



Impact-Driven Research on Software Testing and Analysis Tools: *Getting Real*

Tao Xie

University of Illinois at Urbana-Champaign

taoxie@illinois.edu

Successful Samples: Research 🗲 Practice









...

CAV Awards

- 2016: Separation logic (Facebook Infer)
- 2013: UPPAAL
- 2012: PVS
- 2011: SLAM/SDV
- 2010: Cadence SMV

• ..

http://i-cav.org/cav-award

https://www.research.ibm.com/haifa/conferences/hvc2016/award.shtml ³

ACM Software System Awards

- 2013: Coq
- 2012: LLVM
- 2007: Statemate
- 2006: Eiffel
- 2005: Boyer-Moore Theorem Prover
- 2001: SPIN
- •

NSF Workshop on Formal Methods



- Goal: to identify the future directions in research in formal methods and its transition to industrial practice.
- The workshop brought together researchers and identified primary challenges in the field, both foundational, infrastructural, and in transitioning ideas from research labs to developer tools.

NSF Workshop on Formal Methods



- Successes
- Unrealized Successes
- Obstacles
- Future Directions

http://goto.ucsd.edu/~rjhala/NSFWorkshop/static/Questions.md http://goto.ucsd.edu/~rjhala/NSFWorkshop/static/survey.pdf

Getting Real to Produce Practice Impact

Software testing/analysis tools are naturally tied with software development practice



Industry Academia Collaboration

- Academia (research recognitions, e.g., papers) vs. Industry (company revenues)
- Academia (invention/innovation) vs. Industry (likely involving engineering efforts)
- Academia (long-term/fundamental research or out of box thinking) vs. Industry (short-term research or work)
- Industry: problems, infrastructures, data, evaluation testbeds, ...
- Academia: educating students, ...

Getting Real to Produce Practice Impact

Software testing/analysis tools are naturally tied with software development practice



Google Scholar: "Pointer Analysis"

Google	pointer analysis
Scholar	About 400,000 results (0.07 sec)
Articles	[воок] Efficient context-sensitive pointer analysis for C programs
Case law	Abstract This paper proposes an efficient technique for context-sensitive pointer analysis
My library	using partial transfer functions. A partial transfer function (PTF) describes the behavior of a Cited by 706 Related articles All 36 versions Web of Science: 32 Cite Save More
Any time	Pointer analysis: haven't we solved this problem yet?
Since 2016	M Hind workshop on Program analysis for software tools and, 2001 - dl.acm.org
Since 2015	Abstract During the past twenty-one years, over seventy-five papers and nine Ph. D. theses
Since 2012	wonder "Haven'trdguo: we solved this problem yet?" With input from many researchers in the
Custom range	Cited by 576 Related articles All 23 versions Web of Science: 56 Cite Save

"Pointer Analysis: Haven't We Solved This Problem Yet?" [Hind PASTE'01]

"During the past 21 years, over 75 papers and 9 Ph.D. theses have been published on pointer analysis. Given the tones of work on this topic one may wonder, "Haven't we solved this problem yet?" With input from many researchers in the field, this paper describes issues related to pointer analysis and remaining open problems."

Michael Hind. **Pointer analysis: haven't we solved this problem yet?.** In *Proc.* ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering (PASTE 2001) Source@M. Hind



"Pointer Analysis: Haven't We Solved This Problem Yet?" [Hind PASTE'01]

Section 4.3 Designing an Analysis for a Client's Needs

"Barbara Ryder expands on this topic: "... We can all write an unbounded number of papers that compare different pointer analysis approximations in the abstract. However, this does not accomplish the key goal, which is to design and engineer pointer analyses that are **useful for solving real software problems for realistic programs**."

Michael Hind. **Pointer analysis: haven't we solved this problem yet?.** In *Proc.* ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering (PASTE 2001) Source©M. Hind&B. Ryder



Google Scholar: "Clone Detection"

Google

clone detection

Scholar

About 797,000 results (0.06 sec)

Articles

Case law My library Clone detection using abs <u>ID Baxter</u>, A Yahin, <u>L Moura</u>... - So

Abstract Existing research suggest of large-scale computer programs i such clones promises decreased s Cited by 1109 Related articles A

Any timeCCFinSince 2016scale sSince 2015T KamiySince 2012AbstractSince 2012Since coCustom range...detection

Sort by relevance Sort by date

✓ include patents

✓ include citations

CCFinder: a multilinguistic scale source code

T Kamiya, S Kusumoto, <u>K Inoue</u> - I Abstract–A code **clone** is a code p Since code clones are believed to a **detection** techniques and tools have Cited by 1262 Related articles A

Comparison and evaluation qualitative approach Typically focus/evaluate on intermediate steps (e.g., clone detection) instead of ultimate tasks (e.g., bug detection or refactoring), even when the field already grows mature with *n* years of efforts on intermediate steps

<u>CK Roy</u>, <u>JR Cordy</u>, <u>R Koschke</u> - Science of Computer Programming, 2009 - dl.acm.org Abstract Over the last decade many techniques and tools for software **clone detection** have been proposed. In this paper, we provide a qualitative comparison and evaluation of the current state-of-the-art in **clone detection** techniques and tools, and organize the large Cited by 544 Related articles All 19 versions Web of Science: 151 Cite Save

Some Success Stories of Applying Clone Detection [Focus on Ultimate Tasks]

🖉 patterninsight.com

pattern insight

Bugs fixed in one instance commonly propagate to many other parts of code.

SCM systems and traditional static analysis are not able to catch them.

LEARN MORE FRE



Yingnong Dang, Dongmei Zhang, Song Ge, Chengyun Chu, Yingjun Qiu, and Tao Xie. XIAO: Tuning Code Clones at Hands of Engineers in Practice. In *Proc. ACSAC 2012*,

http://research.microsoft.com/en-us/groups/sa/

Zhenmin Li, Shan Lu, Suvda Myagmar, and Yuanyuan Zhou. **CP-Miner: a tool for finding copy-paste and related bugs in operating system code**. In *Proc. OSDI 2004.*

http://patterninsight.com/



XIAO Usage: MS security bulletin MS12-034

Combined security update for Microsoft Office, Windows, .NET Framework, and Silverlight, published: Tuesday, May 08, 2012

3 publicly disclosed vulnerabilities and seven privately reported involved. Specifically, one is exploited by the <u>Duqu malware</u> to execute arbitrary code when a user opened a malicious Office document.

Insufficient bounds check within the font parsing subsystem of win32k.sys Cloned copy in gdiplus.dll, ogl.dll (office), Silverlight, and Windows Journal viewer

Microsoft Security Research & Defense Blog about this bulletin

"However, we wanted to be sure to address the vulnerable code wherever it appeared across the Microsoft code base. To that end, we have been working with Microsoft Research to develop a "Cloned Code Detection" system that we can run for every MSRC case to find any instance of the vulnerable code in any shipping product. This system is the one that found several of the copies of CVE-2011-3402 that we are now addressing with MS12-034."

Trying Real Tool – Parasoft Jtest 5.0

Back in Year 2004

```
Test 1 (T1):
  BST t1 =
    new BST();
  t1.insert(2);
  t1.insert(1);
  t1.remove(1);
  t1.insert(3);
  t1.size();
```

```
Test 2 (T2):
BST t2 =
    new BST ();
t2.insert(2);
t2.insert(3);
```

```
Test 3 (T3):
  BST t3 =
    new BST ();
  t3.insert(2);
  t3.insert(1);
  t3.size();
```

Approach using method sequence removes **no tests**

Rostra removes **T2 and T3** because T2 and T3 are redundant w.r.t. T1

Redundancy among Jtest-Generated Tests



- 90% of generated tests are redundant!
- Minimized tests preserve the same code (branch) coverage and seeded-bug coverage

Industry Impact — Parasoft Jtest

Back in Year 2004

• People do use Jtest

—Jo

• Recognized with numerous awards, including Jolt Product Excellence Award and JDJ Editor's Choice Award in 2004; adopted by thousands of development teams worldwide.

- businesswire.com

- But didn't seem to love its test generation
 - "I can't i ink of anyone telling me that they love Jtest's testgenerat ture."

r, JUnit book author, 02/05@junit user mailing list

Parasoft VP contacted me informing me that they incorporated our new techniques into Jtest's later versions!

Getting Real to Produce Practice Impact

Software testing/analysis tools are naturally tied with software development practice









NUnit Extension – 13024 downloads

xUnit.net Extension - 8589 downloads

From http://bbcode.codeplex.com/

Why is it as stable as we claim?

We have used Visual Studio IntelliTest (Announcement for Visual Studio 2015: https://www.visualstudio.com/enus/products/vs-2015-product-editions.aspx) to extensively test some important properties of this BBCode-Parser. We used IntelliTest to ensure that the parser never crashes and that it never emits any dangerous tag such as

(Dynamic) Symbolic Execution (DSE): Explosion of Search Space

There are decision procedures for individual path conditions, but...

- Number of potential paths grows exponentially with number of branches
- Reachable code not known initially
- Without guidance, same loop might be unfolded forever

Fitnex search strategy

[Xie et al. DSN 09] http://taoxie.cs.illinois.edu/publications/dsn09-fitnex.pdf



DSE Example

```
public bool TestLoop(int x, int[] y) {
                                                 TestLoop(0, {0})
 if (x == 90) {
       for (int i = 0; i < y.Length; i++)</pre>
              if (y[i] == 15)
                    X++;
        if (x == 110)
                                         Path condition:
              return true;
                                         !(x == 90)
    }
                                               ſ
                                         New path condition:
   return false;
                                         (x == 90)
}
                                         New test input:
```

```
TestLoop(90, {0})
```

DSE Example

```
public bool TestLoop(int x, int[] y) {
                                           TestLoop(90, {0})
 if (x == 90) {
      X++;
       if (x == 110)
            return true;
                               Path condition:
   }
                               (x == 90) \&\& !(y[0] ==15)
  return false;
                               New path condition:
}
                               (x == 90) \&\& (y[0] ==15)
```

```
New test input:
TestLoop(90, {15})
```

Challenge in DSE

```
public bool TestLoop(int x, int[] y) {
                                               TestLoop(90, {15})
 if (x == 90) {
        for (int i = 0; i < y.Length; i++)
              if (y[i] == 15)
                                       Path condition:
                    X++;
       if (x == 110)
                                       (x == 90) \&\& (y[0] == 15)
                                       \&\& !(x+1 == 110)
              return true;
    }
                                       New path condition:
                                       (x == 90) \&\& (y[0] == 15)
   return false;
                                       \&\& (x+1 == 110)
                                              L
                                       New test input:
```

No solution!?

A Closer Look

```
public bool TestLoop(int x, int[] y) {
                                                TestLoop(90, {15})
 if (x == 90) {
        for (int i = 0; i < y.Length; i++)</pre>
              if (y[i] == 15)
                    X++;
                                        Path condition:
     • if (x == 110)
                                        (x == 90) \&\& (y[0] == 15)
                                        && (0 < y.Length)
              return true;
                                        && !(1 < y.Length)
    }
                                        && !(x+1 == 110)
   return false;
                                               ſ
                                        New path condition:
                                        (x == 90) \&\& (y[0] == 15)
                                        && (0 < y.Length)
                                        && (1 < y.Length)
                                        ➔ Expand array size
```

A Closer Look

```
TestLoop(90, {15})
public bool TestLoop(int x, int[] y) {
 if (x == 90) {
        for (int i = 0; i < y.Length; i++)</pre>
               if(y[i] == 15)
                     X++;
                                     We can have infinite paths!
       if (x == 110)
               return true;
                                     Manual analysis \rightarrow need at
    }
                                     least 20 loop iterations to
    return false;
                                     cover the target branch
}
                                     Exploring all paths up to 20
                                     loop iterations is infeasible:
```

2²⁰ paths

Fitnex: Fitness-Guided Exploration

```
TestLoop(90, {15, 0})
TestLoop(90, {15, 15})
```

Key observations: with respect to the coverag target

- not all paths are equally promising for branch-node flipping
- not all branch nodes are equally promising to flip

- Our solution:
 - Prefer to flip branch nodes on the most *promising* paths
 - Prefer to flip the most *promising* branch nodes on paths
 - Fitness function to measure "promising" extents

Fitness Function

- FF computes fitness value (distance between the current state and the goal state)
- Search tries to minimize fitness value

Predicate	Fitness function				
	True	False			
F(a == b)	0	a - b			
F(a > b)	0	(b - a) + K			
$F(a \ge b)$	0	(b-a)			
F(a < b)	0	(a-b)+K			
$F(a \le b)$	0	(a - b)			
$F(P_1 \&\& P_2)$	0	$F(P_1) + F(P_2)$			
$F(P_1 \mid\mid P_2)$	0	$(F(P_1) * F(P_2))/(F(P_1) + F(P_2))$			

[Tracey et al. 98, Liu at al. 05, ...]

Fitness Function for (x == 110)

Compute Fitness Values for Paths

<pre>public bool TestLoop(int x, int[] y) { if (x == 90) {</pre>	(x, y)	Fitness Value
for (int i = 0; i < y.Length; i++)	(90, {0})	20
if (y[i] == 15)	(90, {15})	19
X++;	(90, {15, 0})	19
return true:	(90, {15, 15})	18
}	(90, {15, 15, 0})	18
return false;	(90, {15, 15, 15})	17
}	(90, {15, 15, 15, 0})	17
	(90, {15, 15, 15, 15})	16
Fitness function: 110 – x	(90, {15, 15, 15, 15, 0})	16
	(90, {15, 15, 15, 15, 15})	15

Give preference to flip paths with better fitness values We still need to address which branch node to flip on paths ...

...

Compute Fitness Gains for Branches

```
Fitness Value
public bool TestLoop(int x, int[] y) { (x, y)
 if (x == 90) {
                                             (90, \{0\})
                                                                                      20
    for (int i = 0; i < y.Length; i++)</pre>
                                              (90, {15}) <del><</del> flip b4
                                                                                     19
       if (y[i] == 15)
                                             (90, {15, 0}) ← flip b2
                                                                                      19
          X++;
                                             (90, {15, 15}) ← flip b4
                                                                                      18
    <u>if</u> (x == 110)
                                             (90, {15, 15, 0}) ← flip b2
                                                                                      18
      return true;
                                             (90, {15, 15, 15}) ← flip b4
                                                                                      17
 }
                                             (90, {15, 15, 15, 0}) ← flip b2
                                                                                      17
 return false;
                                             (90, {15, 15, 15, 15}) ← flip b4
                                                                                      16
}
                                             (90, {15, 15, 15, 15, 0}) ← flip b2
        Fitness function: |110 – x |
                                                                                      16
                                             (90, {15, 15, 15, 15, 15}) ← flip b4
                                                                                      15
       Branch b1: i < y.Length
       Branch b2: i >= y.Length
                                        •Flipping Branch b4 (b3) gives us average 1 (-1)
       Branch b3: y[i] == 15
                                        fitness gain (loss)
       Branch b4: y[i] != 15
                                        •Flipping branch b2 (b1) gives us average 0
```

fitness gain (loss)

Compute Fitness Gain for Branches cont.

- For a flipped node leading to *Fnew*, find out the old fitness value *Fold* before flipping
 - Assign Fitness Gain (*Fold Fnew*) for the branch of the flipped node
 - Assign Fitness Gain (*Fnew Fold*) for the other branch of the branch of the flipped node
- Compute the average fitness gain for each branch over time

Search Frontier

- Each branch node candidate for being flipped is prioritized based on its composite fitness value:
 - (Fitness value of node Fitness gain of its branch)
- Select first the one with the best composite fitness value

To avoid local optimal or biases, the fitness-guided strategy is **integrated** with Pex's fairness search strategies

http://taoxie.cs.illinois.edu/publications/dsn09-fitnex.pdf

Automated Test Input Generation for Android: Are We Really There Yet in an Industrial Case?

Xia Zeng₁ **Dengfeng Li**₂ Wujie Zheng₁ Fan Xia₁ **Yuetang Deng**₁ Wing Lam₂ Wei Yang₂ Tao Xie₂ ₁Tencent, Inc., China ₂University of Illinois at Urbana-Champaign, USA





FSE Industry 2016

http://taoxie.cs.illinois.edu/publications/fse16industry-wechat.pdf

Motivation

- Choudhary et al. [ASE'15]: Do we have good-enough tools to test Android apps?
 - Evaluated six research tools and Monkey
 - Monkey tool outperformed all six research tools
- Their study can be further extended
 - No industrial-strength Android app was studied
 - No demonstration on whether and how techniques can further improve Monkey under industrial settings

Challenges on Testing Industrial Mobile Apps

- Challenges for code coverage measurement:
 - Requiring app's source code
 - Industrial-strength Android app can cause 64K reference limit exception during instrumentation
- Challenges for applicability:
 - Scalability on testing apps with large codebases
 - OS compatibility of testing tool

WeChat Overview



- WeChat = WhatsApp+Facebook+Instagram+PayPal+Uber...
- 846 million monthly active users
- Daily number: dozens of billion messages sent, hundreds of million photos uploaded, hundreds of million payment transactions executed
- WeChat backend: 2K+ microservices running on 40k+ servers
 - 10M queries per second during Chinese New Year Eve
- Large codebase on WeChat Android

# of executable Java code lines:	$610,\!629$
# of Java classes:	$8,\!425$
# of Android activities:	607
# of C or C++ code lines:	$\sim \!\! 40,\!000$

Monkey: Experimental setup

- Experiment Setup
 - Set Monkey to fire random events every 500 milliseconds
 - Run Monkey on WeChat 5 times independently
 - Run Monkey for 18 hours each time (2 hours without login)
- Evaluation Metric
 - Line coverage
 - Activity coverage

Monkey: Coverage Result Findings



Finding 1: Monkey has low line coverage (**19.5%**) and low activity coverage (**10.3%**). s

Finding 2: Monkey allocates a lopsided distribution of exploration time on each activity.

Monkey: Exploration time challenges

- Widget obliviousness: It is difficult to generate events at the small-sized GUI element
- State obliviousness: Monkey explores the same two activities repeatedly without contributing to new code coverage



New Approach

- Design goals
 - Have direct access to UI elements on activity under test
 - Allocate more exploration time towards new GUI states
- Techniques
 - Use UIAutomator framework to gain UI layout tree
 - Abstract GUI states to guide firing of state-changing events

New Approach: Coverage result



New approach covers an additional **11.1% p.p.** more lines and **18.4% p.p.** more activities than Monkey does!

Human Intervention



- Provide client-to-client interaction
 - Line coverage increases 1.9% p.p.
 - Activity coverage increases 0.6% p.p..
- Provide additional app events
 - Line coverage increases 0.8% p.p.
 - Activity coverage increases 0.8% p.p.
- Human intervention doesn't help much

Ongoing/Future Work

- Incorporating fine-grained code analysis
 - Construct activity transition graph to guide the tools to explore not-covered activities
- Reusing exploration paths
 - Extracting relevant paths from explored paths and reusing them to boost tools' effectiveness
- Conducting residual coverage analysis
- Getting out of stuck states with manually written rules

Summary

- First industrial case study of applying Monkey on WeChat
- Empirical findings on Monkey's limitations in an industrial setting
- A new approach that addresses the major limitations of Monkey and accomplished substantial code-coverage improvements
- Empirical insights for future enhancements for both Monkey and our approach



Learning for Test Prioritization An Industrial Case Study

Benjamin Busjaeger Salesforce

Tao Xie University of Illinois, Urbana-Champaign

FSE 2016 Industry Track

http://taoxie.cs.illinois.edu/publications/fse16industry-learning.pdf

Large Scale Continuous Integration

Main repository ...





Motivation

	Before Commit	After Commit
Objective	Reject faulty changes	Detect test failures as quickly as possible
Current	Run hand-picked tests	Parallelize (8000 concurrent VMs) & batch (1-1000 changes per run)
Problem	Too complex for human	Feedback time constrained by capacity & batching complicates fault assignment
Desired	Select top-k tests likely to fail	Prioritize all tests by likelihood of failure



Insight: Need Multiple Existing Techniques

- Heterogeneous languages: Java, PL/SQL, JavaScript, Apex, etc.
 - Non-code artifacts: metadata and configuration
 New/recently-failed tests more likely to fail

Test code coverage of change



Textual similarity between test and package configuration; change import org.auraframework.Aura;

/**
 * AuraConfig This is the spring configuration for the aura module.Prov
 * defining runtime implementations here. This class will be loaded by
 */
@AuraConfiguration
public class AuraImplConfig {
 @Impl
 public static FormatAdapter<?> actionJSONFormatAdapter() {
 return new ActionJSONFormatAdapter();
 }
 @Impl
 public static FormatAdapter<?> applicationDefHTMLFormatAdapter() {
 return new ApplicationDefHTMLFormatAdapter();
 }

Test age and recent failures

J												
est Automat	tion -	Remote (Deployments - Other	-								
his is a Flux enabled Autobuild. To find out more about Flux and how it may impact you, read here.												
Autobuile	d Instan	ces for <u>m</u>	ain ftest basic (configs)	FLUX ACTIVE								
Runs for ma	sin_ftest	_basic (T	oggle Admin)									
• Admin A	ictions											
 Yoda Sta 	itus				Bugging Percenta	ige Histo	Res.	alt Proces	sed Result &	Test Fai	ures Proce	ssed Yoda F
Change	Status	Туре	Owner	Start Date	Elapsed	Delay	Ran	Passed	Falled (new)	Errors	Flappers	% Pass
11603030	٩	FULL	pjain	05/27 13:50 POT	32.50%	1h 34m						In Propress
11602880	٩	FULL	rwang	05/27 13:34 807	56.56%	15.38m						In Propress
11602809	٩	FULL	vpola	05/27 13:18 POT	oures	1h 32m	un int	ormati	oni			In Propress
11602690	٩	FULL	msanghavi	05/27 13:01 POT	105.45%	1h 33m						In Propress
11602504	٩	FULL	sarah.lui	05/27 12:44 POT	129.93%	1h 39m						In Propress
11602351	٩	FULL	autointeg (release)	05/27 12:28 POT	153.79%	1h 41m						In Propress
11602121	٩	FULL	jason.zhang	05/27 11:54 POT	205.40%	1h 36m						In Propress
11602011	٩	FULL	vnatarajan	05/27 11:37 POT	229.43%	1h 33m						In Propress
11601827	~	FULL	dbaumann	05/27 11:20 POT	1h 45m	1h 41m	49218	49152	27	36	з	99.86%
11601634	1	8111	0773600	05/22 10:32 807	2h 27m	10.31m	49218	49144	27	16	11	99.84%



Insight: Need Scalable Techniques

- Complex integration tests: concurrent execution
- Large data volume: 500+ changes, 50M+ test runs, 20TB+ results per day

→ Our approach: **Periodically** collect **coarse-grained** measurements

Test code coverage of change

🗎 jacoco > 🌐 com.example.jacoco > 📄 Rectangle.java

Rectangle.java



Textual similarity between test and change

package configuration;

import org.auraframework.Aura;

/**
 * AuraConfig This is the spring configuration for the aura module.Prov
 * defining runtime implementations here. This class will be loaded by
 */
@AuraConfiguration
public class AuraImplConfig {
 @Impl

```
public static FormatAdapter<?> actionJSONFormatAdapter() {
    return new ActionJSONFormatAdapter();
}
```

@Impl

public static FormatAdapter<?> applicationDefHTMLFormatAdapter() {
 return new ApplicationDefHTMLFormatAdapter();
}

Test age and recent failures

This is a Flux enabled Autobuild. To find out more about Flux and how it may impact you, read here.

Runs for main_ftest_basic (Toggie Admin)												
Admin Actions												
Yoda Status Bugging Percentage History Result Processed Result & Test Failures Processed Yoda									ssed Yoda Fi			
Change	Status	Туре	Owner	Start Date	Elapsed	Delay	Ran	Passed	Falled (new)	Errors	Flappers	% Pass
11603030	٩	FULL	pjain	05/27 13:50 POT	32.50%	1h 34m						In Progress
11602880	0	FULL	rwang	05/27 13-34 807	56.56%	1h 38m						In Progress
11602809	0	FULL	vpola	05/27 13:18 PDT	his cell for more	1h 32m	in infi	ormati	oni			In Progress
11602690	(1)	FULL	msanghavi	05/27 13:01 POT	105.45%	1h 33m						In Progress
11602504	(1)	FULL	sarah./ul	05/27 12:44 POT	129.93%	1h 39m						In Progress
11602351	0	FULL	autointeg (release)	05/27 12:28 POT	153.79%	1h 41m						In Propress
11602121	٩	FULL	jason.zhang	05/27 11:54 POT	205.40%	1h 36m						In Propress
11602011	٩	FULL	vnatarajan	05/27 11:37 POT	229.43%	1h 33m						In Progress
11601827	~	FULL	dbaumann	05/27 11:20 POT	1h 45m	1h 41m	49218	49152	27	36	з	99.86%
11601634	1	FULL	pmason	05/27 10:37 PDT	2h 27m	1h 31m	49218	49144	27	36	11	22.84%



New Approach: Learning for Test Prioritization



→ Implementation currently in pilot use at Salesforce



Empirical Study: Setup

Test results of 45K tests for ~3 month period
In this period, 711 changes with ≥1 failure

- 440 for training
- 271 for testing
- Collected **once** for each test:
 - Test code coverage
 - Test text content
- Collected **continuously**:
 - Test age and recent failures



Results: Test selection (before commit)



- New approach achieves highest average recall at all cutoff points
 - 50% failures detected with top 0.2%
 - 75% failures detected with top 3%



Results: Test prioritization (after commit)



- New approach achieves highest APFD with least variance
 - Median: 85%
 - Average: 99%





- Main insights gained in conducting test prioritization in industry
- Novel learning-based approach to test prioritization
- Implementation currently in pilot use at Salesforce
- Empirical evaluation using a large Salesforce dataset



Getting Real to Produce Practice Impact

Software testing/analysis tools are naturally tied with software development practice



Thank you! Questions?

Getting Real to Produce Practice Impact

Software testing/analysis tools are naturally tied with software development practice

