 - 50% | 51 - 100% | 101 - 999%

**Timesheet Status Chart** edit □ ☒
COMPLETED
DISAPPROVED
APPROVING
MISSING

**Project Forecast Comparative Chart** edit □ ☒
Hours Comparison
P1 Actuals | P2 Plans

**Margin Chart** edit □ ☒
Margin
Cost | Bill | Cumulative Margin

**Project Performance Chart** edit □ ☒
Hours Performance
Actuals | Budgets | Plans

**Earned Value Chart** edit □ ☒
BAC $53,501
Actuals Through
BCWP | ACWP | ETC | BCWS

**Project Tracking Chart** edit □ ☒
Hours/Budget
Percentile

**People Performance Chart** edit □ ☒
0 - 50% Booked
101 - 999% Booked
51 - 100% Booked

**People Forecast Chart** edit □ ☒
Hours Forecast
Actuals | Budgets

2

# Current State of Practice



**System Size Chart**                    edit ▢ ✕

Lines of Code

time

# Our Vision

Steve T. ✔

GUI

File IO

XML

XML ✔

Lines of Code

# Our Vision

## Feature Scatter

| Feature | Scattering |
|---|---|
| GUI | +23% |
| SQL | +21% |
| User options | +12% |
| … | |
| File IO | +0% |
| XML | -5% |

## Feature Changes

7.1.0

| Feature | Change |
|---|---|
| Testing | +123% |
| Undo | +87% |
| Menu options | +12% |
| … | |
| Display | +3% |
| Java Swing | -500% |

*Detected by*

*Used in*

# Software Changes
*(bug fixes, refactorings, new features)*

# Project Dashboards

6

# How to **automatically** detect software evolution?

## Hypothesis:

- *Topic evolution models*, borrowed from the text-mining domain, can automatically detect software evolution
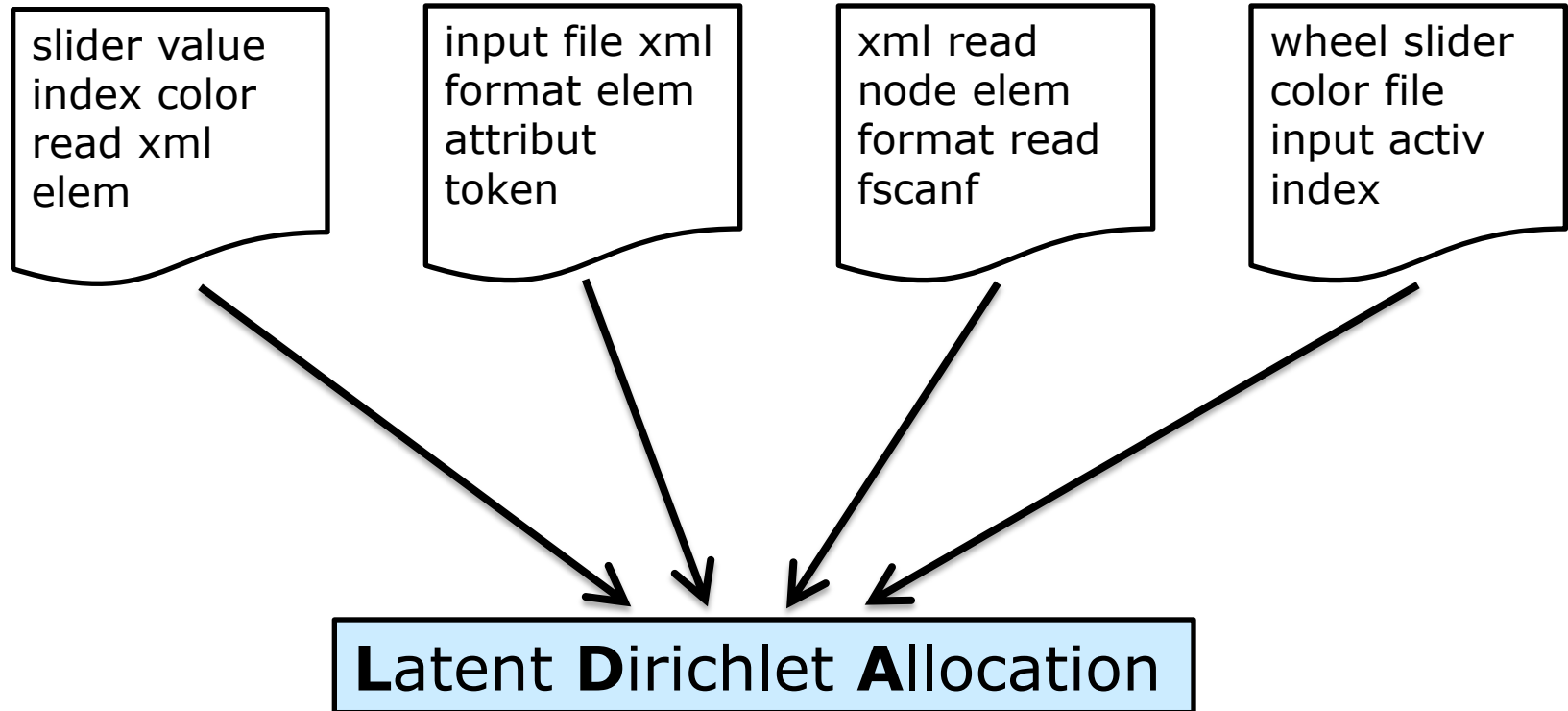
## This paper:

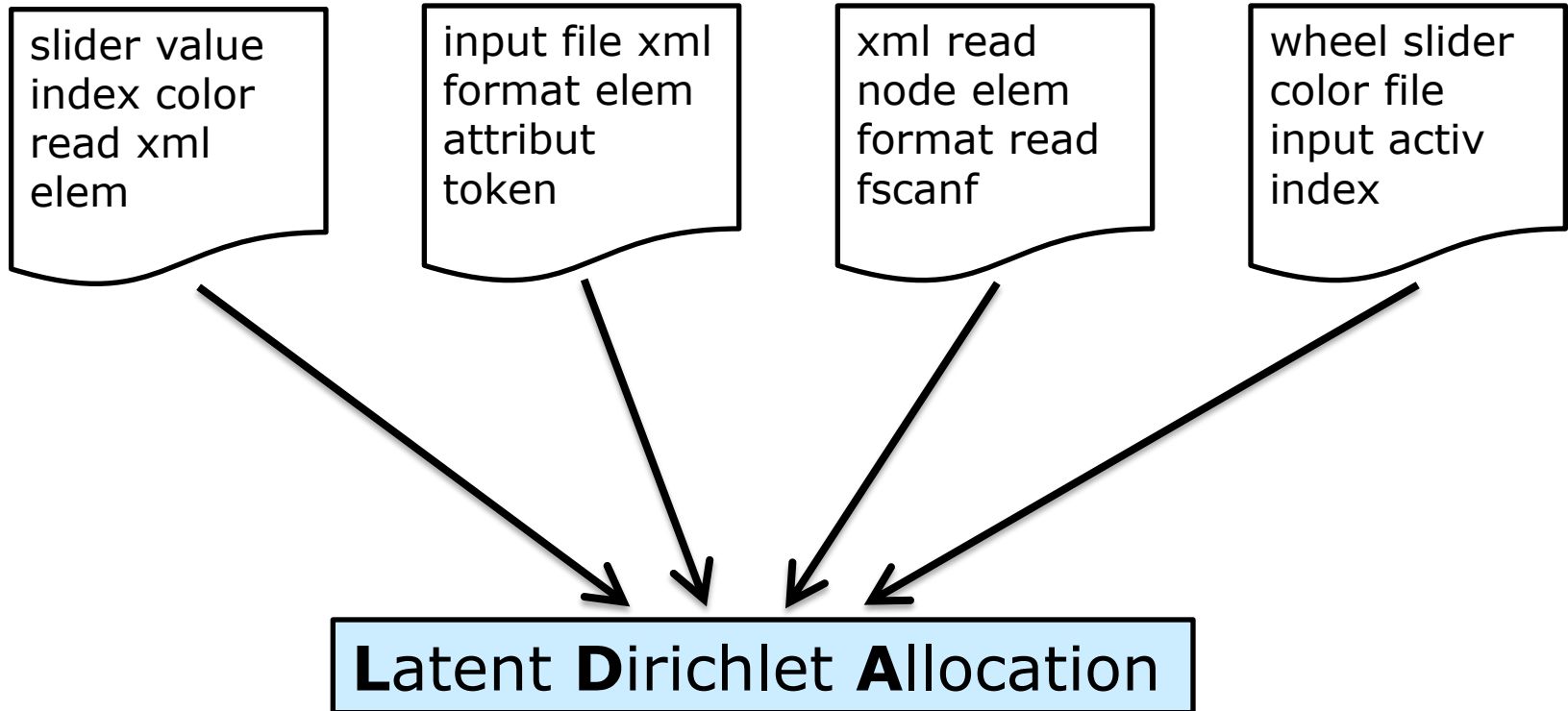- Initial qualitative study to validate hypothesis

# Topic Models

Latent Dirichlet Allocation

# Topic Models

slider value
index color
read xml
elem

input file xml
format elem
attribut
token

xml read
node elem
format read
fscanf

wheel slider
color file
input activ
index

**L**atent **D**irichlet **A**llocation

# Topic Models



slider value index color read xml elem

input file xml format elem attribut token

xml read node elem format read fscanf

wheel slider color file input activ index

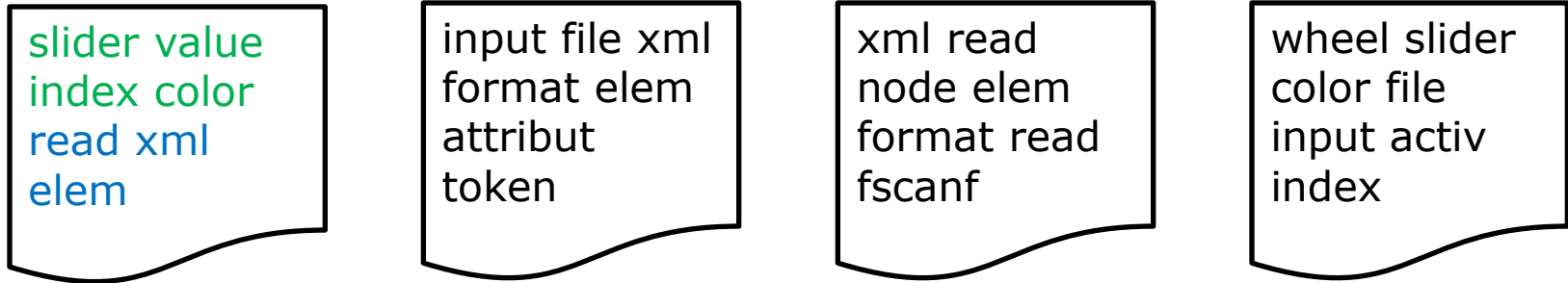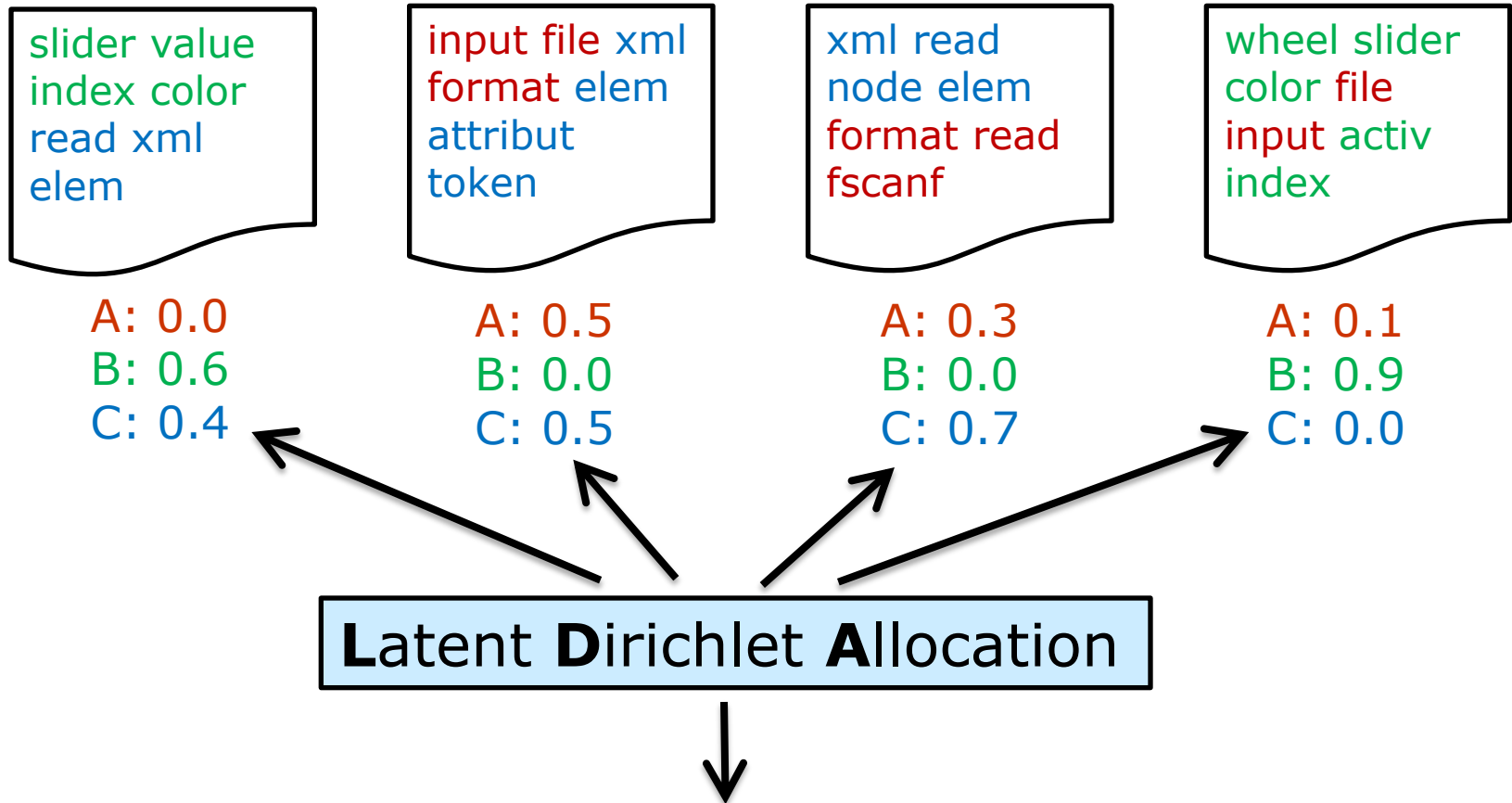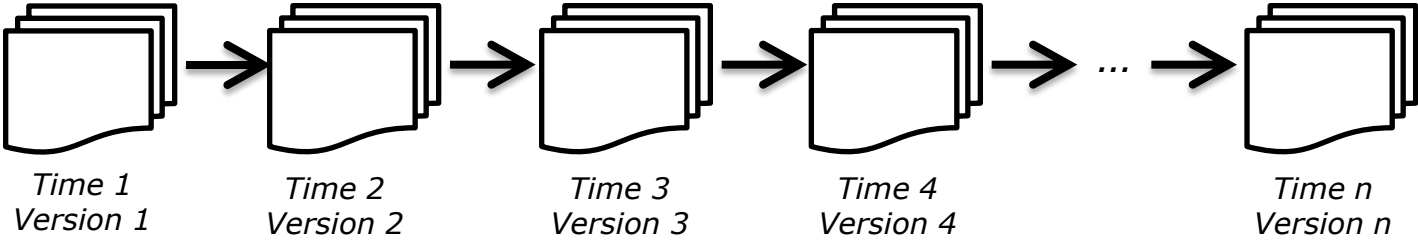**L**atent **D**irichlet **A**llocation

TOPIC A: *{file input read fscanf format ...}*
TOPIC B: *{color index slider wheel ...}*
TOPIC C: *{elem attribut read token xml ...}*

# Topic Models

slider value
index color
read xml
elem

input file xml
format elem
attribut
token

xml read
node elem
format read
fscanf

wheel slider
color file
input activ
index

A: 0.0
B: 0.6
C: 0.4

**L**atent **D**irichlet **A**llocation

TOPIC A: *{file input read fscanf format …}*

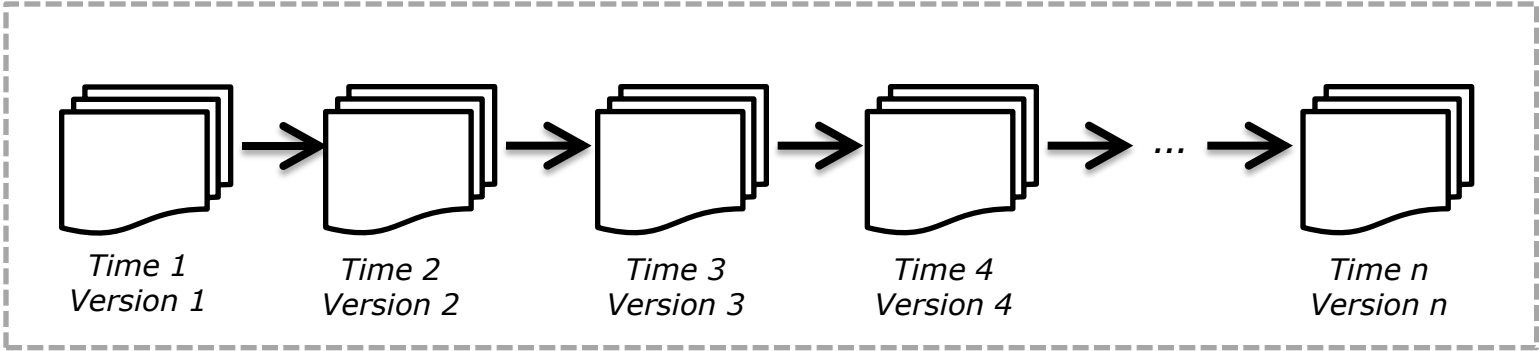TOPIC B: *{color index slider wheel …}*
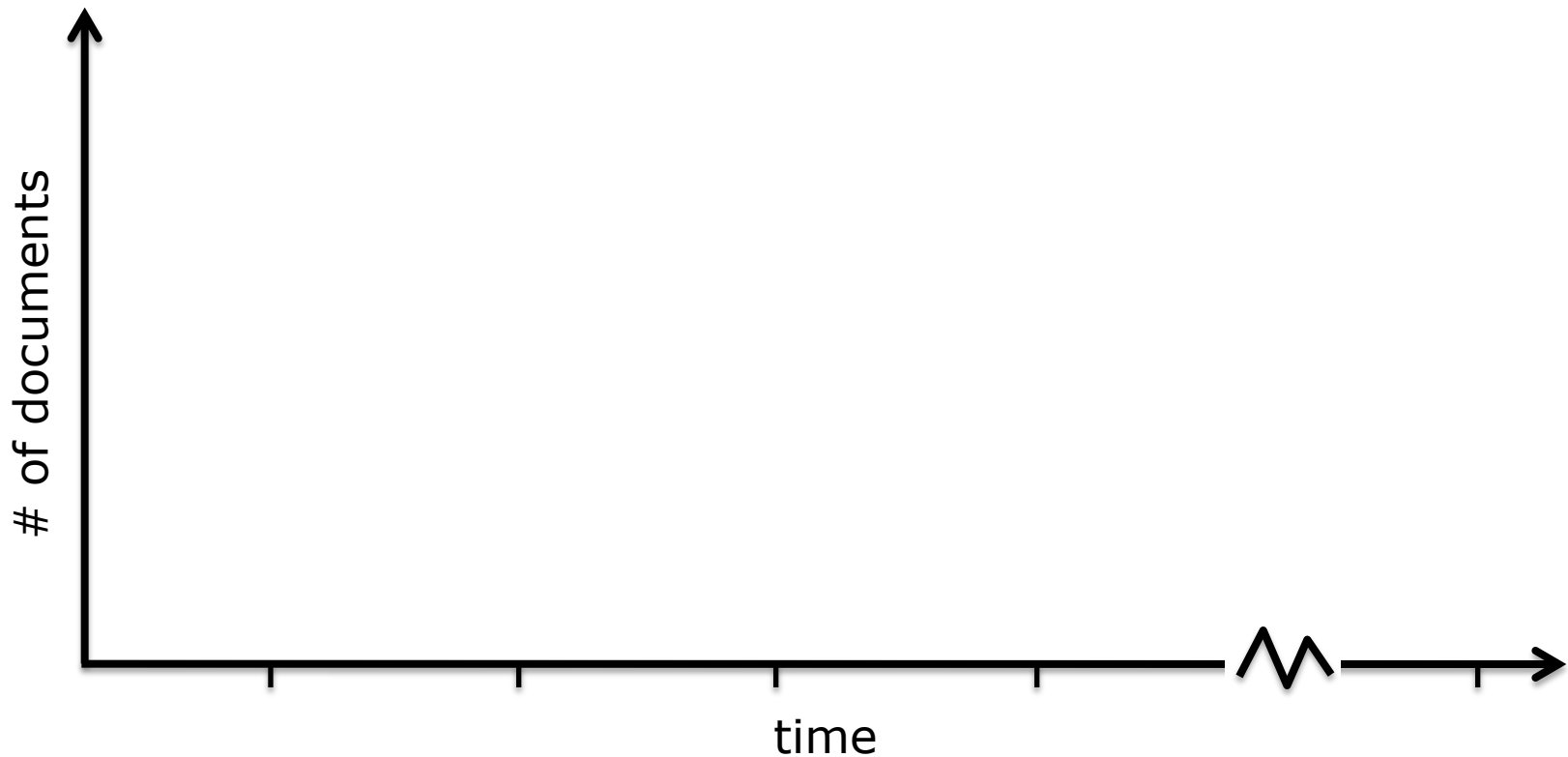
TOPIC C: *{elem attribut read token xml …}*

# Topic Models

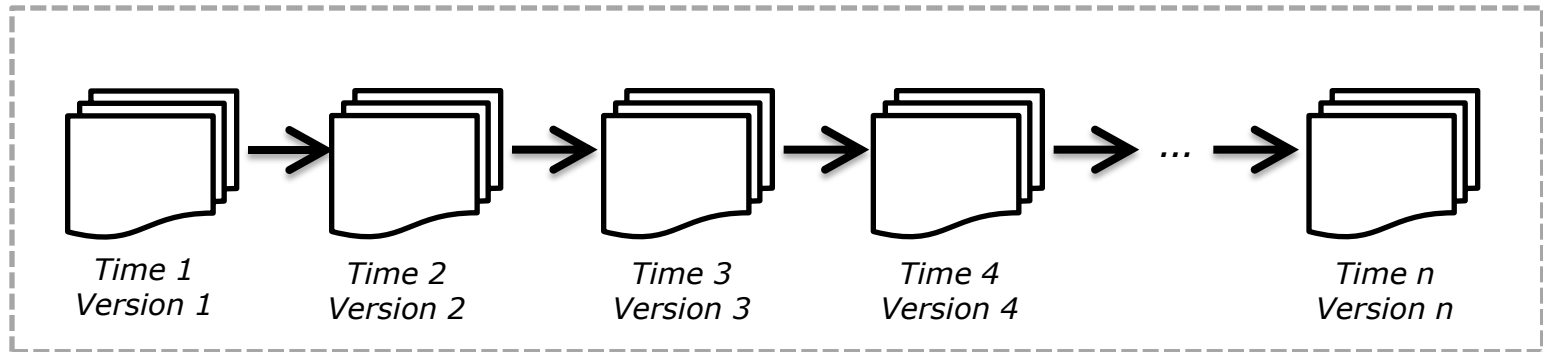

slider value
index color
read xml
elem

A: 0.0
B: 0.6
C: 0.4

input file xml
format elem
attribut
token

A: 0.5
B: 0.0
C: 0.5

xml read
node elem
format read
fscanf

A: 0.3
B: 0.0
C: 0.7

wheel slider
color file
input activ
index

A: 0.1
B: 0.9
C: 0.0

**L**atent **D**irichlet **A**llocation

TOPIC A: *{file input read fscanf format ...}*
TOPIC B: *{color index slider wheel ...}*
TOPIC C: *{elem attribut read token xml ...}*
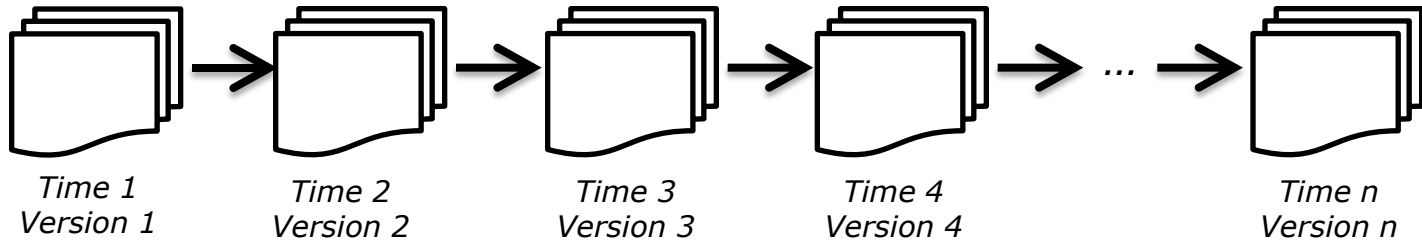
# Topic Evolution Models



Time 1
Version 1

Time 2
Version 2

Time 3
Version 3

Time 4
Version 4

Time n
Version n

# Topic Evolution Models



Time 1
Version 1

Time 2
Version 2

Time 3
Version 3

Time 4
Version 4

...

Time n
Version n

# Topic Evolution Models



Time 1
Version 1

Time 2
Version 2

Time 3
Version 3

Time 4
Version 4

Time n
Version n

# of documents

time

# Topic Evolution Models



Time 1
Version 1

Time 2
Version 2

Time 3
Version 3

Time 4
Version 4

Time n
Version n

# of documents

time

# Topic Evolution Models



Time 1
Version 1

Time 2
Version 2

Time 3
Version 3

Time 4
Version 4

Time n
Version n

# of documents

time

# Topic Evolution Models

Time 1
Version 1

Time 2
Version 2

Time 3
Version 3

Time 4
Version 4

Time n
Version n

# of documents

**XML**

time

# Topic Evolution Models



Time 1 Version 1  Time 2 Version 2  Time 3 Version 3  Time 4 Version 4  Time n Version n

**XML**

**SLIDERS**

# of documents

time

# Topic Evolution Models

## Program Comprehension

- When was the XML feature added?
- At what point did we switch to Java Swing?
- What features were affected by release 7.3?

## Quality Assurance

- What features are becoming too big?
- What features are becoming too scattered?
- What is being worked on right now?

# Topic Models on Code

```
// Add a one if positive
if (sliderValue > 0){
  sliderValue++;
```

# Topic Models on Code

```
// Add a one if positive
if (sliderValue > 0){
  sliderValue++;
```

Parse

Add a one if positive
sliderValue
sliderValue

# Topic Models on Code

```
// Add a one if positive
if (sliderValue > 0){
    sliderValue++;
```
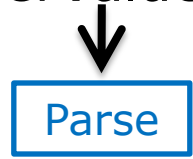
Parse
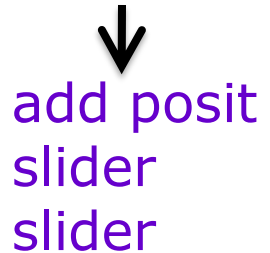
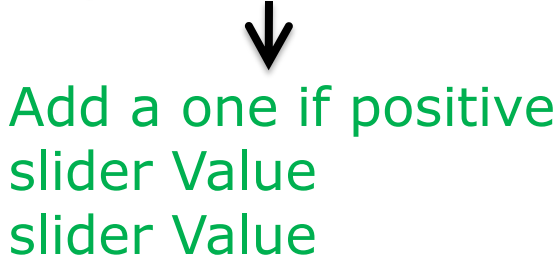Add a one if positive
sliderValue
sliderValue

Tokenize

Add a one if positive
slider Value
slider Value

23

# Topic Models on Code

```
// Add a one if positive
if (sliderValue > 0){
    sliderValue++;
```

Parse

Add a one if positive
sliderValue
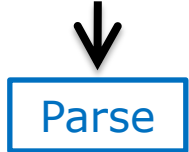sliderValue

Tokenize

Add a one if positive
slider Value
slider Value

Stop

Add positive
slider
slider

# Topic Models on Code

```
// Add a one if positive
if (sliderValue > 0){
    sliderValue++;
```

Parse

Add a one if positive
sliderValue
sliderValue

Tokenize

Add a one if positive
slider Value
slider Value

Stop

Add positive
slider
slider

Stem

add posit
slider
slider

# What causes topics to evolve in source code?

**Developer *change activities*?**
- Corrective evolution (bug fixes)
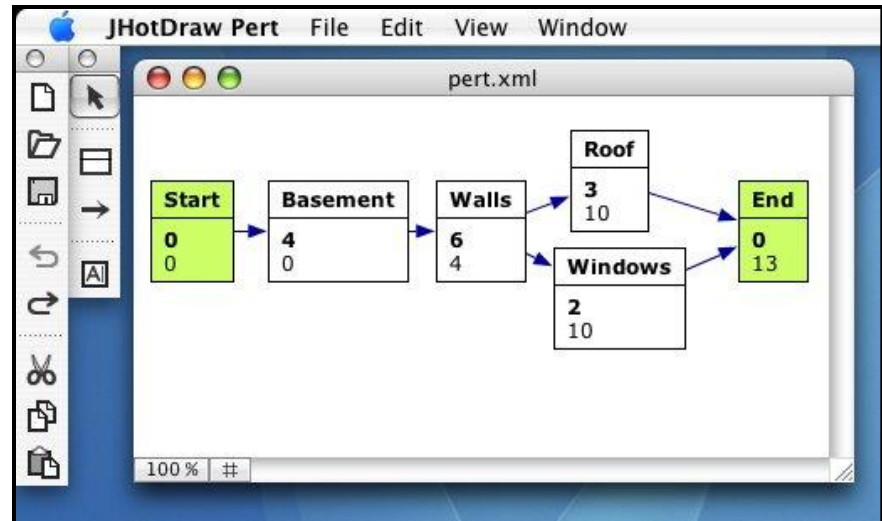- Refactoring (new libraries, structure)
- New functionalities or features

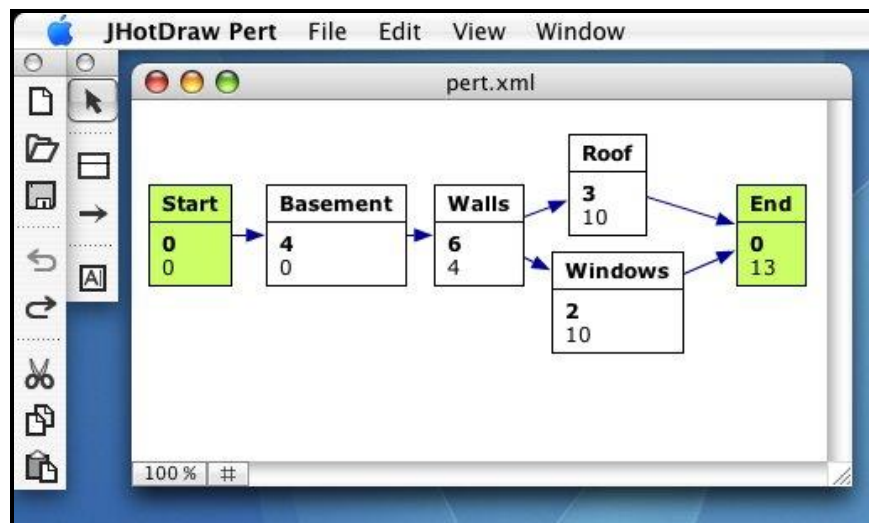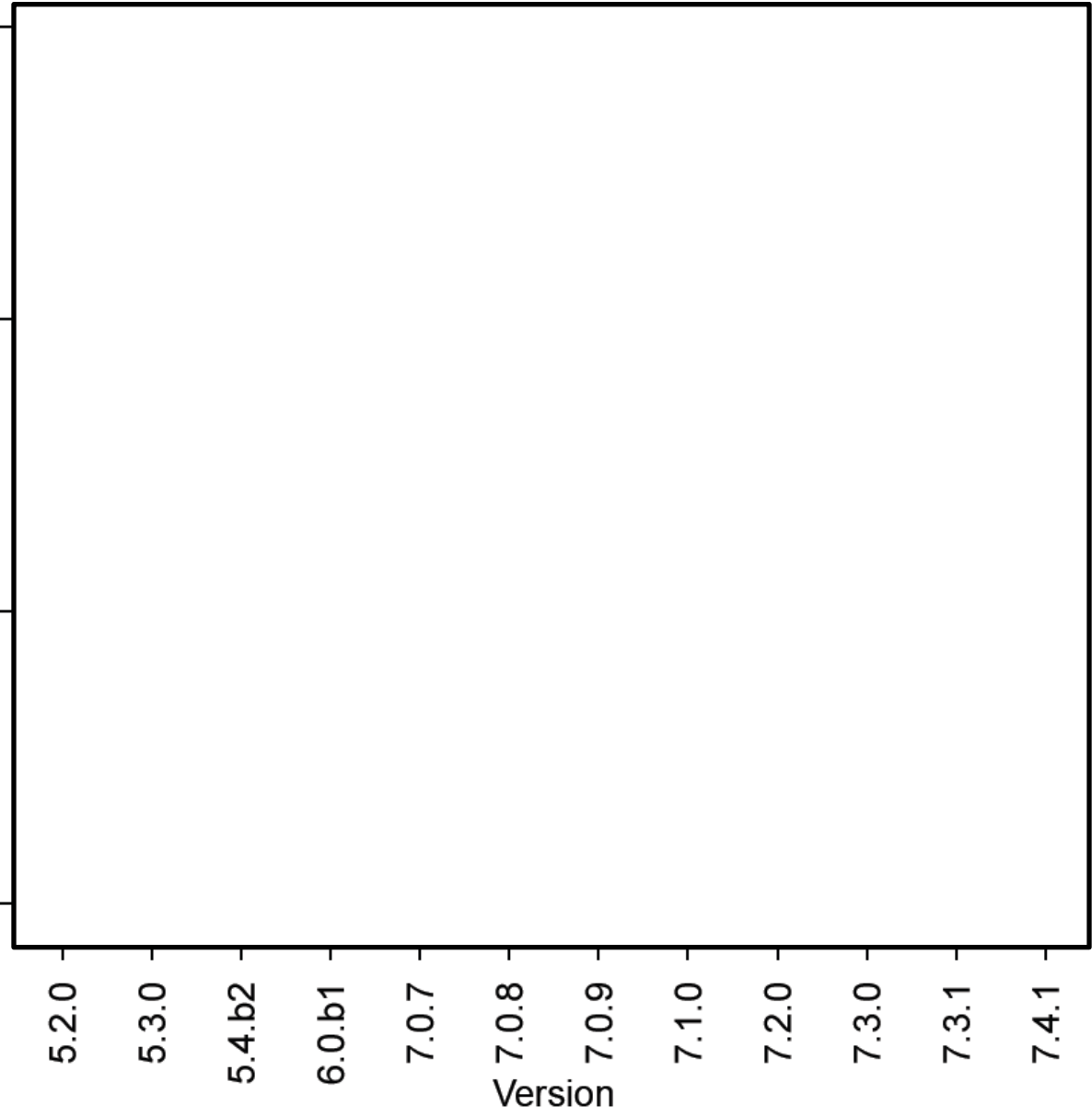**Something else?**
- Noise/error in model
- ???

# Case Study: JHotDraw

- **Medium-sized, open source drawing framework**
  - Well studied and documented
  - Manageable size
- **12 release versions**
  - 9 years
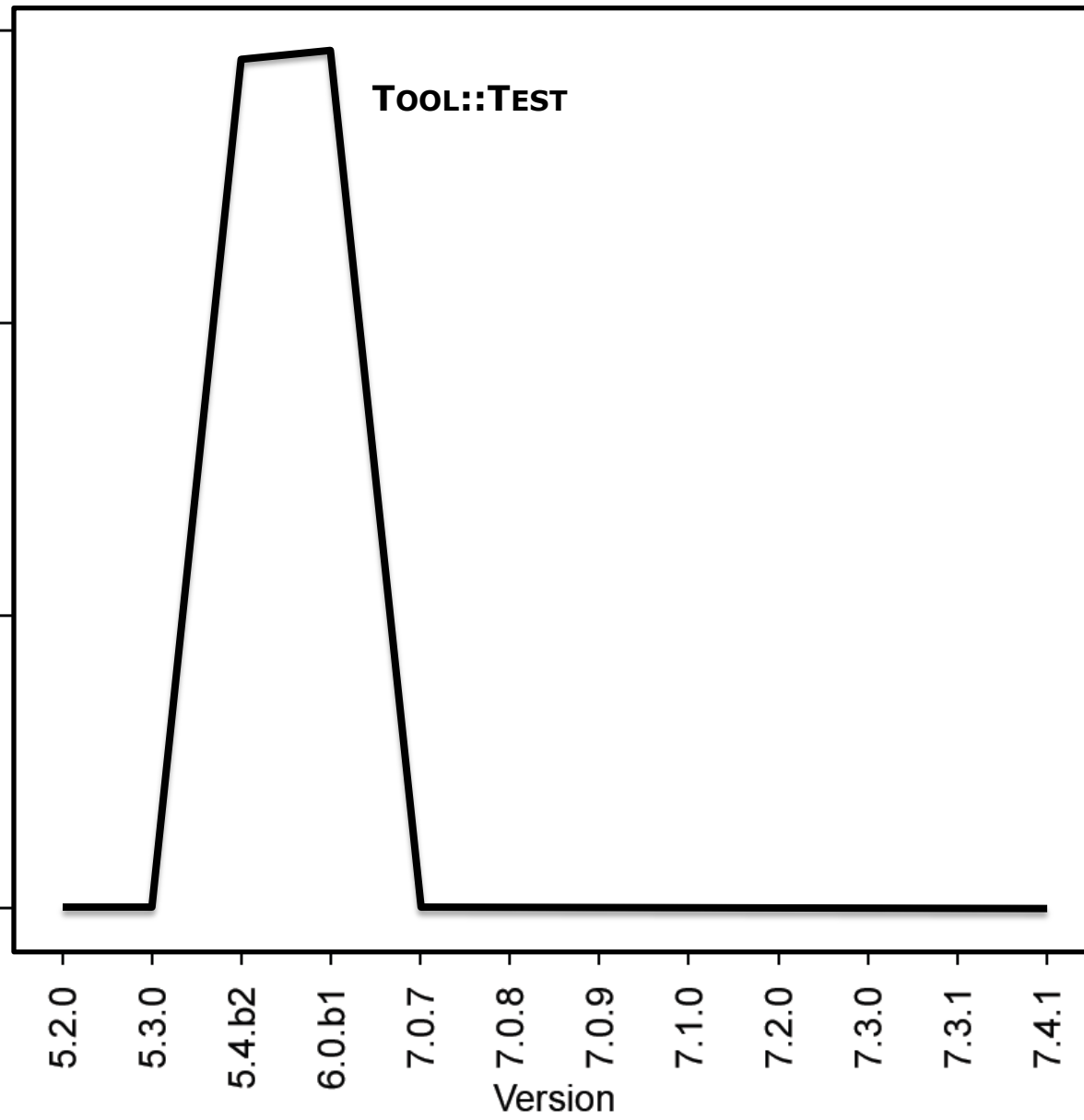  - 17 KLOC → 127 KLOC
- **Used MALLET tool**
  - 45 LDA topics

# Case Study: JHotDraw

- **Medium-sized, open source drawing framework**
  - Well studied and documented
  - Manageable size
- **12 release versions**
  - 9 years
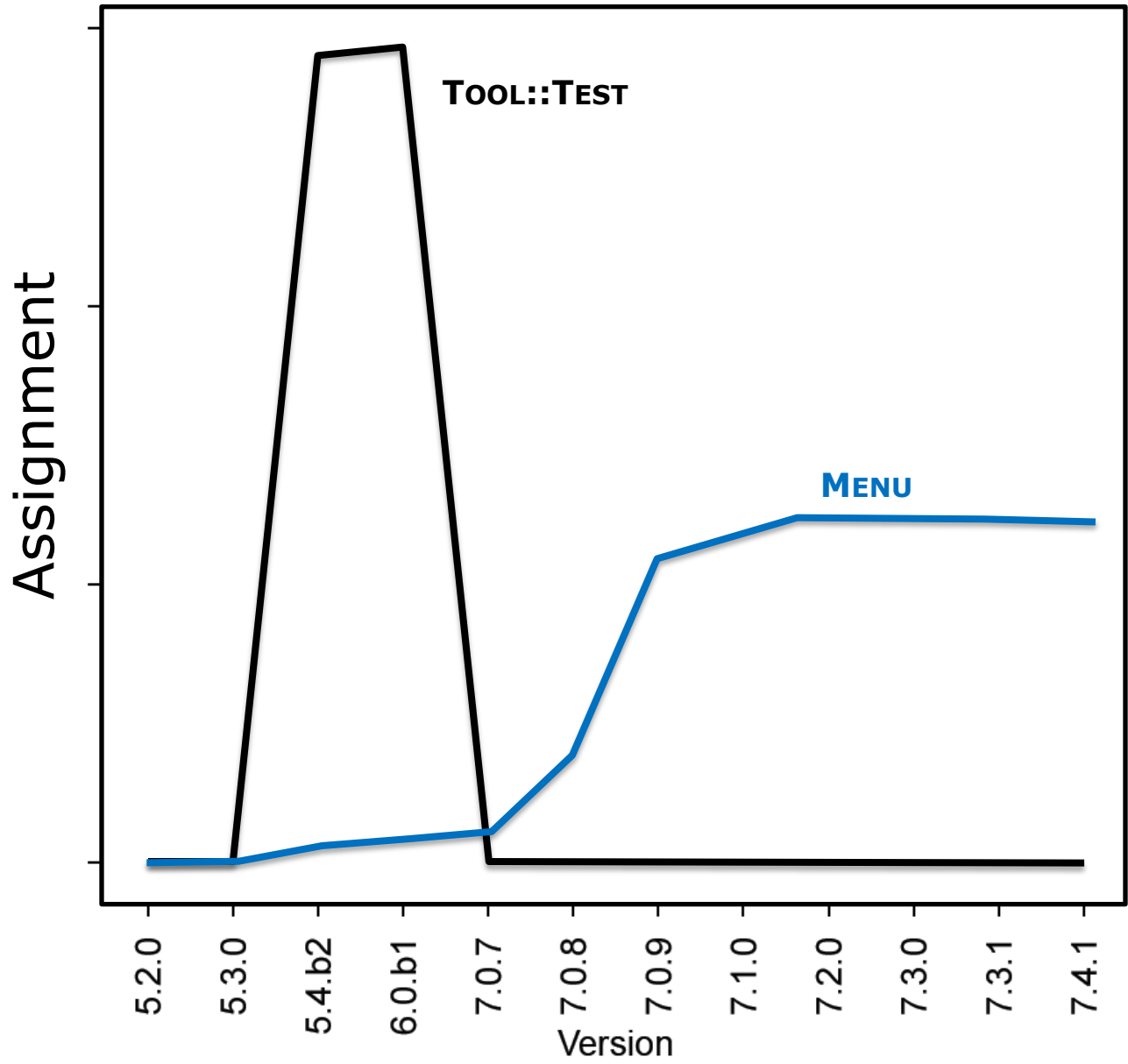  - 17 KLOC → 127 KLOC
- **Used MALLET tool**
  - 45 LDA topics

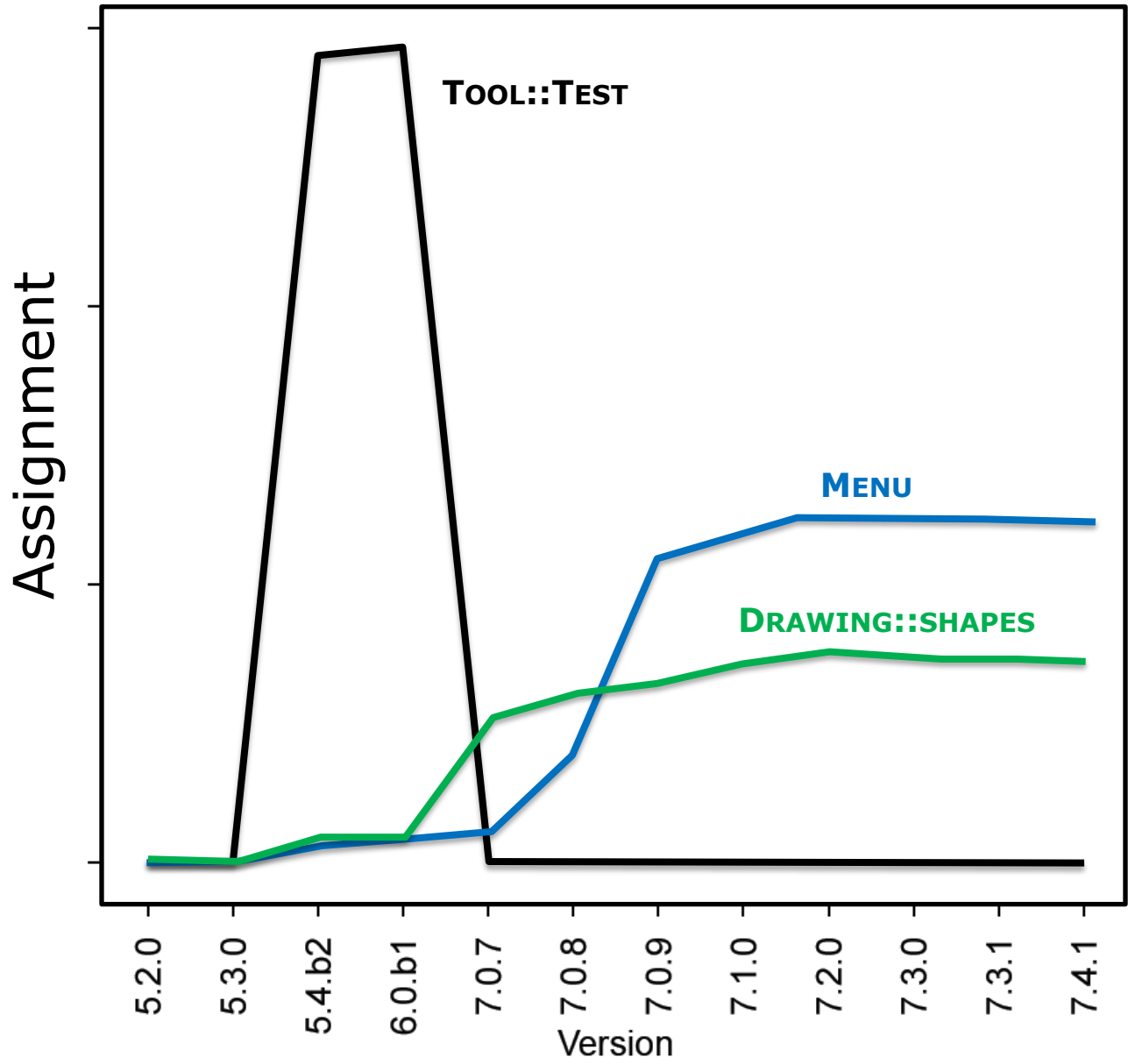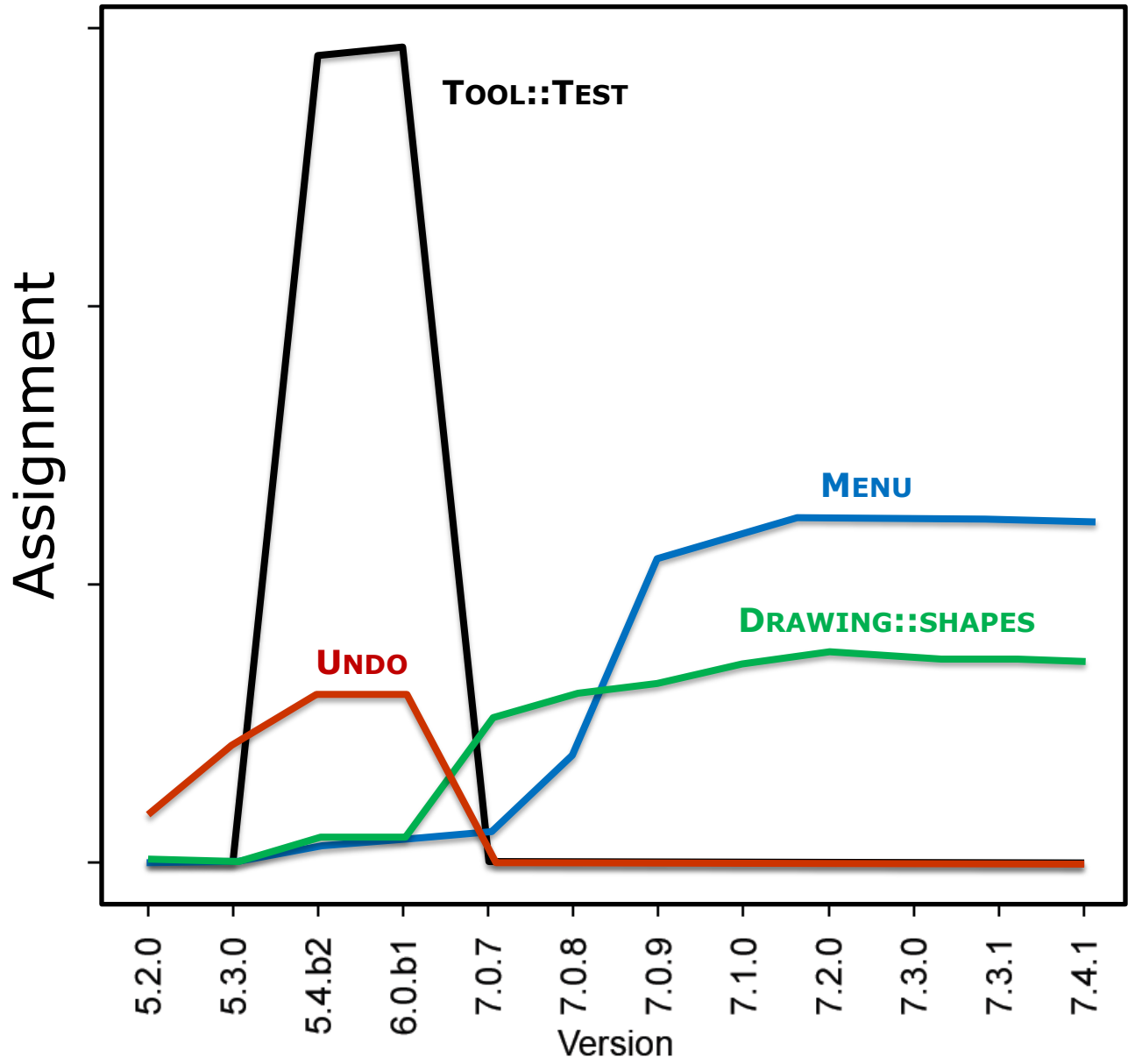Are topic evolutions consistent with the code changes?

Assignment vs Version

Version: 5.2.0, 5.3.0, 5.4.b2, 6.0.b1, 7.0.7, 7.0.8, 7.0.9, 7.1.0, 7.2.0, 7.3.0, 7.3.1, 7.4.1

TOOL::TEST

Assignment

Version

5.2.0 5.3.0 5.4.b2 6.0.b1 7.0.7 7.0.8 7.0.9 7.1.0 7.2.0 7.3.0 7.3.1 7.4.1

TOOL::TEST

MENU

DRAWING::SHAPES

Assignment

Version

5.2.0 5.3.0 5.4.b2 6.0.b1 7.0.7 7.0.8 7.0.9 7.1.0 7.2.0 7.3.0 7.3.1 7.4.1

Assignment

TOOL::TEST

MENU

DRAWING::SHAPES

UNDO

Version

5.2.0  5.3.0  5.4.b2  6.0.b1  7.0.7  7.0.8  7.0.9  7.1.0  7.2.0  7.3.0  7.3.1  7.4.1

# Are topic evolutions consistent with the code changes?

# Are topic evolutions consistent with the code changes?
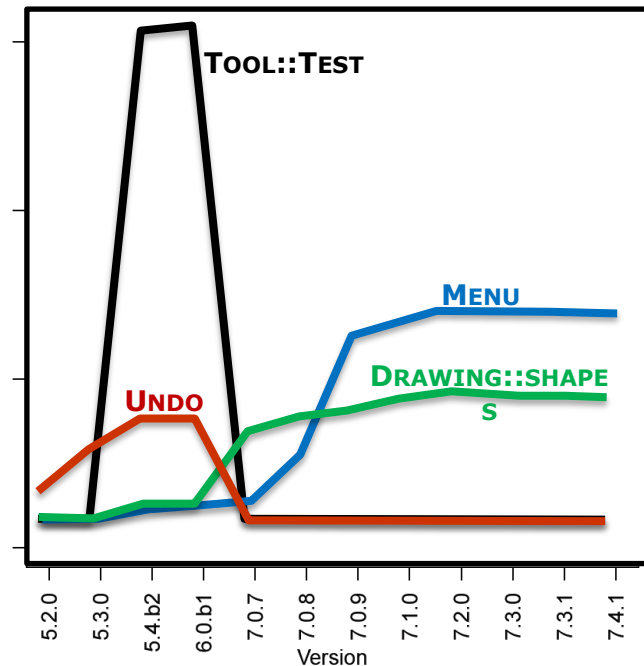
**495 change events**

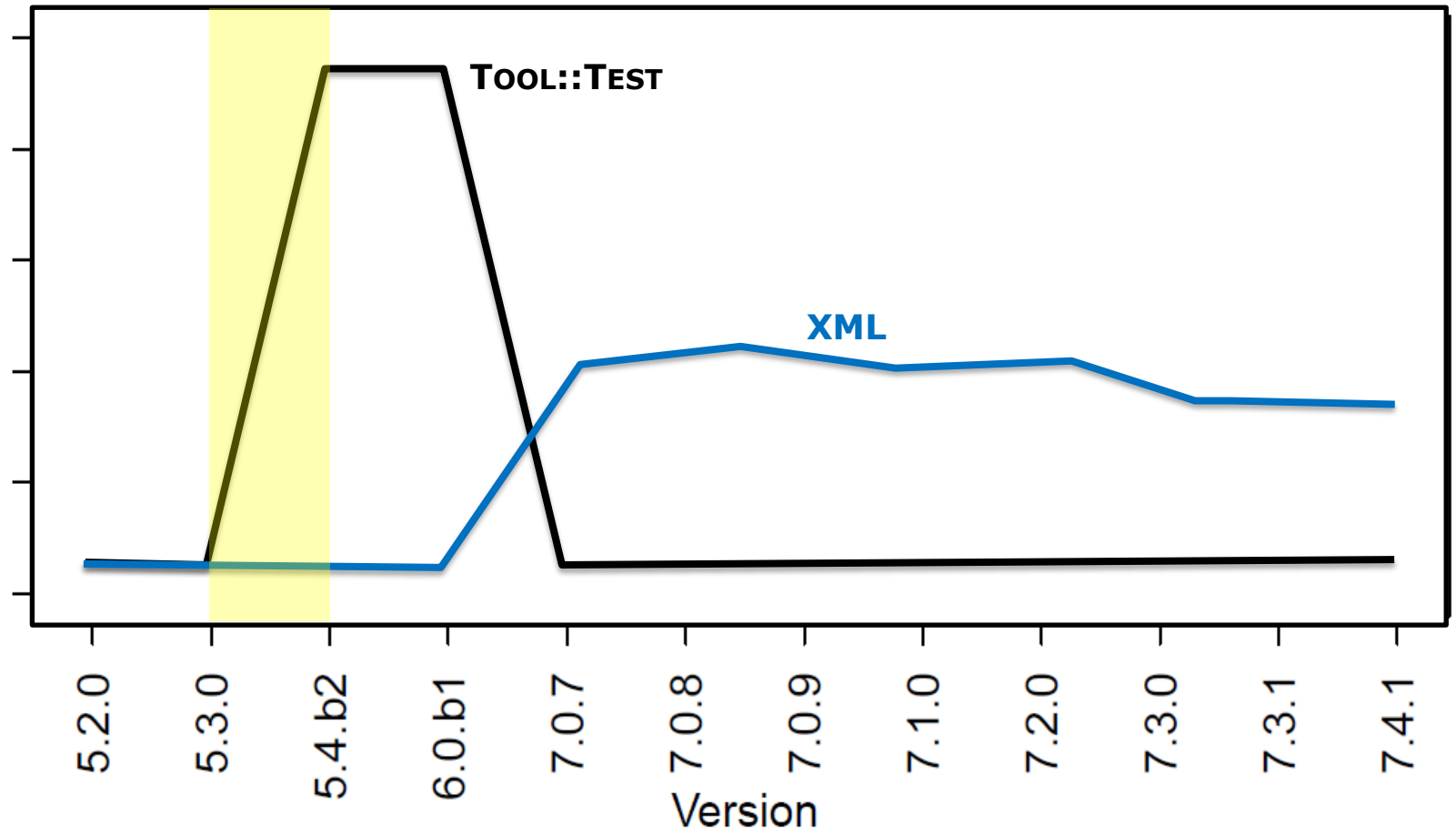*"24% increase in XML topic at version 7.1.0"*
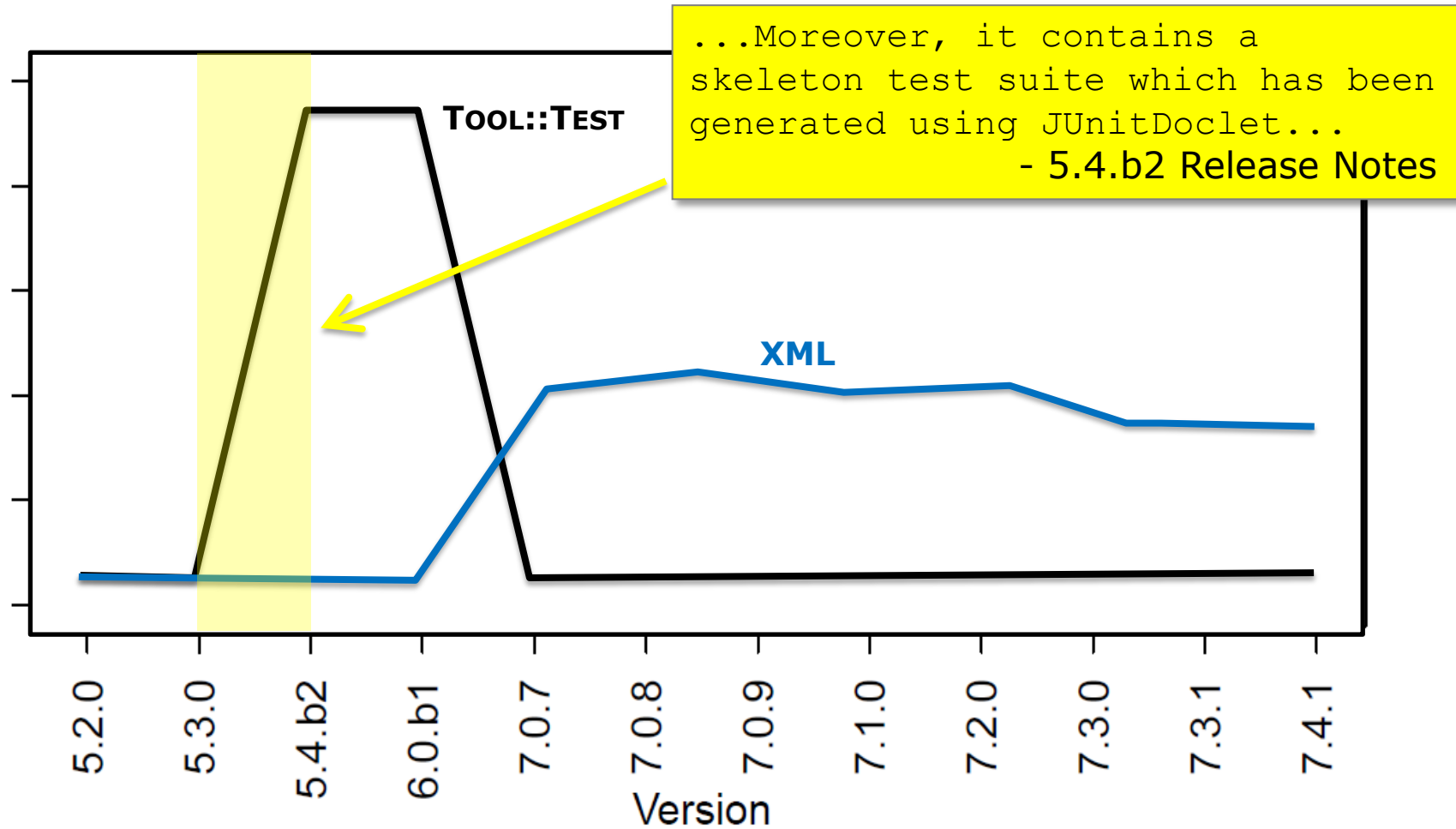
**Selected 80 at random**
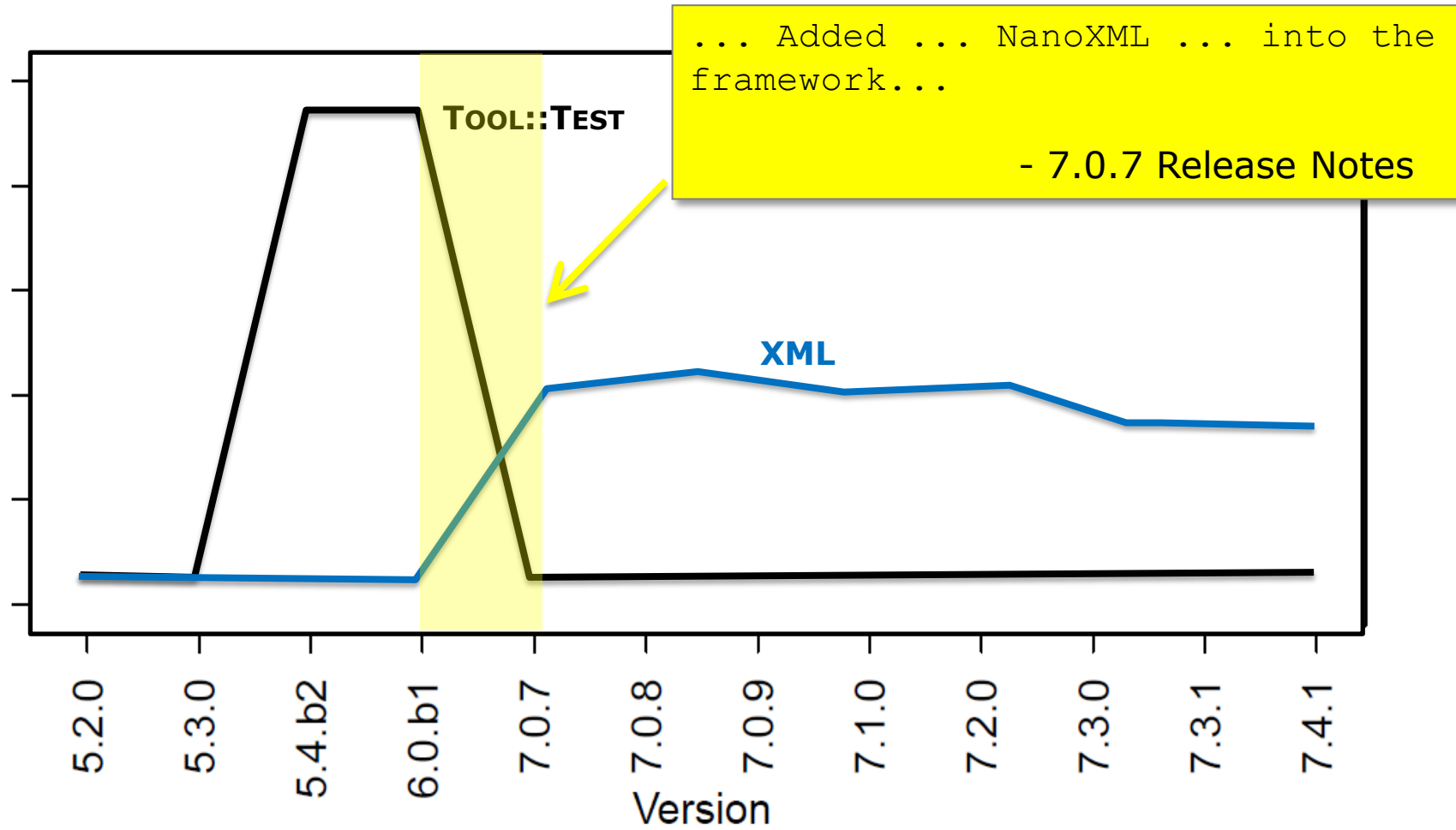
*95% confidence interval, 5% error*
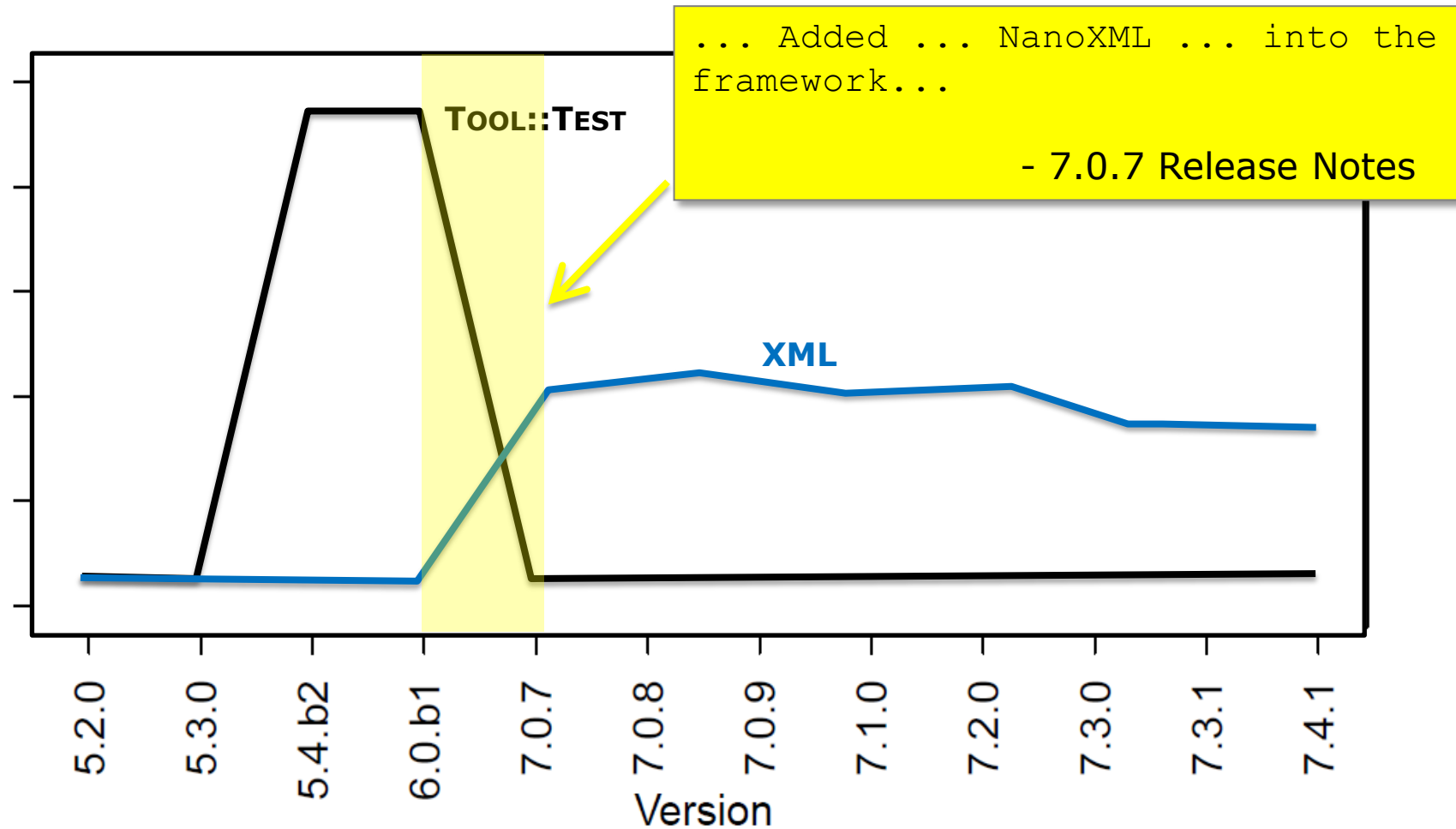
**Manually investigated project documentation**

*Change logs, release notes, commit logs*



TOOL::TEST

MENU

DRAWING::SHAPES

UNDO

Version

5.2.0  5.3.0  5.4.b2  6.0.b1  7.0.7  7.0.8  7.0.9  7.1.0  7.2.0  7.3.0  7.3.1  7.4.1

TOOL::TEST

XML

Version

TOOL::TEST

XML

...Moreover, it contains a skeleton test suite which has been generated using JUnitDoclet...
- 5.4.b2 Release Notes

Version

5.2.0 5.3.0 5.4.b2 6.0.b1 7.0.7 7.0.8 7.0.9 7.1.0 7.2.0 7.3.0 7.3.1 7.4.1

TOOL::TEST

XML

... Added ... NanoXML ... into the framework...

- 7.0.7 Release Notes

Version

5.2.0  5.3.0  5.4.b2  6.0.b1  7.0.7  7.0.8  7.0.9  7.1.0  7.2.0  7.3.0  7.3.1  7.4.1

**Findings:** 92% ± 5% of change events agree with documentation

# Implications



Software Changes — *Detected by* → *Used in* → Dashboards

40

# Implications



Software Changes

Detected by →

L atent
D irichlet
A llocation

Used in →

Topic Models

Dashboards

*Time 1*
*Version 1*

*Time n*
*Version n*

Time 1
Version 1

Time n
Version n

**L**atent
**D**irichlet
**A**llocation

Time 1
Version 1

Time n
Version n

**L**atent
**D**irichlet
**A**llocation

Version

5.2.0  5.3.0  5.4.b2  6.0.b1  7.0.7  7.0.8  7.0.9  7.1.0  7.2.0  7.3.0  7.3.1  7.4.1

44

Time 1
Version 1

Time n
Version n

**L**atent
**D**irichlet
**A**llocation

5.2. 5.3.0 5.4.b2 6.0.b1 7.0.7 7.0.8 7.0.9 7.1.0 7.2.0 7.3.0 7.3.1 7.4.1

Version

45

Time 1
Version 1

Time n
Version n

Latent Dirichlet Allocation

Version

**Time 1
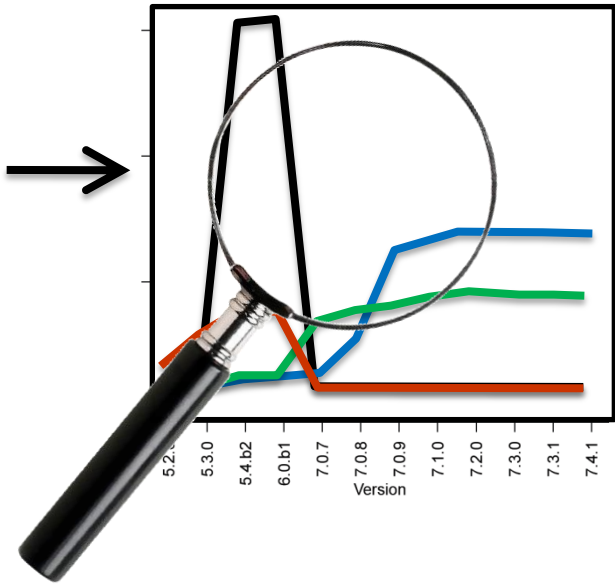Version 1**

**Time n
Version n**

**L**atent
**D**irichlet
**A**llocation

**What do project managers want on the dashboard?**

**How to do recall for topic models?**

# Backups

# Future Work

- **Investigate recall metric**
  - Do changes in source code result in changes in topic metrics?
- **Implement into dashboard**
  - Do managers actually like it?
- **More manual validation**
  - Larger and smaller systems
  - Closed-source systems
  - Systems from other domains