

# From Code Via Builds to Tests and Back to Capture the Full Picture

Kim Herzig (Software engineer & Researcher)

[kimh@microsoft.com](mailto:kimh@microsoft.com)  
[www.research.microsoft.com/people/kimh](http://www.research.microsoft.com/people/kimh)  
[www.kim-herzig.de](http://www.kim-herzig.de)

Tools for  
Software Engineers



Visual Studio  
Team Services



Microsoft



**Build tools**

**Test tools**

**Packaging tools**

**Static analysis tools**

**Code review**

**Data mining tools**

# When I joined Microsoft



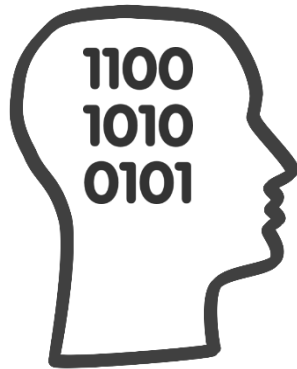
# From Code Via Builds to Tests and back to Capture the Full Picture

Kim Herzig (Software Engineer & Researcher)

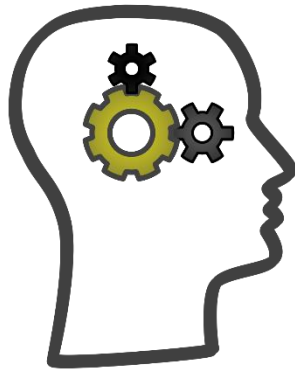
[kimh@microsoft.com](mailto:kimh@microsoft.com)  
[www.research.microsoft.com/people/kimh](http://www.research.microsoft.com/people/kimh)  
[www.kim-herzig.de](http://www.kim-herzig.de)

# MINING SOFTWARE ARCHIVES

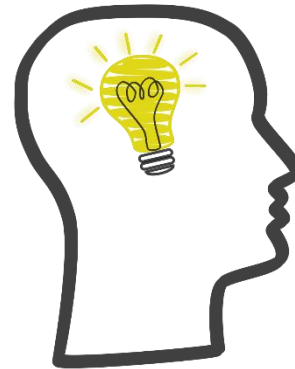
*"The Mining Software Repositories (MSR) field **analyzes** the rich **data** available in software repositories **to uncover interesting and actionable information** about software systems and projects" [MSR website]*



Data



Thinking



Idea



Solution!

# DO YOU KNOW WHAT YOU DO?

- **What is the problem** you are trying to solve?
  - Hypothesis
  - Research Questions
- Is the **data appropriate**?
  - How should the data look like?
  - What do you expect?
- What's the **return of investment** of your solution?

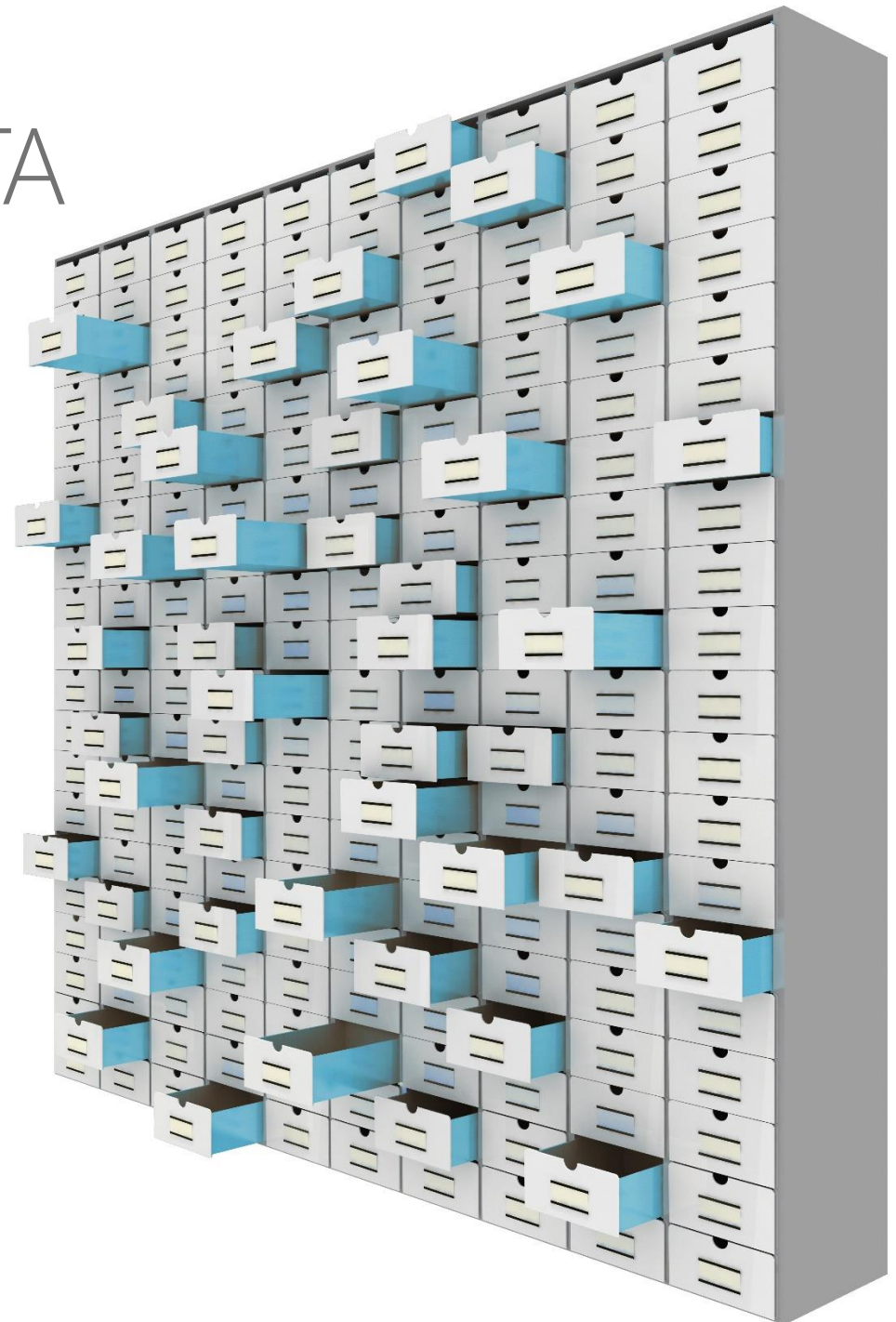




# DON'T GET LOST IN DATA

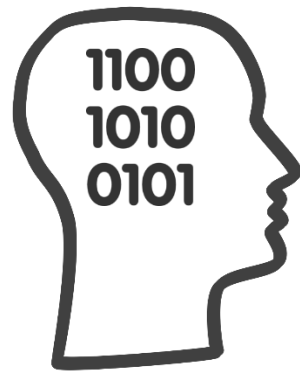
## Mining @Microsoft

- Source code
- Bug report / crashes / stack traces
- Organizational structure
- Releases
- Test executions
- Code reviews
- Builds
- Desktop tracking / telemetry
- ...

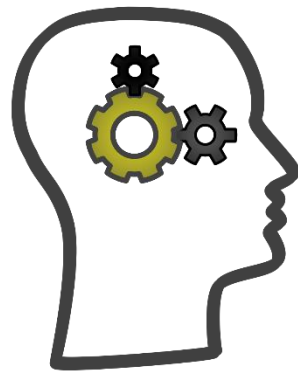




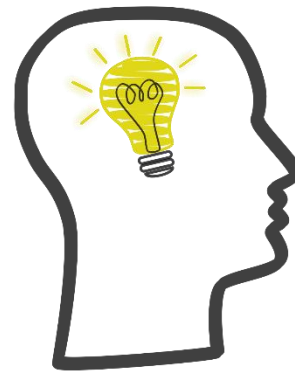
# MINING SOFTWARE ARCHIVES



Data



Thinking



Idea



Solution!

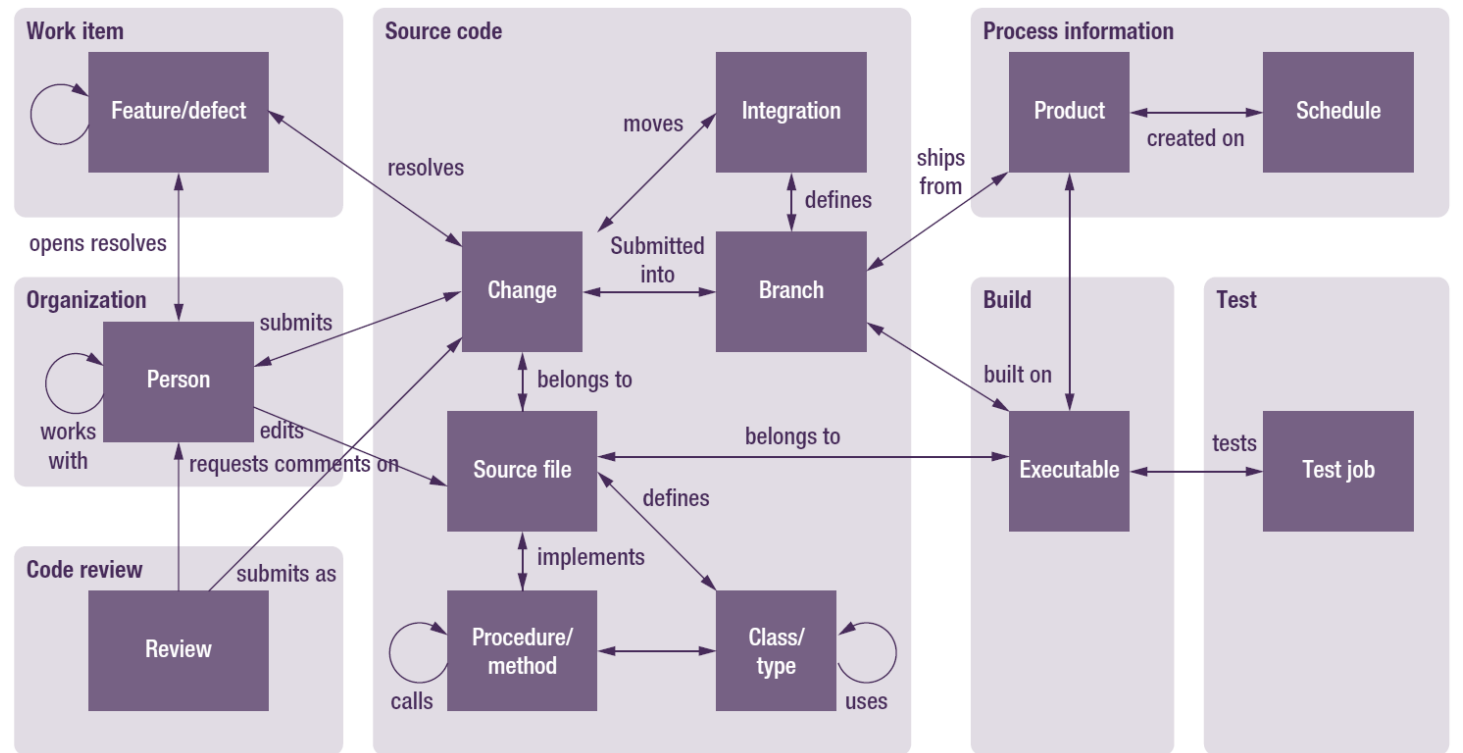
	Complexity of analysis		
	Statistics or machine learning methods required		
	Cross-referencing data from multiple sources		
	Counts and aggregates of existing data		
	<div>Understand usefulness of code review feedback</div> <div>Impact of code ownership on quality</div> <div>Impact of branch tree changes on process</div>	<div>Assess ramp-up time of new hires</div> <div>Engineering effort estimation</div>	
		<div>Churn-based componentization</div>	<div>Confirm beliefs about eng. workflow</div> <div>Measuring integration velocity</div> <div>Monitoring engineering activity</div> <div>Monitoring code ownership</div> <div>Monitoring changes to attack surface</div>
		<div>Assist with design and depl. of tools</div>	<div>Improve code review participation</div>
	New research into process	One time process decision	Ongoing process monitoring
	Usage pattern		

Figure 2. Opportunity vs. complexity of engineering data analysis

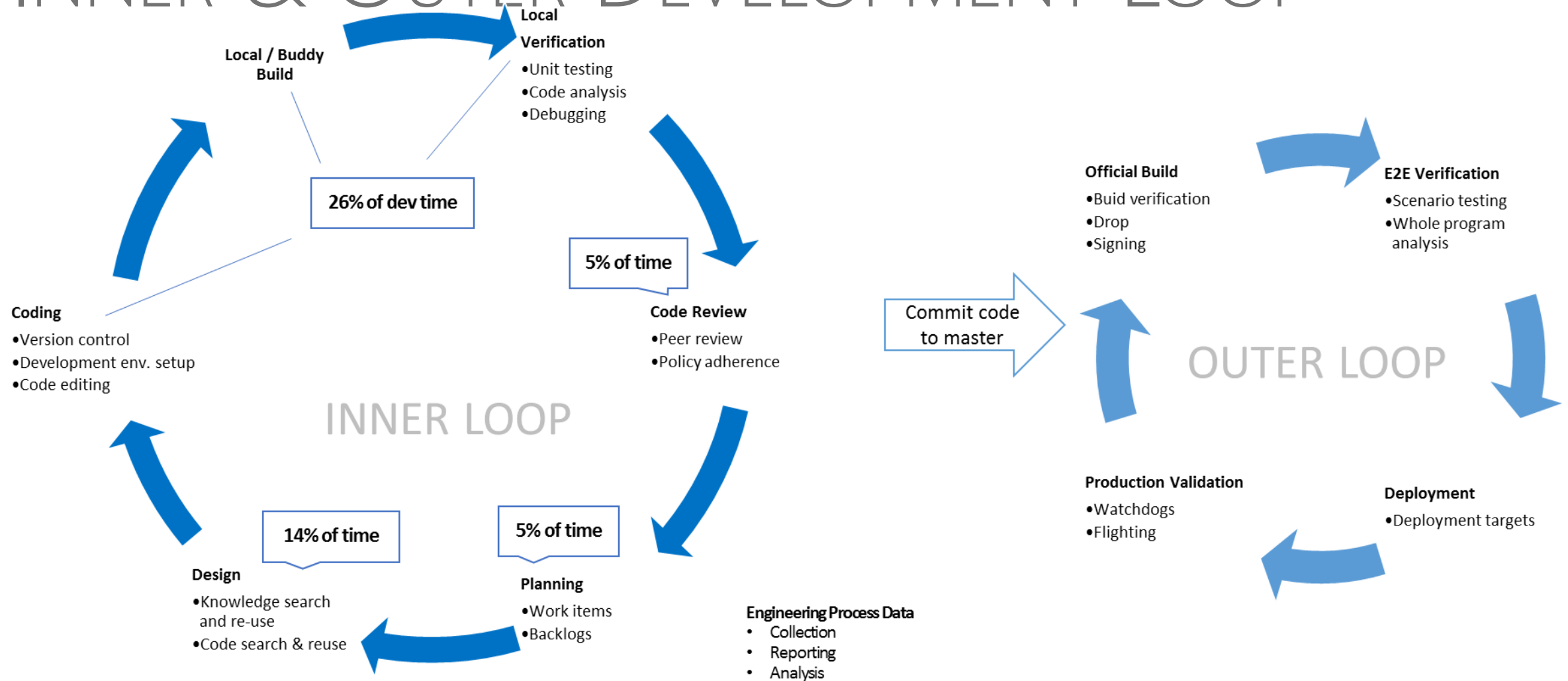
# MSR @ MICROSOFT (CODEMINE / CLOUDMINE)

Measuring & benchmarking the entire development process:  
from change to deploy.

- Code velocity
- Code quality
- Code/test ownership
- Legacy code
- Branching structures
- ...



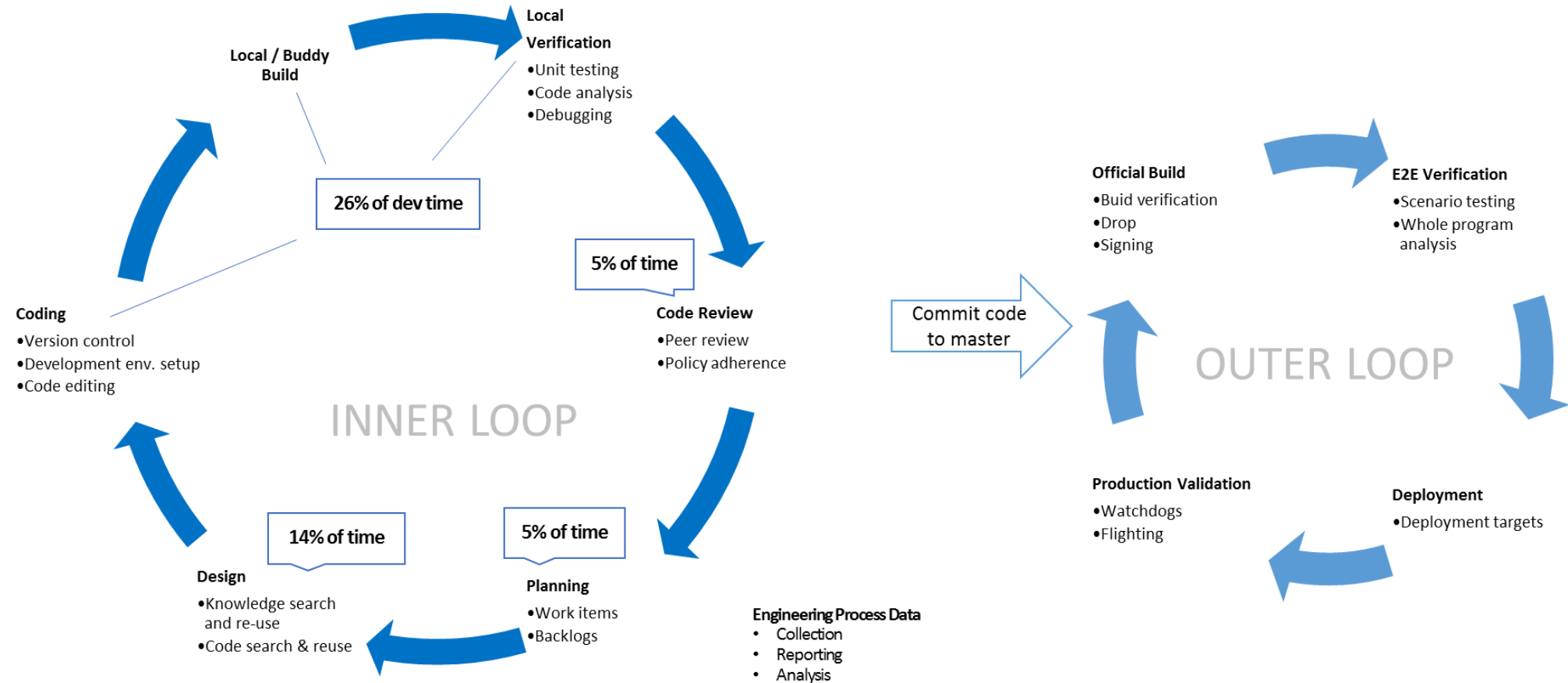
# INNER & OUTER DEVELOPMENT LOOP



Pre-check-in

Integration process

# HIGH EFFORT/COST



# MOST RESEARCH

What you can do with mining ...



<video removed>

0

Pre-Phase

# Data Scientist Play-Book

## ! Know your problem

---

- There is a difference between symptoms and problems

## ! Write down your assumptions

---

- We all make assumptions and it is important to make them explicit
- Share them with engineers.

## ! Engineers are humans

---

- They game the system
- Validate their statements with data and confront them with reality.

## ! Do not trust data

---

- Iterate over data with engineers
- Engineers claiming wrong data: 99% chance of wrong assumptions

## ! You will never get it right

---

- Accuracy of 80% is good. 100% is impossible.
- If you reach 90% do not trust you assumptions / data / results.

1

The source of all evil

# Code Churn



# Version Control <> Version Control

## Centralized

- + Data on server
- + Full data control
- + No loss of information
- Dev needs to be on network
- All commits "public"
- Limited isolation

## Distributed

- + Great work isolation
- + Fully developer controlled
- + More dev satisfaction
- Local dev data lost
- Historic data might be inaccurate
- Hard to mine entire dev process

# The Basics

Thanks to Sascha Just (just@st.cs.uni-saarland.de)

**Getting all changelists and their meta-data in the repository in efficient, machine readable format.**

We are completely avoiding regular expressions. There is no need for it and in some languages (like Java) the implementation tends to be horrendously slow.

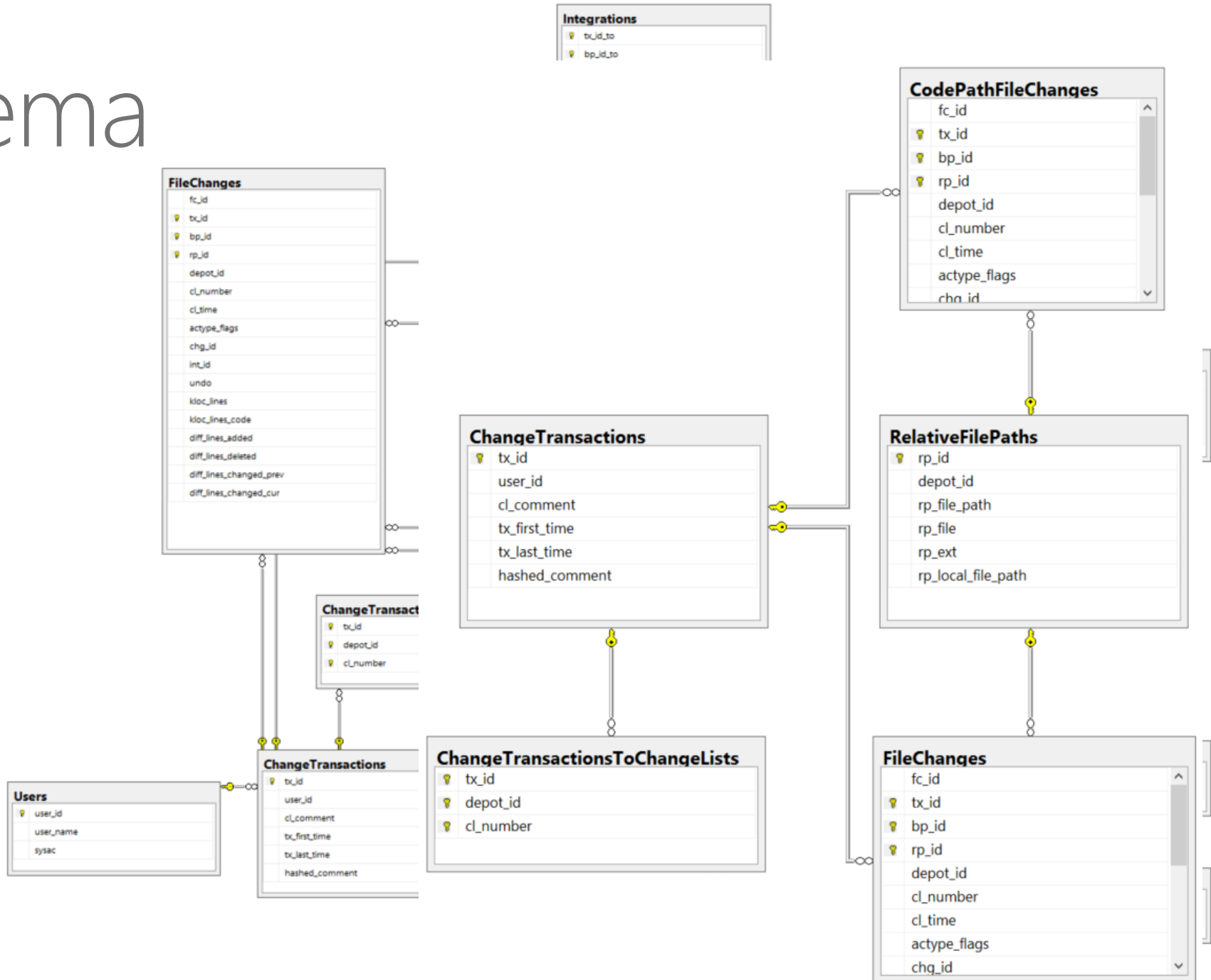
```
> git log --branches --remotes --tags --topo-order --find-copies --raw --numstat --no-abbrev --reverse --format=%n%H%n%T%n%an%n%ae%n%at%n%cn%n%ce%n%ct%n%s%n%bn
```

```
1da198d551de42c41fe96810152274387daed102
914f79e0cc22a4fdcc11a84f4cac4781099bd0aa
Kim Herzig
kimh@microsoft.com
1454546288
Kim Herzig
kimh@microsoft.com
1454546288
Fixing argument conflict
```

```
:100644 100644 6c5a8f04126020b4726821601817230415b6d951 ad0b15c296ec10ae22aabd7c4034a1aef39febae M src/CloudBuildMiner/Configuration/Configuration.cs
:100644 100644 bf53ccf35bf761d166d9a3071fcf51ebe9dc0cf7 125ca8f4dc853060ee8821d443134fb2618c90d8 M src/CloudBuildMiner/Program.cs
59      3      src/CloudBuildMiner/Configuration/Configuration.cs
12      4      src/CloudBuildMiner/Program.cs
```



# DB Schema



# Advanced Mining

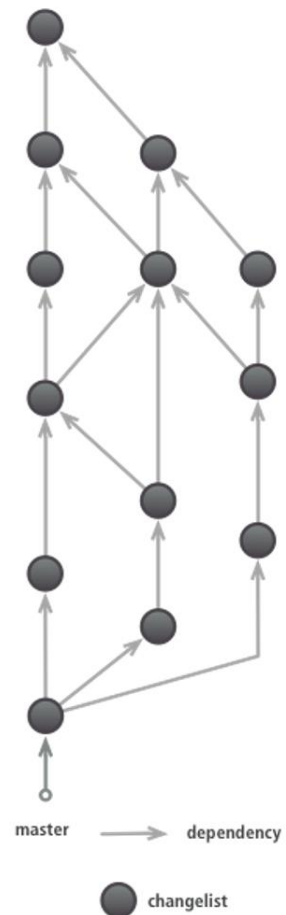
Finding shortest paths

e.g. code integration paths

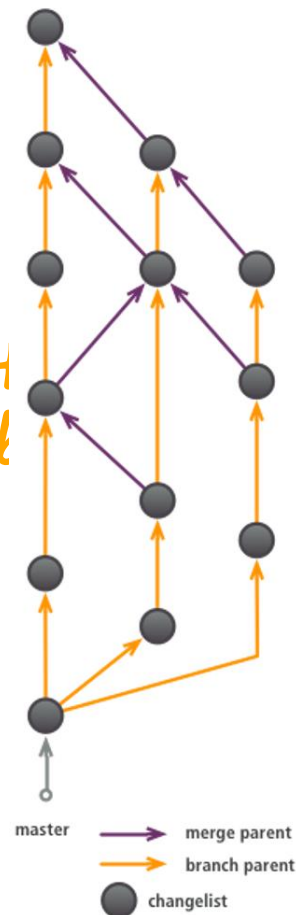
Thanks to Sascha Just (just@st.cs.uni-saarland.de)

Navigation graph usage: which way is the fastest to get to the next branch? and how many commits are in the way? and how many commits are in the way? and how many commits are in the way?

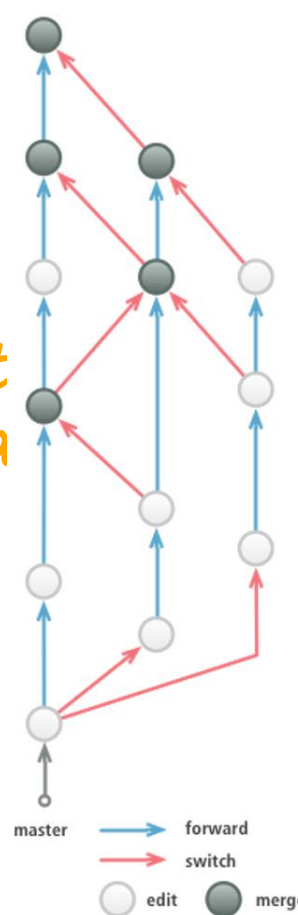
Bare Branch Graph



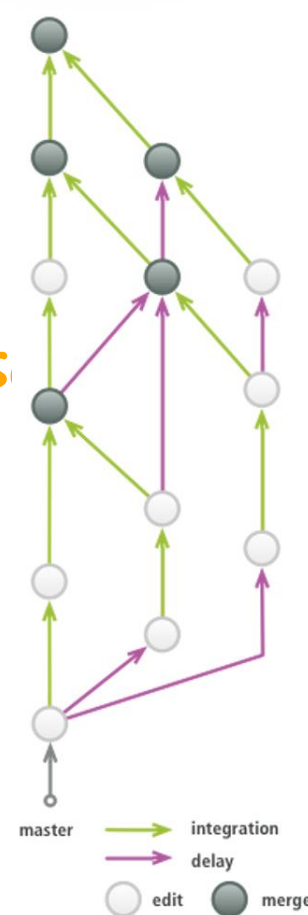
Git Branch Graph



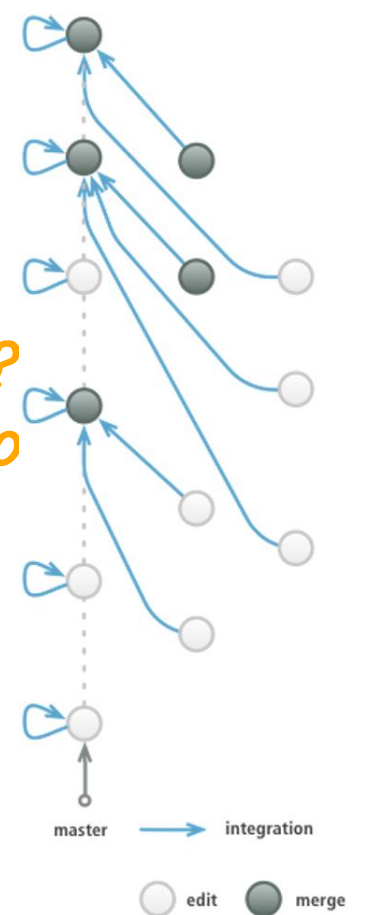
Navigation Branch Graph



Integration Branch Graph



Convergence Branch Graph



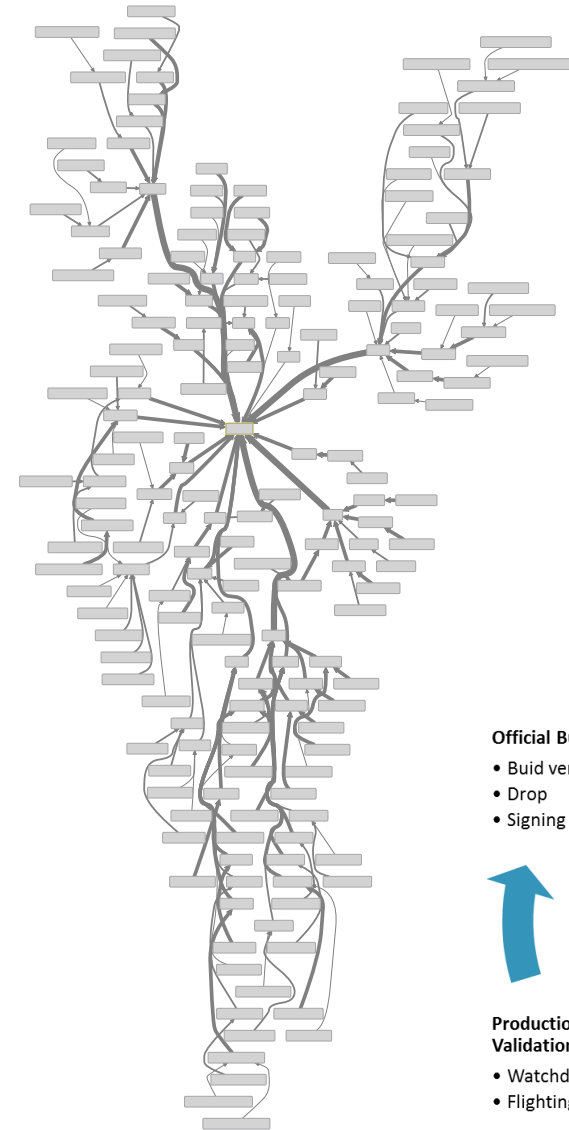
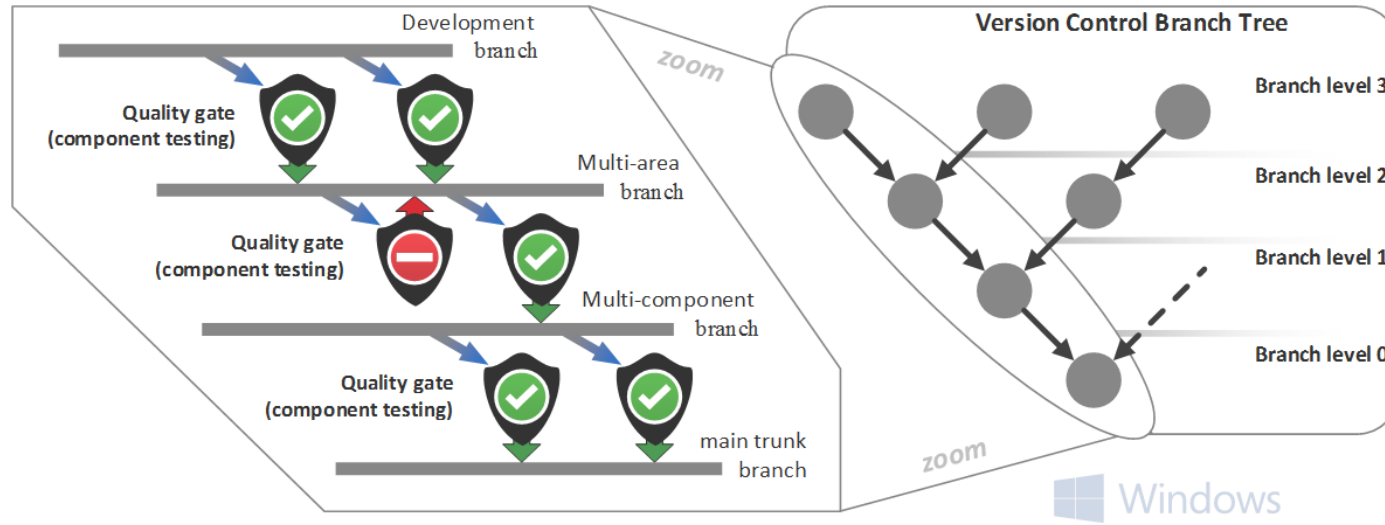
High level

did it / qua

release head

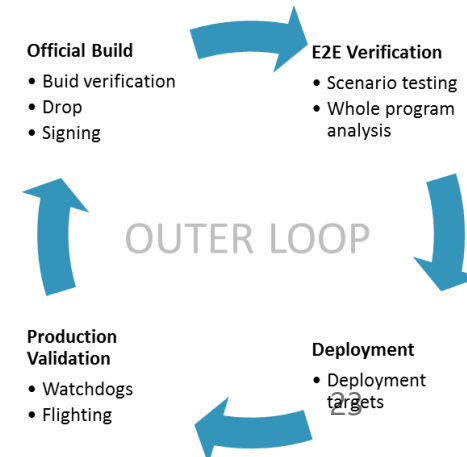
wait? nit to

# Branch structures can be complex ...



## Software testing is expensive

- 10k+ gates executed, 1M+ test cases
- Different branches, architectures, languages, devices, platforms, ...
- Aims to find code issues as early as possible
- Slows down product development



# Data Usage Scenarios

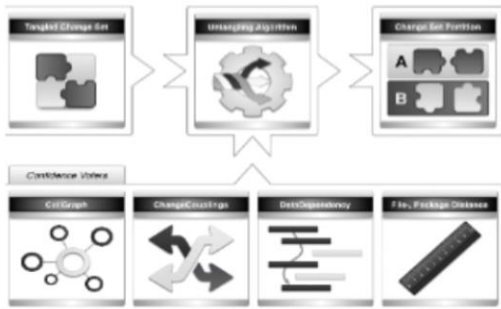
## File level metrics

- A **contributor** (person icon) has made commits/software changes to the software component.
- A **major** (person icon) / **minor** (person icon) contributor(s) applied more/less than 50% of changes to a component.



Ownership: relative number of commits

## Untangling / Measuring Impact



Take a *tangled* change-set as input.

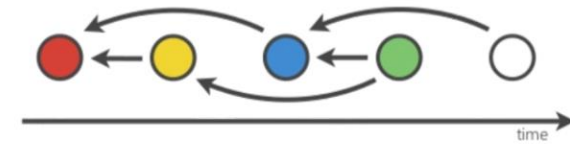
Produce a *change-set partition*

- Each partition holds change operations targeting one and the same development task.
- Partitions do not overlap.

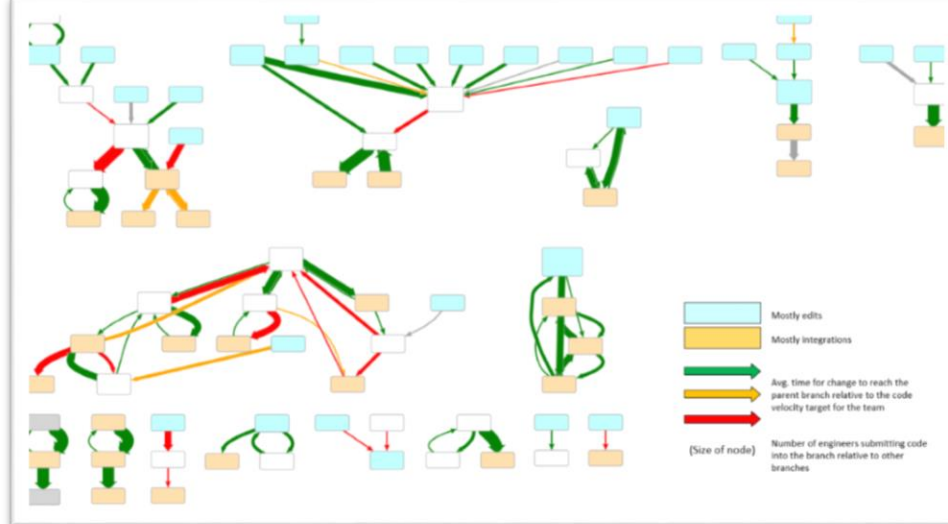
Untangling allows to measure impact on bug count models.

"The Impact of Tangled Code Changes", Herzig, Zeller, MSR 2013

## Change Genealogies



Combine **spatial** and **temporal** dimension of archives.



Code integration paths, Murphy, Czerwinka

# Git is "lying" *a matter of perspective*

## Fast-Forwards



Issue: original branch info lost

Impact: e.g. shared branches

Prevent: collect meta-data

## Rebase



Issue: actual patch changes

Impact: churn metrics, e.g. code quality

Prevent: ban rebases between branches

## Apply

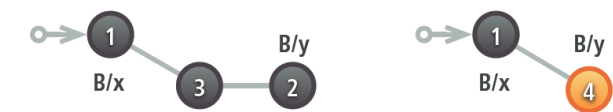


Issue: original info lost, e.g. timestamp & author

Impact: valuable data loss

Prevent: format-patch & am

## Squash



Issue: separate commits tangled

Impact: tangled changes, e.g. code quality

Prevent: ban squash commits

2

It's not a bug it's a feature, isn't it?

# Issue Reports



# The Basics

Create Issue

Configure Fields

Project

mozkito

Issue Type

Bug

Summary

Priority

Major

Due Date

Component/s

Affects Version/s

Fix Version/s

Assignee

Automatic

Reporter

Kim Herzig

Environment

Description

Original Estimate

Remaining Estimate

Attachment

Choose Files

Labels

Create another

Create

Cancel

## Details

Type:	Bug	Status:	Closed ( <a href="#">View Workflow</a> )
Priority:	Critical	Resolution:	Fixed
Affects Version/s:	04-SNAPSHOT	Fix Version/s:	None
Component/s:	mozkito-persistence		
Labels:	None		

## Description

## Activity

All Comments Work Log History Activity Source Reviews

- Kim Herzig added a comment - 07/Dec/12 10:29 AM please pull revision 636b97f24b539c74f880e17205d26906b5baa7e6 to reproduce
- Kim Herzig added a comment - 07/Dec/12 12:04 PM  
Same holds for RepositorySettings. Or I don't get it. A test on GIT runs. A test on HG fails.
- Sascha Just added a comment - 08/Dec/12 12:39 PM  
Well, this is what you get eventually when using Strings as configuration options...  
This was a mess from the beginning and originates from the days we ported the hibernate code.
- Interesting fact: This bug will be solved using refactoring 😊 Now classify this my friend!
- ### Steps to solve this issue (in progress)
- Replace all occurrences of String databaseType by DatabaseType databaseType.
  - Overwrite the toString() in DatabaseType to return name().toLowerCase()
  - Add constructor to DatabaseType(String) that takes the driver string as an argument
  - Add .getDriver() to DatabaseType that returns this string.
  - Remove all occurrences of DatabaseType.name() in the code.
- Sascha Just added a comment - 08/Dec/12 12:43 PM I also added DatabaseType.available() to check if we have the driver class on the classpath.
- Sascha Just added a comment - 08/Dec/12 12:48 PM Out of context: I added 2 static methods to UnrecoverableError: public static UnrecoverableError forma
- Sascha Just added a comment - 08/Dec/12 12:52 PM To avoid messy code, I removed the following method from the PersistenceUtil interface: void createSe
- Sascha Just added a comment - 08/Dec/12 1:33 PM Additionally introducing a new class to wrap around database options. Validity checks are scattered acr
- Sascha Just added a comment - 08/Dec/12 1:37 PM Further, I rename the values in ConnectOptions. These do not reflect the actual actions caused anymor
- Sascha Just added a comment - 08/Dec/12 3:09 PM Additionally: removed driver option from DatabaseTest annotation and replaced type string by enum.
- Sascha Just added a comment - 08/Dec/12 4:12 PM The reported bug will be fixed after the next push. However, there were more problems with the actual t
- Sascha Just <sascha.just@own-hero.net> submitted changeset f4f9b83f9856cb955debc8657448a34f7b950f04 to master in mozkito-research (15 files) - 08/De

Field	Original Value	New Value
Status	Open [ 1 ]	Closed [ 6 ]
Resolution		Fixed [ 1 ]

# API versus HTML

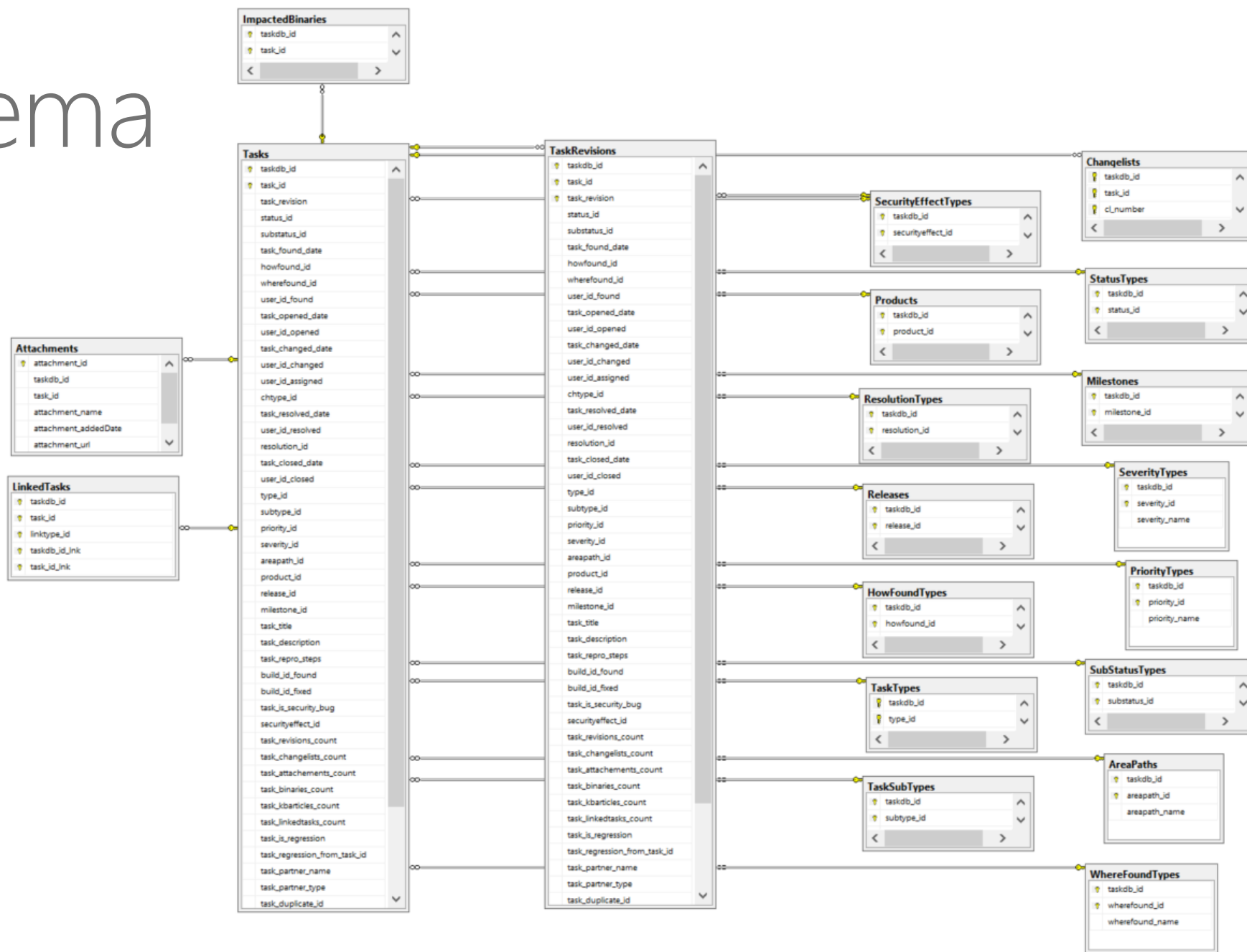
## Using the bug tracker API

- Easier
- Often missing info e.g. history
- APIs tend to change often

## Parse plain HTML / XML

- Access to all info
- More effort
- Will break often. Requires good testing / automation.

# DB Schema

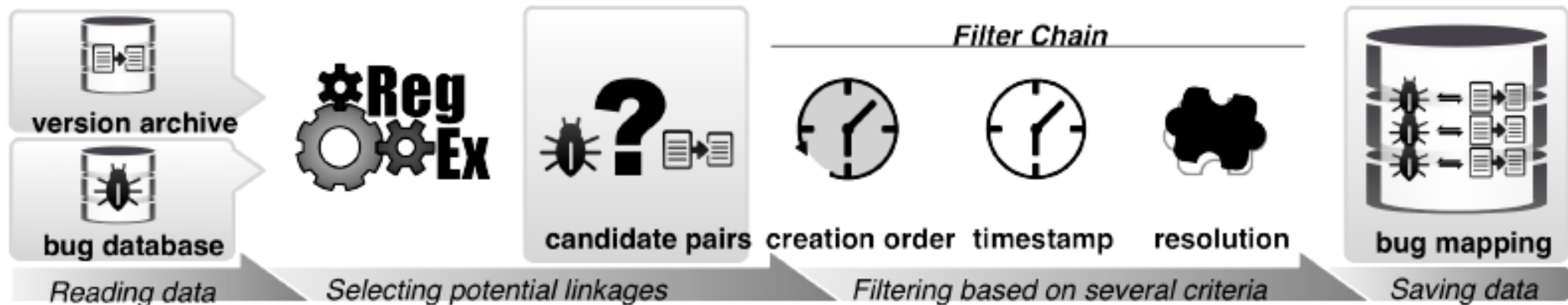


*The problem with issue reports is not the mining part,  
It's the data interpretation part. :-P*

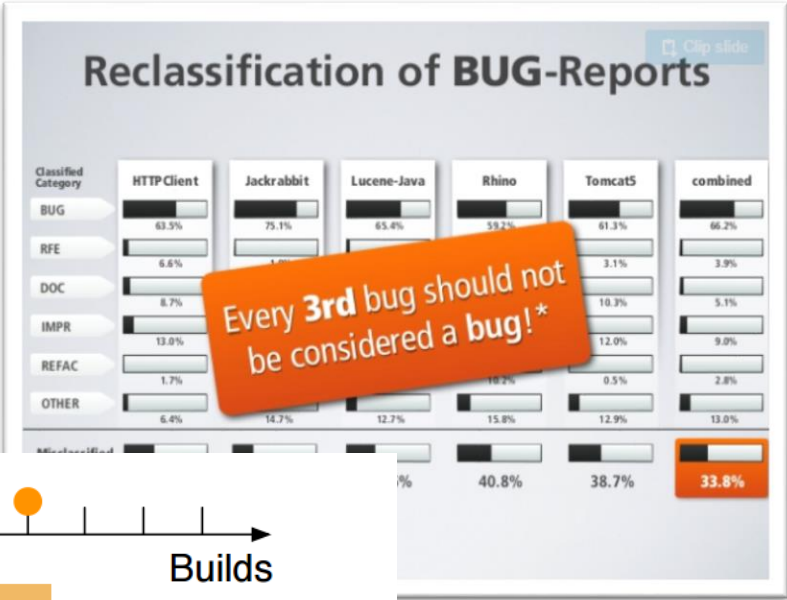
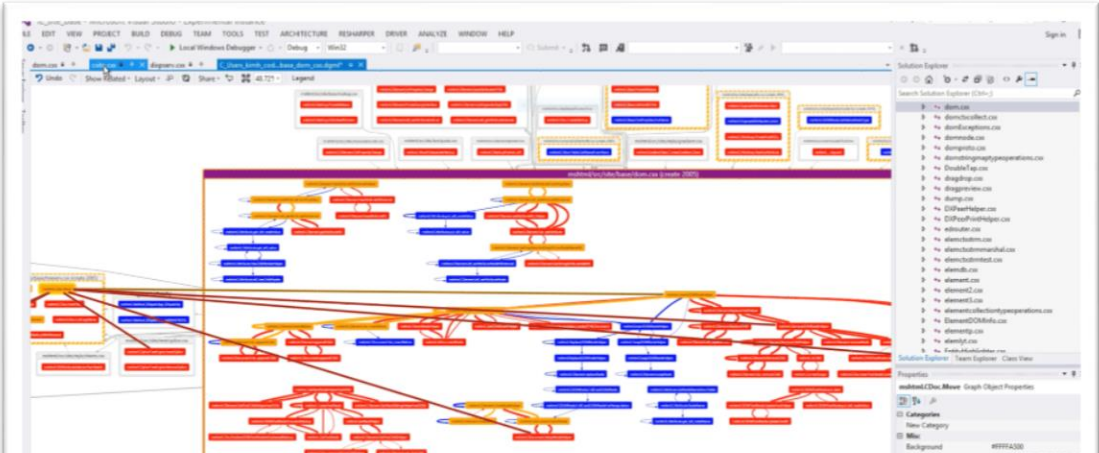


# Mapping bugs to code

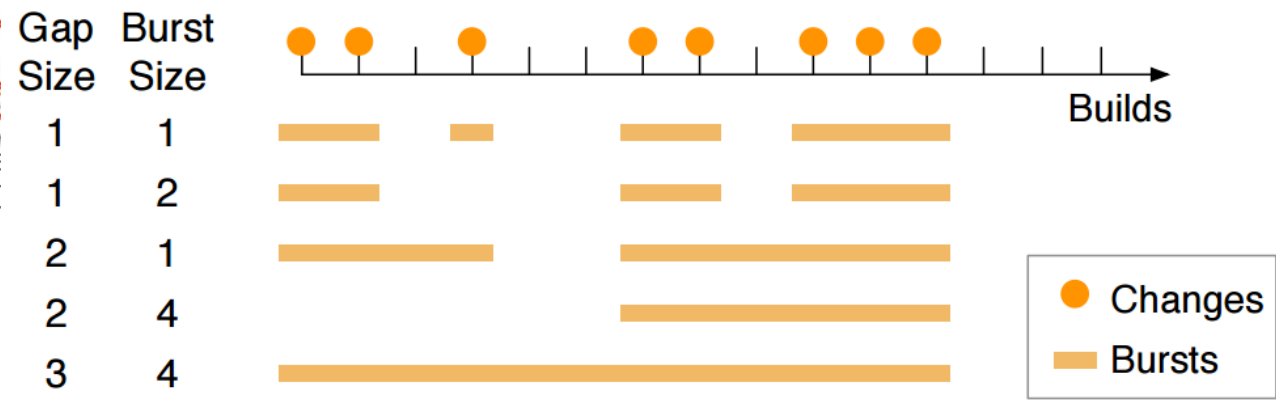
- Check for completeness / correctness of keywords constantly
  - TIP: Monitor mapping rate. Any drop is likely a due to changes in patterns / templates.
- Timestamps are always a challenge.
  - Bugs might be closed after commits (opened after the fact)
  - Different server in different time zones, ...
- Not all bugs get closed
  - There might be hidden patterns for "closed" reports



# Data Usage Scenarios



"Approximating Attack Surface"  
Morrison, Murphy, Williams, IC



"How  
Bug Prediction",  
2015

"Change Bursts as Defect Predictors", Nagappan, Zeller, Zimmermann, Herzig, Murphy  
ISSRE 2010



# Issue Report Pitfalls

- Different tracker => different behavior
  - **Tracker represent data differently.** Users will misuse fields or represent
- Default values
  - **Many issue reporters use default values:** they might not know the right
- What is a bug and what an improvement
  - **Philosophical question!** Be aware that a bug might not be a bug for everyone
- Ambiguous terms
  - **Likely to introduce noise.**
- No standard field set
  - **Labels or tags may replace fields**, e.g. Google tracker or Sourceforge
  - Hard to interpret as there is no common behavior , not even within a te

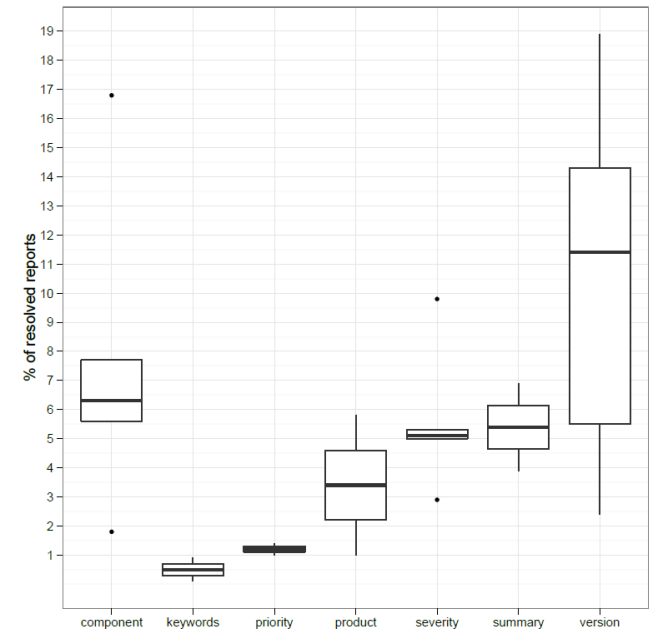
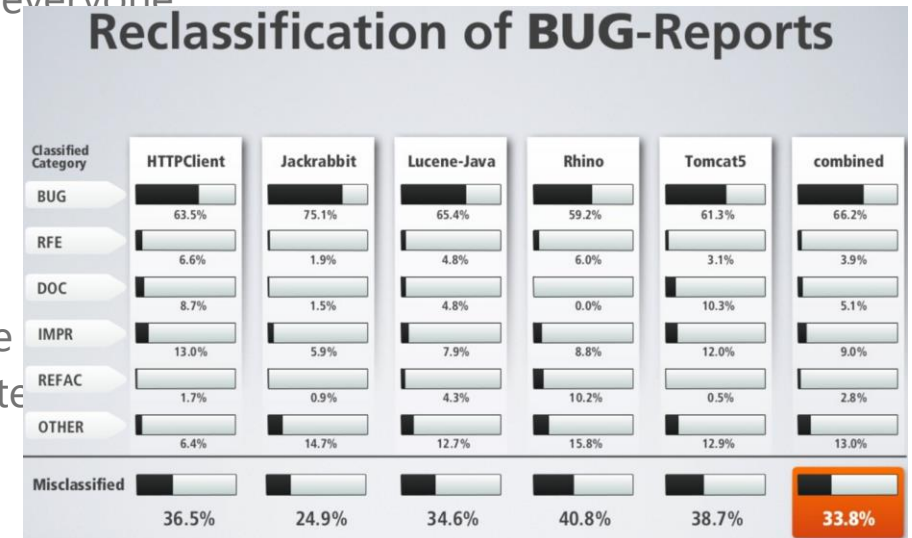


Fig. 2: Percent of resolved issue reports with respect to field changes. Priority, product, and summary only changed in BUGZILLA tracker projects RHINO and TOMCAT.



3

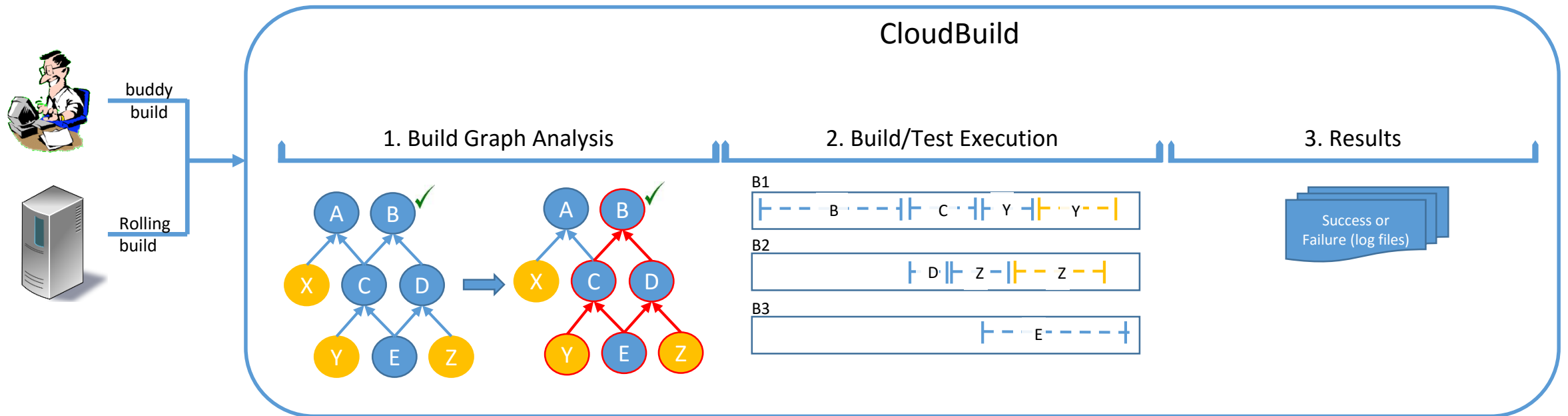
Who cares about build? You should!

# Build Data

# Why Build Data?

- **Lower bound of how fast we can ship software**
- Compiling complex software can take hours.
  - Very expensive!
- Speedup may require componentization / refactoring.
  - Very expensive!
- Catch 22: faster development usually means more builds.
- Each build **failure brings the development pipeline to halt.**
- Modern build systems are extremely complex!

# CloudBuild @ Microsoft



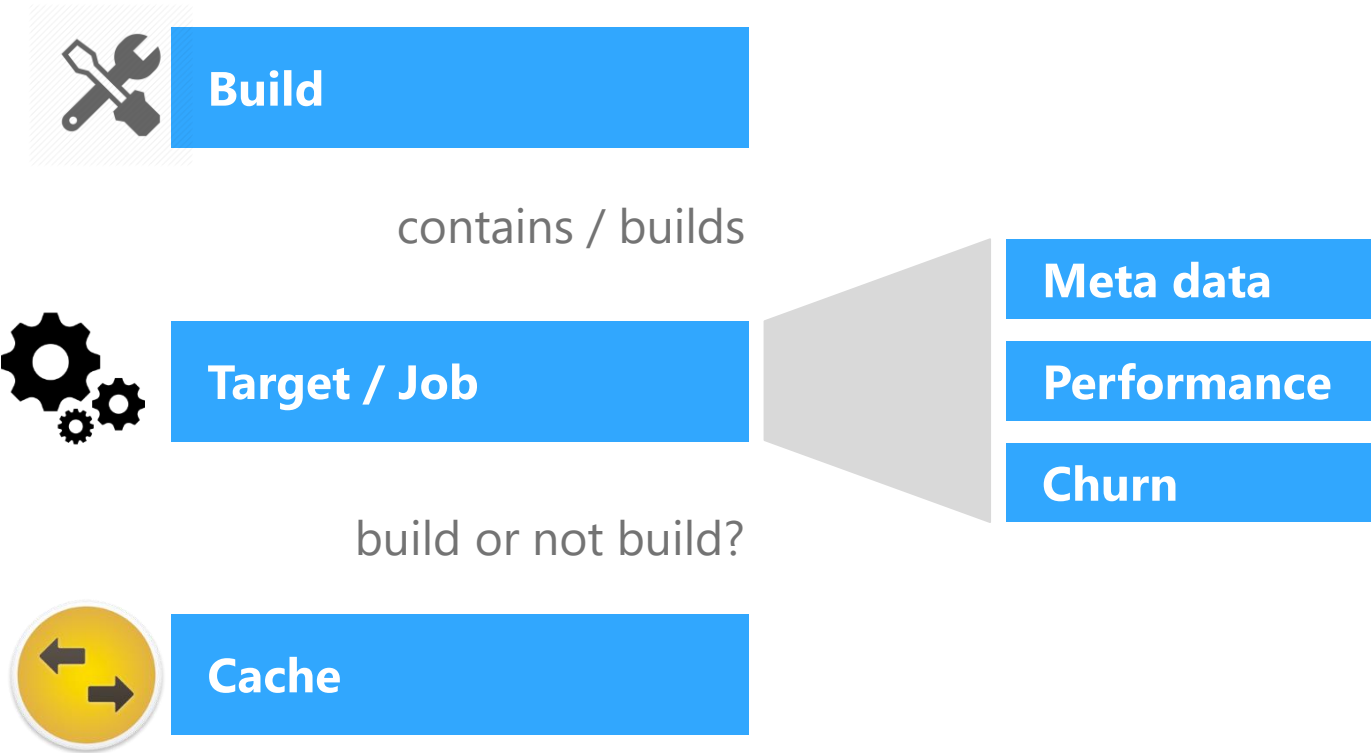
Legend:

- Source module (blue circle)
- Test module (yellow circle)
- Build (blue bar with arrow)
- Test execution (yellow bar with arrow)

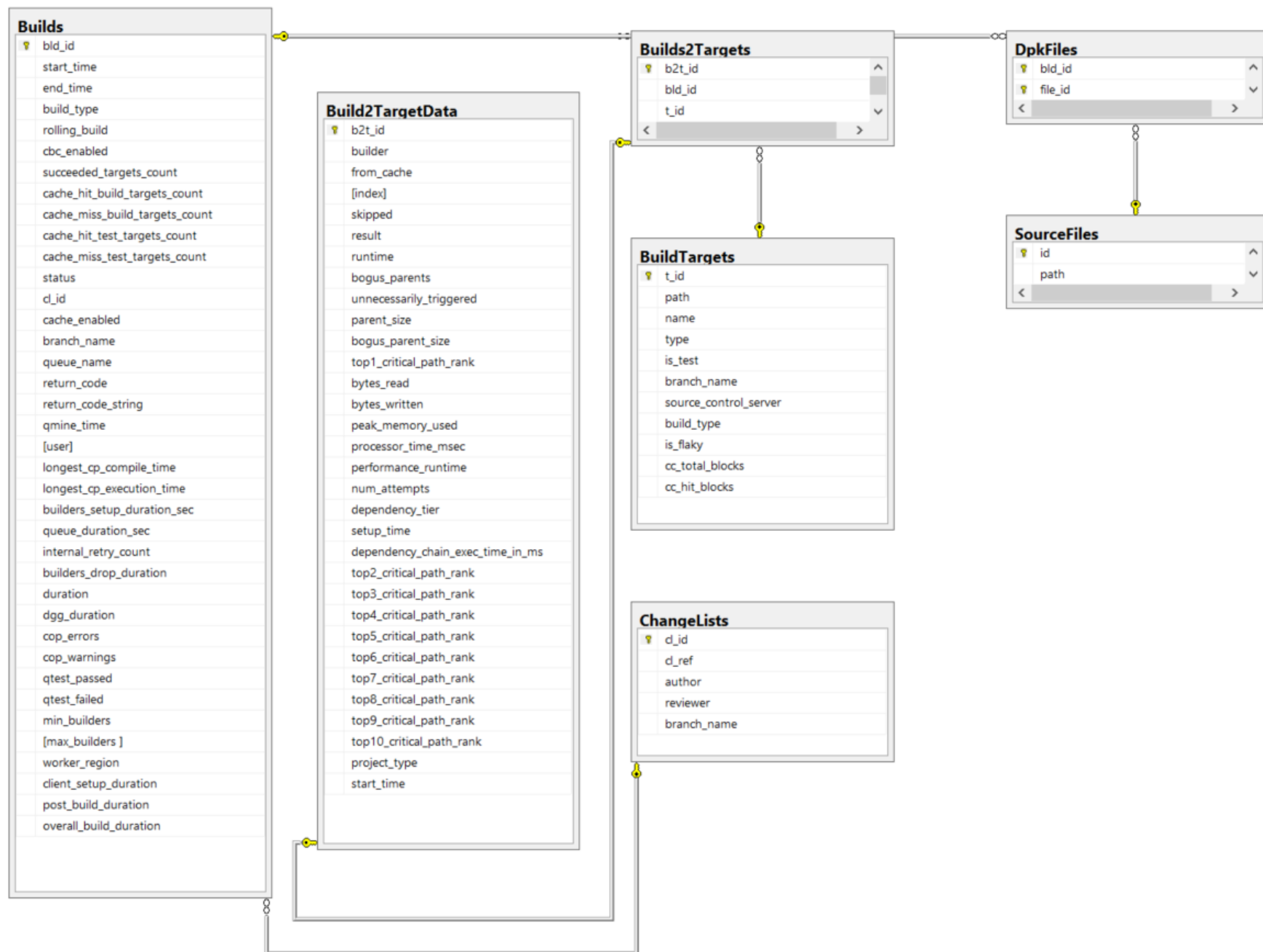
CloudBuild uses module-level test selection  
when any dependencies of a module change, all tests inside that module are executed.

# The Basics

*Many aspects will rely on build telemetry data (white box)*

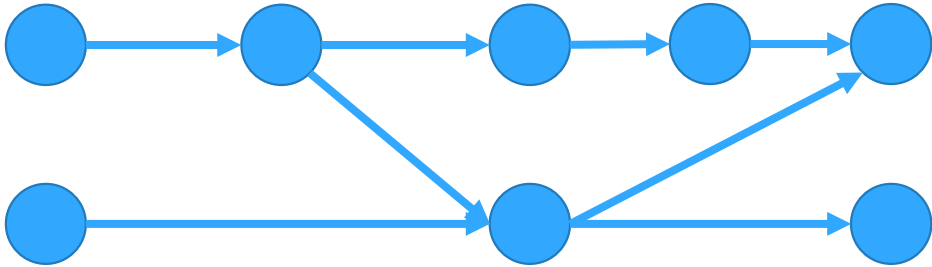


# DB Schema

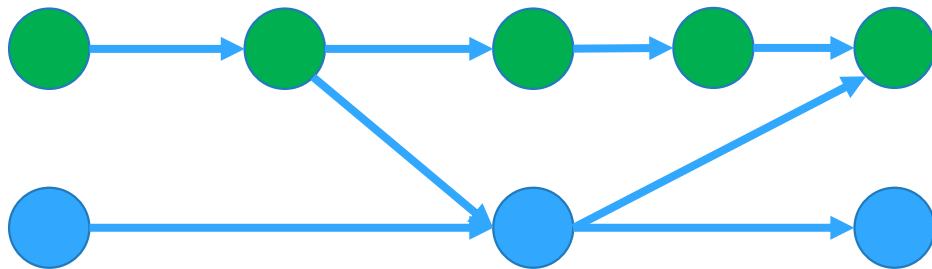


# Longest Critical Build Path (LCP)

- Graph structure: Directed Acyclic Graph (DAG)
  - Edges define dependencies.



- LCP: Defines the **lower bound on how fast we can build** the DAG
  - Defined with respect to build time, not target count.



*You will need to break the LCP to improve build times!*

# Data Usage Scenarios

- Suggesting refactorings to cut build time (*what-if-analysis*)

Tools for  
Software Engineers

Visual Studio  
Team Services

**Branch Health  
Dashboard**

Overview Dashboard

Number of Builds

Build Times

Cache Hits

Build Performance

Longest Critical Paths (LCP)

Target Metrics

Edge Cut Analysis

Config

Full Builds

(Cached) Build Clusters

Parent	Child	#Calls	%Speedup	%Accumulated Speedup	Black-list
...	...	1	0.00 %	0.00 %	<input type="checkbox"/>
...	...	1	4.64 %	4.64 %	<input type="checkbox"/>
...	...	2	0.02 %	4.67 %	<input type="checkbox"/>
...	...	2	0.00 %	4.80 %	<input type="checkbox"/>
...	...	2	0.00 %	4.80 %	<input type="checkbox"/>
...	...	2	0.02 %	4.82 %	<input type="checkbox"/>
...	...	2	0.05 %	4.85 %	<input type="checkbox"/>
...	...	2	0.03 %	4.88 %	<input type="checkbox"/>
...	...	2	0.00 %	4.90 %	<input type="checkbox"/>
...	...	3	0.20 %	4.93 %	<input type="checkbox"/>
...	...	4	1.39 %	6.32 %	<input type="checkbox"/>

- Predict impact of churn on build times.

FilePath	#CLs	Trend?	CL IDs	#Targets churned	#Targets built
...	7		2006795, 2008798, 2016623, 2016702, 2017669, 2017892, 2020499	5	868

-> Great for policy making, e.g. pre-checking



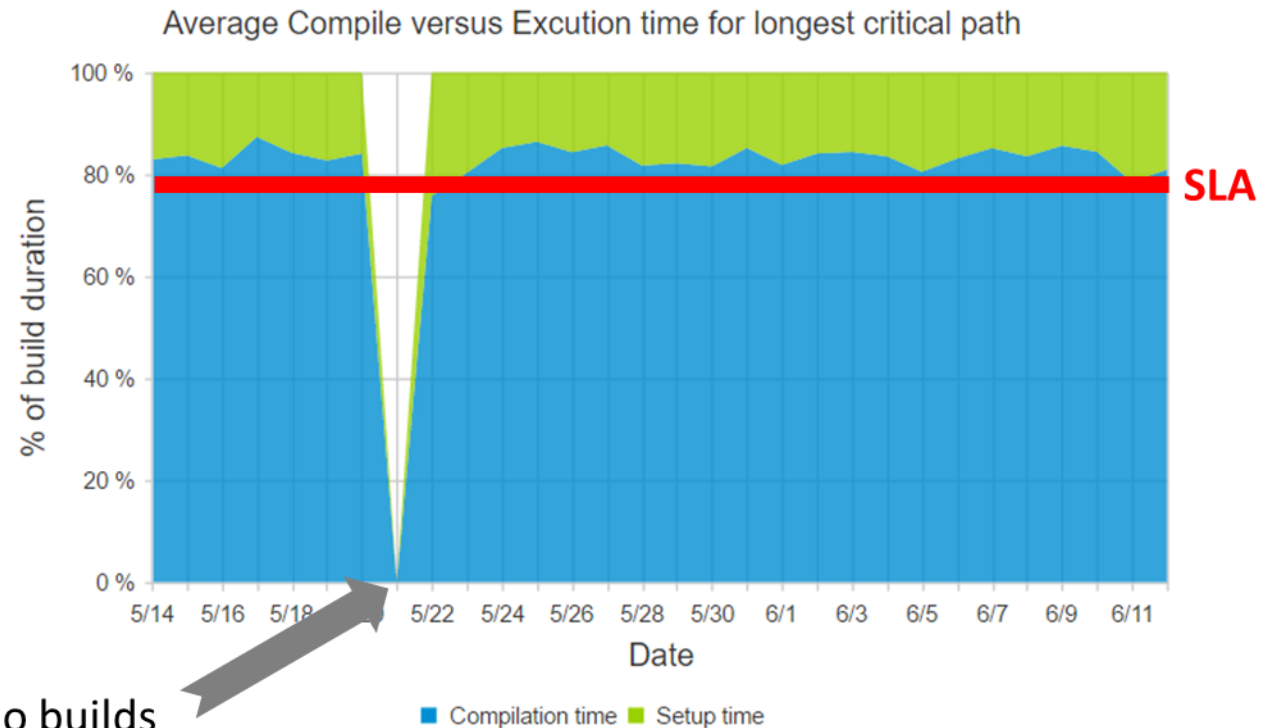
# Data Usage Scenarios (2)

- Detecting architectural build bottlenecks

**Top 25 targets most frequently appearing on longest critical path.**

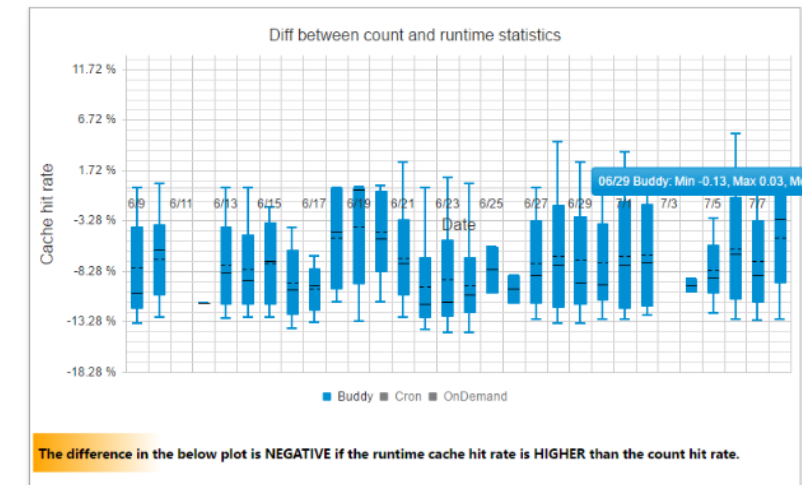
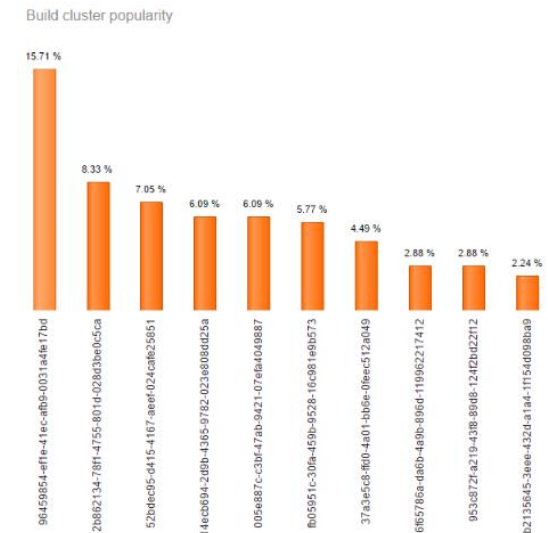
Target Name	%Builds
Microsoft.Exchange.MailboxReplicationService	91.58 %
Microsoft.Exchange.Data.Storage	90.84 %
Microsoft.Exchange.Data.Storage	88.37 %
Microsoft.Exchange.Data.Storage	88.12 %
Microsoft.Exchange.Data.Storage	87.87 %
Microsoft.Exchange.Data.Storage	87.62 %
Microsoft.Exchange.Data.Storage	84.16 %

- Verifying build tool SLAs



# Common Build Data Pitfalls

- No two builds are the same (caching)
  - **Cluster builds based on similarity.** Definition of similarity depends on scenario.
- Only compare comparable builds
  - **Cache hit rate, passing or failing, same code base, build type, etc.**
- Cache hit rates can be misleading
  - **Counts versus build time.**
- Build performance has many factors
  - **Churn, machine capacity, size, I/O, time window, caching, network, ...**



# 4

Where research and industry drift apart

## Test Executions

# The Basics

*Very expensive to mine after the fact. Collect telemetry!*

- **Store all results of all test runs! 24/7 - 365**
  - Data can become large. SQL is no option.
- Record **associations** with builds / users / time / code base
  - Comes in nearly all cases for free.
- Record **details for failures**, e.g. stack traces
- Record **code coverage** *if possible*
  - HUGE data requirements
- Record **performance info**:
  - How long did the test run? CPU usage.

# ~~DB Schema~~

# Azure Data Lake

Test target (suite)

Test case

Column Name	Type
UniqueSessionId	string
FullBranchName	
Queue	
TargetId	
Result	
Path	string
TestRunTimeMS	long
InfrastructureTimeMs	long
Passed	int
Skipped	int

FilesDeployedBytes	long
TestSourceDir	string
TestAdapter	string
DotNetFramework	string
Platform	string
QTestAccountMode	string

Schema

Column Name	Type
UniqueSessionId	string
FullBranchName	string
Queue	string
TargetId	string
Result	string
Path	string
TestRunTimeMS	long
InfrastructureTimeMs	long
Passed	int
Skipped	int
FilesDeployedBytes	long
TestSourceDir	string
TestAdapter	string
DotNetFramework	string
Platform	string
QTestAccountMode	string
Result	string
Attempts	int
ExecutionTimeMs	float
Flaky	string
Message	string
StackTrace	string

> 1 million test case executions  
per day per product.

~ 10GB per day

# Data Usage Scenarios

## Example

Test selection based on code coverage

\$3,350



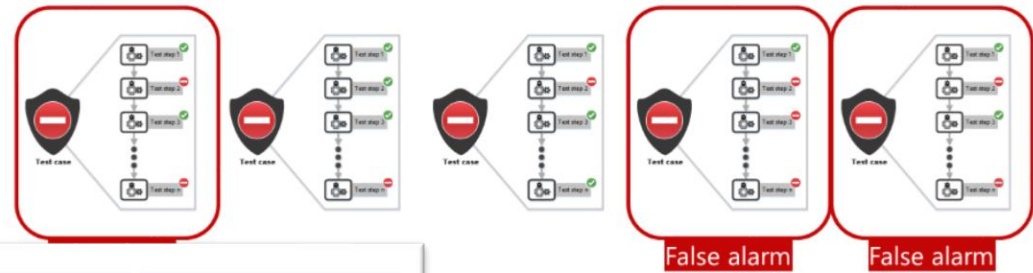
\$722

Slows down test execution by 30%

And this excludes many aspects: e

"Let's assume we had to pay for te

## Historic Patterns as Ground Truth



$p_2 \wedge TestStep_n \Rightarrow \text{False alarm}$

Alarms Using Association Rules", Herzig,

## Current Results

Simulated on Windows 8.1 development period (BVT only)



"The Art of Testing Less Without Sacrificing Code Quality", Herzig, Greiler, Czerwonka, Murphy, ICSE 2015

# Common Test Data Pitfalls

- **Flaky tests**: not every failure is a bug.
  - In fact, most failures are false test alarms.
  - Is it a flaky passing or flaky failing tests?
  - Flaky tests might get re-run.
- **Nobody runs all tests** anymore.
  - Every test event gives you an “incomplete” picture.
- Tests are run in **parallel with build**.
  - Saving test runs may be not as beneficial as you might think.
- Consider the **history** and meta data.
  - Owner, purpose, age, failure rate, run time, expense of a failure, who/how triggered the test ...
- **Code base is moving**. No two test runs are the same.
  - Except if you specifically control for it.

5

“En vogue” but misunderstood

# Code Reviews

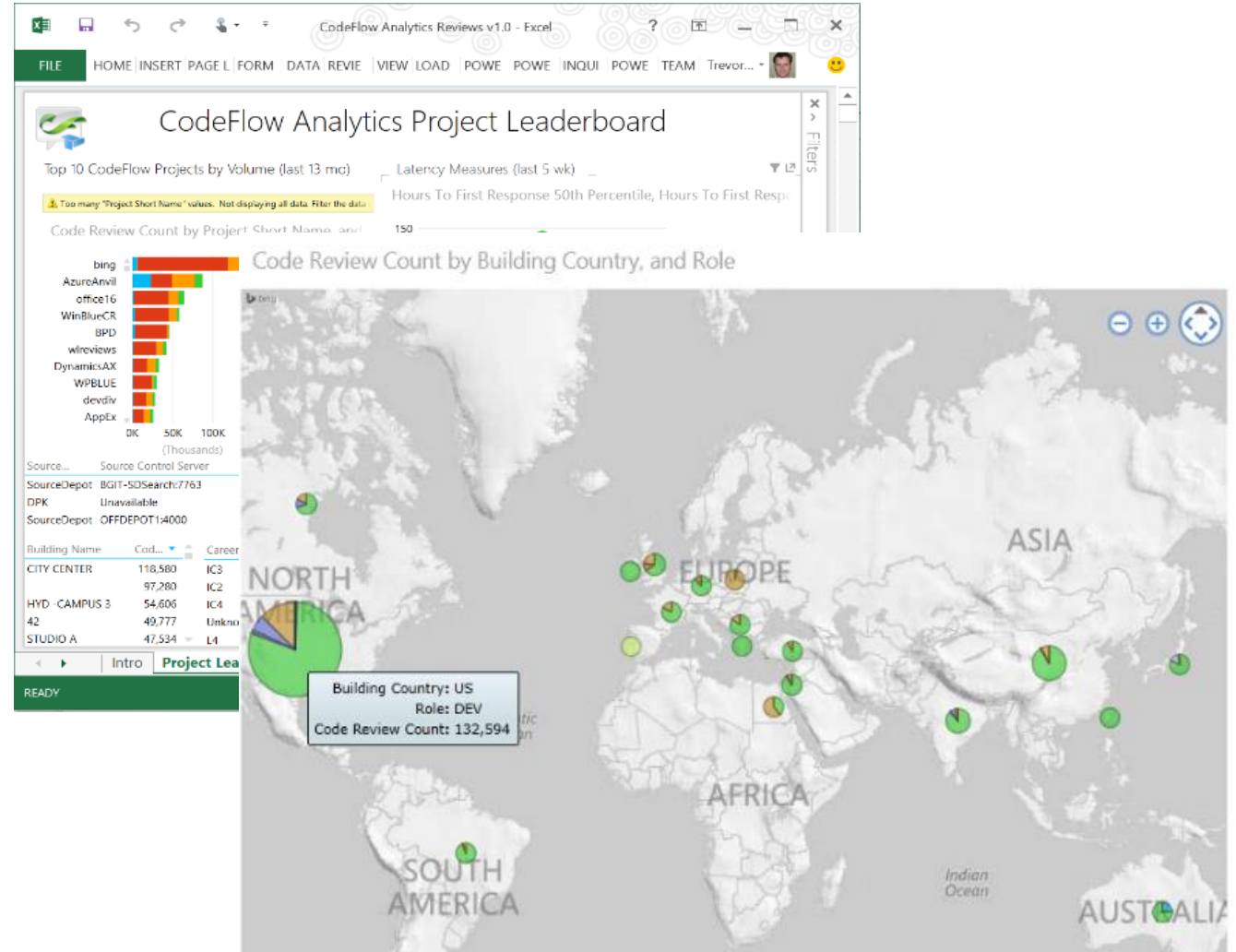




# The Basics

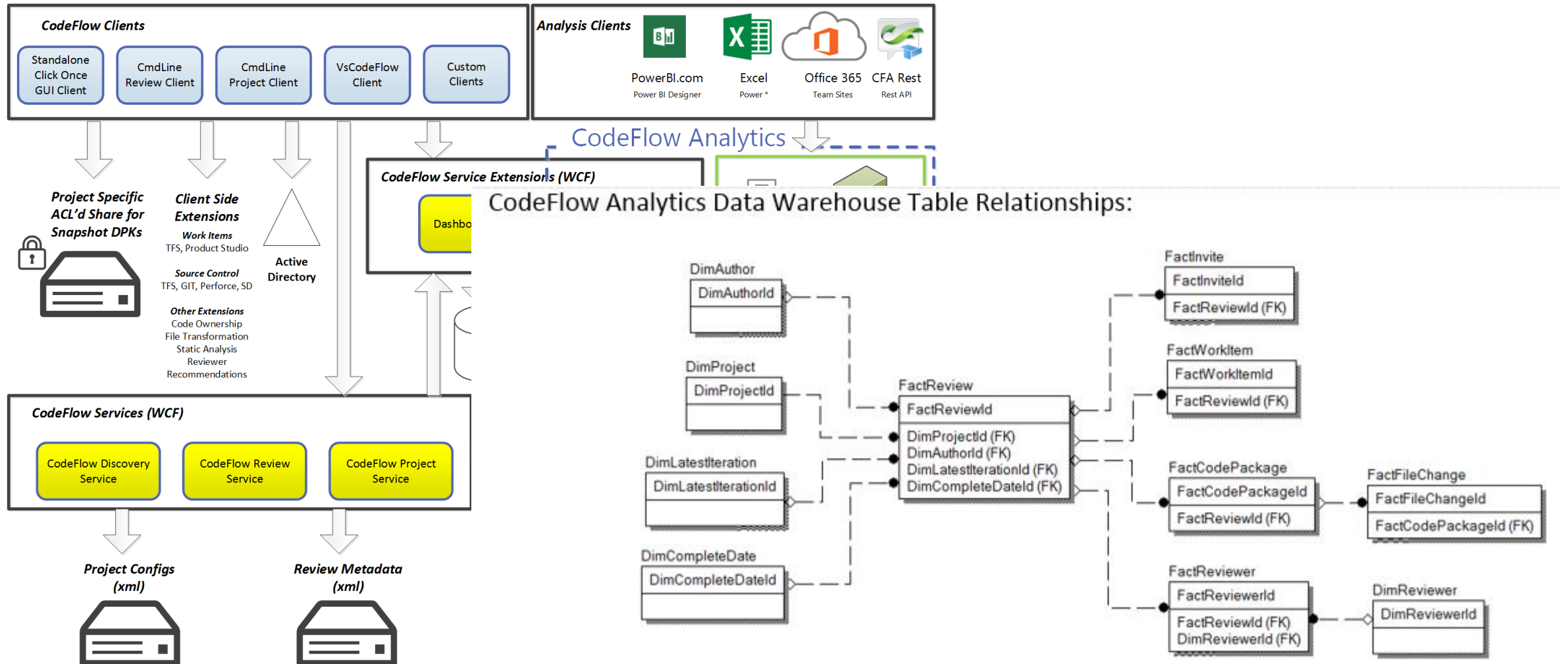
Thanks to Kivanc Muslu <kivancm@microsoft.com>

The screenshot shows the CodeFlow IDE interface. The main window displays a code editor with a file named `ReviewListenerService.svc.cs`. The code contains several lines of C# code, including a comment: `'t die on anything... Consider sending an email here? */`. A comment dialog box is open, showing a discussion between Trevor Carnahan and Christian Bird. The dialog box contains the text: "Does this flush 2 errors for EF exceptions now? From a failure mode analysis PoV, I want exceptions like this to emit only 1 error, with the most actionable message possible." and "changed to output both main and inner exception in one call to TraceError if there is an inner exception (which is often most actionable in my experience)." The dialog box has a "Closed" button and a "Hide" button. The reviewer status panel on the left shows a list of reviewers: Christian Bird (author), Trevor Carnahan (required), CodeFlow Analytics Developers, Ricky Kurniawan, Rock Wang, and Birendra Acharya (optional). The status of the review is "Closed".



# The Basics

Thanks to Kivanc Muslu <kivancm@microsoft.com>



# Where to collect data?

Thanks to Kivanc Muslu <kivancm@microsoft.com>

## Server side (service)

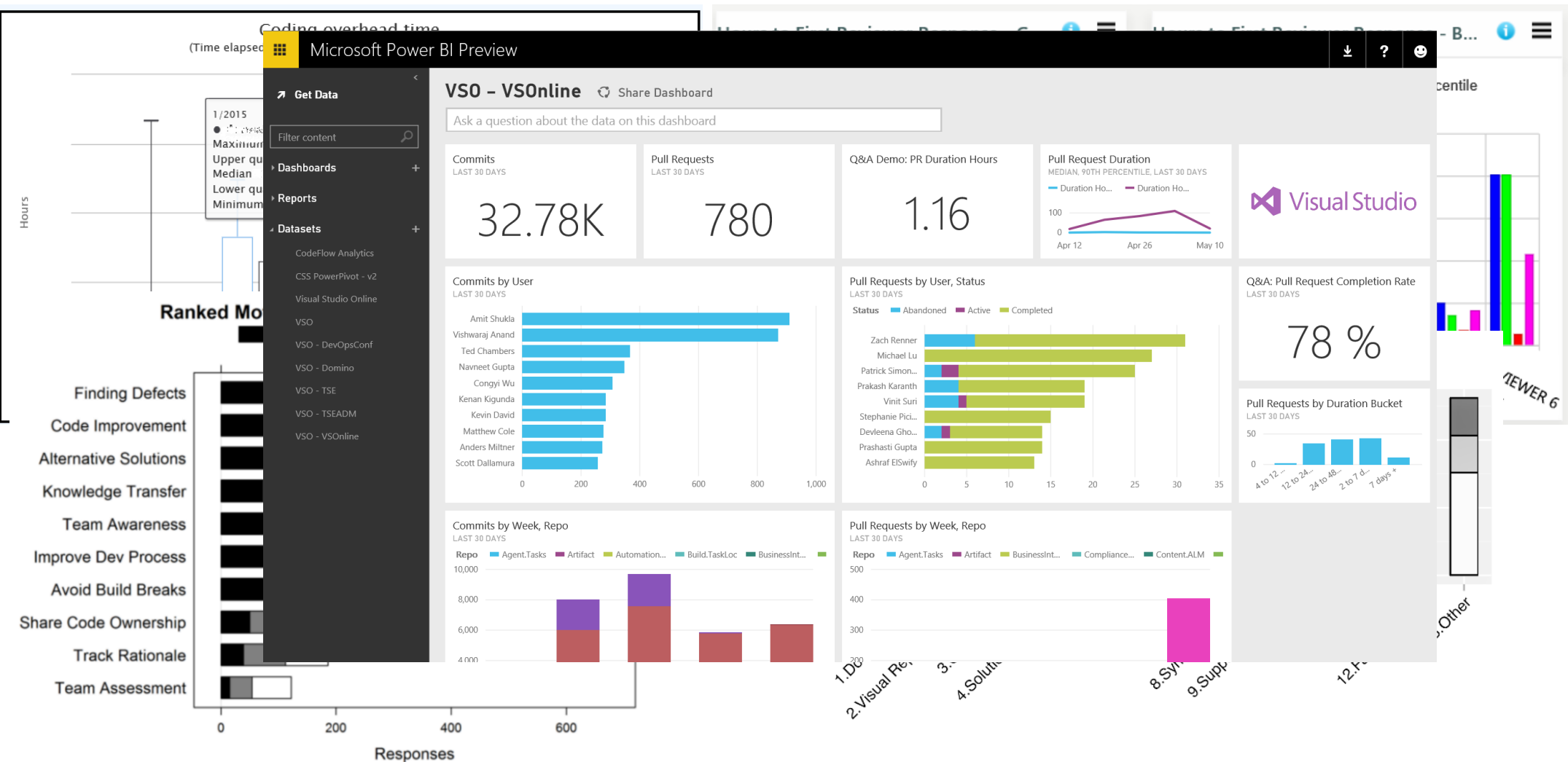
- *Preferred*
- Not all signals/metrics available
- Instant release (no cross-over)
- Full control

## Client side

- May come from different client versions
- Making mistakes can be fatal
  - Runtime, hard to revert, long time to fix
- Data loss due to setups
  - E.g. 4500 users on service, 4000 on client
- Difficult to debug data loss or issues

# Data Usage Scenarios

Thanks to Kivanc Muslu <kivancm@microsoft.com>



# Common Code Review Pitfalls

Thanks to Kivanc Muslu <kivancm@microsoft.com>

- In most cases: **do not use data standalone.**
  - Combing with code, churn, tests, builds, etc.
- Often requires **contextual information**
  - Why are taking the most complex reviews only seconds?
  - Why did the comment exist?
  - Why did no reviewer respond?
- Quality of review depend on **human aspects**
  - E.g. agenda of reviewer, time of submission, engineers block one hour a day to do reviews
- Different **usages** of code reviews
  - Gated check-in
  - FYI
  - Documentation
- **No ground-truth available**

If you take one message from this session ...

# Mining Software Archives

