

Diversity Maximization Speedup for Fault Localization

Liang Gong¹, David Lo², Lingxiao Jiang², and Hongyu Zhang¹

Authored by



Outline

Introduction & Framework

- **Motivation**
 - Fault Localization
 - Test Case Prioritization

- **Diversity Maximization Speedup**
 - Technical Motivation
 - Detailed Approach

- **Experiments**
 - Settings & Results

- **Conclusion & Future work**

Debugging

Problem

- Software errors cost the US economy 59.5 billion dollars (0.6% of 2002's GDP) [1]
- Testing and debugging activities are labor-intensive (30% to 90% of a Project) [2]



[1] National Institute of Standards and Technology (NIST). Software Errors Cost U.S. Economy \$59.5 Billion Annually, June 28, 2002.

[2] B. Beizer. Software Testing Techniques. International Thomson Computer Press, Boston, 2nd edition, 1990.

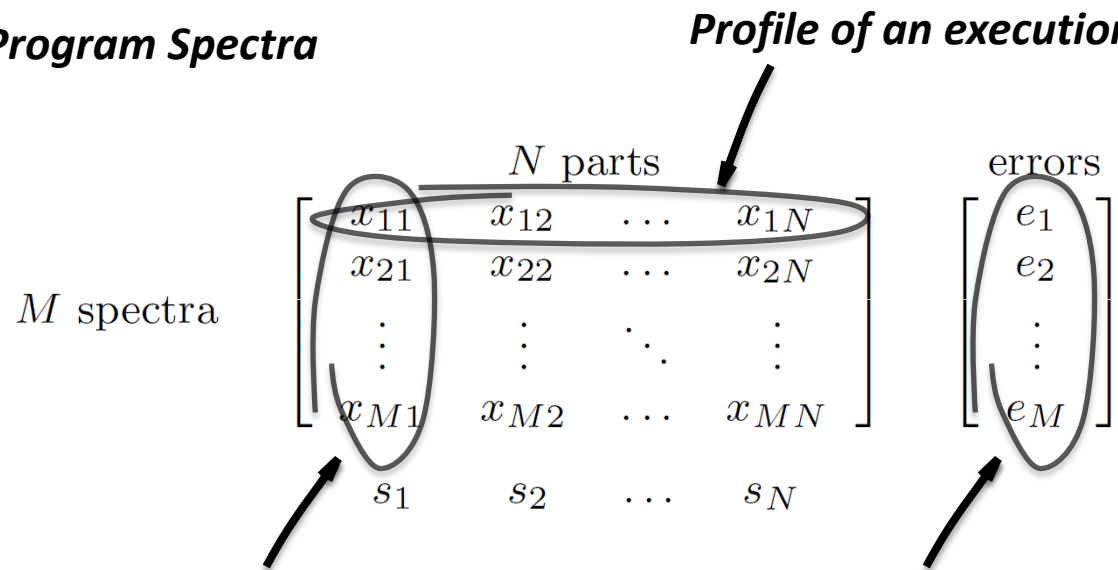
Spectrum-based Fault Localization (abbr. SBFL)

- Automatically recommend a list of suspicious program elements for inspection.
- **Program Spectra** consists of coverage information and execution labels.

SBFL

Introduction

Program Spectra



Coverage information of one element (s_i) in all executions

Correct or incorrect?

Approaches

Fault Localization

For a given statement S

The formula calculates the suspiciousness of S .

No. of failed traces covering S

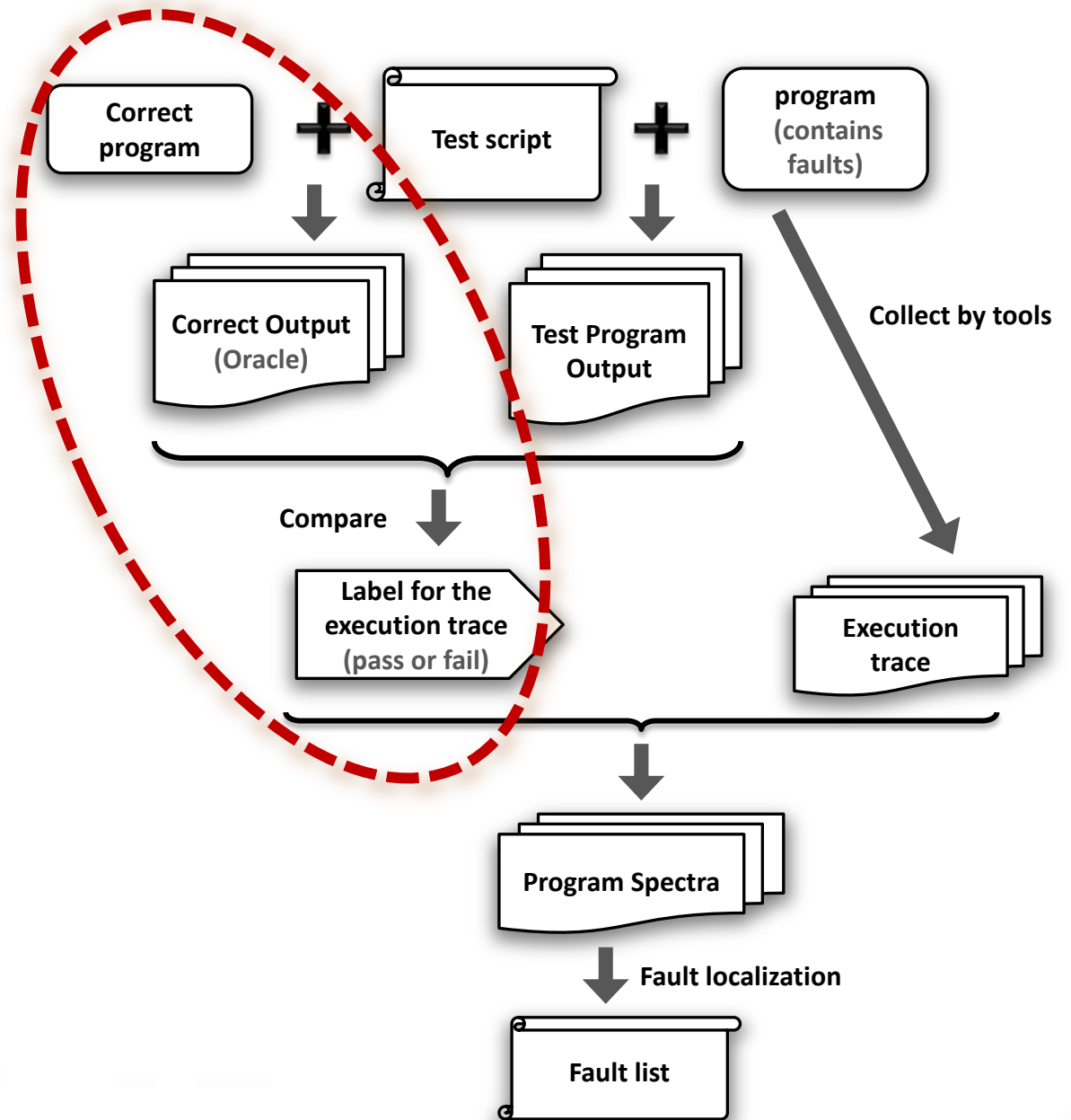
Ochiai

$$\frac{\bar{a}_{ef}}{\sqrt{(a_{ef} + a_{nf}) \cdot (a_{ef} + a_{ep})}}$$

No. of failed traces No. of traces covering S

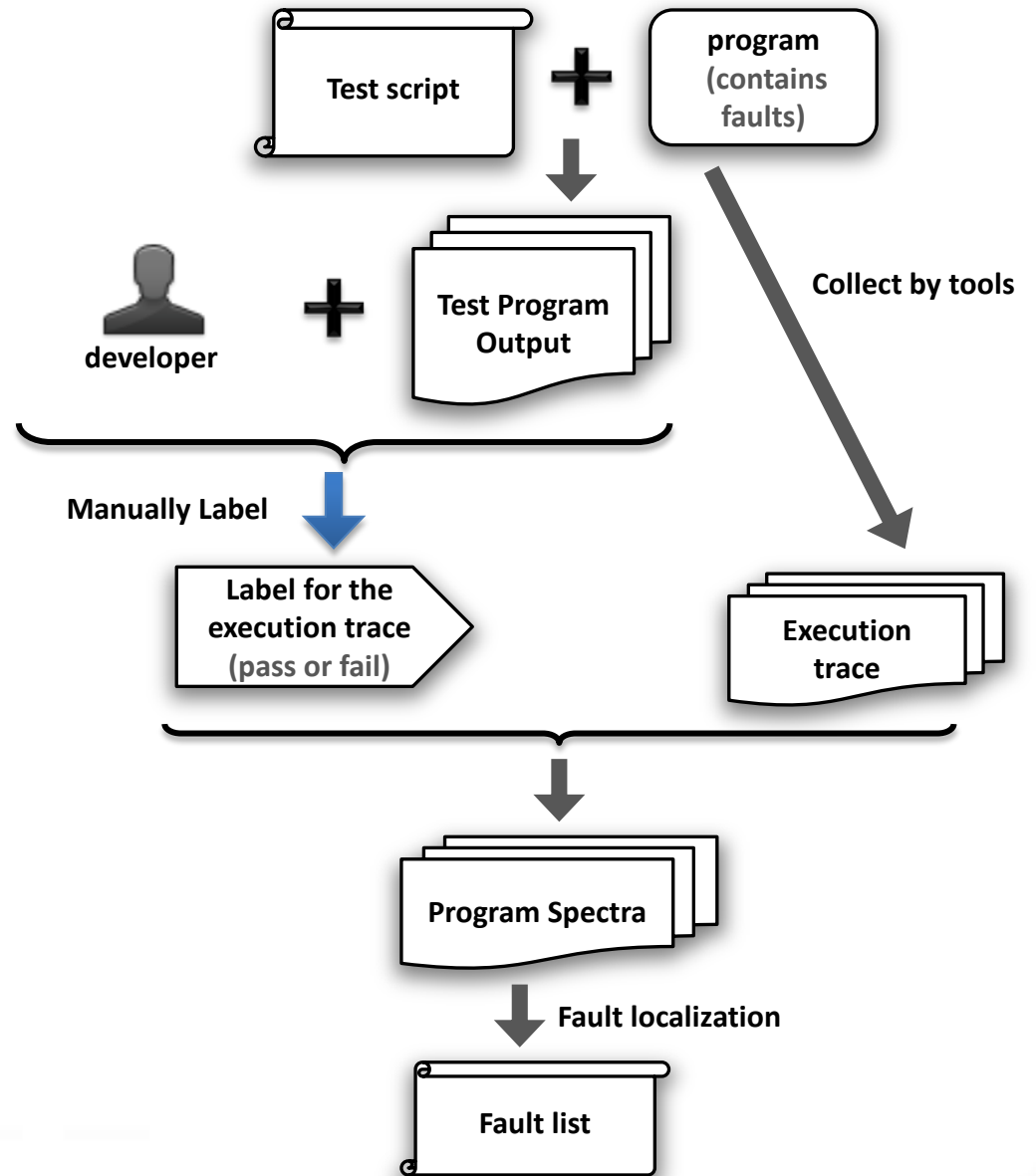
Intuition: If S is covered **more** in **failed** traces and **less** in **passed** traces, it is more likely to contain faults.

Process *In Experiments*



Process

In Practice

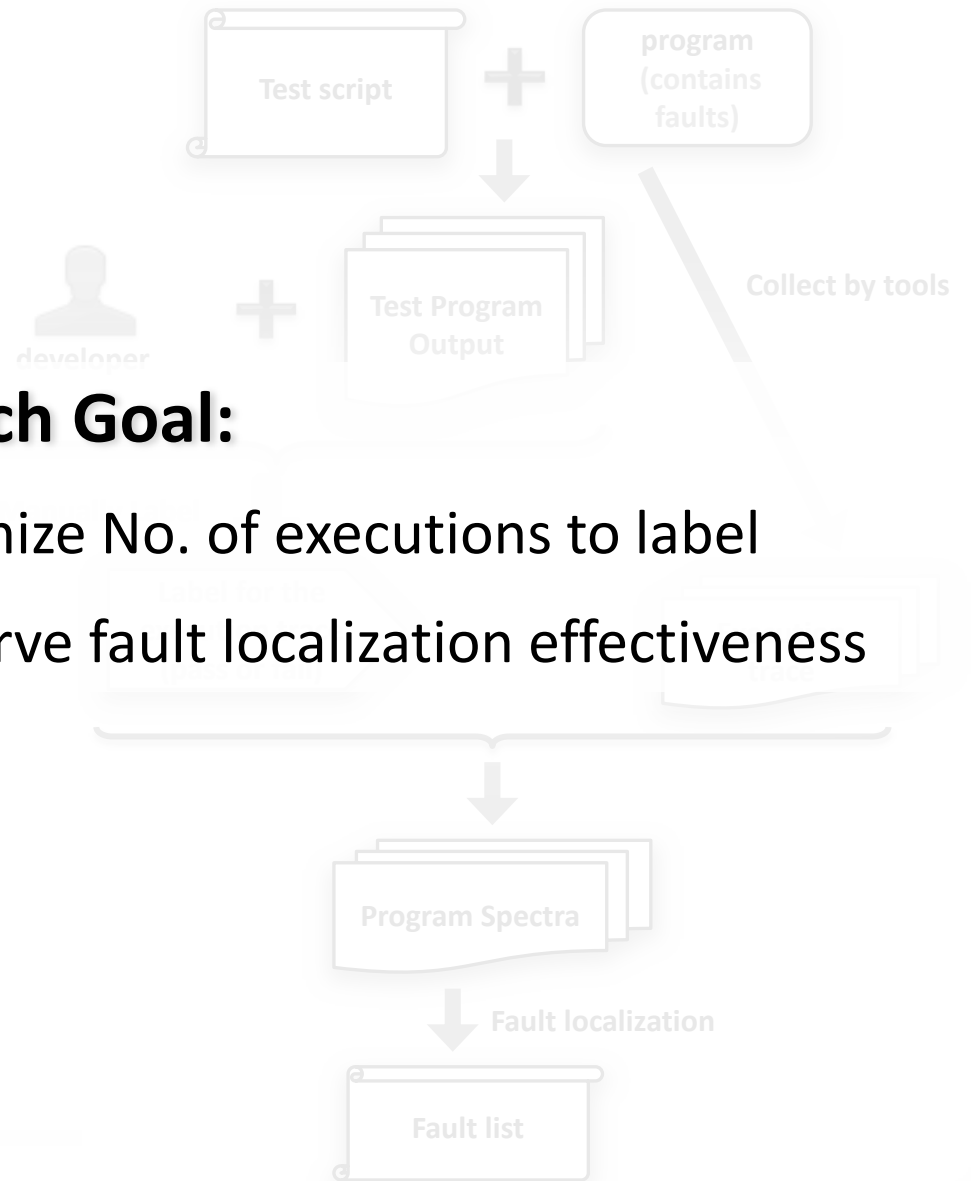


Process

In Practice

Research Goal:

- Minimize No. of executions to label
- Preserve fault localization effectiveness



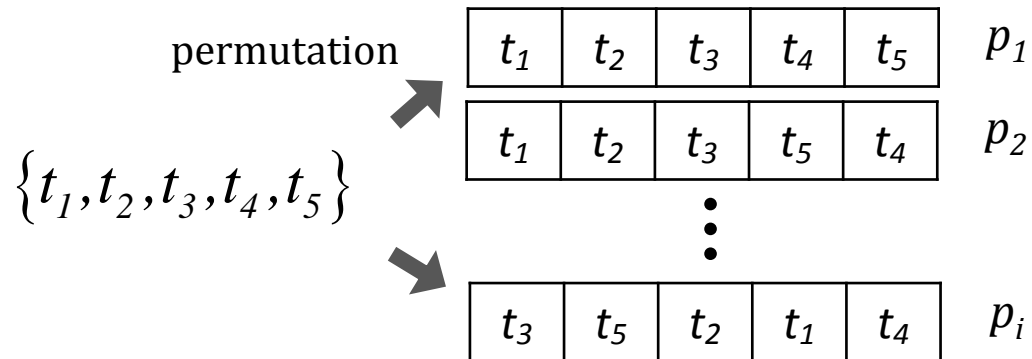
Test Case Prioritization

Introduction

In [3], Rothermel *et al.* define the problem of test case prioritization as follows:

Definition 2.1 (*Test Case Prioritization*). Given

- (1) T , a set of test cases,
 - (2) PT , the set of permutations of T
 - (3) f , a function mapping PT to real numbers,
- the problem is to find a permutation $p \in PT$ such that:
for all $p' \in PT: f(p) \geq f(p')$.



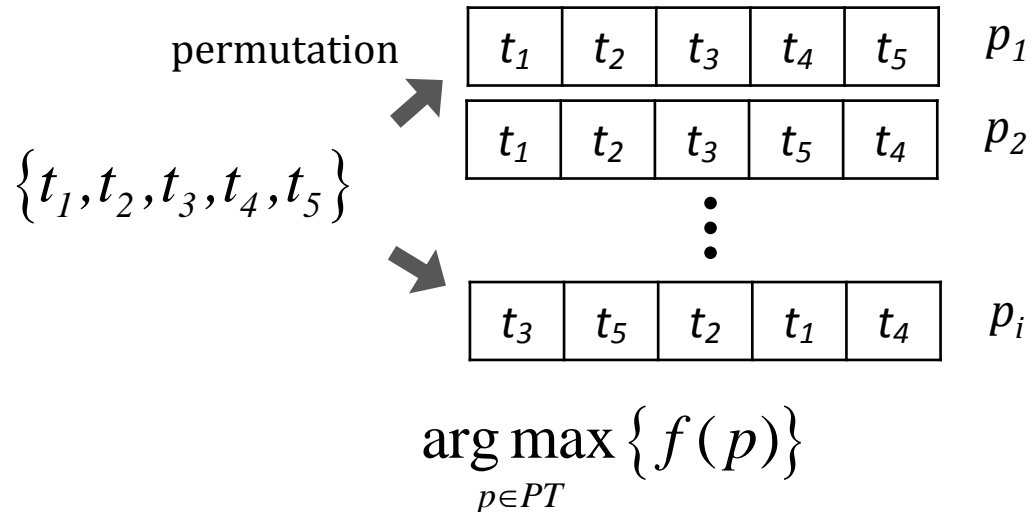
$$\arg \max_{p \in PT} \{f(p)\}$$

Test Case Prioritization

Introduction

In [3], Rothermel *et al.* define the problem of test case prioritization as follows:

Definition 2.1 (Test Case Prioritization) Given $f(p)$ is larger, when permutation p allows the **faulty** program elements to be **ranked higher** meanwhile a **shorter prefix** are considered.
for all $p' \in PT: f(p) \geq f(p')$.



DMS

Introduction

Diversity Maximization Speedup (*abbr.* DMS)

- Greedy algorithm
- Use **diversity of suspiciousness** as the selecting criterion.
- **Speedup** suspiciousness rank changing process of **promising** program elements to further save labeling effort.

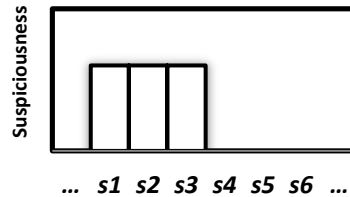
Diversity Maximization speedup (*abbr.* DMS)

- t_0 is the initial failed trace that reveals the fault.
- t_1 and t_2 are candidates to be selected for labeling.

t_0	...	S_1	S_2	S_3	S_4	S_5	S_6	...	<i>p/f</i>
	...	●	●	●	○	○	○	...	<i>fail</i>

Diversity

As Criterion



S_1, S_2, S_3
S_4, S_5, S_6

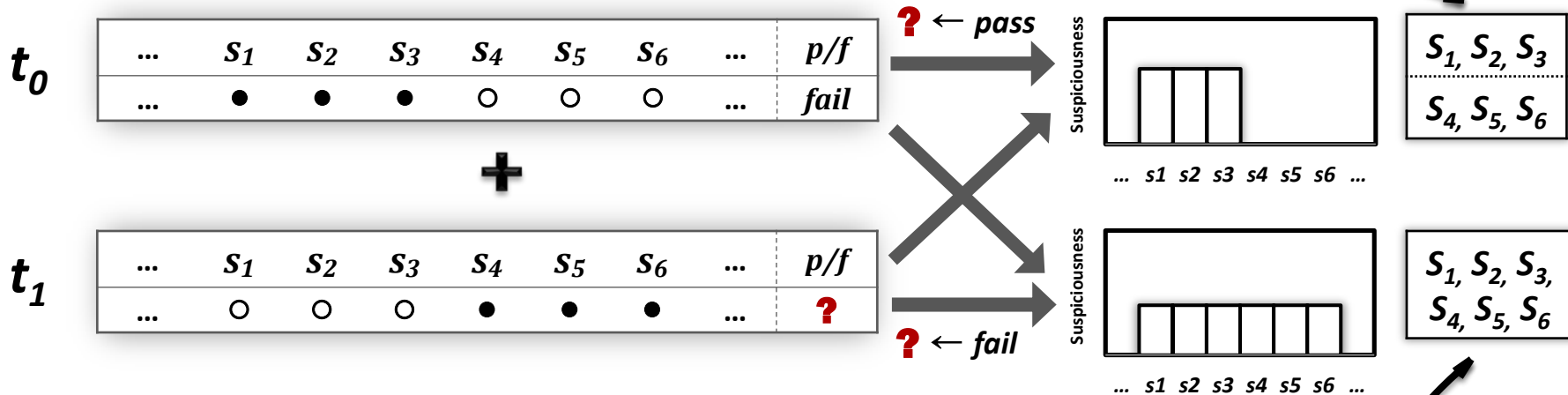
Initial Diagnostic Distribution

Initial Inspection List

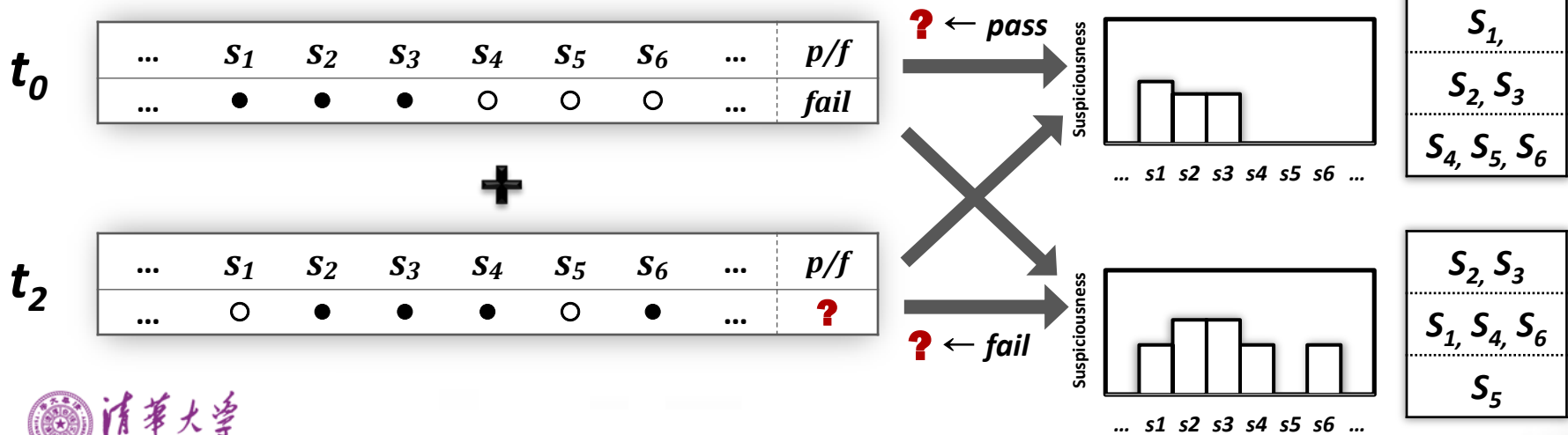
t_1	...	S_1	S_2	S_3	S_4	S_5	S_6	...	<i>p/f</i>
	...	○	○	○	●	●	●	...	?

t_2	...	S_1	S_2	S_3	S_4	S_5	S_6	...	<i>p/f</i>
	...	○	●	●	●	○	●	...	?

- t_1 is preferred by approaches aiming to detect faults

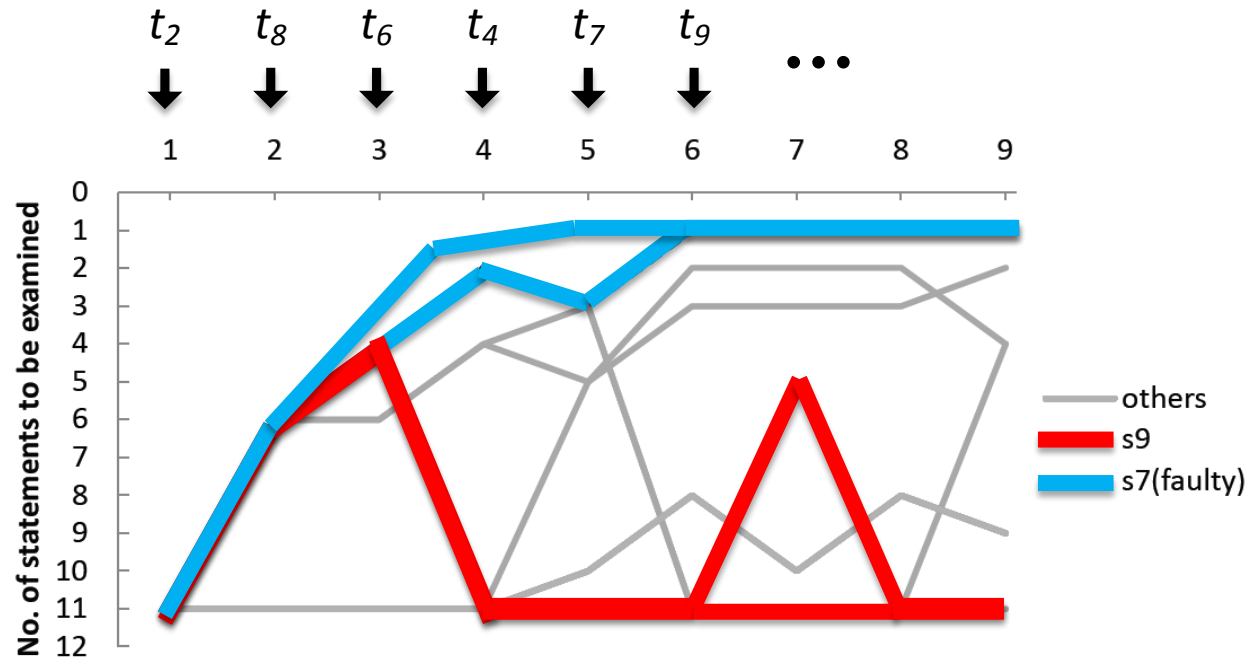


- t_2 is selected according to Diversity Maximization Criterion



☺ **More Diverse, Better Result**

Evolution Trend
Opportunities



Evolution Trend of
Statements' Suspiciousness.

Looking for test cases that could offer more **changing opportunities** to "promising" elements like s7 (with clear trend) — instead of s9 —

Speed up

How to?

Two questions prompt:

- How can we know which statements are “**promising**”?
- With “promising” statements, how can we **speed up** their suspiciousness changing process?

Promising

How to identify?

Representative Time Series

- When the rank of the program element **decreases**, its time series **increases** by 1.
- When the rank of the program element **increases**, its time series **decreases** by 1.
- If the element's rank stays the same, its time series stays the same.

Evolution trend and time series(y_i) of S_8

Iteration (x_i)	1	2	3	4	5	6	7	...
Rank	11	6	4	2	3	11	5	...
Trend (\mathcal{T})		[+]	[+]	[+]	[-]	[-]	[+]	...
y_i	0	1	2	3	2	1	2	...

Promising

How to evaluate?

Evolution trend of S_g

Iteration (x_i)	1	2	3	4	5	6	7	...
Rank	11	6	4	2	3	11	5	...
Trend (\mathcal{T})		[+]	[+]	[+]	[-]	[-]	[+]	...
y_i	0	1	2	3	2	1	2	...

- Linear Regression Analysis:

$$y_i = \beta_1 \cdot x_i + \beta_0 + \epsilon_i$$

- Change-potential Score:

$$W_{\mathcal{T}} = \left| \hat{\beta}_1 \right| \cdot \frac{1}{\hat{\sigma}_{\beta_1} + 1}$$

Changing Rate
Changing Stability

Example trends and their potentials

\mathcal{T}	$\hat{\beta}_1$	$\hat{\sigma}_{\beta_1}$	$W_{\mathcal{T}}$
[+] [+]	1	0	1
[+] [-]	0	0.577	0
[+] [0]	0.5	0.289	0.388
[0] [0]	0	0	0

Speed up

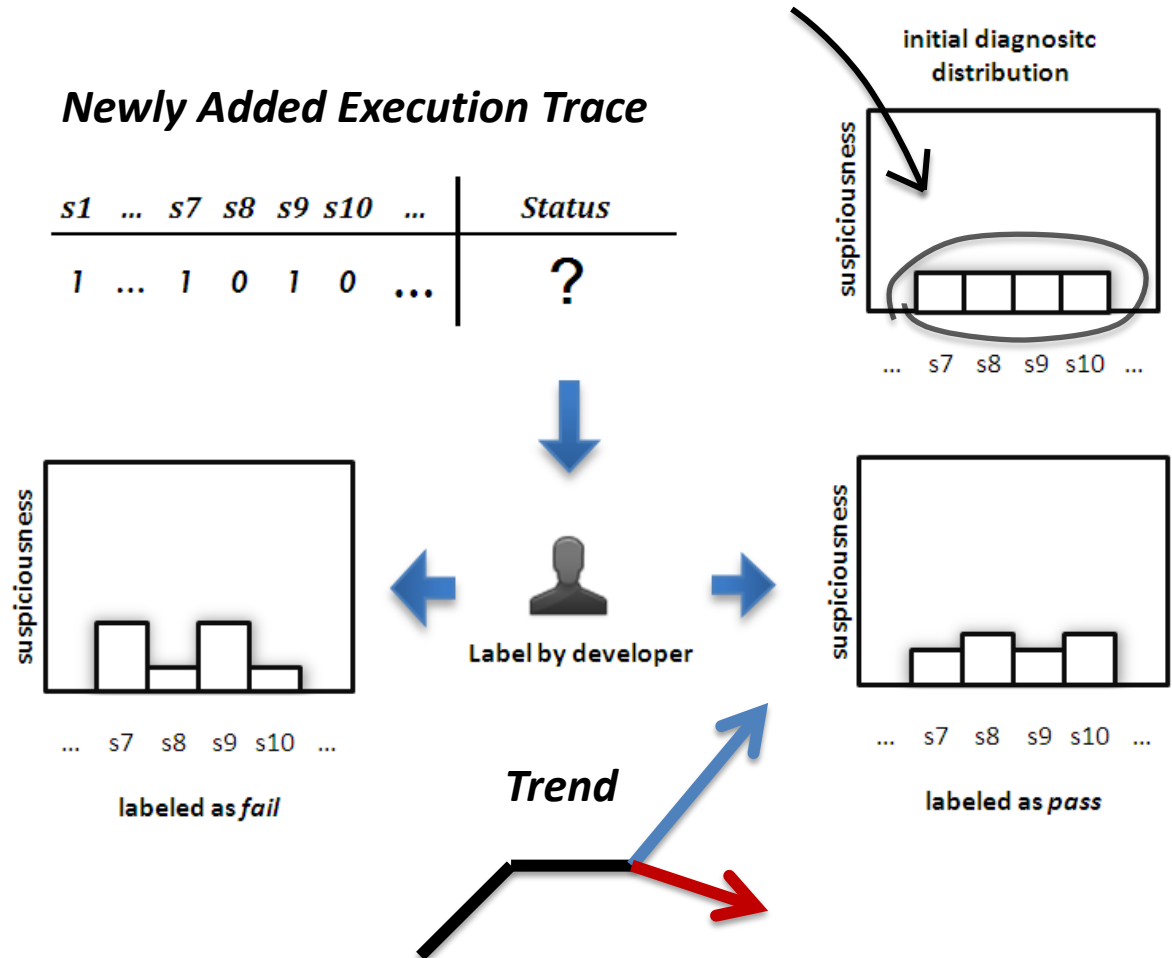
How to?

Two questions prompt:

- How can we know which statements are “**promising**”?
- With “promising” statements, how can we **speed up** their suspiciousness changing process?

- Speed up the suspiciousness ranking changing process by competing in **Suspicious Group**.

Speed up
By Competing



No matter what the label is, ties are broken anyway.
Speedup happens anyway. It could keep “promising” →
or become less “promising” →

Speed up

Our Method

- Change-potential Score of *Suspicious Group*:

$$W_g = \sum_{d \in g} W_{T_d}$$

Change-potential Score of program element d

- Sums of Squares of Change-potential Score of all Groups(G)

$$\mathcal{H}_G = \sum_{g_i \in G} W_{g_i}^2$$

- To choose the next trace t to label,
we use the following formula:

$$\arg \max_{t \in T_U} \{ \mathcal{H}_G - \mathcal{H}_{(G \leftarrow t)} \}$$

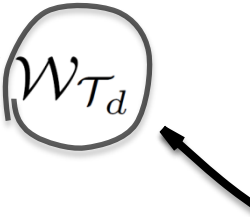
*The Sum of Squares of
change-potential of all groups*

*The Sum of Squares of
change-potential of all groups
when t is added*

Speed up

Our Method

- Change-potential Score of *Suspicious Group*:

$$\mathcal{W}_g = \sum_{d \in g} \mathcal{W}_{T_d}$$


Change-potential Score of program element d

- Sums of Squares of Change-potential Score of all Groups(G)

$$\mathcal{H}_G = \sum_{g_i \in G} \mathcal{W}_{g_i}^2$$

- To choose the next trace t to label,
we use the following formula:

$$\arg \max_{t \in T_U} \{ \mathcal{H}_G - \mathcal{H}_{(G \leftarrow t)} \}$$

Intuition: *When t breaks ties in more promising or larger Suspicious Groups, it is more likely to be selected.*

Test Case Prioritization

Existing Methods

- Coverage Based Prioritization
 - STMT-TOTAL, STMT-ADDTL, and ART.
- Fault-Exposing Potential
 - FEP-TOTAL and FEP-ADDTL.
- Diagnostic Prioritization
 - SEQUIA and RAPTOR.

- Benchmarks for Fault Localization**

Program	Description	LOC	Tests	Faults
tcas	Aircraft Control	173	1609	41
schedule2	Priority Scheduler	374	2710	8
schedule	Priority Scheduler	412	2651	8
replace	Pattern Matcher	564	5543	31
tot_info	Info Measure	565	1052	22
print_tokens2	Lexical Analyzer	570	4055	10
print_tokens	Lexical Analyzer	726	4070	7
space	ADL Compiler	9564	1343	30
flex	Lexical Parser	10124	567	43
sed	Text Processor	9289	371	22
grep	Text Processor	9089	809	17
gzip	Data Compressor	5159	217	15

1

2

1 Siemens Suite

2 UNIX Programs

- Evaluation Metric for Fault Localization:**

$$cost = \frac{|\{j \mid f_{T_S}(d_j) \geq f_{T_S}(d_*)\}|}{|\mathcal{D}|}$$

Experiment
*Dataset &
Evaluation Metric*



Experiment

Settings & Design

Experiments comparing with existing methods:

- Effectiveness on reducing the number of test cases (*i.e.*, labeling effort) needed for a *target cost*
- Effectiveness on reducing cost for a given number of labeled test cases
- Defining *target cost* c_x : *Average fault localization cost when labeling all test cases*

$$c_x = \frac{x}{100} \times c$$

Labeling Effort Needed When Setting c_{101} as Target Cost

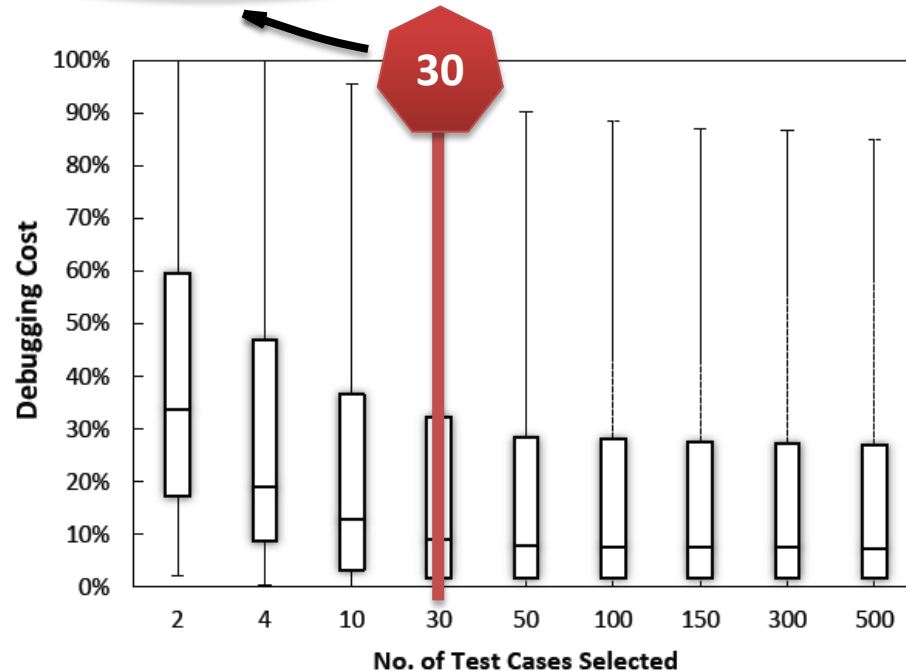
Subject Programs	DMS	RAP-TOR	SEQ-UOIA	STMT-ADDTL	STMT-TOTAL	FEP-ADDTL	ART-MIN
Siemens	18	20	500+	500+	500+	97	150
UNIX	16	48	176	150	500+	98	56

Experiment

Settings & Design

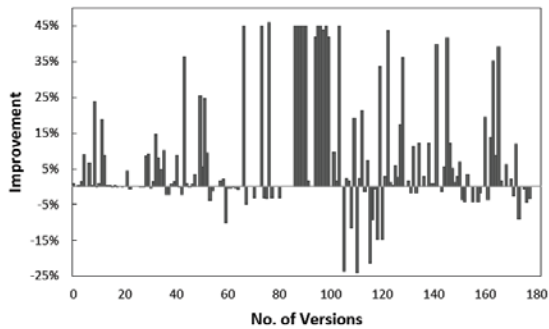
Experiments comparing with existing methods:

- Effectiveness on reducing the number of test cases (*i.e.*, labeling effort) needed for a *target cost*
- Effectiveness on reducing cost for a given number of labeled test cases

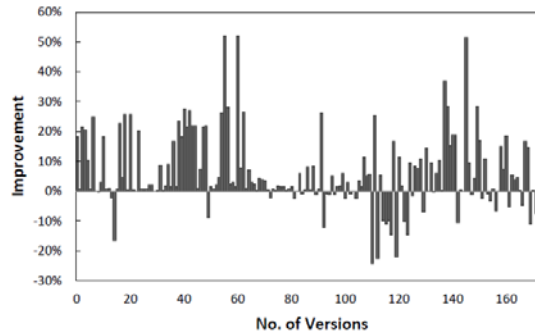


**Average Cost of DMS when Selecting
Different Number of Test Cases**

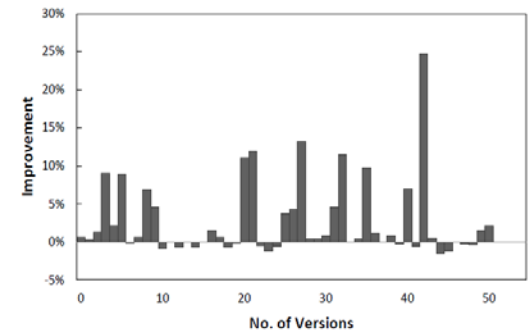
- Effectiveness on Reducing Cost for a Given Number of Labeled Test Cases



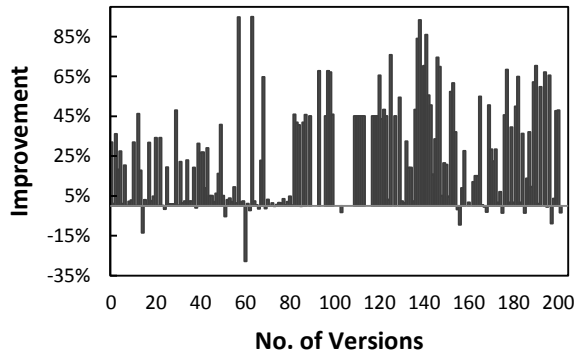
DMS VS FEP-ADDTL



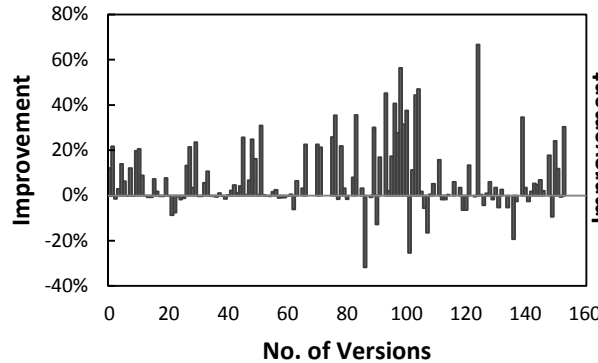
DMS VS ART-MIN



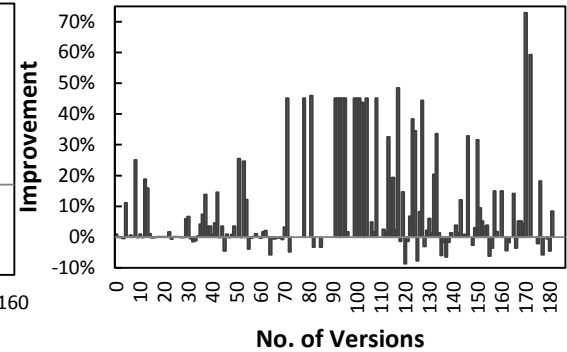
DMS VS RAPTOR



DMS VS STMT-TOTAL



DMS VS SEQUOIA



DMS VS STMT-ADDTL



Pair-wised T-test shows the improvements are statistically significant at 95% interval.



Conclusion

& Future Work

- **Conclusions**

- We propose a new technique aiming to minimize the amount of effort in manual oracle construction, while still permitting effective fault localization.
 - ✓ Given a target fault localization accuracy, our approach can significantly reduce the number of test cases needed to achieve it.
 - ✓ Given a maximum number of test cases that a programmer can manually label, DMS can improve the accuracy of fault localization and thus helps reduce the debugging cost.

- **Future Work**

- Evaluate on more subject programs.
- We will also explore the possibility of adopting more sophisticated trend analysis methods.

Conclusion

& Future Work

GL

Thank you!

- Any questions?

- Evaluate on more subject programs.



- We will also explore the possibility of adopting more sophisticated trend analysis methods.