

Crowdsourced Code Review in Programming Classes

Rob Miller
MIT CSAIL

joint work with Mason Tang, Elena Tatarchenko,
Max Goldman, Joe Henke

Problem: Feedback about Coding Style

- MIT 6.005 Software Construction
 - foundation-level programming course (replaced 6.001/6.170)
 - 400 students per year, mostly sophomores
- Students write lots of code
 - roughly 10kloc in problem sets and projects
- Automatic grading is necessary but not sufficient

```
// compute n! requires n >= 0
int factorial(int n) {
    if (n == 0) return 1;
    else return n * factorial(n-1);
}
```

correct and understandable

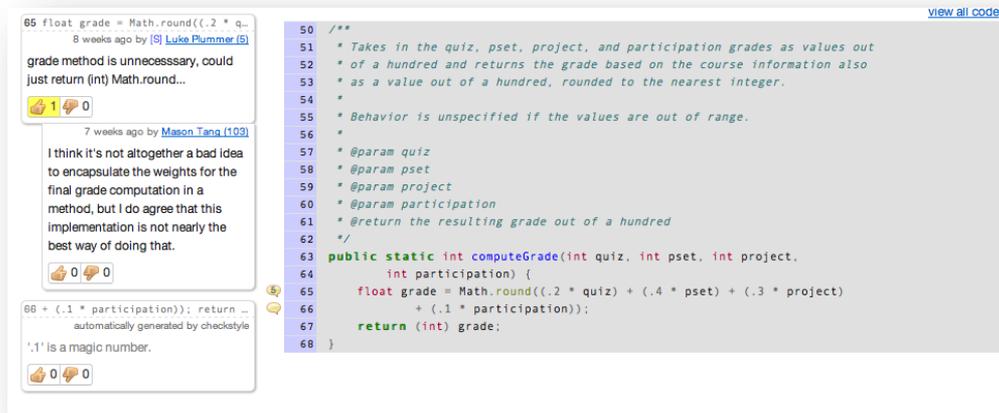
```
int factorial(int n) {
    int i, result=1;
    if (n == 0) result = 1;
    else {
        for (i = 1; i < n; ++i) result *= i;
        result = result*n;
        return result;
    }
    return 1;
}
```

correct but confusing

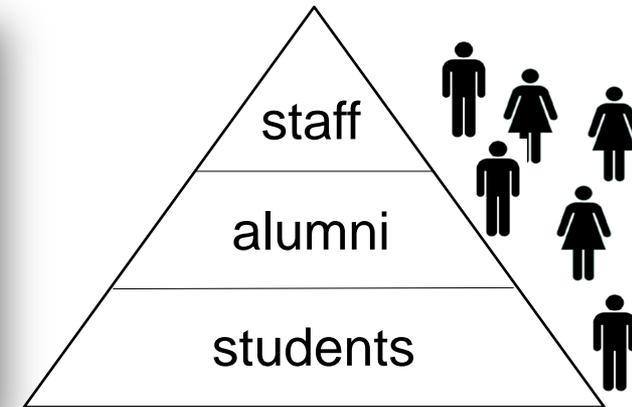
- we need human readers, and we want line-by-line feedback

Approach: Crowd-Driven Code Review

- Chop up student programs into **chunks**
- Review the chunks by a **mixed crowd**: students, staff, alums



The screenshot shows a code review interface. On the left, there are three comment boxes. The top comment, by Luke Plummer (6), says "grade method is unnecessary, could just return (int) Math.round...". The middle comment, by Mason Tang (103), says "I think it's not altogether a bad idea to encapsulate the weights for the final grade computation in a method, but I do agree that this implementation is not nearly the best way of doing that." The bottom comment, by an automatically generated checkstyle, says "'.1' is a magic number." On the right, there is a code snippet for a Java method named `computeGrade` that takes quiz, pset, project, and participation grades as input and returns a rounded float grade.



- Anticipated benefits
 - faster, cheaper, more diverse comments
 - give practice with code reviewing (a widespread industry practice)
 - expose to good and bad solutions
 - reduce workload on teaching staff
 - incorporate alumni back into the course
- Not using for grading... yet

Outline

- ✓ Introduction
- System
 - Caesar code-reviewing system
- Process
 - How we use Caesar in a programming class
- Experience
 - Some results from initial deployment last fall & spring
- Looking ahead
 - Issues to address
 - Opportunities for on-campus and online

Caesar: Divide & Conquer

programs chopped into **chunks**

each chunk assigned to multiple reviewers

```

14 public class RulesOf6005 {
15
16
17 /**
18  * Tests if the string is one of the items in the Course Elements section.
19  *
20  * @param name - the element to be tested
21  * @return true if <name> appears in bold in Course Elements section. Ignores case (capitalization).
22  * Example: "Lectures" and "lectures" will both return true.
23  */
24 public static boolean hasFeature(String name){
25     // TODO: Fill in this method, then remove the RuntimeException
26     String[] elements = { "lectures", "recitations", "laptops required", "text", "problem sets", "i
27     String test = name.toLowerCase();
28     for (int ii = 0; ii < 9; ii++) {
29         if (elements[ii].equals(test)) {
30             return true;
31         }
32     }
33     return false;
34 }
35
36
37 /**
38  * Takes in the quiz, pset, project, and participation grades as values out of a
39  * hundred and returns the grade based on the course information also as a value out
40  * of a hundred, rounded to the nearest integer.
41  *
42  * Behavior is unspecified if the values are out of range.
43  *
44  * @param quiz
45  * @param pset
46  * @param project
47  * @param participation
48  * @return the resulting grade out of a hundred
49  */
50 public static int computeGrade(int quiz, int pset, int project, int participation){
51     return (int)Math.round((quiz*.2) + (pset*.4) + (project*.3) + (participation*.1));
52 }
53
54
55 /**
56  * Based on the slack day policy, returns a date of when the assignment would be due, making sure not
57  * exceed the budget. In the case of the request being more than what's allowed, the latest possible
58  * due date is returned.
59  *
60  * Hint: Take a look at http://download.oracle.com/javase/6/docs/api/java/ut11/GregorianCalendar.html
61  *
62  * Behavior is unspecified if request is neagative or duedate is null.

```

code to review

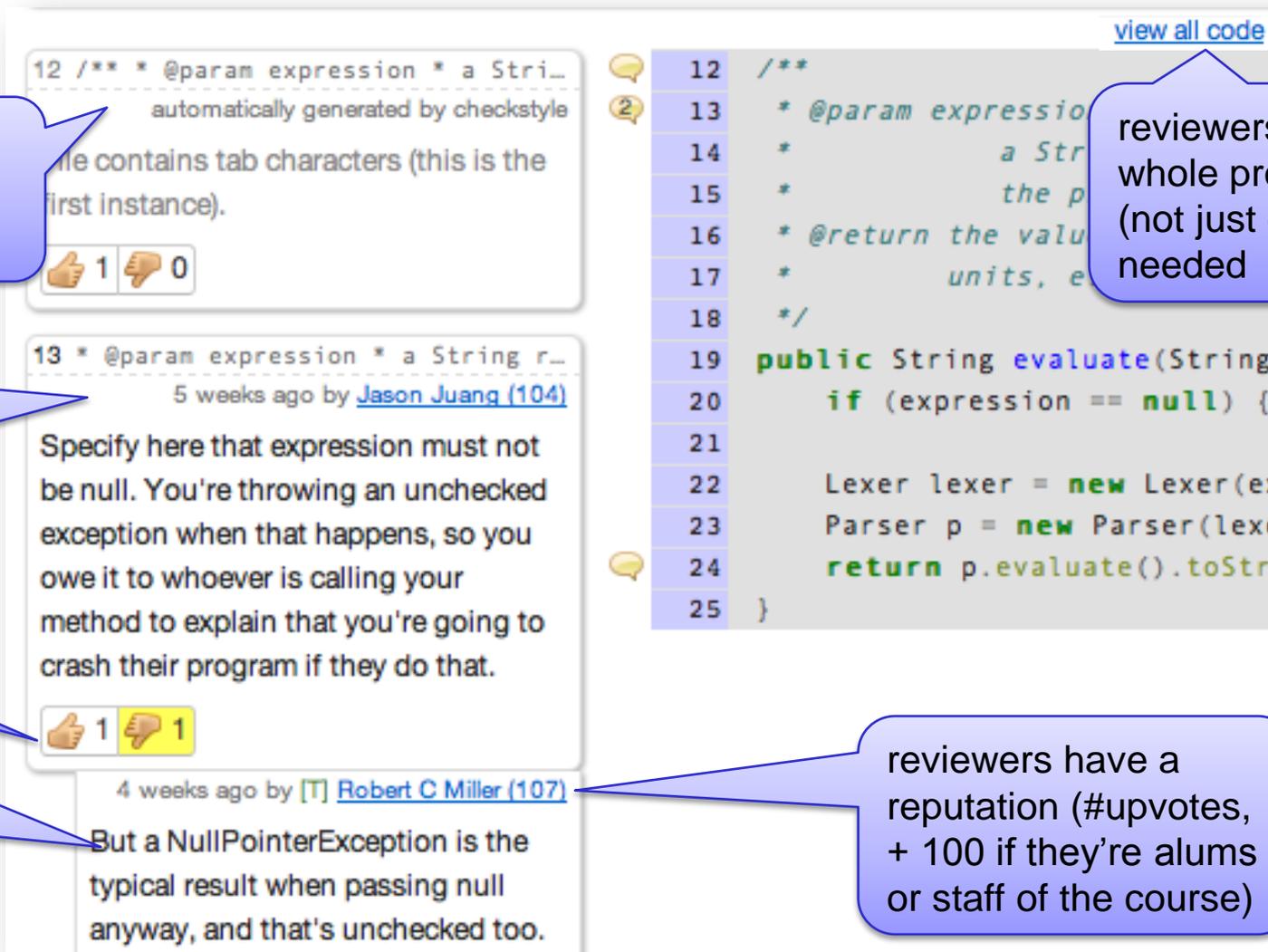
PrimeFactorsServer	5	1
PrimeFactorsServer	5	1
EchoClient	1	1
EchoClient	5	1
EchoClient	3	1
EchoServer	3	1
PrimeFactorsClient	6	1
PrimeFactorsServer	6	1
EchoClient	2	1
EchoClient	2	1

code recently reviewed

RulesOf6005.extendDeadline(..)	18	2
RulesOf6005.extendDeadline(..)	4	3
RulesOf6005.computeGrade(..)	9	3
RulesOf6005.extendDeadline(..)	6	2
RulesOf6005.computeGrade(..)	6	3
RulesOf6005.hasFeature(..)	3	2
RulesOf6005.hasFeature(..)	6	2
RulesOf6005.hasFeature(..)	4	2
RulesOf6005.hasFeature(..)	5	2
RulesOf6005.hasFeature(..)	4	2



Social Reviewing



The screenshot displays a code review interface. On the left, there are two code snippets with associated comments and voting buttons. The first snippet (line 12) is a Javadoc comment: `12 /** * @param expression * a Stri... automatically generated by checkstyle... contains tab characters (this is the first instance).` It has 1 upvote and 0 downvotes. The second snippet (line 13) is another Javadoc comment: `13 * @param expression * a String r... 5 weeks ago by Jason Juang (104)`. The comment text says: "Specify here that expression must not be null. You're throwing an unchecked exception when that happens, so you owe it to whoever is calling your method to explain that you're going to crash their program if they do that." It has 1 upvote and 1 downvote. Below this is a reply from Robert C Miller (107) 4 weeks ago: "But a NullPointerException is the typical result when passing null anyway, and that's unchecked too." On the right, a larger code snippet (lines 12-25) shows a Java method: `12 /** 13 * @param expressio 14 * a Str 15 * the p 16 * @return the valu 17 * units, e 18 */ 19 public String evaluate(String 20 if (expression == null) { 21 22 Lexer lexer = new Lexer(ex 23 Parser p = new Parser(lexe 24 return p.evaluate()).toStri 25 }`. A "view all code" link is visible at the top right of this snippet.

automatic style checker comments

reviewer comments

upvotes & downvotes

replies & discussion

reviewers can see whole program (not just chunk) if needed

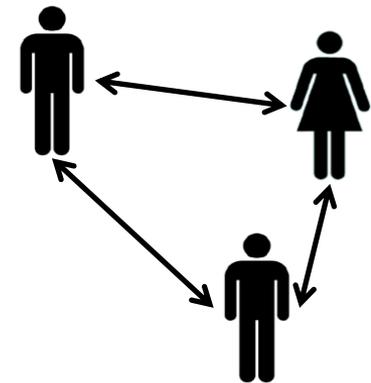
reviewers have a reputation (#upvotes, + 100 if they're alums or staff of the course)

Code Chopping

- Chop each submission into chunks
 - We've tried two sizes: method (~20 loc) and class (~150 loc)
- Sort chunks
 - Prioritize classes marked important by staff
 - Diff against staff-provided code and prioritize chunks with more student-authored lines
 - Deprioritize test code
 - Cluster chunks using MOSS fingerprinting
- Not all chunks get assigned for review
 - Typically 30%
 - So prioritization is important

Task Routing

- Lazy assignment
 - wait until a reviewer logs in before assigning tasks to them, since a given alum/student may or may not show up
- Multiple reviewers per chunk
 - assign 2 students/alums to a chunk to encourage discussion
 - assign 1 staff to the chunk too, for reviewing the reviews
- Fairness to code author
 - balance #reviewed chunks per submission
- Diversity of learning opportunities
 - maximize differences between reviewers on same chunk & same submission
 - role (student, alum, staff)
 - reputation score
 - # times assigned to a chunk together



Privacy vs. Visibility

- Code author is anonymous
- All reviewers are identified & their reviews are browsable

all users

Mason Tang (masont)

Robert C Miller (rcm)

Stephanie Yu (styu)

Jennifer Li (jennyli)

Elena Tatarchenko (elena_11)

Timothy Kaler (tfk)

Tana Wattanawaroon (tana)

Warut Suksompong (warutsuk)

Panupong Pasupat (ppasupat)

Stephen A Freiberg (saf)

Aysylu Biktimirova (aysylu)

Tamara Fleisher (tfleish)

Usman Masood (usmanm)

colleen josephson (cjoseph)

Maksim Kolysh (mkolysh)

Ivan Borsenco (borsenco)

Pedram Razavi (Prazavi)

Robert C Miller

rcm

reputation: 107

ps3

👍 Also, it might be a good idea to append some sort of general error indicator to the except...

👍 I would suggest against catching every error of every possible type; instead, focus on cal...

👍 What if you get a blank string, or a string with improper grammar?

👍 Looks good. Simple and elegant. Are you throwing only RuntimeExceptions? I'm assuming your...

👍 Why not just throw the exception? What if there's an error in the lexer?

👍 Your specs could use a little more detail, but the code looks fine.

👍 Right. Rep exposure!

👍 Looking at your implementation of Parser, this would be okay if result were final. However...

👍 The fact that this can return an error should be in the spec.

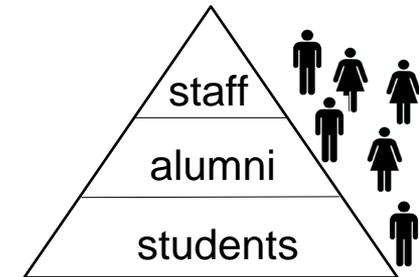
👍 I would argue that this exception should be thrown rather than caught. Otherwise, anyone ...

👍 but would be better to throw as an exception rather than returning as a string. Likewise ..

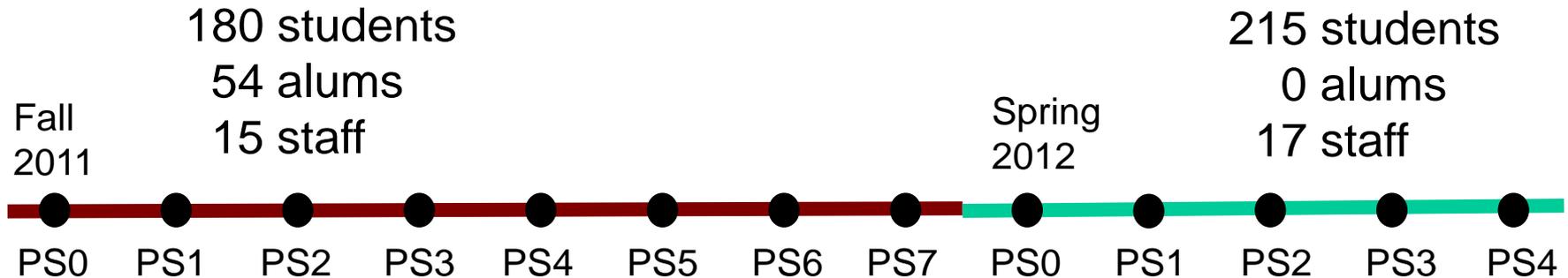
👍 That's fair. Hmm, maybe the course instructors should not put code in Caesar if the author...

Process

- Assignment due on Thurs
- Students review Fri/Sat
 - 10 chunks; must comment or vote on each
 - assignment is still fresh in their minds, curious about other's solutions
- Staff reviews Sun
 - 30 chunks; mostly voting
- Alums review whenever (Fri-Wed)
 - 3-5 chunks; must comment or vote
- Revise & resubmit assignment by Thurs
 - Main motivation is to fix bugs that lost points in autograding, but must address code review comments as well



Experience



13 problem sets, 2.2k submissions

21.5k comments

5% alums

8% staff

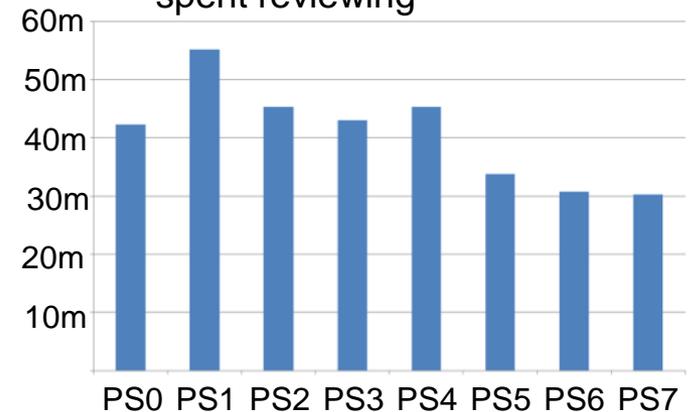
87% students

16.2% upvoted

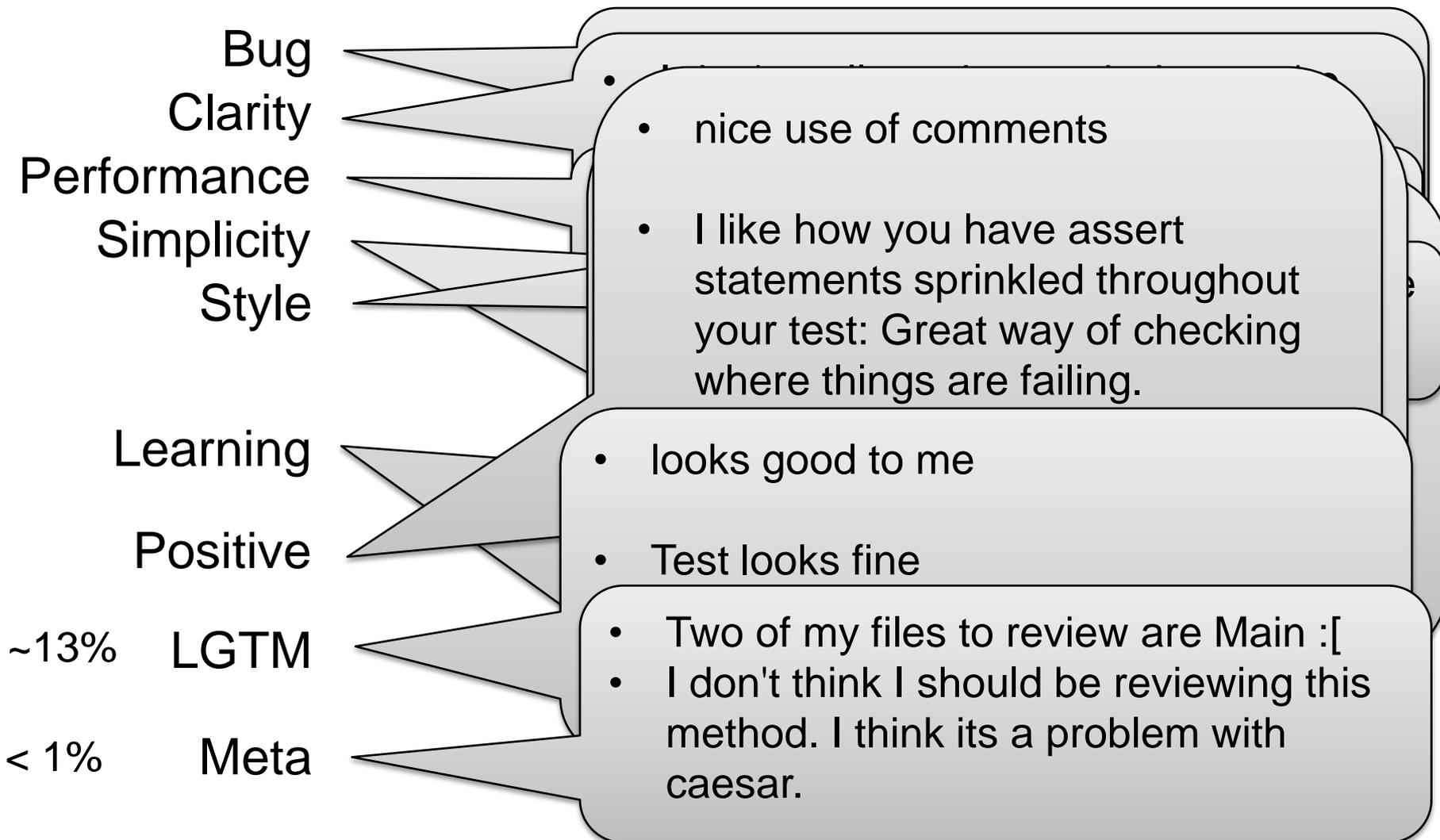
0.7% downvoted

9.6 comments per submission

average time students spent reviewing



Kinds of Comments



Interaction Happened

- 16% of comments were upvoted
- 22% of comments were in conversations (threads with at least 2 human comments)

- Corrections
- Clarifications
- Author responses

- **student:** Prefer equalsIgnoreCase because it says to ignore capitalization
 - **student:** The name = name.toUpperCase line seems to take care of capitalization.

- **student:** Agreed. I think the desired method here is add, not roll

- **student:** what does this method do?
 - **code author:** Essentially assertEquals. As someone else pointed out, this should have been an assert array equals, so that will happen in the future.

Alums Have Different Context

- **alum:** avoid abbreviating variable names... 'hi' would be especially confusing to someone who isn't familiar with english
- **student:** Idk about this one though. I've seen lo and hi at various places. I'd say this is fine. Though, the integer N should be lower case, just to conform with the java naming convention.
- **alum:** **When you're in industry and working on code with people who aren't necessarily familiar with English...** it's a problem to use phonetic abbreviations. There is really no good reason to abbreviate variable names especially since IDEs autocomplete.
- **student:** For that, yes, I agree. Sadly, some IDEs don't autocomplete very well. *cough*
- **code author:** These were the variables that were **given to us**

Reviewers Repeat Themselves

- **staff:** A magic number is a some fixed constant (numbers, Strings, etc.) used in a place where variables would likely add to code clarity or understanding...

- **staff:** Magic numbers are simply primitives that are used without having a descriptive variable name...

- **alum:** Returning null is usually dangerous, especially if your Javadoc doesn't mention that it might return null in some cases...

- **alum:** Don't return null in exception cases, especially if your Javadoc's @return line doesn't mention the possibility of returning null (even though it's mentioned earlier)....

- **student:** Should include more in rep invariant. ie a digit is allowed only once per row, column, and box

- **student:** Should include more in rep invariant. ie a digit is allowed only once per row, column, and box

Some Code Authors Are Self-Conscious

```

91  * You have been warned! DO NOT VIEW THIS FUNCTION
92  */
93  public Lexer.Token recurse(List<Lexer.Token> tokens) {
94
95      //      III  III  IIIIIIIII  IIIIIII  IIIIIIIII  IIIIIII  IIIIIIIII
96      //      III  III  III  III  III  III  III  III  III
97      //      IIIIIIIII  IIIIIIIII  IIIIIII  IIIIIIIII  IIIIIIIII  IIIIIIIII  IIIIIIIII
98      //      III  III  III  III  III  III  III  III  III
99      //      III  III  IIIIIIIII  III  III  IIIIIIIII  IIIIIIIII  IIIIIIIII
100 //
101 // IIIIIII  IIIIIIIII  IIIII  IIIIIII  IIIIIIIII  IIIII  III  IIIIIIIII
102 // III  III
103 // III  III  IIIIIIIII  III  III  III  III  III  III  IIIIIIIII  IIIIIIIII
104 // III  III  III  III  IIIIIIIII  III  III  III  III  III  IIIIIII  III
105 // IIIIIII  III  III  III  III  IIIIIII  IIIIIII  III  III  IIIIIIIII
106 //
107 //
108 //      .  _//_
109 //      / ' ' >
110 //      ) o' _/_ >
111 //      ( / _/ )_ \ >
112 //      ' "_/ / \ >
113 //      _//_//_
114 //      /,---, _//
115 //      "" /_/_//
116 //      /_CC_
117 //      ( _LLLL )\
118 //      |'_LLLLL )_ |
119 //      //__/_L) )_/
120 //      | _ |'_ _'( /'
121 //      L (-'_ _ _ | _ _'_
122 //      _ ) | _LL_ _/_ _'_
123 //      b'ger .((,-_ | ' " _L/_ _'_
124 //      '"/_ _ _'_/
125 //
126 //
127 // image taken from: http://www.chris.com/ascii/index.php?art=creatures%2Fdragons
128

```



Survey Responses

- Learning
 - I thought it was helpful. There are **more things I'll be looking for when I code because I saw them in other people's code and people commented** about it.
 - I like how reviewing the code after we've submitted it **allows us to retain the information** and concepts for a longer time.
- Self-doubt
 - **I think I have more to learn than to offer**, at this point though.
 - I was glad that the **code I was assigned were on the parts that I had understood** in my submission, so I could give better feedback.

Survey Responses

- Skepticism
 - Meh. I **don't think this will have a major impact** on students' coding.
 - No; **I'm fixing code instead of learning.**
- Keen to give and get more
 - Is it **alright to review more** than what is assigned on Caesar? I really like reading and reviewing other people's code, but don't know if over-reviewing would be frowned on.
 - I'm **thrilled** to get 51 human comments, but they are **ALL**, somehow, **on tests.**
- Scary
 - Very good UI. I enjoy commenting on code. I feel like the **Simon Cowell of code review.**

Related Work

- ColorMyGraph [Blank, Sutner, von Ahn]
 - crowd grading of proofs
- Github
 - social commenting on code commits
- Code Review StackExchange
 - threaded Q&A about short code chunks

Looking Ahead

- System changes
 - Capturing common advice in a wiki
 - Increasing reviewer leverage: searching for similar code
 - Promoting good comments for all users to read
 - Alumni engagement: recruiting, social networking, reputation management
- Other possibilities
 - Combining automatic grading (e.g. online tutors) with crowd-driven human feedback
- System is open source
 - <http://github.com/uid/caesar-web>