

On the “Naturalness” of Buggy Code

Baishakhi Ray, Vincent Hellendoorn, Saheel Godhane,
Zhaopeng Tu, Alberto Bacchelli, Prem Devanbu



[illegible]

Real Programs are Natural!!

```
1 public PubnubPublishTest() {  
2     _pubnub = new Pubnub("demo", "demo", false);  
3     _pubnub_enc = new Pubnub("demo", "demo", "demo", "demo", false);  
4  
5     ...  
6  
7     jsarr.put("This is json array");  
8     jsarr.put("One more entry in json array");  
9  
10    ...  
11  
12    jsobj.put("msg1", "Hi");  
13    jsobj.put("msg2", "Java");  
14 }
```

Highly repetitive, predictable,
amenable to large-sample statistical methods

[Hindle Et al., Allamanis et al., Nguyen et al.]



What does it mean when code
is “unnatural” ?

Is “unnatural” code more defect-prone?

What is unnatural
code?

How are they related
to bugs?

Naturalness of Code :

N-Gram Language Model

Given the previous words, learn the *conditional distribution* of the next word.

<i>n-grams</i>	<i>frequencies</i>	<i>probabilities</i>
<code>for (int i = 0 ...</code>	14	0.70
<code>for (int i = start ...</code>	5	0.25
<code>for (int i = end ...</code>	1	0.05

Naturalness of Code : *N-Gram Language Model*

Given the previous words, learn the *conditional distribution* of the next word.

n-grams

frequencies

probabilities

`for (int i = 0 ...`

14

0.70

`for (int i = start ...`

5

0.25

`for (int i = end ...`

1

0.05

Less
seen
token

Un-natural Code

Naturalness of Code :

N-Gram Language Model

<i>n-grams</i>	<i>frequencies</i>	<i>probabilities</i>
----------------	--------------------	----------------------

<code>for (int i = end ...</code>	1	0.05
---	---	-------------

Unnaturalness is measured by **entropy** (“improbability”)

$$Entropy(t) = -P(t)\log_2 P(t)$$

**Higher entropy value = lower probability
= less naturalness**

How does “unnatural” code relate to bugs?

=> Buggy lines of code

Unnatural code **for** Defect prediction

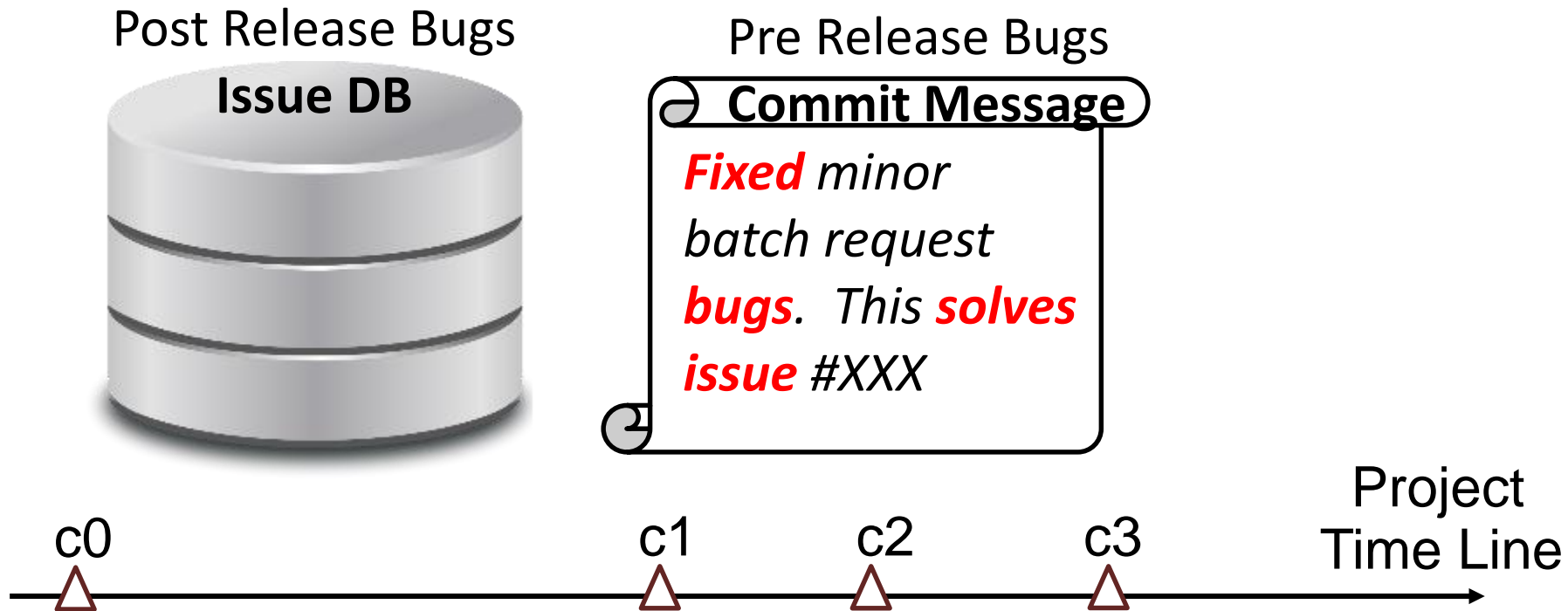
vs. Static analysis tools
in defect prediction

Methodology

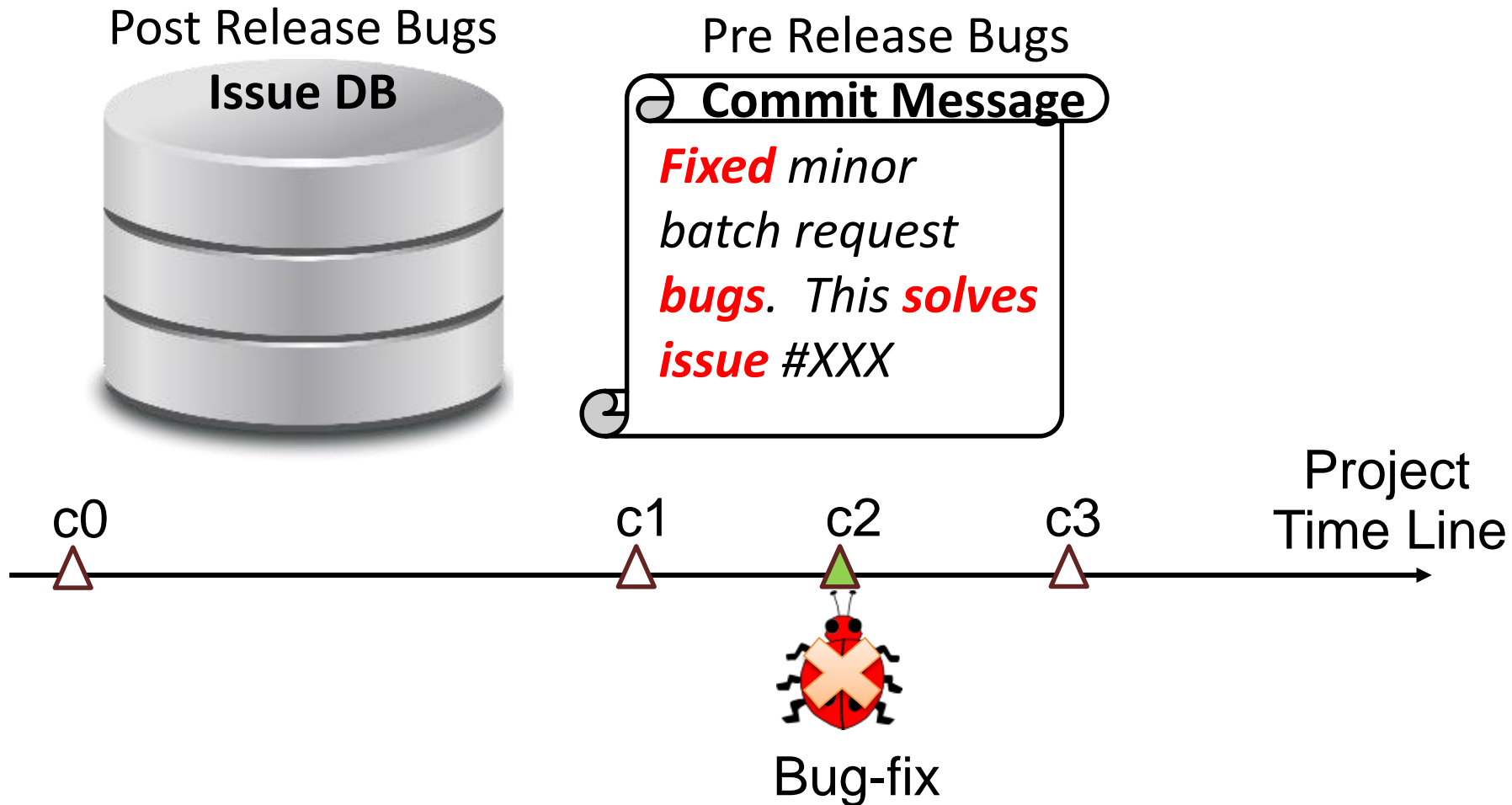
Step 1: Identify buggy lines in each version

Step 2: Measure entropy of each program line

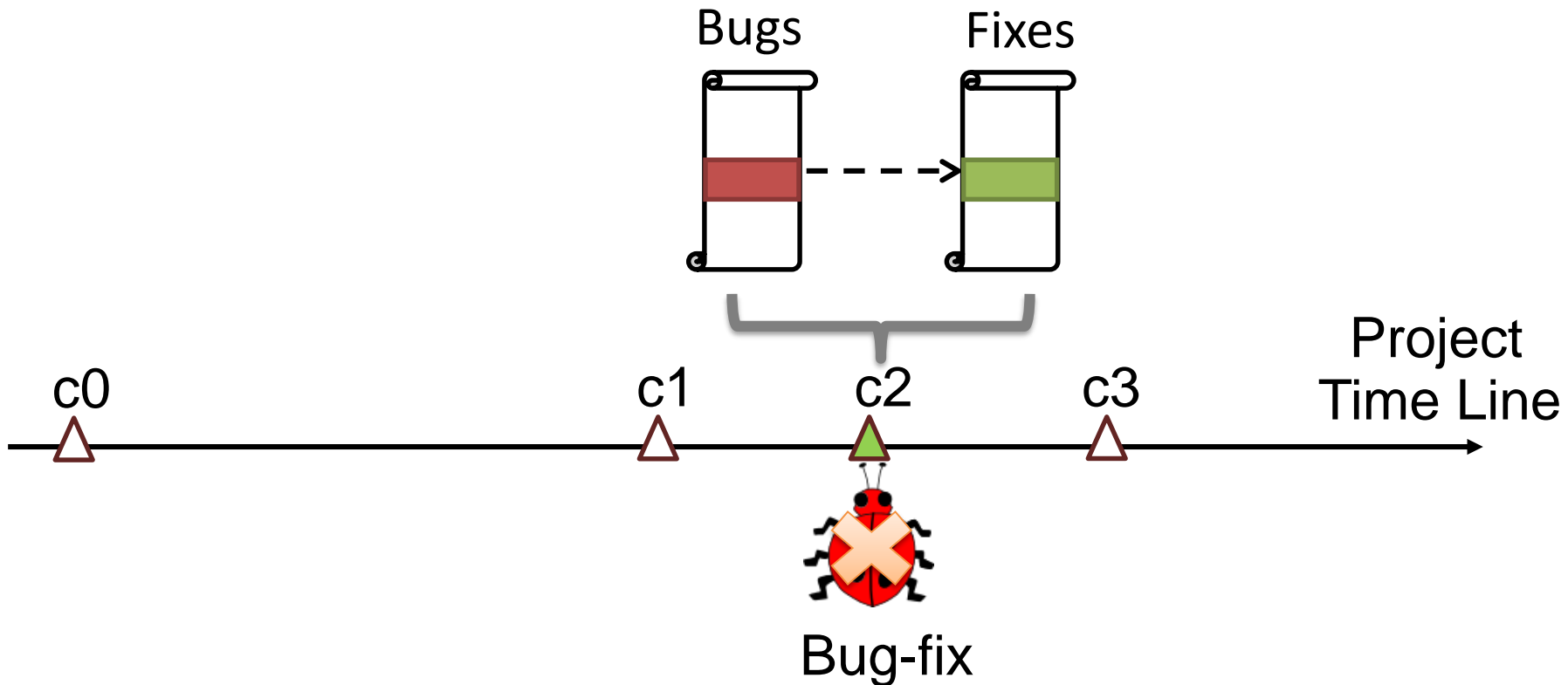
Methodology: Identify Buggy Lines



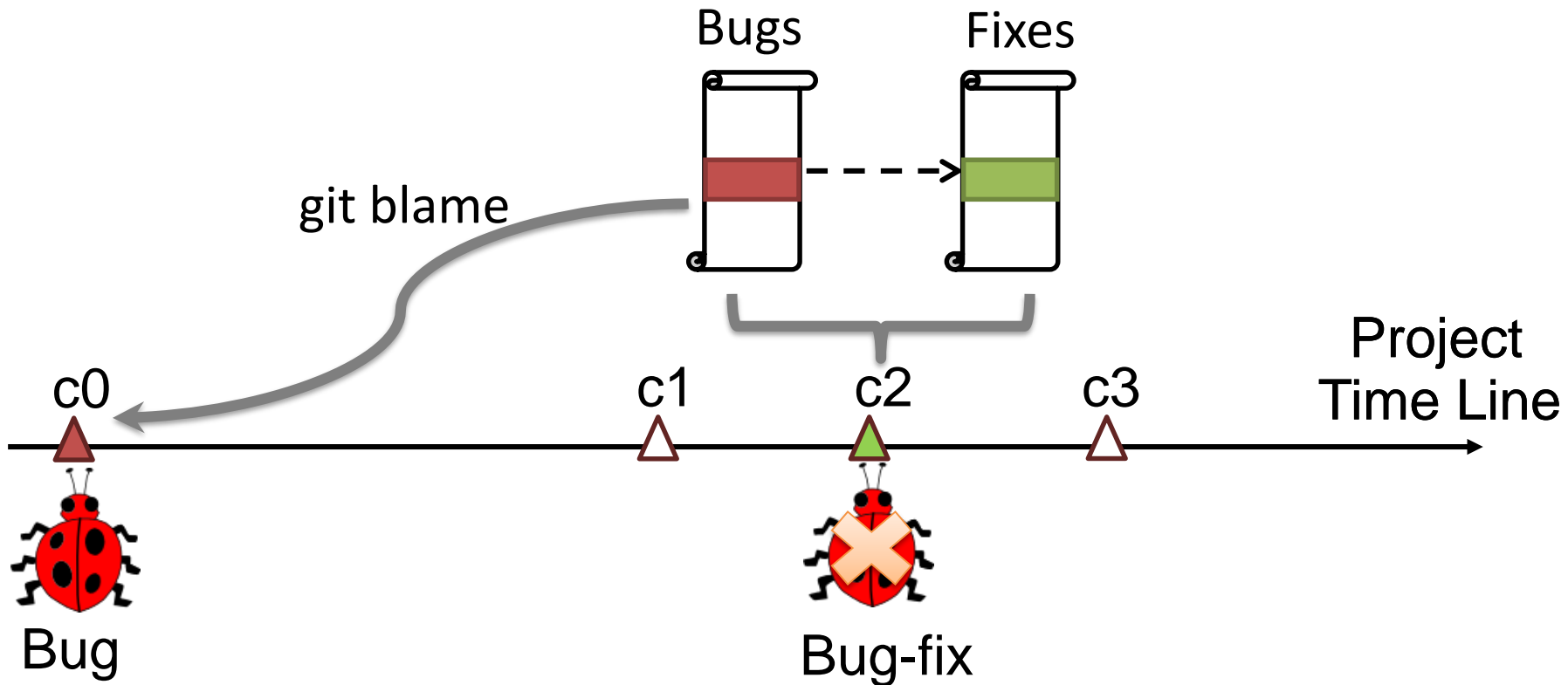
Methodology: Identify Buggy Lines



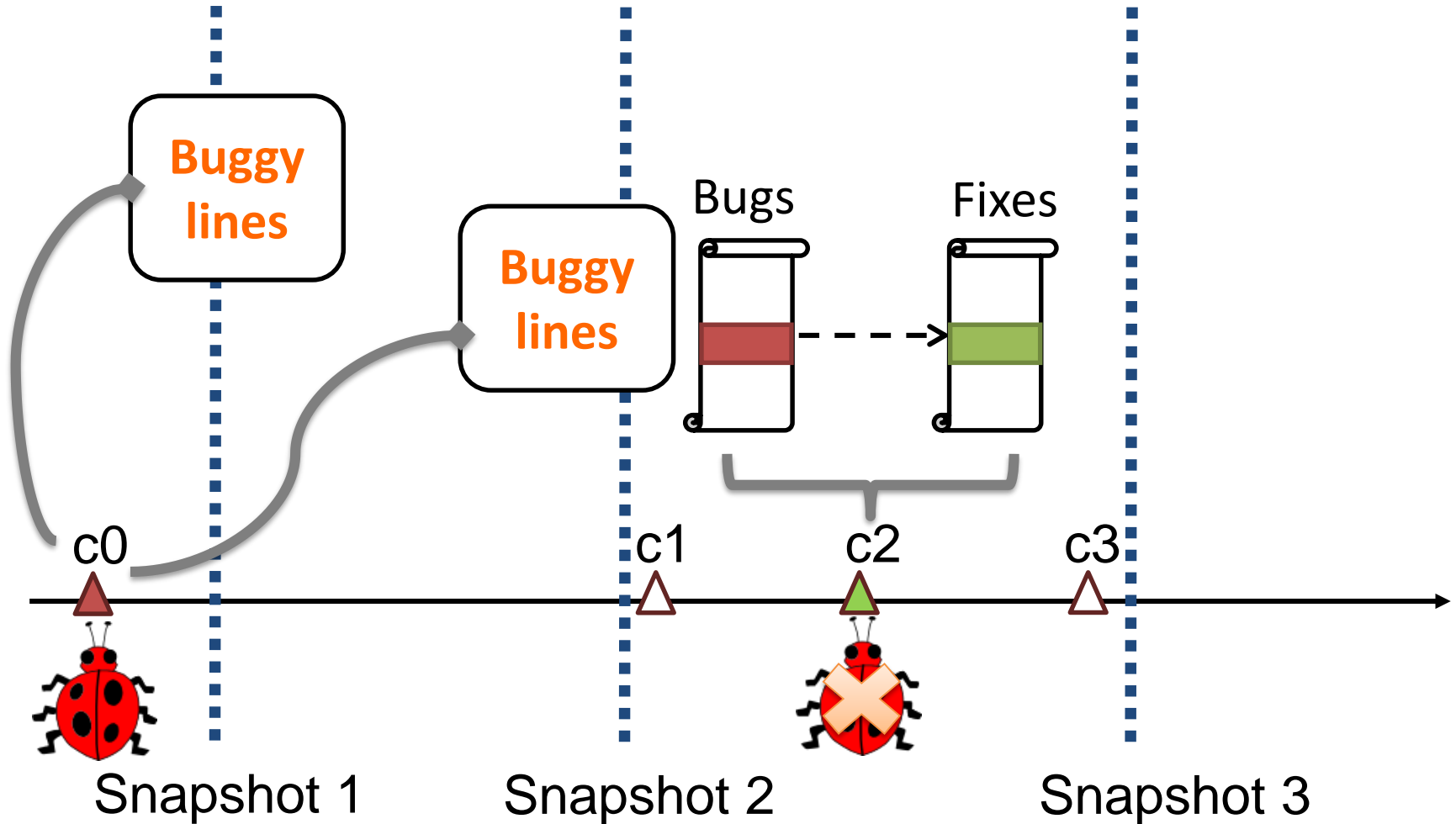
Methodology: Identify Buggy Lines



Methodology: Identify Buggy Lines



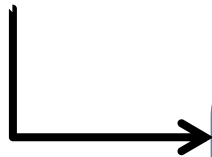
Methodology: Identify Buggy Lines



Methodology: Measure Entropy

Training:

Other files in the snapshot



Test:

Single file in a snapshot

Cache-based n-gram
language model (\$gram)
[Tu et al., FSE 2014]



Output:

Entropy per
program line

Study Subjects

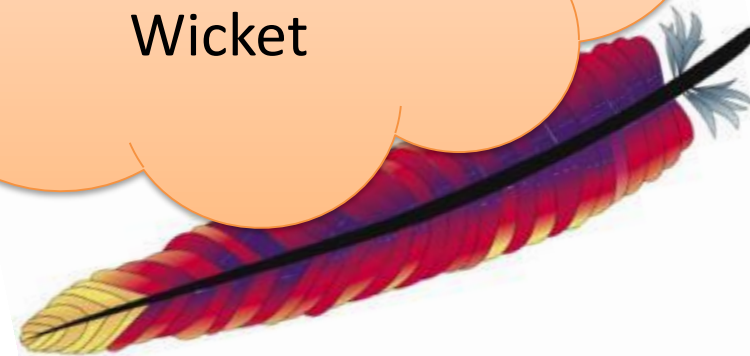


GitHub

Atmosphere, Presto,
Elasticsearch, Netty,
Facebook-android-sdk

Apache

Derby, Lucene,
OpenJPA, Qpid,
Wicket

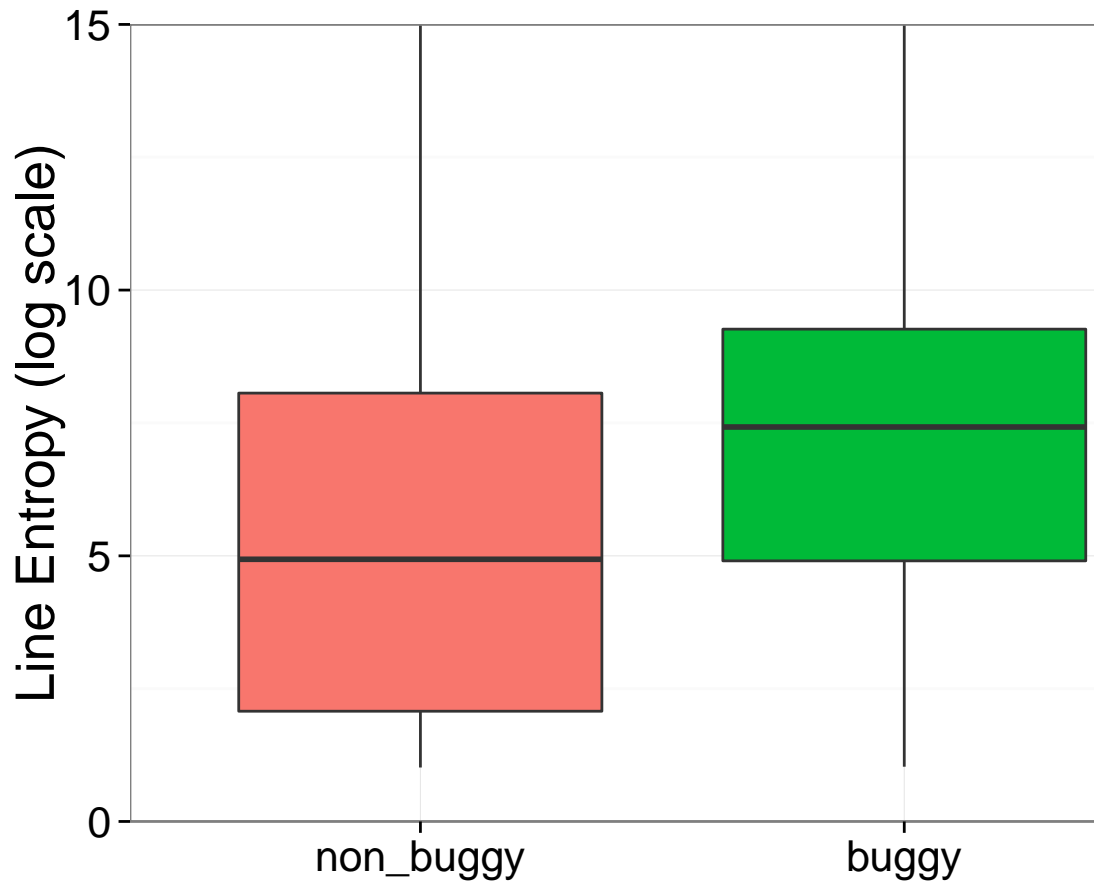


10 projects, 120 versions, 113K Files, 35M Lines, 7K bugs

RESULTS

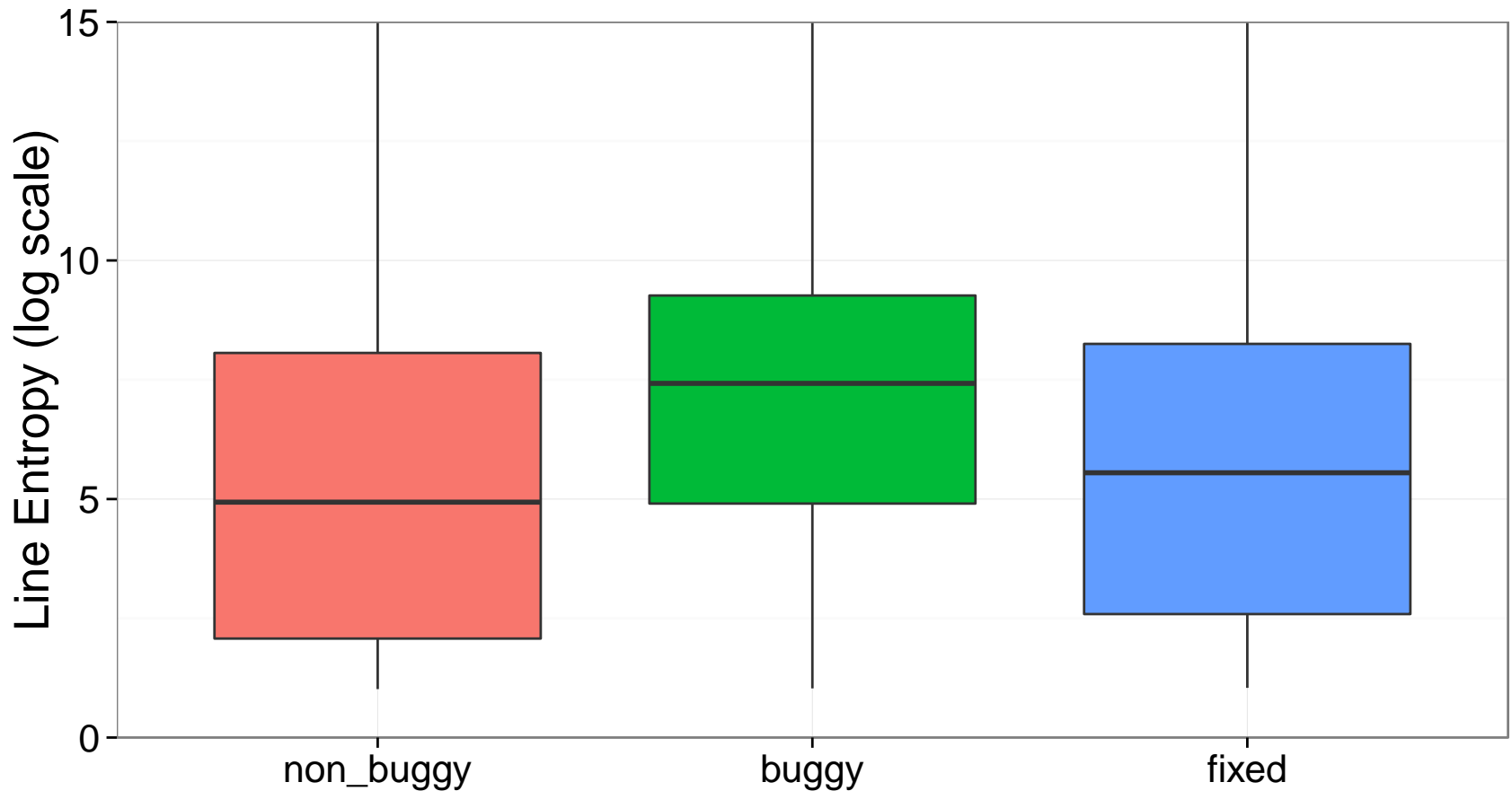
A magnifying glass with a black handle and frame is positioned over the word "RESULTS". The lens of the magnifying glass is centered over the word "RESULTS", which is written in a large, bold, blue, sans-serif font. The word "RESULTS" is slightly tilted upwards from left to right. The background is white. The magnifying glass's handle extends from the bottom right corner towards the center. The lens is a large circle with a black border, and it shows a slight reflection on its surface. The word "RESULTS" is partially obscured by the magnifying glass, with the letters "R", "E", and "S" visible to the left of the lens and "U", "L", "T", and "S" visible to the right. The magnifying glass is held at an angle, with the handle pointing towards the bottom right.

RQ1. Are buggy lines more “unnatural” than non-buggy lines?



The difference is more for low latency & less scattered bugs

RQ2: Do buggy lines become more "natural" after bug-fixes?



RQ2. Do buggy lines become more “natural” after bug-fixes?

Netty: incorrect method call

```
if (isTerminated()) {  
  - terminationFuture.setSuccess(null); // entropy = 5.96  
  + terminationFuture.trySuccess(null); // entropy = 1.34  
}
```

Entropy dropped after
bugfix : 4.63

Lucene: missing conditional check

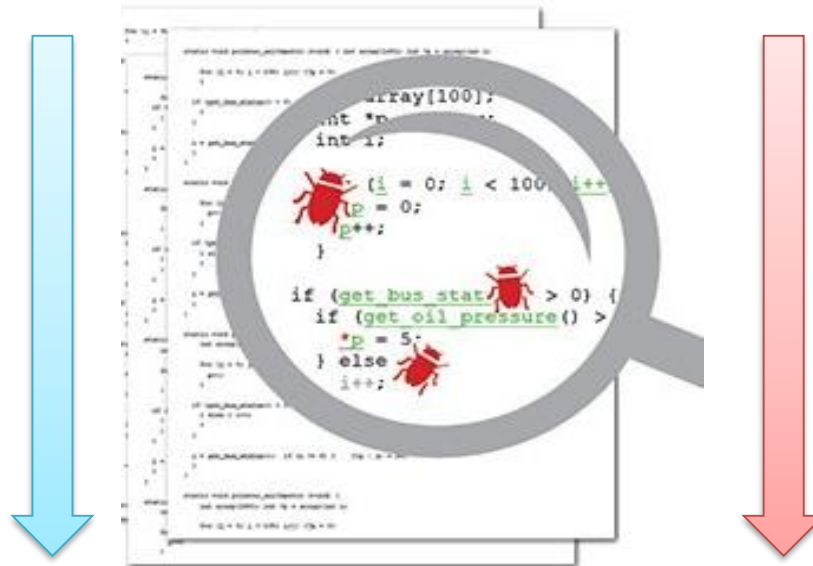
```
if (!directory.exists())  
  - directory.mkdir(); // entropy = 9.21  
  + if (!directory.mkdir()) // entropy = 5.34  
    + throw new IOException (...)
```

Entropy dropped after
bugfix : 3.87

Buggy lines have higher entropies than non-buggy lines;
entropy drops after bug-fixes.

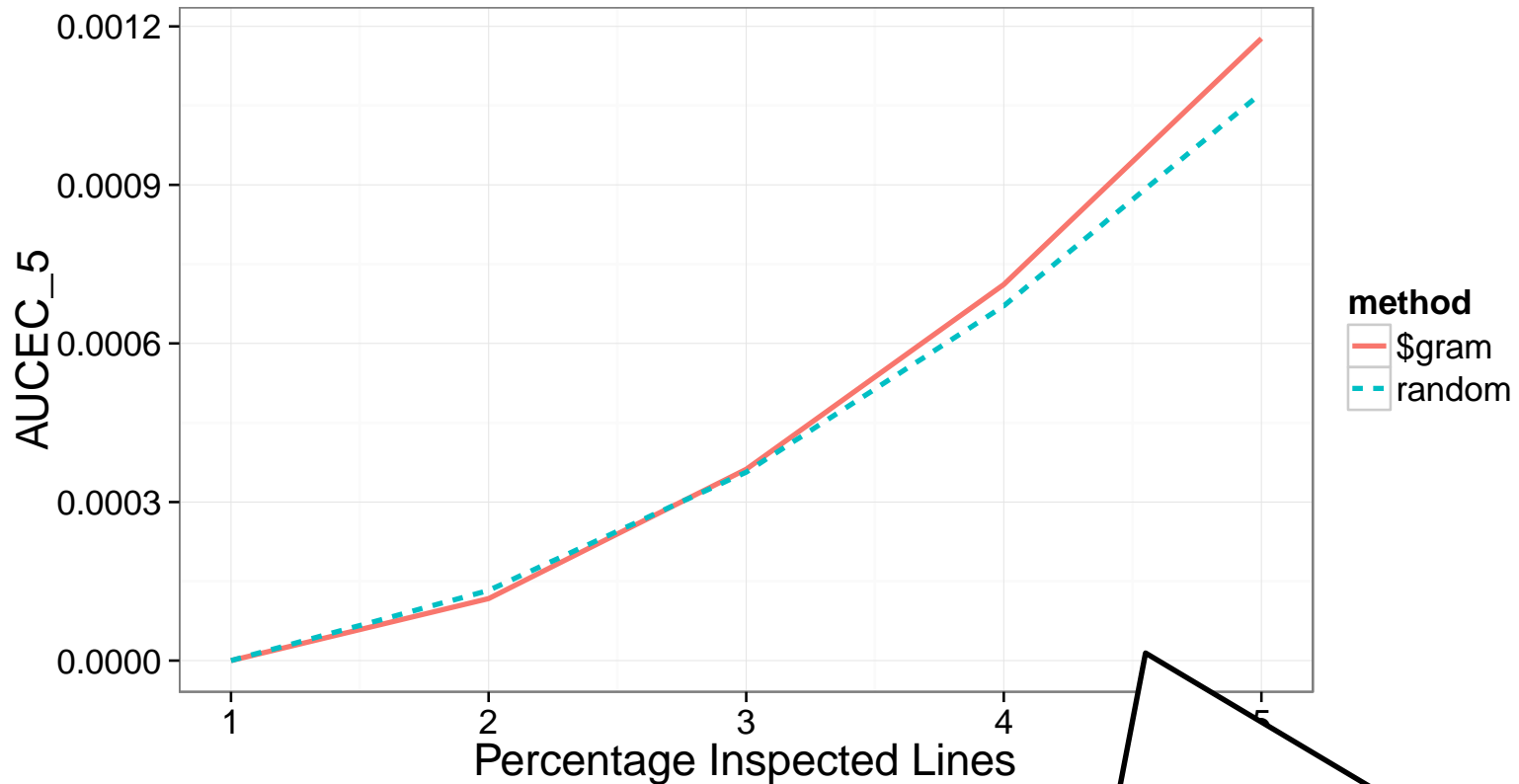
RQ3. Is unnaturalness useful for defect prediction?

Baseline:
Order lines
randomly



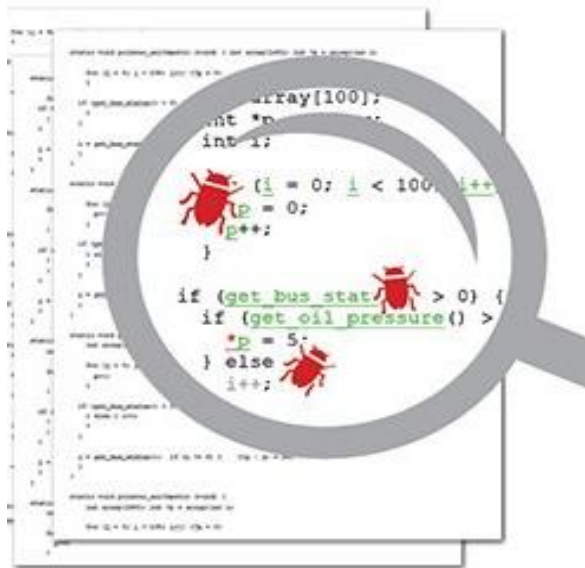
Order lines
by decreasing
entropy

RQ3. Is unnaturalness useful for defect prediction?



AUCEC = Area under the cost-effectiveness curve.
Cost is inspection effort and payoff is number of bugs found

RQ3. Is unnaturalness useful for defect prediction?



Order lines
by decreasing
entropy



Problem: some line types are more entropic than others (eg., import statement vs. for loop)

RQ3. Is unnaturalness useful for defect prediction?



Problem: some line types are more entropic than others (eg., import statement vs. for loop)

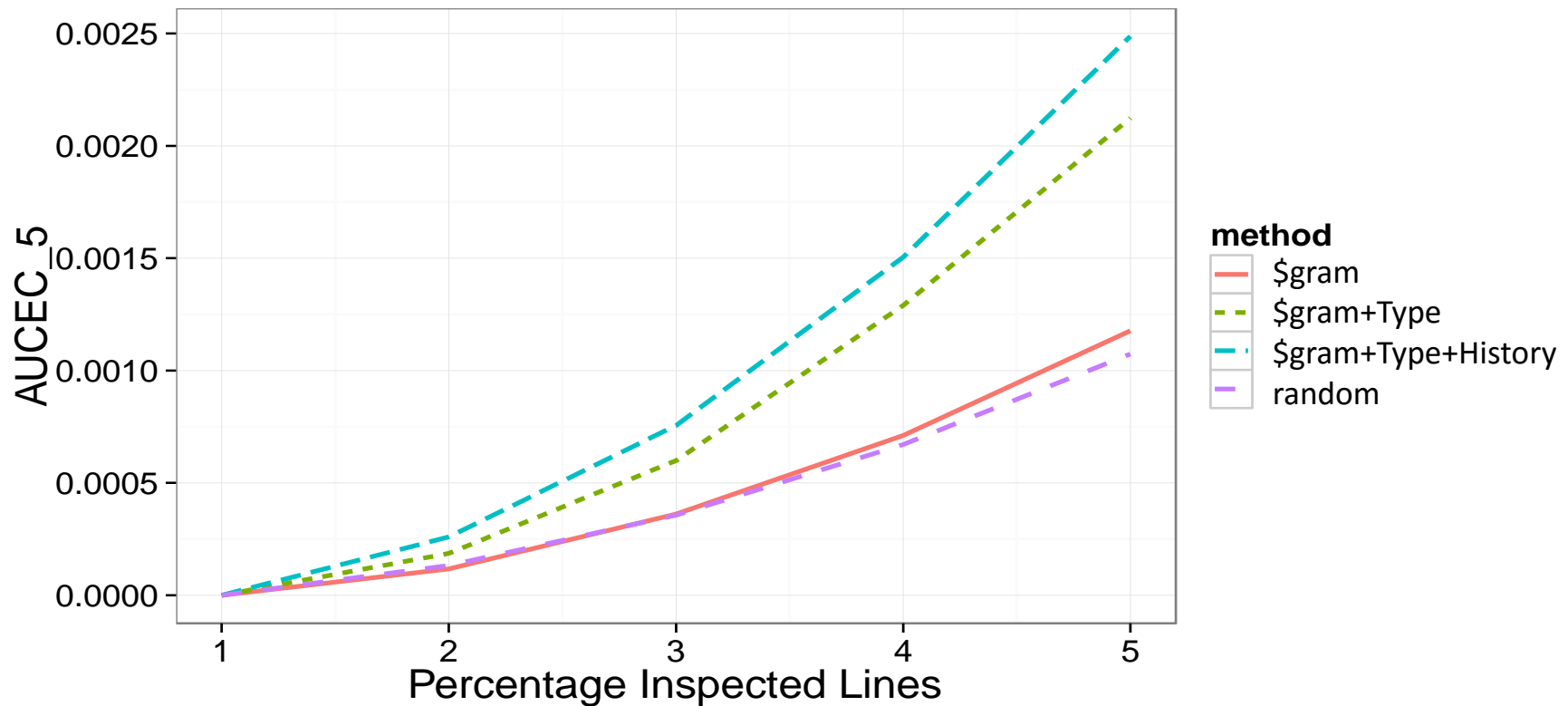
Solution: can we leverage the line types?



$\$gram+Type$: how much a line's entropy deviates from the mean entropy of its own line-type (normalized Z-score)

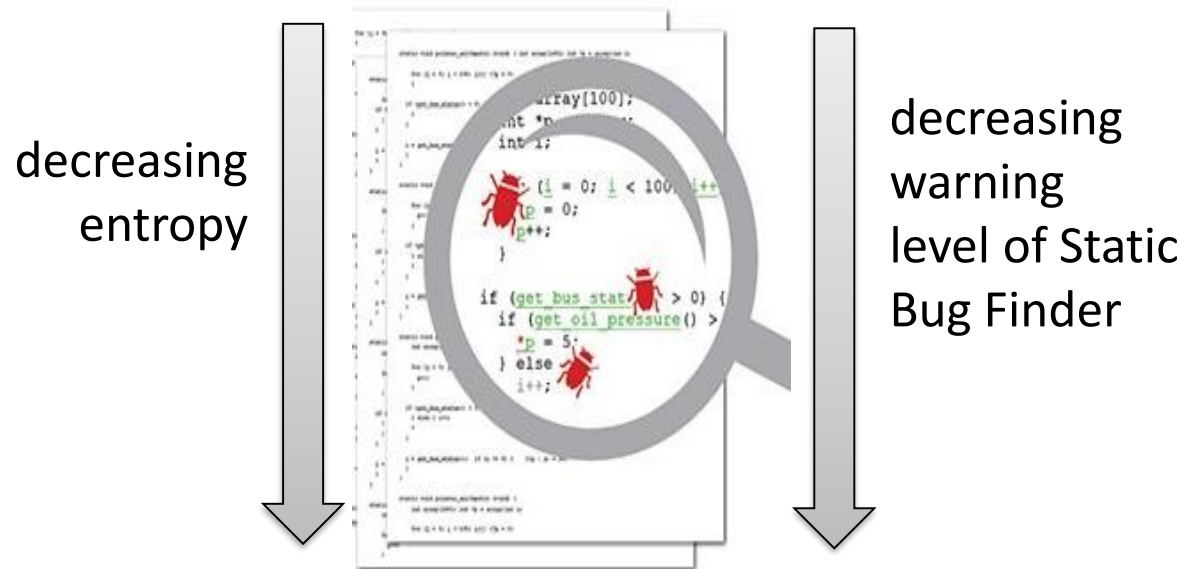
$\$gram+Type+History$: how much a line-type was buggy in the past.

RQ3. Is unnaturalness useful for defect prediction?



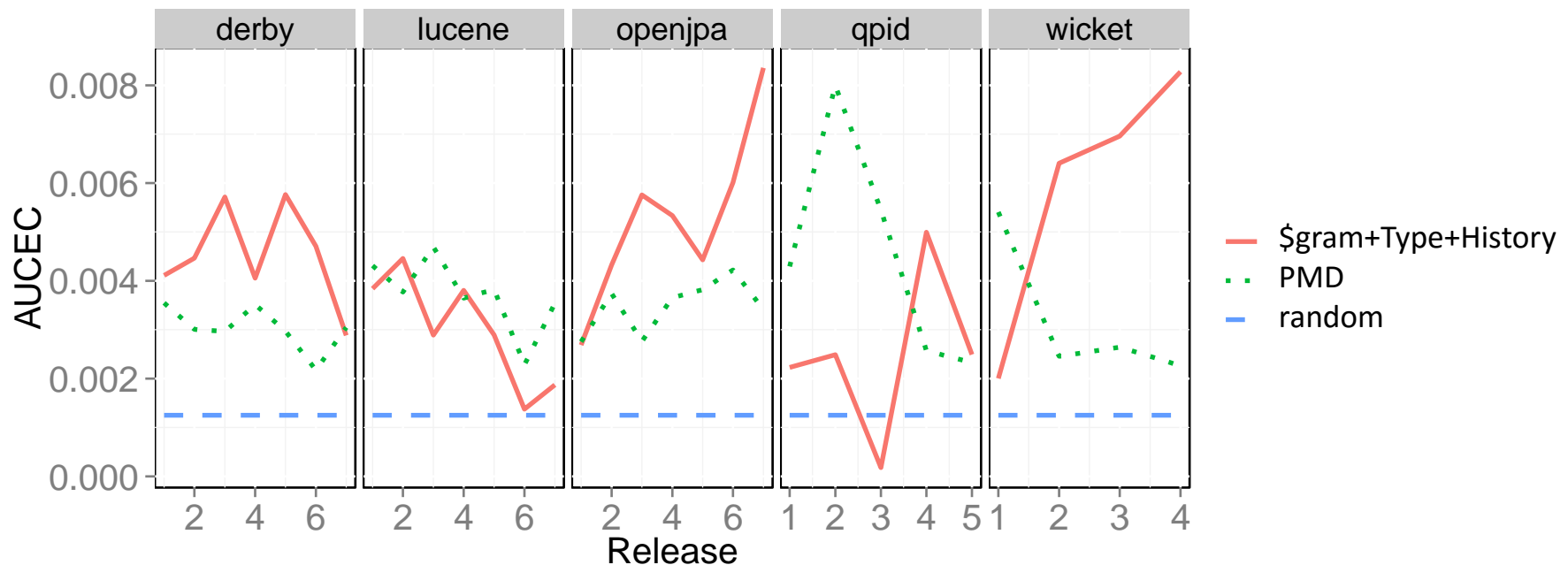
Entropy can be used to guide bug-finding efforts
at **line level**

RQ4. How does unnaturalness perform against static bug finding technique?



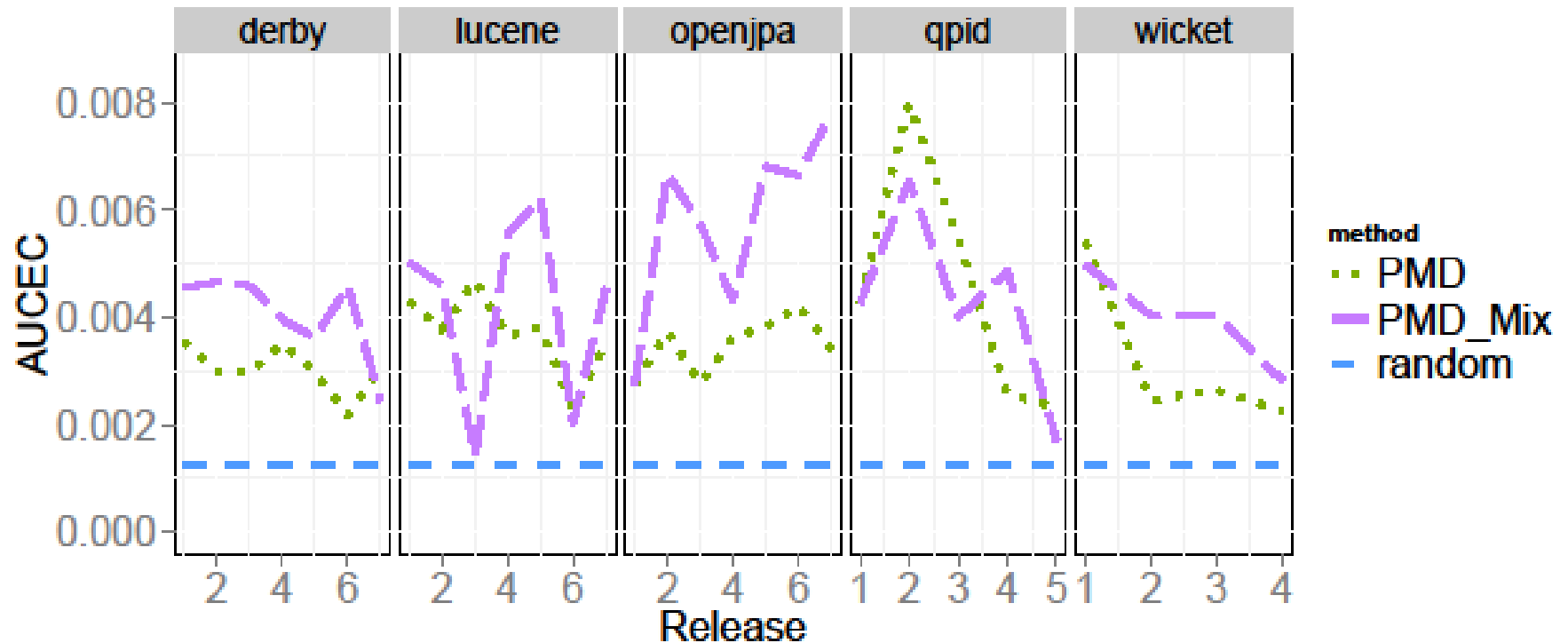
Compared against two popular static bug finding tools:
FindBugs and PMD

RQ4. How does unnaturalness perform against static bug finding technique?



Entropy achieves comparable performance to commonly used Static Bug Finders in defect prediction.

RQ5. Does unnaturalness boost inspection effort on Static Bug Finder's warnings?



Ordering Static Bug Finder's warnings by priority and entropy significantly improves SBF performance.

Summary

=> Buggy lines of code ✓

Unnatural code **for** Defect prediction
at line granularity ✓

vs. Static analysis tools ✓

Implications: Search-based bug repair may benefit for both
fault-localization and searching for fixes

Acknowledgement

National Science Foundation
(Grant 1414172)



Questions?

=> Buggy lines of code ✓

Unnatural code **for** Defect prediction
at line granularity ✓

vs. Static analysis tools ✓

Implications: Search-based bug repair may benefit for both
fault-localization and searching for fixes

Naturalness of Code:

Cache Language Model

Global Context

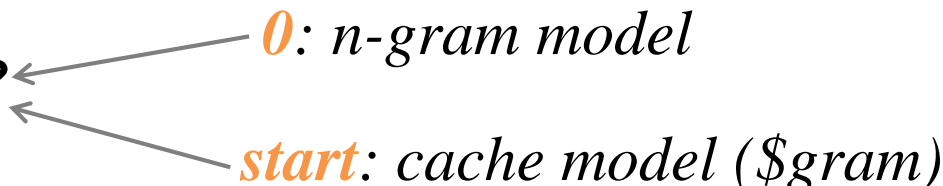
```
for ( int i = 0; i < 10; ++ i )  
...  
for ( int i = 0; i < 10; ++ i )  
...  
for ( int i = 0; i < 10; ++ i )  
...
```

Local Context

```
for ( int i = start; i < end; ++ i )  
...  
for ( int i = start; i < end; ++ i )  
...
```

An additional cache component that memorizes the n-grams in the locality to capture [Tu et al. FSE.14]

for (int i = ?

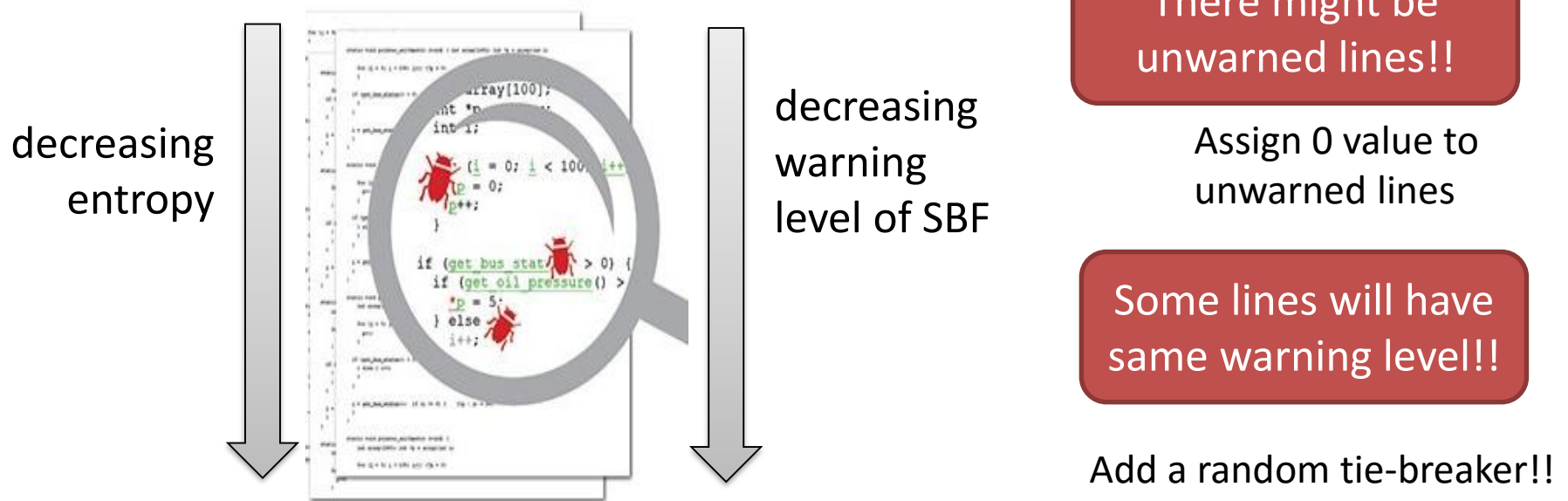


0: n-gram model

start: cache model (\$gram)

\$gram is used measure entropy of each program statement

RQ4. How does “unnaturalness” perform against static bug finding technique?



Summary

- Unnatural code is more likely to be buggy.
- Can be effective for defect prediction, even at line granularity.
- A simple way to complement static bug finding techniques.
- Search-based bug-fixing methods may benefit from using entropy both for fault-localization and searching for fixes.

Check it out @ <http://odd-code.github.io>