

Statistical Debugging for Real-World Performance Problems

Linhai Song¹ and Shan Lu²
¹University of Wisconsin-Madison
²University of Chicago

What are Performance Problems?

- Definition of Performance Problems (PPs):
 - Implementation mistakes causing inefficiency
- An example

```

void ha_partition::start_bulk_insert(int rows) {
    .....
    - if (!rows) //check whether rows is 0
      - DEBUG_VOID_RETURN;
      - rows= rows/m_tot_parts + 1;
      + rows= rows ? rows/m_tot_parts + 1 : 0;
      ..... // fast path using caches
}
                
```

MySQL Bug 26527

Performance Diagnosis

Performance Diagnosis

- Identifying the causes of performance failures
 - In-house and on-line diagnosis

Performance Diagnosis is Challenging

- The state of the art is **preliminary**
 - Profilers
 - Only tools mentioned in bug reports we studied
 - Output time-consuming functions, not root causes
- More effective tools are necessary

```

void ha_partition::start_bulk_insert(int rows) {
    .....
    - if (!rows)
      - DEBUG_VOID_RETURN;
      - rows= rows/m_tot_parts + 1;
      + rows= rows ? rows/m_tot_parts + 1 : 0;
      ..... // fast path using caches
}
                
```

MySQL Bug 26527

How Can We Do Better?

Can we learn from functional bug diagnosis?

How to Diagnose Functional Bugs

- The state of the art is **mature**
- Statistical Debugging(SD)
 - Among the most effective

The diagram illustrates the process of statistical debugging. It starts with 'Input' consisting of 'Good' and 'Bad' examples. These are fed into a 'Program' (represented by a code snippet: `int main(int argc, char ** argv) { if () { } }`). The program's output is categorized into 'Predicates' (e.g., 'taken', 'not-taken') and 'Symptoms' (represented by smile and frowny faces). These predicates and symptoms are then processed by a 'Statistical Model' to produce a 'Rank' table with columns for 'Rank', 'Predicates', and 'Score'. The table shows a 'taken' predicate with a score of 1.

Remember these examples?

```

1 // Print_tokens2 v7
2 if(ch == '\n')
3     return (TRUE);
4 else if(ch == ' ')
5     // Bug: should return FALSE
6     return (TRUE);
7 else
8     return (FALSE);
    
```

How to Diagnose Functional Bugs

- The state of the art is **mature**
- Statistical Debugging(SD)
 - Among the most effective

This diagram is identical to the one in the top-left slide, showing the flow from input through a program to a statistical model and a table of results.

Apply SD to Performance Diagnosis?

Since statistical debugging is effective for functional diagnosis, maybe it will also be effective for performance diagnosis.

What are the challenges?

Is it Feasible?

- Q1: How to tell success runs from failure runs?
- Q2: How to obtain good and bad inputs?

The diagram is similar to the top-left slide, but with red boxes highlighting the 'Q1?' and 'Q2?' questions. 'Q1?' is placed near the 'Symptom' icons, and 'Q2?' is placed near the 'Input' icons.

How to Apply?

- Q3: What predicates to collect?
- Q4: What statistical model to use?

Input: Good, Bad

Program: `int main(int argc, char ** argv) { ... }`

Symptom: (Smiley face), (Sad face)

Q3? Predicates: taken, not-taken

Q4? Statistical Model

Rank	Predicates	Score
1	taken	...

How to Apply?

- Q3: What predicates to collect?
- Q4: What statistical model to use?
- Q5: How to do on-line performance diagnosis?

Input: Good, Bad

Program: (Group of people icon)

Symptom: (Smiley face), (Sad face)

Q2? (Input)

Q3? Predicates: taken, not-taken

Q4? Statistical Model

Rank	Predicates	Score
1	taken	...

Contributions (I)

We answer all these questions by studying 65 user-reported performance bugs.

Input: Good, Bad

Program: (Group of people icon)

Symptom: (Smiley face), (Sad face)

Q1? (Input/Symptom)

Q5? (Program)

Q3? Predicates: taken, not-taken

Q4? Statistical Model

Rank	Predicates	Score
1	taken	...

Contributions (II)

- Is it feasible to apply SD to PPs?
 - Easy to tell failure runs based on users' reports (Q1)
 - Inputs are provided during reporting (Q2)
- How to apply SD to PPs?
 - In-house diagnosis
 - 3 predicates (Q3)
 - 2 statistical models (Q4)
 - On-line diagnosis (Q5)
 - Same diagnosis capability with <10% overhead
 - Not sacrifice diagnosis latency (**Unique**)

Outline

- Overview
- Is it feasible to apply SD for PPs?
- How to conduct SD for PPs?
 - In-house diagnosis scenario
 - On-line diagnosis scenario
- Conclusion

Methodology

- Application and Bug Source

App.	# Bugs ^[1]	# Bug User Perceived
Apache	25	16
Chrome	10	5
GCC	11	9
Mozilla	36	19
MySQL	28	17
Total: 110		65

[1] Guoliang Jin, Linhai Song, Xiaoming Shi, Joel Scherpelz, and Shan Lu. Understanding and Detecting Real-World Performance Bugs. In PLDI'2012.

Is It Feasible? (Part I)

- Q1: How to tell success runs from failure runs?**
 - A large workload? Or inefficient implementation?

Input:

Good

Bad

Program:

```
int main(int argc, char ** argv) {
    #if 0
    // ...
    #endif
}
```

Predicates:

taken 😊

Statistical Model

Symptom:

😊

😞

Q1?

Predicates:

not-taken 😞

Rank	Predicates	Score
1	taken 😊

19

Q1: How to Identify Failure Runs?

- The majority is observed through comparison

Dominating

Category	MySQL	Mozilla	GCC	Chrome	Apache
within one code base	10	10	10	10	10
cross multiple code bases	5	5	5	5	5
Not using comparison	10	10	10	10	10

20

Comparison within One Code Base

- the same input with different **configuration**
- inputs with different **sizes**
- inputs with slightly different **functionality**

Category	MySQL	Mozilla	GCC	Chrome	Apache
within one code base	10	10	10	10	10
cross multiple code bases	5	5	5	5	5
Not using comparison	10	10	10	10	10

21

Comparison across Multiple Code Bases

- same applications' different **versions**
- different **applications**

Category	MySQL	Mozilla	GCC	Chrome	Apache
within one code base	10	10	10	10	10
cross multiple code bases	5	5	5	5	5
Not using comparison	10	10	10	10	10

22

Not Using Comparison

Mozilla#299742: "it frozen the GUI to crawl!"

Mozilla#299742: "Easy to tell failures from successes!" in the page"

MySQL#46461: "causing the test suit to fail due to timeout"

Non-tolerable

Category	MySQL	Mozilla	GCC	Chrome	Apache
within one code base	10	10	10	10	10
cross multiple code bases	5	5	5	5	5
Not using comparison	10	10	10	10	10

23

Is It Feasible? (Part II)

- Q1: How to tell success runs from failure runs?**
- Q2: How to obtain good and bad inputs?**

Input:

Good

Bad

Q2?

Program:

```
int main(int argc, char ** argv) {
    #if 0
    // ...
    #endif
}
```

Predicates:

taken 😊

Statistical Model

Symptom:

😊

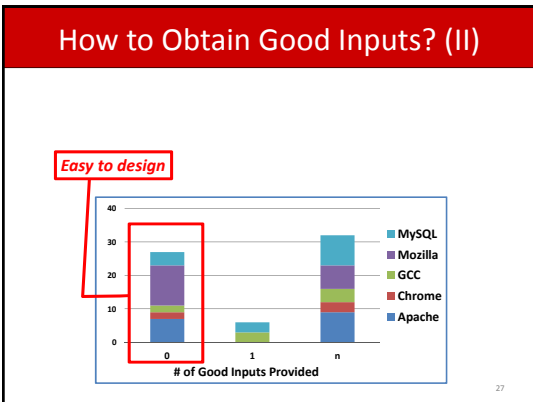
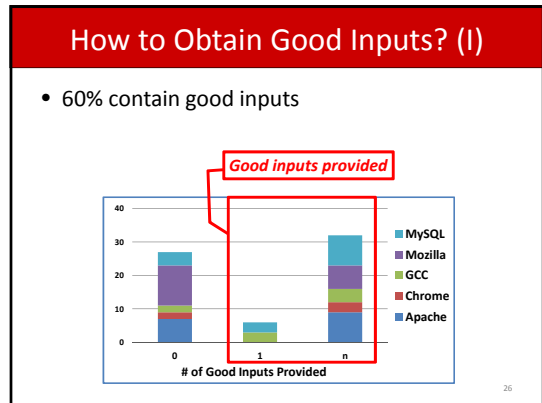
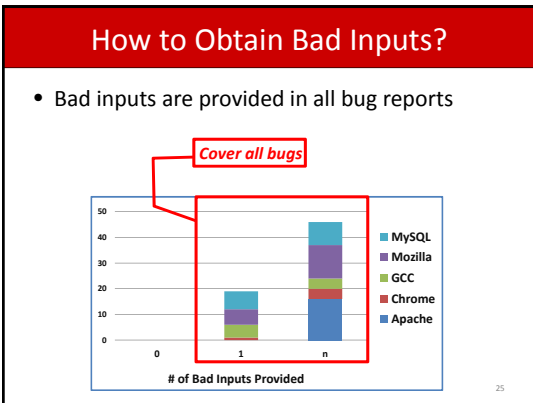
😞

Predicates:

not-taken 😞

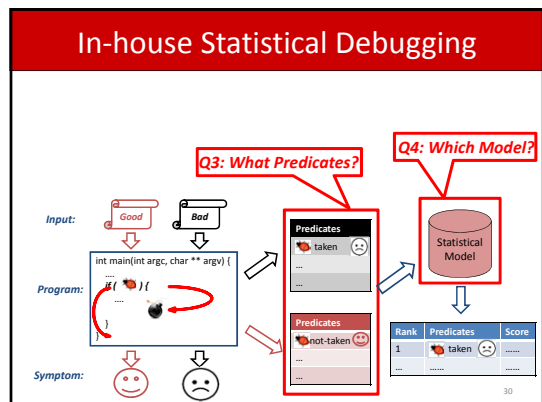
Rank	Predicates	Score
1	taken 😊

24



- ### Other Findings and Implications
- Compared with functional bugs
 - More PPs observed through comparison
 - More PPs reported with good inputs
 - Implications for SD
 - Easy to tell success runs from failure runs
 - Similar good inputs are provided
 - SD is a nature fit for PPs

- ### Outline
- Overview
 - Is it feasible to apply SD for PPs?
 - How to conduct SD for PPs?
 - In-house diagnosis scenario
 - On-line diagnosis scenario
 - Conclusion



Design

- Study widely used predicates and models

```
int x, y;
...
x = ...;
n=fopen(...);
if (p) ...
else ...
```

31

Experimental Methodology

- Experimental setting
 - 10 success runs vs. 10 failure runs
 - 20 unique inputs
- Techniques under comparison
 - CBI[2, 3] for C programs
 - Pin for C++ programs
 - Compared with profiling results from OProfile

[2] Ben Liblit, Alex Aiken, Alice X Zhen, and Michael I Jordan. Bug Isolation via Remote Program Sampling. In PLDI'2003.
 [3] Ben Liblit, Mayur Naik, Alice X Zheng, Alex Aiken, and Michael I Jordan. Scalable Statistical Bug Isolation. In PLDI'2005.

32

Experimental Results

BugID	Basic Model			ΔLDA			Profiler
	Branch	Return	S-pair	Branch	Return	S-pair	
Mozilla258793	v1	-	/	-	-	/	-
Mozilla299742	v1	-	/	-	-	/	-
Mozilla347306	-	-	-	v1	v1	v1	v1
Mozilla416628	-	-	-	v1	-	v1	v1
MySQL15811	-	-	/	v1	v1	/	v1
MySQL26527	v1	-	/	-	-	/	-
MySQL27287	-	-	/	v1	-	/	v1
MySQL40337	v1	-	/	-	-	/	-
MySQL42649	v1	-	/	-	-	/	-
MySQL44723	v1	-	/	-	-	/	-
Apache3278	v1	v1	v1	-	-	-	-
Apache34464	-	-	-	v3	v1	-	v5
...

33

Experimental Results

BugID	Basic Model			ΔLDA			Profiler
	Branch	Return	S-pair	Branch	Return	S-pair	
Mozilla258793	v1	-	/	-	-	/	-
Mozilla299742	v1	-	/	-	-	/	-
Mozilla347306	-	-	-	v1	v1	v1	v1
Mozilla416628	-	-	-	v1	-	v1	v1
MySQL15811	-	-	/	v1	v1	/	v1
MySQL26527	v1	-	/	-	-	/	-
MySQL27287	-	-	/	v1	-	/	v1
MySQL40337	v1	-	/	-	-	/	-
MySQL42649	v1	-	/	-	-	/	-
MySQL44723	v1	-	/	-	-	/	-
Apache3278	v1	v1	v1	-	-	-	-
Apache34464	-	-	-	v3	v1	-	v5
...

34

Experimental Results

BugID	Basic Model			ΔLDA			Profiler
	Branch	Return	S-pair	Branch	Return	S-pair	
Mozilla258793	v1	-	/	-	-	/	-
Mozilla299742	v1	-	/	-	-	/	-
Mozilla347306	-	-	-	v1	v1	v1	v1
Mozilla416628	-	-	-	v1	-	v1	v1
MySQL15811	-	-	/	v1	v1	/	v1
MySQL26527	v1	-	/	-	-	/	-
MySQL27287	-	-	/	v1	-	/	v1
MySQL40337	v1	-	/	-	-	/	-
MySQL42649	v1	-	/	-	-	/	-
MySQL44723	v1	-	/	-	-	/	-
Apache3278	v1	v1	v1	-	-	-	-
Apache34464	-	-	-	v3	v1	-	v5
...

35

Experimental Results

BugID	Basic Model			ΔLDA			Profiler
	Branch	Return	S-pair	Branch	Return	S-pair	
Mozilla258793	v1	-	/	-	-	/	-
Mozilla299742	v1	-	/	-	-	/	-
Mozilla347306	-	-	-	v1	v1	v1	v1
Mozilla416628	-	-	-	v1	-	v1	v1
MySQL15811	-	-	/	v1	v1	/	v1
MySQL26527	v1	-	/	-	-	/	-
MySQL27287	-	-	/	v1	-	/	v1
MySQL40337	v1	-	/	-	-	/	-
MySQL42649	v1	-	/	-	-	/	-
MySQL44723	v1	-	/	-	-	/	-
Apache3278	v1	v1	v1	-	-	-	-
Apache34464	-	-	-	v3	v1	-	v5
...

36

Experimental Results

BugID	Basic Model			ALDA			Profiler
	Branch	Return	S-pair	Branch	Return	S-pair	
Mozilla258793	v1	-	/	-	-	/	-
Mozilla299742	v1	-	/	-	-	/	-
Mozilla347306	-	-	-	v1	v1	v1	v1
Mozilla416628	-	-	-	v1	-	v1	v1
MySQL15811	-	-	/	v1	v1	/	v1
MySQL26527	v1	-	/	-	-	/	-
MySQL27287	-	-	/	v1	-	/	v1
MySQL40337	v1	-	/	-	-	/	-
MySQL42649	v1	-	/	-	-	/	-
MySQL44723	v1	-	/	-	-	/	-
Apache3278	v1	v1	v1	-	-	-	-
Apache34464	-	-	-	v3	v1	-	v5
...

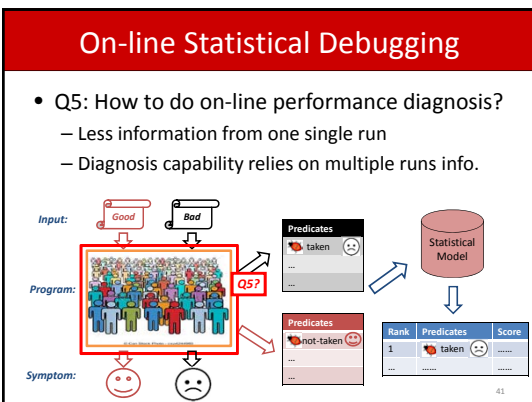
Experimental Results

BugID	Basic Model			ALDA			Profiler
	Branch	Return	S-pair	Branch	Return	S-pair	
Mozilla258793	v1	-	/	-	-	/	-
Mozilla299742	v1	-	/	-	-	/	-
Mozilla347306	-	-	-	v1	v1	v1	v1
Mozilla416628	-	-	-	v1	-	v1	v1
MySQL15811	-	-	/	v1	v1	/	v1
MySQL26527	v1	-	/	-	-	/	-
MySQL27287	-	-	/	v1	-	/	v1
MySQL40337	v1	-	/	-	-	/	-
MySQL42649	v1	-	/	-	-	/	-
MySQL44723	v1	-	/	-	-	/	-
Apache3278	v1	v1	v1	-	-	-	-
Apache34464	-	-	-	v3	v1	-	v5
...

Experimental Results

BugID	Basic Model			ALDA			Profiler
	Branch	Return	S-pair	Branch	Return	S-pair	
Mozilla258793	v1	-	/	-	-	/	-
Mozilla299742	v1	-	/	-	-	/	-
Mozilla347306	-	-	-	v1	v1	v1	v1
Mozilla416628	-	-	-	v1	-	v1	v1
MySQL15811	-	-	/	v1	v1	/	v1
MySQL26527	v1	-	/	-	-	/	-
MySQL27287	-	-	/	v1	-	/	v1
MySQL40337	v1	-	/	-	-	/	-
MySQL42649	v1	-	/	-	-	/	-
MySQL44723	v1	-	/	-	-	/	-
Apache3278	v1	v1	v1	-	-	-	-
Apache34464	-	-	-	v3	v1	-	v5
...

- ### Outline
- Overview
 - Is it feasible to apply SD for PPs?
 - How to conduct SD for PPs?
 - In-house diagnosis scenario
 - On-line diagnosis scenario
 - Conclusion



- ### Experimental Methodology
- Tool design
 - CBI in sampling mode for C benchmarks
 - LBR for C++ benchmarks
 - Benchmarks
 - Reuse benchmarks from in-house experiments
 - Most effective predicate and model
 - Experiment design
 - Default sampling rate is roughly 1/10000
 - 1000 success runs and 1000 failure runs

Experimental Results

BugID	Diagnosis Capability	Run-time Overhead	Requested Failure Runs
Mozilla258793	v1	2.39%	100
Mozilla299742	v1	4.27%	500
Mozilla347306	v1	1.42%	10
Mozilla416628	v1	2.03%	10
MySQL15811	v1	2.25%	10
MySQL26527	v1	6.05%	500
MySQL27287	v1	3.02%	10
MySQL40337	v1	2.69%	100
MySQL42649	v2	6.10%	500
MySQL44723	v1	3.16%	100
Apache3278	-	0.23%	>1000
Apache34464	v3	0.18%	10
...

Experimental Results

BugID	Diagnosis Capability	Run-time Overhead	Requested Failure Runs
Mozilla258793	v1	2.39%	100
Mozilla299742	v1	4.27%	500
Mozilla347306	v1	1.42%	10
Mozilla416628	v1	2.03%	10
MySQL15811	v1	2.25%	10
MySQL26527	v1	6.05%	500
MySQL27287	v1	3.02%	10
MySQL40337	v1	2.69%	100
MySQL42649	v2	6.10%	500
MySQL44723	v1	3.16%	100
Apache3278	-	0.23%	>1000
Apache34464	v3	0.18%	10
...

Same Diagnosis Capability

< 10%

Experimental Results

BugID	Diagnosis Capability	Run-time Overhead	Requested Failure Runs
Mozilla258793	v1	2.39%	100
Mozilla299742	v1	4.27%	500
Mozilla347306	v1	1.42%	10
Mozilla416628	v1	2.03%	10
MySQL15811	v1	2.25%	10
MySQL26527	v1	6.05%	500
MySQL27287	v1	3.02%	10
MySQL40337	v1	2.69%	100
MySQL42649	v2	6.10%	500
MySQL44723	v1	3.16%	100
Apache3278	-	0.23%	>1000
Apache34464	v3	0.18%	10
...

10 X 10000 ??

Experimental Results

BugID	Diagnosis Capability	Run-time Overhead	Requested Failure Runs
Mozilla258793	v1	2.39%	100
Mozilla299742	v1	4.27%	500
Mozilla347306	v1	1.42%	10
Mozilla416628	v1	2.03%	10
MySQL15811	v1	2.25%	10
MySQL26527	v1	6.05%	500
MySQL27287	v1	3.02%	10
MySQL40337	v1	2.69%	100
MySQL42649	v2	6.10%	500
MySQL44723	v1	3.16%	100
Apache3278	-	0.23%	>1000
Apache34464	v3	0.18%	10
...

Experimental Results

BugID	Diagnosis Capability	Run-time Overhead	Requested Failure Runs
Mozilla258793	v1	2.39%	100
Mozilla299742	v1	4.27%	500
Mozilla347306	v1	1.42%	10
Mozilla416628	v1	2.03%	10
MySQL15811	v1	2.25%	10
MySQL26527	v1	6.05%	500
MySQL27287	v1	3.02%	10
MySQL40337	v1	2.69%	100
MySQL42649	v2	6.10%	500
MySQL44723	v1	3.16%	100
Apache3278	-	0.23%	>1000
Apache34464	v3	0.18%	10
...

Why?

- ### Conclusion and Future Works
- Study diagnosis process for PPs
 - Statistical debugging is a natural fit
 - Study statistical debugging on PPs
 - Branch predicates + two statistical models
 - Future works
 - Analyze inefficient loops
 - Provide detailed fix hints

Break



49