

Conceptual Mutation Testing for Student Programming Misconceptions



Siddhartha Prasad
Ben Greenman
Tim Nelson
Shriram Krishnamurthi



Understanding CS Problems

Understanding CS Problems

Q. median



Understanding CS Problems

Q. median



Understanding CS Problems

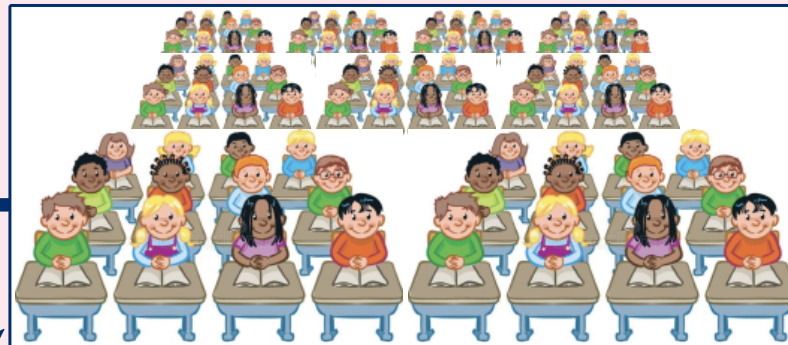
Q. median



feedback?

Understanding CS Problems

Q. median

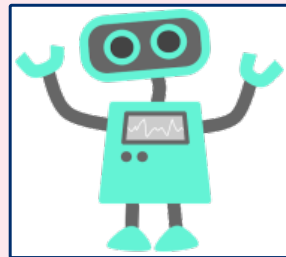


~~feedback?~~



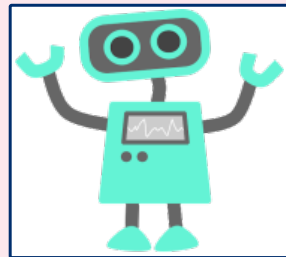
Understanding CS Problems

Q. median

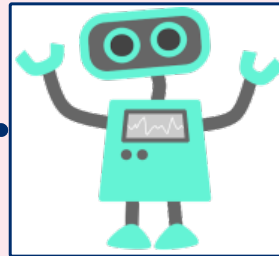


Understanding CS Problems

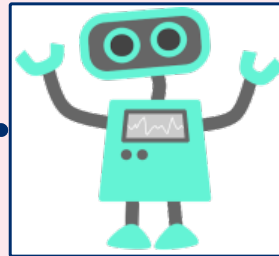
Q. median



How to communicate?

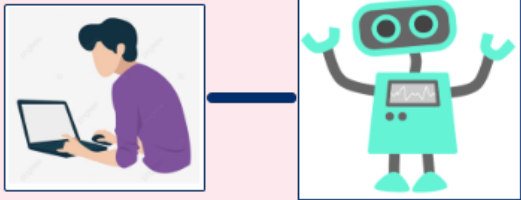


How to communicate?



Test Cases

Q. median



Q. median

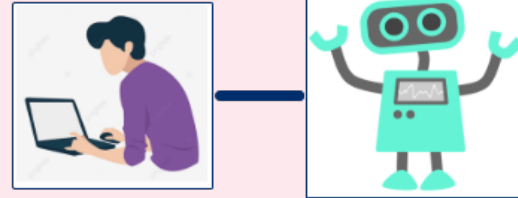


median [1] is 1

median [1, 2, 3] is 3

median [3, 3, 3] is 3

Q. median



```
median [1] is 1
```

```
median [1, 2, 3] is 3
```

```
median [3, 3, 3] is 3
```

What's **wrong** with these tests?

```
1 include my-gdrive("median-code.arr")
2 # DO NOT CHANGE ANYTHING ABOVE THIS LINE
3
4 check:
5 median([list: 1]) is 1
6 median([list: 1, 2, 3]) is 3
7 median([list: 3, 3, 3, 3]) is 3
8
9
10 end
```

The image shows a code editor window with a file named `median-tests.arr`. The code contains an `include` statement for `median-code.arr`, a comment `# DO NOT CHANGE ANYTHING ABOVE THIS LINE`, and a `check` block with three test cases for the `median` function. The tests are: `median([list: 1]) is 1`, `median([list: 1, 2, 3]) is 3`, and `median([list: 3, 3, 3, 3]) is 3`. The `end` keyword is at the bottom of the `check` block.

On the right side of the editor, a test runner window is open, displaying the results of the tests. It shows a large **INCORRECT** message and the text `CONSEQUENTLY, THOROUGHNESS IS UNKNOWN`. Below this, a message states: `These tests do not match the behavior described by the assignment:` with a link to `definitions://5:2-5:30`. A specific test failure is highlighted in a red box: `6 median([list: 1, 2, 3]) is 3`.

```
1 include my-gdrive("median-code.arr")
2 # DO NOT CHANGE ANYTHING ABOVE THIS LINE
3
4 check:
5   median([list: 1]) is 1
6   median([list: 1, 2, 3]) is 3
7   median([list: 3, 3, 3, 3]) is 3
8
9
10 end
```

median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment: [definitions://5:2-5:30](#)

6 median([list: 1, 2, 3]) is 3

```
1 ▾ include my-gdrive("median-code.arr")
2 # DO NOT CHANGE ANYTHING ABOVE THIS LINE
3
4 ▾ check:
5 ▾ median([list: 1]) is 1
6 ▾ median([list: 1, 2, 3]) is 2
7 ▾ median([list: 3, 3, 3, 3]) is 3
8
9 #Shows that Median is not Mean
10 ▾ median([list: 1, 1, 3]) is 1
11
12 # Shows that Median is not Mode
13 ▾ median([list: 1, 1, 3, 4, 4]) is 3
14 end
```



```
1 include my-gdrive("median-code.arr")
2 # DO NOT CHANGE ANYTHING ABOVE THIS LINE
3
4 check:
5 median([list: 1]) is 1
6 median([list: 1, 2, 3]) is 2
7 median([list: 3, 3, 3, 3]) is 3
8
9 #Shows that Median is not Mean
10 median([list: 1, 1, 3]) is 1
11
12 # Shows that Median is not Mode
13 median([list: 1, 1, 3, 4, 4]) is 3
14 end
```

VALID

These tests are valid and consistent with the assignment handout. They caught 2 of 4 sample buggy programs. Add more test cases to improve this test suite's thoroughness.

```
1 include my-gdrive("median-code.arr")
2 # DO NOT CHANGE ANYTHING ABOVE THIS LINE
3
4 check:
5 median([list: 1]) is 1
6 median([list: 1, 2, 3]) is 2
7 median([list: 3, 3, 3, 3]) is 3
8
9 #Shows that Median is not Mean
10 median([list: 1, 1, 3]) is 1
11
12 # Shows that Median is not Mode
13 median([list: 1, 1, 3, 4, 4]) is 3
14 end
```

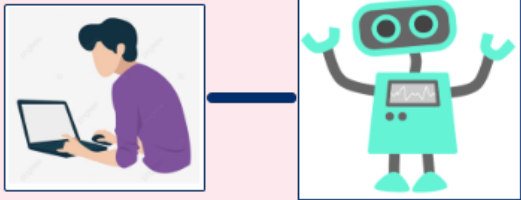
median-tests.arr

VALID

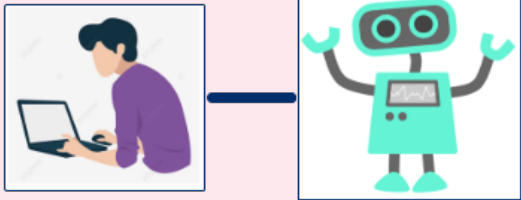
These tests are valid and consistent with the assignment handout. They caught 2 of 4 sample buggy programs. Add more test cases to improve this test suite's thoroughness.

What's **wrong** with these tests?

Q. median



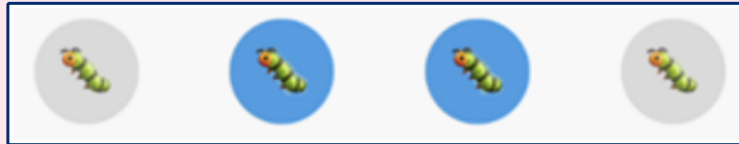
Q. median



Tests must distinguish:
mean
median vs. mode
middle ...

Valid & Thorough

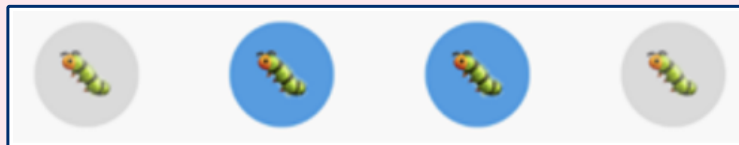
How to check thoroughness?



Buggy solutions
(mutation testing)



RQ. How to design buggies?



RQ. How to design buggies?

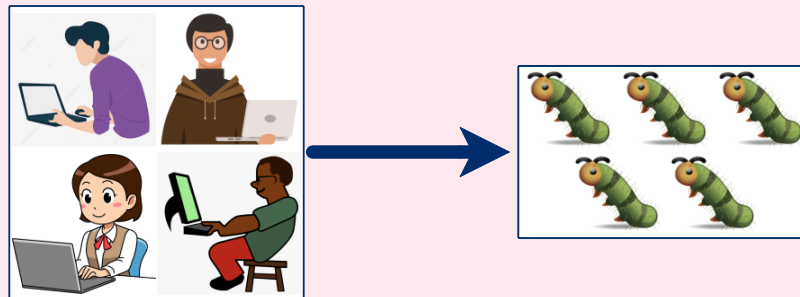


Need to **discover** misconceptions

Prior Work:
Expert-Driven



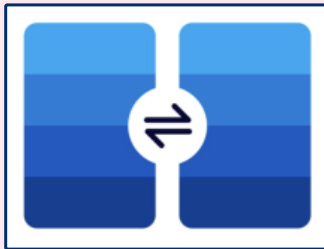
Today, a recipe:
Buggies from Data



1. Design problem

1. Design problem

Running example:
Doc Diff



```
docdiff ['a'] ['A'] is 1
```

```
docdiff ['one', 'two'] ['one'] is 1/2
```

```
docdiff ['hello'] ['world'] is 0
```

2. Collect invalid tests

2. Collect invalid tests



median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment:

median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment:

median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment:

median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment:

2. Collect invalid tests

Doc Diff ==> 1,500 invalids in ~1 week



median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment:

median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment:

median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment:

median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment:

median-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN

These tests do not match the behavior described by the assignment:

3. Cluster tests by feature vector

3. Cluster tests by feature vector

median-tests.arr
INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN
These tests do not match the behavior described by the assignment.

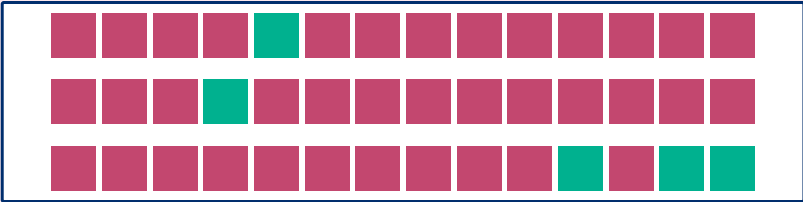
median-tests.arr
INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN
These tests do not match the behavior described by the assignment.

median-tests.arr
INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN
These tests do not match the behavior described by the assignment.

median-tests.arr
INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN
These tests do not match the behavior described by the assignment.

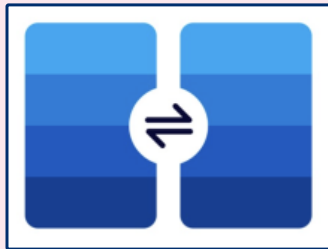
median-tests.arr
INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN
These tests do not match the behavior described by the assignment.

median-tests.arr
INCORRECT CONSEQUENTLY, THOROUGHNESS IS UNKNOWN
These tests do not match the behavior described by the assignment.



Feature vectors \leq problem characteristics

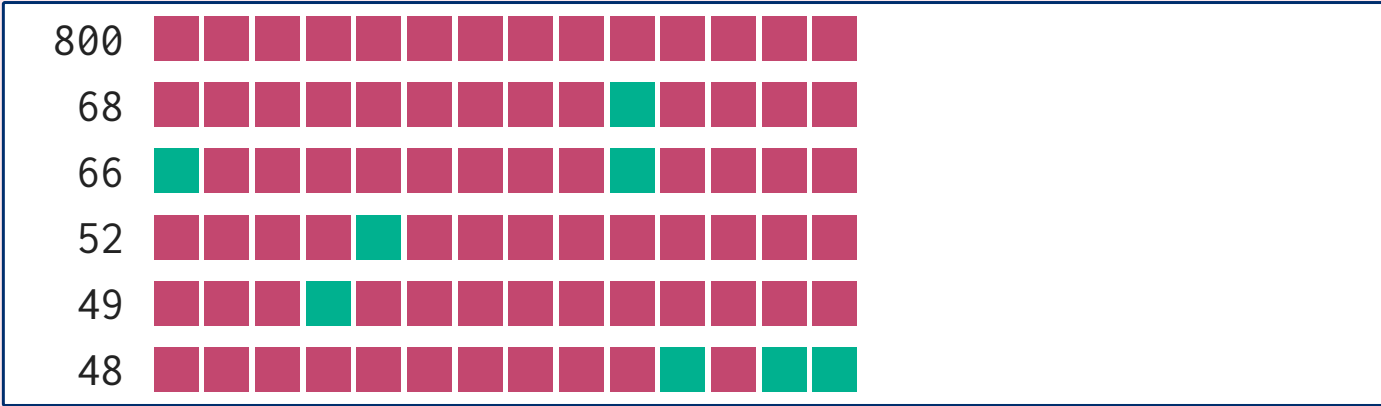
Feature vectors \leq problem characteristics









- Case-insensitive
- Words may repeat
- Diff may be a fraction
- ... [14 in total]

5. Sort clusters

5. Sort clusters

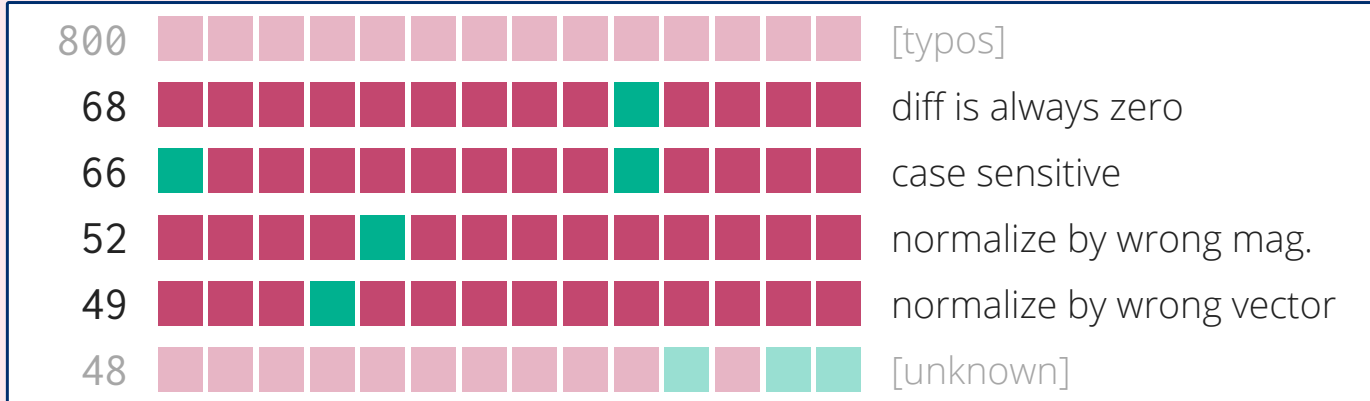


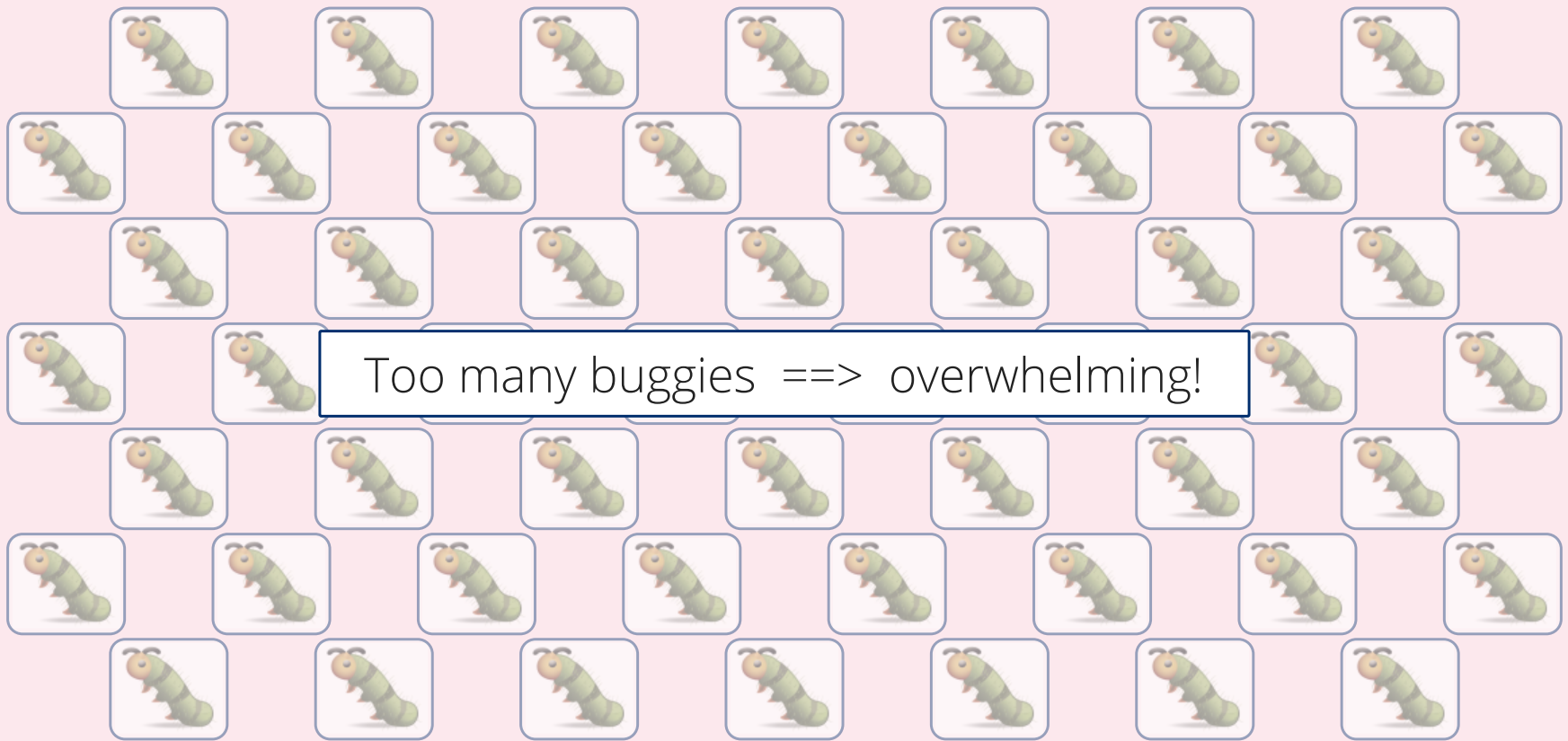
5. Sort clusters

800		[typos]
68		diff is always zero
66		case sensitive
52		normalize by wrong mag.
49		normalize by wrong vector
48		[unknown]

5. Sort clusters

6. Make buggies





Too many buggies ==> overwhelming!



6. Make buggies

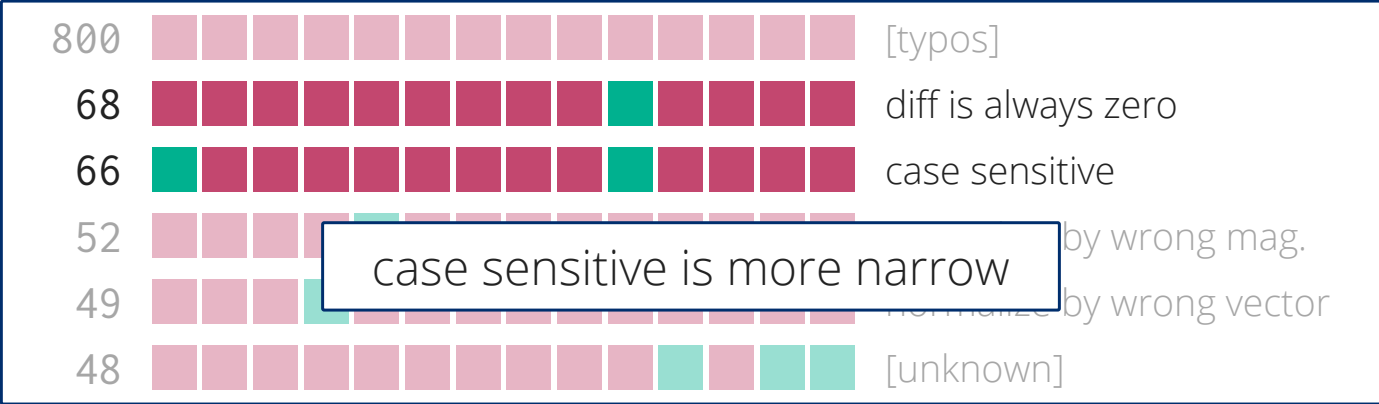
- * Focus on 1-2 ■
- * Favor narrow characteristics
- * Maximize subproblem coverage

6. Make buggies

- * Focus on 1-2 ■
- * Favor narrow characteristics
- * Maximize subproblem coverage

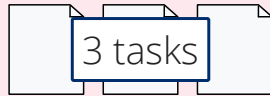
6. Make buggies

* Focus on 1-2 ■
* Favor narrow characteristics



Evaluation

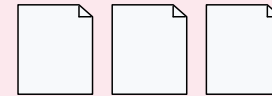
2020



2021

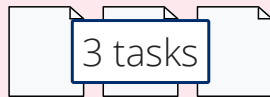


2022



Evaluation

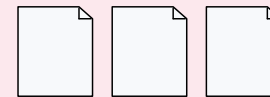
2020



2021

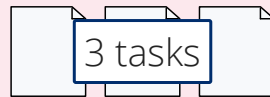


2022



Evaluation

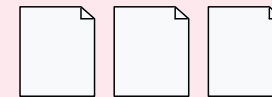
2020



2021

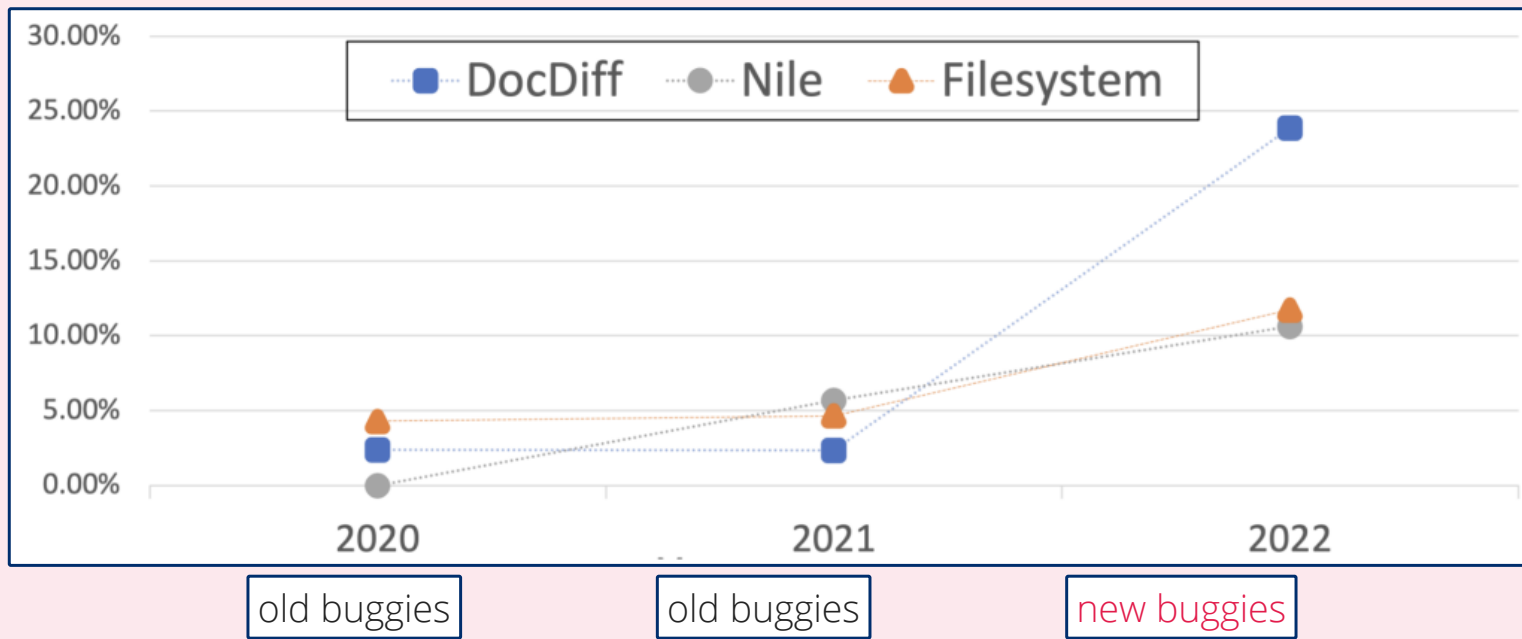




2022

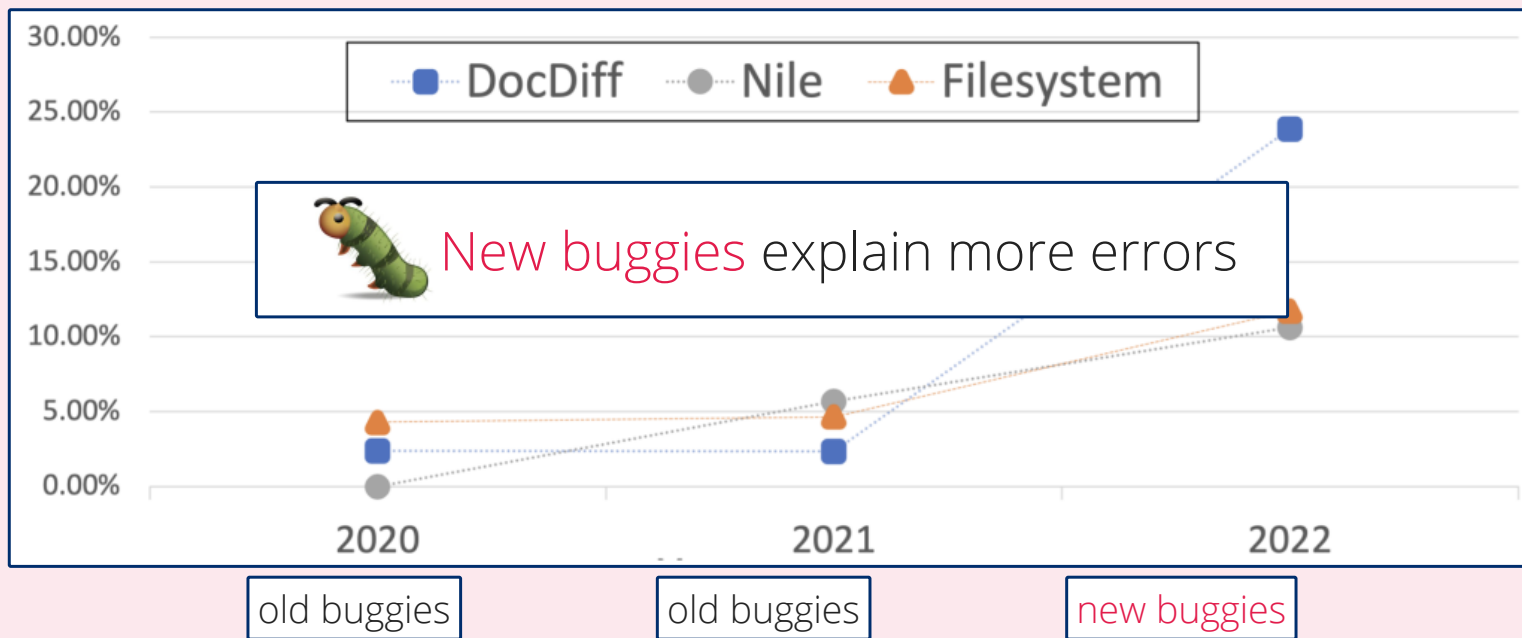


2020, 2021 ==> test
2022 ==> deploy

% explainable invalid tests
explainable = 1-■ or 2-■



% explainable invalid tests
explainable = 1- or 2-



High Effect Sizes for 2022



Matchup	Problem	95% CI	p value
2022 vs 2020	DocDiff	[-0.75, -0.57]	1.35E-29
	Nile	[-0.55, -0.26]	9.07E-14
	FileSys	[-0.35, -0.21]	2.35E-10
2022 vs 2021	DocDiff	[-0.70, -0.51]	6.87E-29
	Nile	[-0.27, -0.07]	1.82E-3
	FileSys	[-0.33, -0.19]	2.32E-9
2021 vs 2020	DocDiff	[-0.07, 0.08]	4.60E-1
	Nile	[-0.39, -0.13]	1.15E-17
	FileSys	[-0.06, 0.03]	2.52E-1

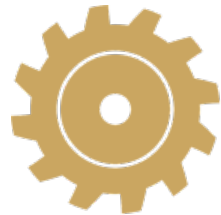
Weeks of
Data > Years of
Tuning



Promising approach for new problems



Recipe to uncover misconceptions
semi-automatic



Recipe to uncover misconceptions
semi-automatic

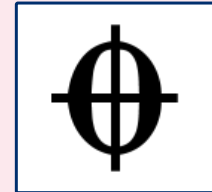


Data ==> better teaching

What's next? Hinting



What's next? Hinting





docdiff-tests.arr

INCORRECT CONSEQUENTLY, THOROUGHNESS UNKNOWN

These tests do not match the behavior described by the assignment.

```
13 overlap([list: "a", "b"], [list: "b"]) is 0
```

The assignment says:
Overlap must be proportional to the dot product of two vectors.





Deep Goal:

Rigorous methods for CS Ed research





Let's talk!

1. design problem

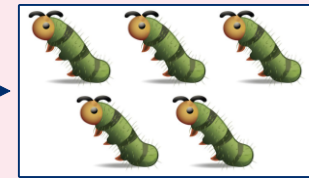
2. identify characteristics

3. collect invalid tests

4. cluster by feature vector

5. analyze top clusters

6. select buggies



Future

Data collection is a bottleneck
~1 semester ramp-up

+70% typos! How to reduce?
D4 / Data Druid

■ **Table 8** Our 2022 chaffs gave 1-m/2-m outcomes significantly more often than prior chaffs. The 2021 vs. 2020 results are similar except for Nile, which used D4 in 2021.

Matchup	Assignment	p value	Z score	Effect Size [95% CI] (Cohen's D)
2022 vs 2020	DocDiff	1.35E-29	-11.24	0.66 [-0.75, -0.57]
	Nile	9.07E-14	-7.36	-0.41 [-0.55, -0.26]
	Filesystem	2.35E-10	-6.22	-0.28 [-0.35, -0.21]
2022 vs 2021	DocDiff	6.87E-29	-11.09	-0.61 [-0.70, -0.51]
	Nile	1.82E-03	-2.91	-0.17 [-0.27, -0.07]
	Filesystem	2.32E-09	-5.86	-0.26 [-0.33, -0.19]
2021 vs 2020	DocDiff	4.60E-01	0.1	0 [-0.07, 0.08]
	Nile	1.15E-17	-8.48	-0.26 [-0.39, -0.13]
	Filesystem	2.52E-01	-0.67	-0.02 [-0.06, 0.03]

