# **In-Flow Peer Review**

Dave Clarke
Uppsala Universitet
dave.clarke@it.uu.se

Tony Clear

Auckland University of Technology tony.clear@aut.ac.nz

Kathi Fisler WPI kfisler@cs.wpi.edu

Matthias Hauswirth
University of Lugano
Matthias.Hauswirth@usi.ch

Shriram Krishnamurthi Brown University sk@cs.brown.edu Joe Gibbs Politz Brown University joe@cs.brown.edu

Ville Tirronen University of Jyväskylä

ville.e.t.tirronen@jyu.fi

Tobias Wrigstad

Uppsala Unviersitet
tobias.wrigstad@it.uu.se

#### **Abstract**

Peer-review is a valuable tool that helps both the reviewee, who receives feedback about his work, and the reviewer, who sees different potential solutions and improves her ability to critique work. *In-flow* peer-review (IFPR) is peer-review done while an assignment is in progress. Peer-review done during this time is likely to result in greater motivation for both reviewer and reviewee. This working-group report summarizes IFPR and discusses numerous dimensions of the process, each of which alleviates some problems while raising associated concerns.

## 1. In-Flow Peer-Review

Peer-review has been employed for various reasons in Computer Science courses [61]. It is a mechanism for having students read each others' work, learn how to give feedback, and even to help with assessment. Indeed, of the six major computational thinking skills listed in the current draft of the AP Computer Science Principles curriculum [13], the fourth is:

P4: Analyzing problems and artifacts

The results and artifacts of computation and the computational techniques and strategies that generate them can be understood both intrinsically for what they are as well as for what they produce. They can also be analyzed and evaluated by applying aesthetic, mathematical, pragmatic, and other criteria. Students in this course design and produce solutions, models, and artifacts, and they evaluate and analyze their own computational work as well as the computational work that others have produced.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Copyright © ACM [to be supplied]...\$15.00. http://dx.doi.org/10.1145/

Students are expected to:

- Evaluate a proposed solution to a problem;
- Locate and correct errors;
- · Explain how an artifact functions; and
- Justify appropriateness and correctness.

Peer review clearly has a role to play in developing each of these skills. Students must read and evaluate proposed (partial) solutions, try to at least locate (if not offer corrections to) errors, offer their explanations for what an artifact is doing (especially if it does not match the expectations set by the problem), and justify their views on the appropriateness and correctness of presented solutions. Giving authors the ability to respond to reviews further reinforces the quoted principles.

Peer review has uses beyond merely evaluating programs. Writing benefits from peer review (Topping's review lists several exampes [61, page 261]), as do other artifacts that aren't just programs, like design diagrams, test suites, formal models, documentation, and presentations. All of these artifacts are also fair game for peer review in computer science courses and more, and peer review addresses similar underlying learning goals of evaluation and explanation of existing work.

This working group explored a particular variant of peer-review called *in-flow peer review* [47] (IFPR). In this model, peer review occurs while an assignment is in progress, before students submit their work for final grading. Performing peer-review in-flow has several potential benefits:

- It helps students better understand the problem specification. If the work they see others doing is inconsistent with their understanding, one or the other (or both!) might be confused. It is better to discover this while the assignment is in progress rather than after it is over
- Students are motivated to read feedback they get since it can affect their performance on the current assignment. In contrast, feedback given when the assignment is over may get less attention if students have moved on to other assignments.
- Students can apply what they learn from seeing examples of one another's work, and also learn to exercise judgment when

evaluating existing solutions. When a student sees another's work, she does not know the quality of the work she sees: it could be better than her own work, but it could also be worse. This takes some potential problems with plagiarism and turns them into a part of the learning process.

- It further emphasizes the comparative examination of work against a student's own.
- It helps students develop skills in a standard component of industrial software development.

Several challenges arise with this model, including figuring out how to decompose assignments for meaningful reviews, how to prevent students from gaming the process to avoid doing their own work, how to minimize the extra time this takes to complete homeworks, and how to help students not be led astray by weak or inaccurate reviews. Considering the potential learning objectives of IFPR, these challenges seem worth tackling.

This report summarizes activities of a working group on the promises and pitfalls of in-flow peer-review in computer science classes. The group members represented several countries and taught various courses at different levels (though the majority taught courses related to programming, programming languages, or other aspects of software development). We arrived at several different learning objectives an instructor might have for using peer review in a course (section 3) that drove our discussion. In addition, prior to the group's in-person meeting, each group member created two assignments for in-flow peer-review. These case studies, which are summarized in two figures (figure 2, figure 3), also helped to form the basis of many of our discussions.

Several assumptions and decisions scoped our work. We viewed IFPR as a mechanism for achieving certain learning goals, not as a way to scale grading (which is a use of peer-review in MOOCs [36]). We focused on person-to-person reviewing, rather than consider automated assessment tools that also provide a form of in-flow feedback to students. The members wanted to understand the benefits of writing reviews, a task which automated feedback tools eclipse. We also focused on *peer feedback* rather than *peer assessment*; Liu and Carless [38] describe the latter as targeting grading while the former targets collaboration.

# 2. An IFPR Roadmap

IFPR is a mechanism open to many policies. These policies are a function of a course's goals, student maturity, cultural context, and more. Therefore, an instructor who chooses to use IFPR will have to make several decisions about exactly what form they will employ. This section briefly outlines some of these decision points, with references to the rest of the document for more details.

#### 2.1 The IFPR Process

IFPR follows a particular process for assignments. In order to have an in-flow component, the assignment requires at least one reviewable submission that occurs before the final deadline of the assignment. This requires thinking through a few procedural questions that all in-flow assignments must address:

• The choice of submissions. How should an assignment be broken down into multiple stages? Even in a programming assignment, there are many choices: tests before code; data structures before tests; design documents and architectural specifications before code; multiple iterations of code; and so on. Some choices raise more concerns regarding plagiarism, while others only work under certain assumptions about software development methodologies. (Section 5.1)

- The distribution of reviewing. Should reviewing be distributed in order of submission? Randomly? Between students of similar or opposite attainment levels? Synchronizing review across students enables more policies on reviewer assignments, but incurs overhead for students and staff through more course deadlines. (Section 5.2)
- The manner in which reviews are conducted. This includes the choice of review "technology": should reviewing be mediated by a computer application or should it be done face-to-face (perhaps as a small group meeting around a table)? This also includes the use of rubrics: On the one hand, rubrics for reviewing guide the reviewer and may result in more concrete, actionable outcomes. On the other hand, a rubric can result in less constructive engagement and may result in important issues being missed. (Section 5.3.1)

#### 2.2 Issues Surrounding IFPR

There are a number of other cross-cutting issues that inform the choices made in the in-flow process, and affect the appropriateness and effectiveness of IFPR in particular contexts:

- The role of anonymity. When, if ever, should authors and reviewers know about each others' identity? Using single- and double-blind reviewing systems introduces trade-offs between protecting students' identity, creating the potential for abuse, and introducing students to norms of professional behavior, all of which need to be taken into account. Anonymity may enable more students to participate comfortably, at the cost of missed opportunities for creating cultures of collaboration and professional working behavior. (Section 6.3)
- The role of experts. Experts can play any number of roles: being entirely hands-off and treating this as an entirely student-run process (presumably after introducing the process and its purpose); intervening periodically; or constantly monitoring and even grading the responses. These roles set different tones between students regarding expertise and authority, but also define standards while students are new to the process. (Section 6.4)
- Suitability for non-majors. While IFPR is easy to justify for majors because of its correspondence to industrial practice (codereviews), the industrial argument makes less sense for non-majors. Objectives around collaboration and creating standards of evaluation, however, seem to apply to both majors and non-majors. (Section 6.5)

#### 2.3 Terminology

Throughout this report, we use the following terminology for the various artifacts, roles, and aspects of IFPR (figure 1 illustrates these terms and their dependencies graphically):

- An exercise is a problem set or assignment associated with a course; it may consist of multiple independent subproblems.
- A *stage* is a problem or task within an exercise that will be sent out for peer review.
- A piece of work on which a student will be reviewed is called a *submission*: this could be a piece of code, a paper, a presentation, or any other work on which peers will provide feedback.
- The *author* of a submission is the student who did the work associated with the submission (and who presumably would receive any grade associated with the submission, as well as any reviews about the submission).
- A review is written by a reviewer and responds to a specific submission.

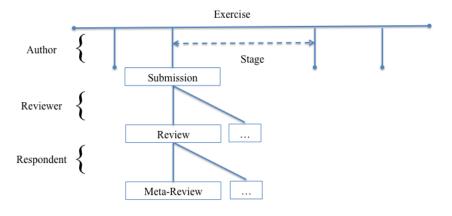


Figure 1: Illustration of Key Terminology

- A meta-review contains feedback on a review: this could be a
  grade of a review produced by course staff, feedback from the
  original submission author back to the reviewer, or any other
  commentary on the contents of a review. The provider of a
  meta-review is called a respondent.
- A reviewing assignment indicates which reviewers are expected to produce reviews for each submission.

The figure loosely shows their temporal ordering, with time increasing both downward and rightward. The diagram is intentionally ambiguous about the overlap of some events because different configurations of IFPR engender different temporal orders of reviewing events across students. The boxes labeled "..." mean there are one or more of the box immediately to the left.

## 3. Educational Goals of IFPR

Both peer review and the in-flow variant target a complex and interesting set of educational goals, some student-focused and some instructor-focused. Working group members were surprised at the subtleties that these goals brought to questions about how to configure IFPR. Indeed, many found discussions of the educational goals (and their impacts) the most thought-provoking aspect of our discussions. We lay out the goals here, referring back to them as we discuss configurations of IFPR throughout the report.

# 3.1 Student Learning Objectives

Fundamentally, IFPR fosters collaborative learning in which students can practice several critical skills:

- Assessing whether another's work satisfies problem requirements
- Providing actionable, useful, and appropriate feedback to others
- · Extracting high-level design choices from another's work
- Comparing others' high-level design choices and practices to one's own
- Deciding whether to adopt or ignore particular feedback or ideas
- Learning to value and grant authority to feedback from peers

The first two items arise primarily in students' role as reviewers, and are common to all forms of peer-review. The last two tasks arise more in students' roles as recipients of reviews, and have

more urgency in an in-flow context. The middle two arise in both roles. In the context of Bloom's taxonomy [5], these skills move students beyond "remember," "understand," and "apply" to "evaluate". They engage students in reflection and meta-cognitive thinking about their own work, while also requiring students to be able to communicate technical issues clearly to others.

Beyond these goals, regular comparison of one's own work to that of others can help students calibrate their abilities. In particular, it should provide means for students to gain confidence and selfefficacy in their work, and in discussing the works of others.

The extent to which IFPR targets these goals depends significantly on the artifacts students are asked to review, the criteria or rubrics through which they produce reviews, the means through which students are expected to respond to reviews, and the feedback students receive on their reviews. Some configurations of peer review, particularly those designed simply to scale grading, naturally and necessarily de-emphasize some of these goals. Section 6 explores these tradeoffs in detail.

The emphasis on collaboration in these goals illustrates that IFPR is an example of a Contributing Student Pedagogy (CSP), a pedagogy in which students (1) contribute to the learning of others and (2) value the contributions of other students. A 2008 ITiCSE working group report describes various facets of such pedagogies [24]. IFPR targets the second criterion (valuing the contributions of other students) more than traditional, post-submission peer review. Various parameters in implementations of IFPR affect the extent to which students contribute to the learning of others in practice: shallow reviews, for example, arguably meet the letter but not the intent of a CSP. Separately, IFPR has goals beyond CSP: writing reviews offers benefits to the reviewer as much as the reviewee, and often the learning goals that inspire IFPR (and peer-review in general) are more focused on the reviewer than the reviewee. Nonetheless, many of the theoretical underpinnings of CSPs also apply to IFPR, and thus affect the ideas in this report.

Of course, IFPR also has the potential to interfere with student learning. Reviewing asks students to switch between very different tasks (programming and reviewing); depending on the timing of reviewing, this could affect students' cognitive load. Careful design of exercises is important: allowing students to share parts of their solutions through reviewing can, for example, affect whether students stay in their zones of proximal development [67].

#### 3.2 Instructor Goals

From an instructor's perspective, IFPR can target several objectives, including:

- 1. Emphasizing the importance of writing in technical contexts
- Providing human feedback more scalably and more timely than with only expert assessment
- Providing an additional perspective on how students perceive course material, since students express their understanding in a different medium than their assignment submissions
- Increasing social interaction within computing and programming, addressing a common misconception about working in the discipline
- Fostering engagement of and interaction between students from different cultures
- Helping students improve performance and learning while actively engaged with course material
- Re-examining plagiarism issues by casting the re-use of classmates' ideas and code in a positive light, and including grading mechanisms that take this into account

The first three objectives arise in most forms of peer-review. The next two arise in general peer-review, though the immediacy of the in-flow context likely enhances their impact. The last two are more directly associated with IFPR.

All of the working group members were interested in IFPR more as a way to enhance students' learning than as a way to scale grading. Indeed, most members were open to (if not already) investing additional staff resources in making sure students were learning reviewing skills from a peer-review process. The members were interested in the insights they could gain as instructors from reading students' reviews (as per the third objective), though none believed that grading all of the reviews was scalable or cost-effective.

The group also coalesced around the social benefits of peer review, seeing this as an important aspect of developing competent professionals. Many of our discussions about giving review feedback and whether reviewing should be anonymous revolved around the impacts these issues could have on collaboration and socialization through peer review.

# 4. Examples of IFPR

The case studies from working group members covered a variety of student levels and course types. More interestingly, they varied widely in the kinds of artifacts and processes that they suggested for IFPR. Figure 2 and figure 3 summarize the key parameters of the case studies: the former covers introductory-level courses and the latter covers upper-level courses. The table lists the name and level of the course, describes the assignment in which IFPR was used, describes the submissions that were reviewed, and describes the review criteria for submissions. Assignments with [\*] after their descriptions have been used in actual courses; the rest are hypothetical uses proposed by working group members, based on exercises in their current courses. The descriptions of the case studies are available at https://github.com/brownplt/iticse-in-flow-2014/tree/master/in-flow-assignments.

## 5. The In-Flow Process

Several steps are part of any in-flow assignment, as illustrated in figure 1. This section lays out the process of submissions, reviews, and meta-reviews in more detail, along with the design decisions that the group identified for each activity. Section 6 continues by discussing issues that cross-cut the process.

#### 5.1 Stages and Submissions

As the case studies in figure 2 and figure 3 show, IFPR can be used with many different kinds of submissions. Even once an instructor has identified the general class of artifacts to review (such as papers versus code), she can choose different ways to use IFPR to build up to final versions. The group identified four broad choices in the artifacts to be reviewed:

• Multiple iterations of the same specific deliverable: This approach is the most similar to existing peer grading approaches, where an entire deliverable is presented for review. This mirrors common practice in courses where students do peer review of written work, which is well-studied in contexts other than computer science [8, 61]. Little extra work in assignment design is required to have students review drafts, so this provides a low-friction way to adapt an existing assignment for IFPR. One consideration is that plagiarism can be more of a problem in programming tasks that are the same across students than in writing tasks where goals are less objective and more variance is expected. We discuss plagiarism (and mitigations) more in section 6.1.

Case studies 6, 9, and 12 included submission steps that were prototypes or drafts of the final product.

• Multiple iterations of an evolving deliverable: Some projects don't have drafts as much as an evolving set of specifications and deliverables. For example, in a long-term software engineering project, the demands of the system may change over time as new requirements are discovered. This is distinct from multiple iterations of the same deliverable because the goal itself changes along with the submission. Both high-level feedback on the direction of the project and low-level implementation feedback can be helpful in this setting.

Case studies 10, 12, 13, and 15 have elements of this approach, where the deliverable's nature changes more over time, and in response to review.

• Separate deliverables that build on each other: Often, programming assignments can be broken down into several subproblems, often expressed via decomposition into helper functions or separate classes. If the assignment naturally fits this kind of breakdown, a natural strategy for using in-flow review is to review individual pieces, to catch mistakes in the components before composing them into a final solution. This approach lends itself well to detailed evaluation, because the separate components can be assessed in isolation, in both their prereview state and their state after final submission. Also, if the instructor (rather than the students) decides on the decomposition, they have a lot of control over the path students take through the assignment, which can inform decisions about rubrics and feedback guidance.

Case studies 1, 2, 5, 7, 8, and 15 use this approach.

• Incremental views of the same deliverable: Programming problems can have a number of distinct artifacts associated with them beyond code: They have documentation, tests, specifications, underlying data representations, measurable behavior like time and memory complexity, visualizations, and more. Having students produce different views on the same problem, with review in between, is another way to break up the assignment. For example, an assignment might proceed with a test-first or test-last approach, with review before or after a test suite is written. This focuses students on different aspects of the problem in their different stages of review.

#	Course and Level	Exercise	Submissions	Review Criteria
1	Computing for	Cluster data on voting records	Code and tests for	Provide scores from 0 to 100 on each of
	Social Sciences and	(US Senate) to identify	instructor-defined subsets of	(a) whether tests meaningfully capture
	Humanities	senators with similar	overall functionality	the assignment purpose, and (b)
	(undergrad	ideology [link]		whether code performs the
	non-majors)			corresponding computation correctly
2	CS1	Write code and assertions for	Work so far on subset of functions	Rate readability and correctness;
		various components of a	designated by instructor	additional free-form comments
		pinball game [link]		
3	Advanced CS1 with	8	Datatype definition with instances	Indicate whether (a) data structure can
3	Data Structures	incremental and functional	of the data, test cases, complete	support required operations within time
		updates on trees [*] [link]	programs	bounds, (b) interesting examples of data
				are missing, (c) tests offer good
				coverage and are correct
4	Programming	In-class clicker assignment to	CFGs for code snippets (drawn	Provide yes/no assessment of whether
-	Fundamentals 2	explain control-flow through	through custom software package)	CFG is accurate
	(2nd semester	if-statements [*] [link]		
	undergrad)			
5	Imperative and OO	Implement a program that	Description of learning goals	Instructor-provided template on choice
	Programming	satisfies a student-selected set	covered by program, program	of goals, whether program satisfied
	Methodology (2nd	of learning goals [link]	code, and give presentation on how	them, and presentation quality
	year)		program achieves goals	
6	Imperative and OO	Implement simple Pong game	Tests, two draft implementations,	Questions about whether key
	Programming	in model-view-controller	and a final implementation	components are present; whether tests
	Methodology (2nd	style [link]		are reasonably complete and motivated;
	year)			whether good code practice followed
				(i.e., naming, structure, indentation)

Figure 2: Summary of case studies: first-year, second-year and non-majors courses ([\*] indicates use in an actual course)

Case studies 3, 6, 10, 13, 14, and 15 take this approach, where submissions are different kinds of artifacts, but all contribute to the understanding of a single artifact.

The nature of the reviewed submission naturally affects the time required for review, as well as the amount of expert guidance required (students will have more experience with evaluating some artifacts over others). In turn, the course's learning objectives should guide the choice of artifacts: preparing students to participate in industrial code review, for example, will be better served by using IFPR on code-based artifacts rather than written papers.

### 5.2 Assigning and Scheduling Reviews

The IFPR process requires building time into assignments for performing review, and deciding how to assign reviewers to artifacts for review. We discuss each in turn.

#### **5.2.1** Scheduling Decisions

The major scheduling decision in IFPR is whether reviewing happens *synchronously* or *asynchronously*. Synchronous reviewing occurs when all authors submit their work for review (and reviewing commences) at the same time. Asynchronous reviewing occurs when authors submit their work for review when it is ready, and different students are in different stages of review at the same time. The two scheduling modes have several tradeoffs:

- If all students are forced to submit before reviewing starts, there
  is the full pool of reviews to draw from in any review assignment strategy (section 5.2.2). In the asynchronous setting, students can only review submissions that occurred before theirs,
  which can skew the reviewing process if, for example, highachieving students tend to submit early.
- With intermediate deadlines, all students have the same time to use review information. If there are no intermediate deadlines,

students who submit later have less time to use the information from the review process before the deadline.

- In the asynchronous setting, students who want to work at their own pace can, and the process doesn't discourage them from continuing with their work while they await reviews. With intermediate deadlines, a student cannot progress on her own schedule
- If reviews are available and presented to students after they submit, the problem is more likely to be fresh in their mind. In the synchronous setting, there can be a longer gap between submission and review. It's not clear if one is particularly better than the other: coming back to a problem after not thinking about it for a while can be beneficial, but it also takes time to recall the problem and re-load it into working memory in order to perform review.
- Synchronous reviewing requires extra scheduling overhead that
  is likely to lengthen assignments for purely logistic reasons.
  Asynchronous reviewing doesn't require extra scheduling in the
  assignment, it just changes the workflow that students follow.

With synchronous reviewing, there is an extra question of whether there should be separate time set aside for reviewing in between submissions (with a review submission deadline), or if reviews can be completely in parallel with the next submission step. This can affect the timeliness of reviews, which can affect how useful the review is to the reviewee as they move forward with the assignment.

# 5.2.2 Assigning Reviewers to Submissions

Whether reviewing is synchronous or asynchronous, there needs to be a strategy for assigning students to review submissions. We identified both a number of methods for assigning submissions to students to review, and several miscellaneous modifiers that could

#	Course and Level	Exercise	Submissions	Review Criteria
7	Introduction to Functional Programming (upper undergrad/MS)	Implement Boggle (find all valid words in 4x4 grid)  [link]	Decomposition of overall problem into tasks (with QuickCheck assertions), tests, code	Check decomposition makes sense, presenting alternative if own differs from reviewed one; try own test suite on the code being reviewed
8	Advanced Software Design (upper undergrad/MS)	Design and implement subset of a mobile app+server for a game using iterative development [link]	Design documents so far	Free-form comments on comprehensibility, quality of documentation, coverage of use cases, adherence to design principles, and choice of subsystem to implement; concrete examples required to illustrate each point
9	Collaborative Computing (MS)	Collaboratively produce a research article [*] [link]	Drafts of article	Conference-paper reviewing rubric: questions on suitability for audience, originality and demonstrated knowledge in contribution, evidence for arguments, methods, presentation, etc.
10	Software Security (upper undergrad/MS)	Find ways to attack a web-based application (black-box, then white-box)	Description of strategy to use in attacking the application in black-box fashion.	Free-form comments on comprehensiveness and appropriateness of attack strategy
11	Software Modeling and Verification (upper undergrad/MS)	Use model checking to find flaws in a protocol [link]	Proposed model of the system environment and desired properties that should (not) hold under this model	Assess whether model conforms to problem and whether model supports/masks the properties provided with the model; comment on good/bad features of this model
12	Software Performance (MS)	Develop an extension to the Jikes visual debugger [link]	Proposed extension, prototypes, final artifact	Comment on one thing they particularly like and one aspect that could be improved; evaluate prototypes following in-class presentations by each team; review final artifact for usability, extensibility, and documentation
13	Logic for System Modelling (upper undergrad/MS)	Write a relational (Alloy) model of an elevator [*]	Model of data components, description of desired properties of model, initial model of elevator operations	Comment on whether components/properties are missing, whether they are reasonable, and whether model is suitably operational or too declarative.
14	Programming Languages (upper undergrad/grad)	Provide a test suite and implementation for a type checker [*] [link]	Tests first, implementations later (submission deadlines not synchronized across students, but must occur in order per student)	Set of ~10 specific questions about test coverage, plus free-form comments on style or organization of test suite; no peer review on implementations
15	Software Security (MS)	Implement simple online web-app on a strict timetable, then create attack trees for it [link]	Initial program, attack trees, and secured application along with review of differences between original and secured application and results of using static analysis and fuzzing tools on the implementation	Free-form comparison to what was done in own solution

Figure 3: Summary of case studies: upper-level undergraduate and graduate courses ([\*] indicates use in an actual course)

apply. We also note when a particular strategy is more or less appropriate in synchronous or asynchronous settings.

- Random Assignment: Perhaps the most obvious and simple method for reviewer assignment is random: each reviewer is assigned one or more submissions at random to review. There are of course many types of randomness; it is probably useful to ensure that all submissions get the same number of reviews, for example. In the asynchronous setting, the pool of reviewable submissions will necessarily be smaller (since only a subset of submissions have already come in); this skews the selection. In this setting, random assignment also lacks temporal fairness: the most recent submission isn't guaranteed to be reviewed first, which can weaken the benefit of quick feedback in the asynchronous model.
- Temporal Assignment: Reviews can also be assigned in the order submissions were received. It's not clear that this makes much sense for synchronous review, where temporal order is somewhat unrelated to motivations for assigning reviews. However, in the asynchronous case, assigning reviews in the order submissions are received helps ensure that feedback happens as quickly as possible (assuming that students complete reviews at around the same rate). We also note that if students are aware that temporal ordering is occurring, they can collude to time their submissions in order to ensure or avoid particular review pairings. Using a mix of randomness and temporal ordering could alleviate this somewhat, at some minor cost to review turnaround time.
- By Metric: There are a number of metrics that could be used to assign reviews with the goal of getting more effective feedback for students. We identified:
  - Achievement: Reviewers could be matched to authors with similar or different levels of achievement on past assignments (or, if the assignment can be evaluated automatically, even on the current assignment). Existing research shows that on group work, pairing weak and strong students can help the weak students (though the strong students don't do as well either) [30]. The effects of such an assignment are certainly course- and assignment-specific.
  - Prior Review Quality: If the course tracks review quality through meta-reviewing (section 5.4), the system could assign consistently strong reviewers to weak work in order to maximize improvement (again, "weak work" could be predicted by past achievement of students, or by an automatic grading system).
  - Similarity Between Solutions: Reviewers may learn more from reviewing a variety of solutions that are different from their own. They may also be more able to review solutions that are similar to their own. Depending on the assignment and learning goals, it could be valuable to group students based on solution approach or code similarity.

We expect that by-metric assignments work best synchronously, because it is difficult to perform the assignment until the metric can be measured for all students. Doing a by-metric assignment of reviews asynchronously is possible if the metric is known before submission (e.g. if only using past submission performance), but it would result in some students waiting for their assigned-by-metric reviewee to submit.

• Student-chosen: Reviewers could also be involved in choosing which submissions (including by which authors) they review. For example, a simple model could have all submissions go into a publicly-visible pool of submissions, from which students choose submissions to review. The assignment could re-

quire that students perform some number of reviews, and remove submissions from the pool once they have been reviewed enough times in order to avoid a small number of submissions getting more reviews than others. This works with both asynchronous and synchronous-style scheduling, but can be problematic if the limits are too rigid and some students submit very late (leaving no submissions available to review for students who submit early). A solution is to pre-seed the set of submissions with some instructor-provided submissions in order to ensure enough supply.

Students could also choose partners to review independent of particular submissions. It might be reasonable to switch to a student-chosen partner approach after having assignments with other review assignment strategies, once students have decided they enjoy collaborating. This also works fine in both synchronous and asynchronous styles: students may even arrange to complete their work at a time convenient for their reviewer in order to get the most prompt feedback.

Papadopoulos et al. explored different strategies for assigning students in peer review in a computer science course on networking [45]. They find that students who select their peers freely perform better (according to experts rating the utility and clarity of reviews) than pairs where the students were assigned (randomly) by the instructors beforehand.

All student-chosen strategies for review interact heavily with choices about anonymity of the review process, which we discuss in more detail in section 6.3.

In addition to these strategies, there are a few other factors for educators to consider in assigning reviews:

- **Groups vs. individual**: The review assignment strategies in this section aren't limited to only pairs of reviewers and reviewees. It would be perfectly reasonable for a group of reviewers to create a review together through discussion, whether online or in person. The same parameters of randomness, temporality, and so on apply. Anonymity is possible but more difficult when a discussion among multiple students is involved. We discuss the contents of reviews and review discussions more in section 5.3.
- Class-wide review: At the extreme end of group review is a
  review that involves (potentially) the whole class. This could
  be, for example, a presentation that the whole class comments
  on immediately, or a collective review process as in a studio
  art course, where work is presented and discussed publicly.
  Using an online tool, work can even be published publicly but
  anonymously, and allow for any interested class member to
  comment.
- Mutual reviews: Reviewer-reviewee pairs can be mutual or disjoint – students may form pairs (or groups where everyone reviews everyone else) that review one another, or there may be only a one-way connection between reviewer and reviewee. Mutual reviewing could provide more concrete motivation (other than abstract altruism or a grade), if students are helping someone who is actively helping them in return. Mutual reviewing can still retain anonymity.
- Persistent review assignments: In an assignment with two or more reviewed stages, reviewers can change at each stage, or continue to be the same throughout the process. One study reports that, for assignments with more than 4-5 stages, switching authors at each stage made the reviewing burden onerous [47]. Continuing to review the same author's work may have lessened this burden, since the comprehension effort from earlier stages would carry over.

Other strategies for assigning reviewers to submissions exist. We have assumed that the reviewers are drawn from within the class (which may not be the case in a peer-mentoring situation [41]). In a class where there are communication or language barriers between students, it may also make sense to assign reviewers so that communication is maximized or the challenge of communicating in another language is maximized. It may also be useful to secretly or not secretly assign instructors or TAs as reviewers and reviewees sometimes in order to guide or monitor the process. We discuss instructor and TA participation in reviews more in section 6.4.

### 5.3 Performing Review

Much of our discussion of how to conduct reviewing focused on *review rubrics*, which can be used to focus student feedback on specific features of the assignment. We considered whether information beyond submissions would help reviewers. We also discussed several forms that reviews could take, noting that technology often guides programming courses towards text-based feedback.

# 5.3.1 Review Rubrics

Rubrics serve two important goals in any form of peer review: they communicate expectations to reviewers (serving as a form of scaffolding), and they help foster a baseline of quality in all reviews. While these goals suggest highly-structured rubrics, overly structured rubrics can limit reviewers' and authors' attention to the questions on the rubric. They can also provide too much scaffolding, especially once students need to practice evaluating work from scratch. The tradeoffs around designing rubrics must balance these tensions.

Rubric design must consider the rubrics' utility for the reviewer, the author, and the instructor seeking to understand how students are performing. These goals are not necessarily at odds with one another, but may conflict incidentally when picking a particular configuration of rubrics.

The working group identified several potential roles for rubrics:

- Rubrics as scaffolding for reviews: Rubrics help students learn how to construct good reviews, especially for students new to the process. A beginning student who is learning to both read and write code might not know where to start in critiquing a program. Prompting with specific questions helps in situations where students don't yet know how to structure a review from scratch.
- Rubrics for focus: A rubric can focus reviewers' attention on different questions that reflect the goals of an assignment. For example, it could prompt code-specific questions ("Is this code well-documented?", "Are all the type annotations correct?", etc.), problem-specific questions ("Is there a test for a list with duplicate elements?", "Does this program meet the problem specification for input X?", etc.), or questions that encourage actionable feedback ("Provide a test case that this solution does not pass."). Students may also optionally ask for reviews with a certain focus when offering submissions for review. This might help ensure relevance of the review for the author.
- Rubrics as an alibi: Rubrics can be used as alibis for reviewers who fear criticizing works of others because of cultural values, self-image, or other factors. For example, being asked to point out one part which could be done better, or to identify errors will shift the blame from the reviewer who found the bugs to the instructor who provided the rubric.
- Rubrics for reviews of good solutions: In at least one case of using IFPR in the classroom, reviewers reported not knowing what to write when reviewing good solutions [47]. A rubric could explicitly prompt for feedback even on good work (e.g.

"What did you like about this submission?", "List one thing you would change, regardless of correctness", "What should the author *not* change in this solution?"), so that reviewers don't simply sign off on a solution as good enough without reflecting and providing some useful feedback.

- Rubrics for conduct: Rubrics can guide reviewers towards a professional and appropriate tone for giving feedback, and help frame negative feedback in a constructive way. For example, forcing a review to contain comments on the strengths of the submission under review can soften other criticism. When appropriate, rubrics can guide reviewers towards more constructive language for example, "This could be done differently" vs. "This is wrong".
- Rubrics for time management: Open-ended review tasks don't make it clear how much time reviewers should spend on them. Just having a specific rubric can make it easier for a reviewer to identify when they are done (and estimate the time themselves). A rubric could even specify the amount of time that a reviewer should spend, and how much on each part of the assignment, to ensure that reviewing does not take up more time than intended.

Evolving rubrics across a course or curriculum may offer a good balance between initial scaffolding (for reviewers and authors) and eventual opportunities for both groups to demonstrate critical-thinking skills. One model would evolve rubrics from having fairly targeted questions to asking broad questions: this model gradually removes scaffolding. Another model starts with concrete questions (such as "Do these tests look correct") and progresses to questions on more abstract issues (such as "Do these tests cover the space of possible inputs") as students master more of the subject material.

A variation on evolving rubrics would allow different students to work with different review forms, depending on their ability as reviewers. This comment arose from the working group members' experience as conference program-committee members: members often found overly structured forms to be annoying, feeling they interfered with how they wanted to convey issues with a work. However, they also noted that early on in their paper reviewing career, they appreciated the rubrics' ability to get them past the initial blank form.

A different form of variation might pose more questions to reviewers than are conveyed to authors. This situation could make sense when the review is used to assess the reviewer's understanding of a work, or when too much information in a review might distract the author from the critical information in a review.

Structure enables certain comparisons between reviews. Inexperienced authors may benefit from structure when aggregating the feedback of multiple reviews: for example, structure could help authors understanding that two or more reviews give contradictory advice. Two students discussing reviews (that they are making or have received) may be similarly helped by an imposed structure. Certain kinds of structure can enable automated analysis of reviews, which can provide useful diagnostics to both instructors and students. Similarly, software tools have the potential to provide richer dashboards when review comments are structured.

Discussing reviews and rubrics with the entire class is another good example of using rubrics for communication. Students or experts might see common problems which should be communicated to all, either by sharing sufficiently general comments with the entire class or even adding an entry to a rubric which brings attention to the issue in subsequent reviews.

#### 5.3.2 Information Provided to Reviewers

In some cases, reviewers can be provided with information beyond the submission. When submissions are source code, for example, reviewers could be given both the submission and information about how the submission held up against an instructor-defined test suite (whether or not that information is available to the submission's author). On the one hand, information such as a test-suite score may reduce the time burden of reviewing; on the other hand, it could have the downside of reviewers only focusing on the issues that auto-grading revealed, masking situations in which the auto-grading missed something important (Politz et al. observed cases in which reviewers were more negative than grades from an instructor-provided test suite [46]).

Additional information for reviewers provides an implicit rubric, subject to the same tradeoffs we discussed regarding rubric structure. Instructors should bear this in mind when considering whether additional information is actually helpful to the overall process.

#### 5.3.3 Forms of Reviews

Reviews can take various forms, from written documents to verbal feedback, from paragraphs to small comments associated with particular fragments of prose or code, and from individual to groupwide feedback. For written reviews, the group noted the general applicability of plain text, but modern software tools (such as Github and other graphical version-control tools) enable targeted comments and conversations between authors and reviewers down to the line number in a particular revision. These conversations have more structure than untargeted comments about the entire submission

In some situations, non-text artifacts can be effective, as not all submissions need to be code. Code-architecture diagrams can be critiqued and marked up with freehand annotations, pictures of the state of a running program can be drawn, and code patches can be used to convey comments. In these cases, however, technology choices can become a limitation.

The Informa tool allows students to give live feedback on problems with several interfaces that could also be useful for review [28]. For example, during a Java program comprehension task, students use a drawing tool to create a graphical representation of the heap at particular program points. Another example had students highlight portions of code that exhibited certain behavior or had a certain feature. Both of these interfaces go beyond simple text or scalar feedback, and can be used to provide richer information in reviews.

Reviews can also be conducted face-to-face, whether solely between students or moderated by TAs or instructors. Moderation can make arguments more constructive, guide discussion towards relevant points, and make a face-to-face meeting less intimidating. Moderated review moves the process more towards a studio-like setting, and may be appropriate especially for teaching students what is involved in a constructive code review process. Hundhausen, Agrawal, and Agarwal discuss this kind of in-person review, dubbed *pedagogical code review*, in early courses [31]. In pedagogical code review, a small group led by a moderator use a set of predefined coding practice guidelines to guide a group review of student programs.

## 5.4 Review Feedback (Meta-Reviewing)

Any instructor using peer-review must choose whether to include grading or feedback on the contents of reviews themselves. We use the term *meta-review* to refer to any feedback on a review (because feedback can be considered a review of a review). Feedback can take many forms: the author who received a review could report on whether the review was constructive or led to changes, course staff could formally grade reviews and return comments to the reviewer, or third parties could comment on the relative merits across a set of reviews. Which model makes sense depends on factors including the learning objectives for IFPR, features of peer-

review software, and course logistics (such as staff size relative to student population). Many of the issues here apply to peer-review in general, rather than only to IFPR.

According to Ramachandran and Gehringer [49], reviews consist of (1) summative, (2) problem detection, and (3) advisory content. Meta-reviews can report on each of these three types of contents, each of which is valuable in its own way. While summative contents can reflect a reviewer's understanding, problem-detection content directly helps a student identify opportunities for improvement, and advisory content points out ways in which students might improve. Meta-reviews can include information on which parts of a review were constructive, and which led to actual changes. Metareviews written by authors of submissions can also include rebuttals to aspects of a review; in IFPR, such rebuttals can arise when students are debating the requirements of an exercise through the review process (a healthy outcome relative to the goals of IFPR). With enough iteration of this form, IFPR more closely resembles traditional collaboration rather than peer-reviewing of each others' work.

#### 5.4.1 Types of Meta-Reviewing

Around Ramachandran and Gehringer's framework, there are several ways to structure the information in meta-reviews, and provide useful feedback to reviewers.

- Direct feedback from course staff: Feedback from instructors or TAs can repair incorrect advice and reinforce good behavior (case studies 1, 11, and 13 call out the importance of correcting faulty reviews explicitly). Class size is clearly a factor here. If someone wants to use peer-review to help scale human feedback in large courses, then giving expert feedback on reviewing might not be feasible.
- Feedback based on assessment of submissions: In situations where an expert evaluation of the *assignment* is available (whether through auto-grading or by human TA) and reviews are quantitative, it should be possible to automate a meta-review that tells a reviewer something about the quality of work they reviewed. For instance, if a student indicates in a Likert scale that they "strongly agree" that a solution is correct, but the grade for the assignment they reviewed is low, an automated meta-review can indicate that this review likely mis-evaluated the work under review.
- · Reporting correspondence among reviews: Reviewers could be told about the correspondence between their evaluation of a submission and those of other students. For example, the SWoRD tool for peer review of writing tells student reviewers, on each criterion they reviewed, how they did relative to the average of other students' scores [8]. An example of feedback that they show says "Your ratings were too nice for this set of papers. Your average rating was 6.50 and the group average was 5.23." This hints to the reviewer that he may have missed something in his review. This is related to Hamer et al.'s work on identifying "rogue" students in peer assessment [26], which is focused on identifying outliers' impact on grades. This does run into issues of calibration and opinion; just because a student disagrees with the average, it doesn't mean they are wrong! The outlying reviewer may have understood something the other reviewers didn't, in which case comparing his review to an expert's, or to a trusted automated process, may be more useful feedback.
- Having students review submissions of known quality: In CaptainTeach programming assignments, half the time students are asked to review a known-good or known-bad solution (implemented by the course staff) [47]. Students use a Likert scale

in each review to indicate whether they think the submission under review is correct. If a reviewer gives a strong score to a known-bad solution, or a weak score to a known-good solution, she gets immediate feedback telling her of the discrepancy.

Existing research has explored ways to provide or assess meta reviews. Nelson and Schunn describe a rubric for evaluating peer feedback in writing assignments which includes criteria like the concreteness and actionability of the review, and whether it was generally positive or negative [43]. Swan, Shen, and Hiltz study assessment strategies for comments in online discussion forums used to discuss class content [59]. Though the discussions are not necessarily critiques of student work—they are simply prompts for questions and comments—they do have similar requirements to reviews in relevance, accuracy, and focus.

The Expertiza peer review process contains an explicit reviewof-review phase for collaborative work [48], and a related Expertiza tool attempts to give some more qualitative feedback automatically by a natural-language analysis of student work [49].

In Aropä, each review is "an assignment in itself... Reviews can thus be reviewed using all the facilities for normal assignments" [25]. Based on this observation, Aropä's main workflow suggests two explicit kinds of meta-review. First reviewees can *rate* their reviews, which gives feedback on the perceived review quality. Second, there is an explicit *dispute* phase, in which the reviewee can disagree with the content of a review and request that the reviewer reconsider, after which the reviewer can submit revised feedback. This review-dispute-revise loop is, in effect, an in-flow review of the review, since the student writing the review gets feedback on multiple versions of the review.

#### 5.4.2 Using Meta-Reviews

While one generally may prefer to eliminate low-quality contents in reviews, in a pedagogical context receiving some low-quality review contents can be beneficial. While in traditional educational settings authors may trust all the feedback they receive from the instructor, in IFPR authors have to learn to assess the value of the reviews they receive. They will have to learn to separate review comments into those they will act upon and those they will ignore, then triage those they wish to act upon. Moreover, having a diversity of reviews, maybe even contradictory ones, can be a starting point for valuable discussions in class. Having to wade through reviews can implicitly train reviewers that they, in turn, should not submit "brain dumps" of everything they think of, but instead provide valuable and concise reviews. The important metric is actionability, not volume.

Instructors may seek to use meta-reviews to monitor the IFPR process. Given the quicker turn-around times inherent to IFPR, such monitoring benefits from tool support and structural elements of meta-reviews. For example, asking authors to rate the reviews they receive on a simple Likert scale makes it easy for an instructor to focus on potentially problematic reviews without imposing undue burden on the students. In some IFPR configurations, software tools that include automatic grading could report partial information on whether student performance improves following the review phase. Such information would be most useful for identifying cases in which poor work did not improve, prompting the instructor to check on whether the author had received useful and actionable advice through reviews.

Meta-reviewing incurs a cost. Whether meta-reviews are worth that cost depends on the learning goals. If teaching how to review is important, meta-reviews are essential; however if the learning goals focus on artifact production or performance, and if the reviewers are experienced, meta-reviews may be less essential. An alternative to providing meta-reviews for each review is to provide a few example reviews and their meta-reviews. To not tempt students to

simply reuse the best example review comments, these exemplar reviews can come from an assignment that is different from the current assignment.

A live demonstration of how to do a code review is a form of scaffolding on the process-level, but does not drive content as specifically as rubrics. Regardless of how reviewing is introduced and scaffolded, it is important to allot time to deal with misconceptions on how to create a review as part of the course design.

## 6. Parameters and Issues

Several issues and parameters cross-cut the stages of the IFPR process discussed in section 5. Questions about preventing plagiarism, integrating IFPR with course-level grading, deciding where to use anonymity, involving experts, making IFPR relevant for non-majors, engaging students in the process, and identifying software needs all guide one's particular configuration of IFPR. We discuss each of these issues in turn.

#### 6.1 IFPR and Plagiarism

IFPR, like many course and assignment structures, requires careful mechanism design to ensure that students aren't incentivized towards detrimental behavior that lets them get a good grade at the cost of their (or others') education.

One of the most immediate problems with IFPR is that, by definition, students are shown one another's work while in the middle of an assignment. Since the final submission happens after students have been exposed to other students' work, the IFPR educator must determine how to account for this exposure when assigning a grade to the final submission.

At the extreme, a student could submit an empty initial submission, copy what he sees during the reviewing phase, and submit the copied solution as his own final solution. In less extreme cases, a student may copy all or part of another solution into her own after submitting an initial first try that she becomes convinced is incorrect. There are a number of course- and grading-design decisions that can affect the degree to which copying is a problem:

• Variation in Assignments: One major factor in determining whether copying is even a problem is how similar students' submissions are expected to be. In many programming courses, students implement to the exact same algorithmic specification; other than coding-style issues, one implementation is just as good as another. This is in contrast to other domains where peer review is often used, like creative or critical writing, in which students often write on different topics or choose different positions to represent on the same topic.

One approach is to provide variants of a programming problem to different students. Zeller [74] gives each student a variation on a theme to avoid students reviewing another who is working on exactly the same problem. Indeed, it is often possible to generate large numbers of different problems automatically from a specification, as Gulwani et al. have done for algebra problems and more [1, 56].

A drawback of variation in assignments is that it weakens one of the benefits of IFPR – having students review the same problem they are already thinking about! Especially for beginning students, where program comprehension skills are still being learned, one goal of IFPR is to lessen the cognitive load of the comprehension task by having the student review code for a problem they already understand. If they have to internalize an additional problem description along with new code, this puts significantly more overhead into the reviewing process.

Depending on the learning goals, it may be good for the reviewer to learn to incorporate ideas from different solutions into

her own, since it requires a more abstract understanding of the techniques. For novices, it may be enough of a challenge to recognize a good solution and apply it to her own.

• Weighted Submission Grading: There is often value in having reviewers copy parts of other submissions that they see in order to improve their own work. It happens all the time in professional software development, and the act of recognizing a good solution demonstrates understanding that is far beyond blind plagiarism. Reviewers should take things from the examples they see and demonstrate that they learned from them; however, a student has no guarantee that what he is seeing is correct, so blindly copying can hurt!

However, wholesale copying (where a student submits an empty file then copies the best of what they see) should be discouraged (to say the least!). In order to mitigate this, Politz et al. [47] grade IFPR assignments by assigning heavier weights to initial submissions than to post-review submissions: an initial program submission counts for 75% of the grade. Students can still improve the 25%-weighted part of their score based on review feedback and copying others' solutions, but they can also hurt their score if they make incorrect changes. Different weightings put different emphases on the importance of review. Having the post-review score count for more might be acceptable in some classroom settings, and ultimately comes down to a choice about student maturity, class culture, and other course-specific factors.

• Alternative or Supplemental Grading: Another solution to the grading problem is to supplement assignment grading with other techniques that cannot be copied. For example, in an inperson code review of a student's solution, an instructor can quickly ascertain whether the student has simply copied something or actually understands the code they have submitted. This can be done by, for example, asking the student to change his program to match a new specification, or asking her to understand a proposed change to her submitted code.

## 6.2 Interaction with Course-Level Grading

Instructors must determine the extent to which IFPR activities impact course grades and the mechanisms through which they do so. Section 1 noted Liu and Carless' distinction between *peer feedback* and *peer assessment* [38], where the latter's goal is grading. Most of the working group discussion focused on peer-feedback (which fit the course contexts of the participants), though we gave some attention to peer-assessment (more often proposed to address grading at scale in large courses).

## 6.2.1 Peer Assessment and IFPR

Kulkarni et al. have shown that, with careful rubric and mechanism design, peer assessment can produce similar results to TA grading in MOOCs [36]. Reily et al. report on the accuracy of a combination of peer reviews at assessing programming assignments [50]. Hamer et al. describe a technique for deriving grades from weighted averages of peer assessments to identify "rogue" reviewers and generate grades that weight more apparently accurate students' assessments more heavily [26]. A common theme across these studies is that a combination of peers' assessments provides an accurate enough assessment, even if some particular students or reviews are inaccurate.

Peer assessment changes the motivation structure of IFPR. For example, a student who is afraid of affecting his peers' grades with negative feedback may be more hesitant to give that feedback. In contexts where students are still learning to review and give feedback, inaccurate reviews are expected and an important part of the learning process; in this case, reviews probably should not be

used for grading purposes. Using peer review for grading should be adopted with care, and practitioners should carefully consider its effects on the other design decisions discussed in this report.

## 6.2.2 Should Reviews Be Graded?

Although section 5.4 discussed various forms of feedback on reviews, it did not discuss whether reviews should be assigned scores that affect students' course grades. Grading schemes can range from checkbox-style points for submitting reviews (without grading content), to more detailed assessments. Of our case studies, six (3, 5, 6, 8, 9, and 11) explicitly tie reviews to the course or assignment grade in some form. No case study discussed grading metareviews.

In general, the group members were reluctant to ascribe grades to reviews (as opposed to giving meta-reviews, which the group strongly endorsed). The group shared concerns that having reviews influence course grades (beyond required participation) would misdirect student motivation for reviewing [33].

Nonetheless, the group did discuss options for having reviews figure into grades. We discussed basing assignment grades on reviews rather than on the work submitted (on the grounds that reviews reflect students' understanding of the assignment): this would allocate staff grading time to meta-reviews rather than to code, which could be more valuable (as some, though not all, aspects of code can be assessed automatically). We also discussed basing students' grades on the improvements that their reviews inspired in the work of others, but felt the nuances (submissions with little room for improvement, students who chose not to act on reviews) made this infeasible.

The group found relatively little in the literature on grading reviews. Sims [55] proposes grading the reviews according to compliance to the review writing guidelines. The Review Quality Instrument (RQI) of van Rooyen et al. [64] is a simple, reliable, and valid scale for studying scientific peer review processes. The authors claim high internal consistency for RQI. Trautman et al. propose using this framework for educational peer reviews [62], though note a significant limitation: the instrument gauges how well the reviewer has considered the key aspects of the work and less whether the review is accurate or correct. They also question whether grading reviews might diminish the less tangible benefits of peer reviews such as increased motivation, ownership and increased interest in learning. This question is put forth as topic for research instead of a claim.

Wessa et al. [70] identify statistical measures of peer review process and demonstrate that [71] these measurements can be used to built statistical models, and therefore automatic evaluations of reviews. They prescribe that objective measures of review quality, such as word count and number of received and given reviews, can be used as basis for assessing overall student performance. Similar automated review process is proposed by Ramachandran et al. [49], who suggest metrics such as content, tone and quantity of feedback to suitably represent a review. The group did not discuss automated assessment of reviews, though some members have looked informally at word counts in student reviews and found low correlation with the value of reviews.

#### 6.2.3 Interaction with Relative and Curve Grading

The group noted that IFPR (like other collaborative course structures) interacts poorly with grading strategies that evaluate students relative to one another. Such strategies conflict with students' motivations to help one another improve their work. The working group identified three distinct ways in which the conflict could manifest itself:

- **Demotivative**: There is disincentive to do good review, because it can push others past oneself in achievement, adversely affecting one's own grade.
- **Destructive**: Instead of just being apathetic about review, students could even sabotage one another with bad feedback, hoping to reduce others' scores to actively improve their own.
- Unmotivative: Since the curve puts a limit on how much one
  can achieve, there is a disincentive to respond to feedback or
  reflect (e.g. especially for high-achieving students, there's little
  reason to take feedback seriously).

No one in the group used relative grading in their own courses, so we lacked first-hand experience in mitigating these problems in that context. Boud, Cohen and Sampson [6] discuss various tensions between standard assessment practices and learning from peers. Rick and Guzdial [51] discuss the impact of curve-based grading on collaboration and peer learning.

#### 6.3 Anonymity

Several issues relating to anonymity and privacy come to the fore with peer review. Developing a culture of positive and constructive critique where students can both give and take feedback appropriately, can take time and require a degree of practice. While the broader aim may be to develop a positive, supportive and professional approach to peer review, there may be a need to provide some initial shielding from scrutiny for students who are new to the institution or the practice.

## **6.3.1** Types of Anonymity

Each of the IFPR roles—authors and reviewers—can be anonymous to other students or faculty. Reasonable arguments can be made for each configuration, and different configurations have been used in practice.

In PeerWise [16], a system to which students submit proposed study questions on course material, student submissions are ranked for quality and correctness and the content of the contributed questions is public. In that case, not revealing contributor identities to peers is important: this creates a degree of safety for the novice student contributor and helps identify the public ranking with the work and not with the student. However the contributor identities are visible to the instructors as student contributions may be summatively graded. A variant on this (similar to some conference reviewing models) could also see reviews made visible to all reviewers, which may help in establishing and reinforcing norms and standards and mitigate the risks of unduly harsh or abusive reviews.

In case study 10, anonymity is staged: at first students review one another anonymously, then groups are formed and groups review one another, revealing identities and adding a social element. This lets students do a first round of reviewing to get comfortable with the process of feedback, and after has the benefits of encouraging professional collaboration.

In total in our case studies, only two (4 and 11) explicitly stated that reviews were anonymous. Several (studies 5, 6, 8, 12, and 15) had reviews that were in-person, and therefore cannot be anonymous. The others left it unstated, and in different course contexts the assignments could be administered either way.

#### 6.3.2 Upsides of Anonymity

At the introductory level, students suffer from both confidence and maturity problems, making anonymous review an attractive option (at least initially). Authors will not face personal embarrassment if others see work that they are not comfortable showing publicly, and anonymity gives reviewers the freedom to be more candid.

In addition, reviewers who know the author under review might make assumptions based on the author rather than the submission. This could even cause authors who consider themselves weaker to not question incorrect work that comes from a supposedly smart reviewer. Anonymity helps level the playing field in the face of such preconceptions.

In general the group considers anonymity to be a less jarring initial option for IFPR that is more likely to protect students who are new to the peer-review process. However, sharing identity during review has significant benefits in the right contexts, and anonymity isn't without its own problems.

## 6.3.3 Downsides of Anonymity

Anonymity can unwittingly enable a culture of excessive criticality, or even of "flaming" and online harassment, to develop among some students. The former can probably only be monitored by course staff. The latter would require policies and techniques for reporting abuse and inappropriate behavior: for example, a "Flag Abuse" button in a Web interface that allows students to bring offensive or inappropriate content to the attention of the staff. In general, these issues are similar to the unsatisfactory aspects of peer assessment schemes in group work that ask group members to evaluate relative contributions [11].

The benefits of non-anonymity center mainly around creating collaborative cultures and helping students learn professional behavior. Non-anonymity creates opportunities for students to acknowledge each others' contributions. As a broader educational goal, ethics and professionalism are meant to be covered as part of our curricula [21]; an open model for peer review gives a clear opportunity for enforcing appropriate behaviour. Industrial code reviews are not done anonymously [52], so students gain relevant skills from learning to give and receive non-anonymous feedback.

## 6.3.4 Anonymity and Cultural Considerations

Anonymity may have cultural connections at several levels. For students who come from a consensus-based national culture [29] where preserving harmony is a strong value and overt criticism can be considered offensive or cause a loss of face, a greater level of anonymity may be required at first for students to feel more at ease in speaking their minds. A converse cultural aspect may be in operation with non-anonymous contributions; if students know the other person, they may be less or more critical a priori, based upon judgments of the peer's relative status or perceived expertise, rather than reviewing the material in its own right.

At some institutions such as Brown University (where three of the group members have studied IFPR), total anonymity is already hard to establish, because a robust undergraduate TA program means that students often act as TAs for one another independent of peer review, and know about one another's performance. At Brown, the institutional and student culture makes it commonplace to know who is reading your code.

The group debated whether giving students the choice to anonymize was a good idea, and concluded that it was not desirable because it encouraged hiding attitudes, and may make people justify being more objectionably critical because they were allowed to be anonymous. A further negative was that it would not bring the shy students out (which is sometimes the educator's goal). It was concluded that the preferable approach was to make anonymity or identifiability a matter of policy, and educate students about the importance of professionalism in either case.

Overall the group doesn't recommend that IFPR practitioners adopt anonymity by default, but rather that they take course and culture into account. Overly stressing anonymity could unwittingly give the impression in students that peer review is dangerous. We do recommend that in cases where work is identifiable, the underlying goals and expectations related to professionalism and community-building be consciously introduced to students from the outset. This

could be talked about as explicitly and strongly as the discussions about plagiarism, with potential penalties for abuse of the system through lack of respect for one another.

## 6.4 The Role of Experts

Peer review can sometimes usefully be complemented with expert review. There are several arguments for and against combining the two. Potential benefits of expert review include: moderating conflicts between reviewers and reviewees; facilitating contribution and collaboration; overcoming cultural adjustment problems; and providing exemplars of good work and good reviewing.

• Experts as Moderators and Facilitators: Experts can act as moderators to make sure that issues and conflicts that arise, whether in a live situation or asynchronously, can be dealt with by an authority figure. As moderators, experts do not take on the role of reviewers, which keeps students in charge of the feedback itself. Moderation is thus a form of process- rather than content-expertise (akin to pedagogical code reviews [31], which are led by an expert moderator). Experts may also act as facilitators of group discussion of work or reviews. Reviewers who are not sure of themselves may not contribute much; experts can assist in getting students to contribute, and push the idea of review as a learning process, in addition to acting as figures of objective authority on grading. Facilitation is likely harder to integrate into IFPR configurations in which students submit work and reviews asychnronously relative to one another.

Whether acting as moderators or facilitators, expert reviewers can help address cultural issues where students devalue the opinions of their peers and overvalue the opinions of instructors. They can also offset deficiencies in the knowledge and insight of reviewers who are themselves adjusting to the process of giving informed critiques, and ensure that some knowledgeable, high quality feedback is received.

• Perception and Quality of Expert Review: Expert review (especially instructor or TA provided) may have unwarranted special status in students' minds: feedback from experts may be interpreted as "more relevant for my grade," and hence more likely to be acted upon. This last concern is called out explicitly by case study 9, in which students produce a peer-reviewed research article, and in the past often discounted peer feedback in favor of instructor-provided comments.

Cho and MacArthur hide the provenance of expert- vs. peerprovided reviews in order to study whether expert review is in fact more effective at improving students' grades [7]. They compare three approaches to giving feedback on written assignments in a psychology course: feedback from a single peer, feedback from a single topic expert, and feedback from multiple peers [7]. The results indicate that feedback from multiple peers results in better quality revisions than feedback from an expert, with feedback from a single peer being the worst. The hypothesised reason for this was that peers gave feedback that was phrased in terms that students could more easily comprehend.

Hamer et al. investigate differences between reviews provided by peers and expert tutors on Matlab programming assignments in a peer assessment context [27]. The main significant difference they find is that tutors not only write longer feedback, but provide more *specific negative* comments about inaccuracies in the program under review. In addition, there were two kinds of feedback in the review rubric—*correctness* and *style*—the longer, more specific negative feedback tended to be in the correctness comments rather than style comments. There was no significant difference in other features of review comments, like

whether there was concrete, actionable advice for improving the program, if the review contained positive encouragement or reinforcement, or the final mark given by the tutor.

• Expert-provided Exemplars and Models: In IFPR the goal is to avoid attaching the idea of constructive review to an expert being present. However, an expert can provide models for review that students can follow. So one strategy would have experts give examples, or be present for some (early) sessions and not others. These exemplars and supports represent a form of scaffolding reviewing.

The working group members discussed an "Editor's picks" option, where instructors and TAs highlight good examples of reviews for others to learn from. TAs might monitor reviews and post a handful of high-quality ones for the entire class. As a tweak to avoid simply externally rewarding good examples (which comes with the downsides of extrinsic motivation [33]), the reviews picked don't necessarily all need to be objectively good. Instead, the editor's picks could highlight interesting reviews, and what is interesting (or could be better!) about them. Students then learn specifically why something is considered good or bad in a review, and can apply that knowledge to their own reviewing. The process of breaking down the review shifts the emphasis from quality of reviewer to important aspects of the review itself.

#### 6.5 Does IFPR Make Sense for Non-Majors?

When instructing computer science majors, reviews of code can be motivated by their resemblance to code reviews and similar activities in the software profession, even though this may not be the actual driver behind their inclusion in a course. When using inflow peer review with non-majors, such motivations might work less well due to a disconnect between the students and the IT profession. Regardless of the future profession, we view the reviewing skill as an important one and note that reviewing comes in naturally in a number of fields, from academic writing to zoology, but also in the students' everyday lives in reviews of movies and restaurants on websites, etc.

We frame in-flow peer reviewing as a technique for improving learning based on timely feedback throughout a process. This can be different from a means of quality assurance of a final artifact, or a technique for learning to produce and consume reviews, although these concepts arise naturally in a setting where in-flow peer review is used. Knowing how to produce and consume reviews is a useful skill in its own right, regardless of the profession or the context in which it is used. On this note, using reviews in classes with a diverse student population will have the additional benefit of producing more diverse reviews. In many cases, such as for GUI mockups, this is highly desirable. Therefore, whether or not a student will face a review professionally is not necessary for IFPR to be useful.

Test cases are useful for non-majors to review as they are often less tangled with complicated "non-functional aspects" of the problem such as memory management, performance, or code elegance. Case study 1 features an assignment from a course designed explicitly for non-majors, with review of tests. We also envision that non-majors reviewing code might reduce "code fear" by making code less magical and establishing that code can be wrong. This is also related to the "community of practice" of programming in that students are seeking to join the broader community of programmers by taking a CS course in the first place. Giving non-majors a sense of what code—and programmers—are doing might teach respect for code and programming, while not necessarily trying to make students into professional programmers. Many programmers are managed by non-programmer experts in other fields who could

benefit from an understanding of practices surrounding code. There are also cases where non-programmers must be able to read or relate to code to make sure that code follows complicated legal practices or financial algorithms.

Finally, as an important aside, reviewing might help with retention especially by removing the stigma of computer science being an anti-social, "inhuman" subject. Humanities students generally read more than they write – as reviewing is more like reading, it might feel more comfortable for them than programming. This could be a component of a CS course in the style of fine arts, which Barker et al. found to create a better community culture and encouraged retention of female students [4].

#### 6.6 Bringing Students Along

One challenge of peer review in general lies in getting students to value the contributions of their peers. Peer evaluations potentially contradict a model students have of the instructor as the sole authority. In addition, many students aren't initially comfortable acting as reviewers (and giving criticism to one another!), so a desire for harmony may hinder their ability to produce effective reviews at first. These two dimensions—taking the reviews seriously and taking the act of reviewing seriously—differ in force and mitigations.

When students are in the role of reviewers, instructors may need to consider how to balance anonymity in reviewing with the need to engage students in the professional practice of reviewing (section 6.3). When addressing students coming from cultures that don't emphasize review and criticism, explicit reminders that review is part of the working programmer's life in western culture can be helpful. Chung and Chow [9] discuss such culture factors when discussing peer- and problem-based learning in a class in Hong Kong.

When students are in the role of receiving reviews, instructors should require students to demonstrate engagement with the reviews. Asking students to indicate how they used reviews in revising their work is one option. Authors may need explicit instructions on how to read reviews; this might also help them cope with criticisms. For some students, peer review may be their first time getting negative feedback; instructors should make response to review a positive activity. For authors whose work was strong and did not yield actionable reviews, an instructor could ask "what have you learned from looking at others' solutions?" Hopefully, such activities on a few early assignments would help students develop self-reflection skills.

Specific advice one could give students includes:

- When you receive criticism, go back to other reviews are you making the same mistake over and over?
- Have a cooling-off period if you get negative reviews you may view them more positively with a little passage of time.
- Do you see a contradiction in your reviews? Maybe one of the reviewers is wrong, or maybe your work isn't clear enough.
- How will you avoid making these mistakes again in the future?
   Look at the review, break out the actionable items, prioritize.
- Are you taking the review personally? Remember that the review is about the work (and improving it!), not about you the ability to recognize your weaknesses and improve them is incredibly valuable. This is an attitude encouraged by *egoless programming* [69].

If students remain skeptical of the value of peer-review, a midcourse survey on the value of reviewing might help convey the experiences of others. This also fosters a culture of peer-review as a collective effort rather than one of criticism, judgment, or assessment

#### 6.7 Software and Analytics for IFPR

Good software tools are key to making IFPR manageable and informative for both students and course staff. For students, tools that integrate reviewing with the IDE used for programming assignments, for example, mitigates some of the context switching that the process otherwise requires. For staff, good tools not only manage the logistics of the process, but can also be instrumental in producing meta-reviews and in monitoring the effectiveness of IFPR.

Software systems could (1) give basic feedback on review quality by comparing reviews between students [8] or through use of machine learning methods [49], (2) aggregate student responses for discussion [28], (3) flag students who consistently write certain kinds of reviews (weak, strong, superficial, etc.), or (4) maximize variety of review tasks by using static analysis tools to appraise source code similarity. Some of these features are more sophisticated than those included in current peer-review platforms.

Features for organizing submissions or reviews in various ways could help instructors run effective IFPR processes. Software could allow instructors to associate arbitrary tags with each submission or review, such as "sloppy" or "discuss in class," "ignore," "insufficient test coverage," or "misunderstood requirement B." Instructors could then search or group submissions and reviews by tags, or use tags as categories when producing analytical visualizations. Tags could be visible only to instructors, or instructors could make tags visible to authors or reviewers. Another way to organize submissions or reviews would be to cluster them automatically by various criteria, similar to Expertiza's automatic meta-reviewing categories [49]. Example criteria include the length of the reviews or their tone ("positive," "neutral," or "negative"), which may be detectable automatically using natural language processing techniques. Based on this automatic and manual organization, instructors could identify recurring misconceptions, or they could select representative exemplars of reviews to be presented in class. Finally, the information in reviews (e.g., feedback in each rubric, or scores on Likert scales) could be used to organize submissions, and information in meta-reviews could be used to organize reviews.

Instructors as well as students could benefit from analytic visualizations or dashboards. Visualizations presenting activity over time could help instructors understand the timeliness of reviews or a student's progress over time in terms of the quality of their submissions or their reviews. Visualizations of reviewer-author relationships, which could include various historic aspects, such as the quality of their past submissions or the usefulness of their past reviews, could help with reviewer assignment. Visualizations superimposing reviews on top of submissions or meta-reviews on top of reviews could provide a compact picture of a certain artifact to authors, reviewers, or instructors.

Luxton-Reilly [39] provides a systematic survey of tools for peer-assessment. The survey uses several key parameters to compare tools, such as the flexibilities of work flows, rubric designs, and the nature of evaluation supported. The summary data in the survey is not sufficient to determine whether a given tool supports IFPR, but the descriptions of individual tools contain details that may be useful in checking for IFPR support.

## 6.7.1 Are Conference Managers Suitable?

Software tools for managing conference-paper submissions are designed to support peer review. It thus makes sense to ask whether these tools are suitable for peer review (IFPR or otherwise) in pedagogic contexts. Popular conference managers in computer science include EasyChair [65], HotCRP [32], CyberChair [63], START [23], Linklings [42], and CONTINUE [34]. Some working group members had tried using conference managers for peer review, but found the fit to be problematic unless the process matched that of conference reviews.

Peer-review systems must support two distinct tasks: handling submissions and handling reviews. For IFPR, students must submit revised work; in some peer-review configurations, students may submit multiple artifacts at different stages. While some conference managers support revisions, they typically don't track multiple submissions from the same student over time. Conference managers assume that submission deadlines are fixed and synchronous; this assumption is reflected in the synchronous assignment of papers to reviewers. Reviewers are typically given read access to all of the submitted papers, not just the subset that they are due to review (this sets aside conflict-of-interest, which is less of an issue in the pedagogic context). These are two examples of configurations that instructors may wish to make in a pedagogic context that are inconsistent with the design of conference managers. HotCRP [32] does support more flexibility anonymity handling than most other tools: "selective reviewer anonymity" allows reviewers to explicitly decide on whether or not to keep their reviews anonymous.

#### 7. Industrial Code-Review Practice

Much of our discussion of IFPR has revolved around code review, as programs are one of the most common artifacts in computer science courses. Given the rich history of code-review in industry, it is worth asking what IFPR can learn from this practice. Code review is an essential component of industrial software-development. Industrial peer review is inherently in-flow, as reviews are conducted on a regular basis during development – it would be odd indeed to only perform code review after a product had shipped to customers! The industrial product-development life cycle is longer than that in many courses, but best practices in industrial peer review are still useful context for this report.

#### 7.1 Motivations for Industrial Code Review

Industrial code reviews differ in motivation from pedagogic code reviews. The goal is often to reduce the defect rate before releasing software or committing to a design. This is the primary measure of an effective review process in Fagan's seminal work [18], in followup work to it [66], and also in some modern surveys tied to large case studies [12]. The goals of in-flow peer review in pedagogic settings are much broader than just finding bugs in peers' code, though finding problems is certainly one worthy cause for review

However, in one modern study that studies attitudes about code review in both developers and managers at Microsoft, researchers found that defect finding, while important, was only one of several top motivations developers saw for review [3]. Also scoring high in developer surveys as motivations for code review are general code improvement, suggesting and finding alternative solutions, knowledge transfer, and team awareness. While performing reviews, developers indicated specific cases where they learned about a new API that they could use in their own code, or providing links to the code author with documentation of other alternatives, suggesting that these activities do indeed occur.

Activities like knowledge transfer and dissemination of ideas are certainly goals of in-flow peer review, and in-flow strategies should consider ways to foster them. Bachelli and Bird [3] note that these metrics are harder to measure than defect rate, but observe them coming up spontaneously in interviews and in observations of reviewers.

## 7.2 Staged Code Inspections

Fagan's seminal work on code inspections in an industrial setting [18] finds that putting inspections at carefully-delineated points throughout a product's life cycle can save time by fixing faults earlier, before other work builds on the buggy code. In Fagan's experiments, there are three inspections: one after an initial

design phase, one after initial coding, and one after unit testing and before system-wide testing. Experiments show that maximal productivity is reached by including only the first two inspection steps due to the high cost in developer time relative to the time saved by early detection, but that using the first two steps increases programmer productivity by 23%, according to their metrics.

This result mirrors our intuitions about the value of staging assignments for review at points that are useful for catching and fixing misconceptions about the assignment. Our primary goal is not simply to improve the programming output of students, however. We care about what they learn from seeing other examples, teaching them to effectively review others' code, and more. Still, it is useful to consider that the in-flow experience is similar to effective industry practices, and note that professional developers benefit from the staging process.

#### 7.3 Meetings vs. Asynchronous Code Review

Fagan's original results are for *formal code inspections* [18], which consist of a meeting of several developers (including the original author), conducted with prior preparation and with a separate *reader*, separate from the author, who presents the work. Defects' cause and detection are documented in detail, which acts as a sort of "rubric" for the code review.

While formal code inspections demonstrably find valuable defects, it is not clear that the organization of a meeting is required in order to have a comparable effect. Votta studied the necessity of meetings for code inspection, and found that the majority of defects—over 90%—were found in the *preparation* for the meeting, rather than in the meeting itself [66]. Votta concludes that much of the benefit of code review can be had without the overhead of scheduling in-person meetings.

There is further research on this debate, but the only clear conclusion is that significant benefits of review remain even without inperson review. In a pedagogic setting, in-person reviews may serve other goals, like training students to review, encouraging productive feedback, and supporting the social aspects of review. However, industry research suggests that the overhead of scheduling and holding meetings isn't a prerequisite of effective reviews in professional settings.

#### 7.4 What and How to Review

A large industrial case study on code review [12], which also documents a survey of code review in industry (including a longer discussion of formal vs. lightweight review), identifies guidelines for effective reviewing practices. By measuring defect rates found against the number of lines of code under review and the length of the review session, their study recommends "the single best piece of advice we can give is to review between 100 and 300 lines of code at a time and spend 30-60 minutes to review it." While this advice may be appropriate for peer-review in upper-level or graduate courses, this much code would likely overwhelm lower-level students.

#### 8. Additional Related Work

Pedagogic uses of peer review have a long history that predates in-flow reviewing or peer review in computing courses. A survey of this history is beyond the scope of this document, though some relevant citations appear in other surveys [61]. Here, we focus on more recent work in computing education and cognitive aspects of education that bear on in-flow peer review.

#### 8.1 Pair Programming

Pair programming (henceforth PP) is a software-development technique in which two programmers work together on one computer.

PP involves significant (and continuous) in-flow peer feedback, though coding together is a rather different activity than writing and responding to reviews. One surface-level difference lies in the number of reviewers: IFPR students may receive multiple reviews, whereas comments in PP come from a dedicated programming partner. Other interesting differences arise along three dimensions: responsibility, skills developed, and dynamics. To better align the practices, we contrast PP with an IFPR model called *mutual review* (discussed in section 5.3) in which a pair of students are tasked with reviewing one another's submissions.

• **Responsibility**: The key difference between mutual review and PP is that in PP both students are responsible for the quality of a *single* artifact, whereas in IFPR each student is responsible for her own artifact. This can naturally lead to a significant difference in motivation and responsibility.

The differences in shared responsibility can also be seen when pairing students of different strengths. In PP, if one student is stronger, that student may drive the production of the artifact, or even take over the work, without the weaker student having the opportunity to participate fully; thus, the weaker student may not benefit from the experience. In IFPR, a strong student's assignment (and hence grade) is not affected by the weak student. The strong student is in a better position to help the weaker student; although the strong student may receive no beneficial comments, no harm will be done. With IFPR, any problems due to mismatched pairings can be alleviated by assigning multiple reviewers. In practice (in several authors' experience), in PP students often work alone even when expected to work collaboratively.

- Skills: The skills required and developed differ between IFPR and in PP. PP deals primarily with the collaborative creation of an artifact. IFPR focuses also on the creation of an artifact, though less collaboratively as students do not work face-to-face. IFPR also focuses on the high-level skill of performing reviews, which requires both program comprehension and judgment. These activities require students to take a higher-level perspective on the task at hand, including the difference between understanding one's own code and understanding the code of others. IFPR also forces students to trade off work, namely, time spent on their own assignment versus time spent on reviewing. This trade-off is not present in PP. To continue the analogy with academic paper writing and reviewing: PP is more like working as coauthors, whereas IFPR is more like the relationship between a journal paper reviewer and an author (especially if there are multiple IFPR stages).
- **Dynamics**: The difference between solo and shared responsibility for the artifact yields different dynamics between PP and IFPR. Due to the continuous communication required in PP, students are often immediately alerted to problems in program comprehension. In contrast, IFPR is separated by both space and time. As a result, programs can acquire significant accidental complexity. Students realize this, and thus learn about the difficulties of producing and reading code, both by trying to make sense out of the submissions of others and by seeing the feedback their own submissions receive.

Finally, with IFPR, students are compelled to review other students' code: they cannot ignore problems and let their partner do the work. On the other hand, in IFPR it may be easier to ignore the advice given than in PP, because of the difference in ownership; it is harder to argue that someone not make suggested changes to a shared program than to ignore the feedback on one's own program.

#### 8.2 Intrinsic vs. Extrinsic Motivation

Ways in which particular assignment designs tradeoff between intrinsic and extrinsic motivation have been a theme in this report. The group frequently considered arguments from Kohn's "Punished by Rewards" [33], which argues that praise and rewards (such as grades) are ineffective and even demotivating. Kohn claims that rewards effectively punish students who aren't rewarded, discourages risk-taking in students, devalues the reasoning behind a low or high grade since the grade is emphasized, and can create a bad relationship between teacher and student. IFPR is designed to increase intrinsic motivation for reviewing relative to traditional post-submission peer review. The extent to which this shift occurs likely interacts with how the grading system treats reviewing.

## 8.3 Metacognitive Reflection

One goal of in-flow review is to encourage reflection while in the middle of an assignment. Meta-cognitive reflection has been studied as an important part of the learning process. Indeed, it has been indicated that one difference between experts in program comprehension and novices is the focus on metacognition [17]. Others in different contexts have found effective reviewing and prompting strategies for encouraging reflection that IFPR can learn from.

Palinscar and Brown study *reciprocal teaching*, in which a teacher alternates with a student in a dialog that prompts for reflective activities, like generating summaries or clarifying confusing elements [44]. They used reciprocal teaching with seventh graders struggling with reading comprehension, with an emphasis on letting students take over as the session progresses. The entire point of the exercise is to encourage reflective activities in students, and similar prompts—for summarization or clarification—may lead students in IFPR contexts to give reviews that cause more reflection, or more directly reflect themselves.

Davis and Linn [15] use explicit self-review prompts at different stages of assignments given to eighth graders. For example, in one assignment, students had to perform a repeated task (designing clothing and environments to help cold-blooded aliens survive). They compared responses to direct prompts submitted along with a design, like "Our design will work well because...", to prompts designed to encourage reflection after the fact, like "Our design could be better if we...", and plan-ahead prompts designed to cause reflection during the assignment, like "In thinking about doing our design, we need to think about..." Their sample size was small, and they did not find a significant difference in design quality between the direct and reflective prompts. They did find that students gave better explanations when given the reflective prompts, but the difference could easily be attributed to the small sample size.

Frederiksen and White have done a series of studies on reflective assessment [72] and reflective collaboration [20] in middle school science classes. In an online environment, students work on mock experiments using a scientific-method like flow for a project: they start with an initial inquiry, form hypotheses, analyze mock data, and draw conclusions. In between steps, they are asked questions that urge them to reflect on their work: why they think a hypothesis is true, if they are being meticulous in analyzing their results, and more. In addition, the environment contains simple autonomous agents, called advisors, that give automated feedback and suggestions to students. They are also given the opportunity to assess the work of other students in a few ways: they can simply rate the contents, or they can make specific suggestions, like telling a student that they should pay attention to a particular advisor. Finally, the course is complemented with role-playing activities where students take on the role of advisors, and give specific feedback - the advisors have specific flavors of feedback; an example in the paper is

that a student acting as the "'Skeptic' might [be asked to] say 'I disagree' or 'Prove it'" to another student's submission.

#### 8.4 Learning From Examples

In IFPR, students have the opportunity to take what they learn from examples of others' work and apply it to their own. There is existing research in how students take what they learn from existing examples (whether from peers or experts) and use it to learn principles they subsequently apply to new problems.

There is a large body of work that studies learning from worked examples [2], in which students learn solely or primarily from example solutions to problems with accompanying descriptions of the problem solving process used. Worked examples and peer solutions differ significantly. Typically, worked examples are carefully crafted by experts to help students learn a problem-solving process. In peer review, students see solutions produced by peers after attempting the problem themselves, not as part of learning how to do the problem. In addition, worked examples generally have explanations at a finer granularity than we would consider presenting peer review at. Nevertheless, some of the recommendations of the worked examples literature may be relevant in the IFPR setting: in a broad survey of worked examples research, Atkinson et al. recommend that students "experience a variety of examples per problem type," and that examples be presented in "close proximity to matched problems" [2]. Both of these recommendations are consonant with strategies we have proposed for IFPR assignments.

Kulkarni et al. [35] discuss changes in creative output between subjects who saw varying numbers of examples, and diversity in examples, prior to creating their own artwork. Subjects seeing more diverse examples created artwork with more unique features than subjects seeing a less diverse set or fewer examples. In an IFPR setting, students who see examples from others, especially when already primed to think about the same problem, may similarly have more options to draw on in their solution, rather than only using whatever techniques they would have tried in their initial submission.

In PeerWise [16], students created and reviewed one another's multiple-choice questions, which has elements both of learning by example and of review. Denny et al. find that students who engaged with the system more—by contributing and exploring more example questions than others (and more than they were required to by the course)—performed better than those who did not. However, it's not clear that the exposure to more examples caused the difference in performance.

## 8.5 Peer Instruction

Peer Instruction (PI), which is a specific form of student-centered pedagogy [40], has been shown to be a promising way to improve student performance [14] and engagement [54] both in introductory courses [14] and upper-division courses [37]. Peer instruction, as defined by Crouch et al. [14], focuses on engaging students in activities that require them to apply the core concepts under study and to explain these concepts to their peers. Concretely, a class taught using PI principles can consist of short presentations, each of which focuses on a particular core concept, which is then tested by presenting students a conceptual question, which the students first solve individually and then discuss in groups.

Related to our undertaking, the most interesting component of PI is the peer discussions that occur after presentation of each concept. This differs from (in-flow) peer review in that there isn't a submission in question that students are trying to improve. The exercises are a vehicle for discussion, and the discussion is the end goal, not producing a quality submission. Neither of these goals is necessarily more helpful than the other, but depending on expected outcomes, one can be more effective. In programming and writ-

ing disciplines, for example, one explicit goal is to train students to produce quality programs and written work. In mathematics or physics, it may be more important that students understand concepts and know how and where to apply them, rather than producing any particular artifact.

#### 8.6 Comprehending Program Structure

Program comprehension is at the same time a prerequisite and a learning goal of in-flow peer review of programming assignments. Students need some ability to read code in order to provide a meaningful review to one another, but at the same time, IFPR can lessen the cognitive burden of comprehension by having students review problems that are conceptually close (or identical) to something the reader has just encountered. The degree to which it tends toward one direction or another is a function of the experience level of the students and the goals of the particular course.

In addition, for in-flow assignments centered around programming, one goal of peer-review is to help students reflect on their own code structure. There is a rich literature on program comprehension (including contrasting experts and novices), but much of that focuses on understanding the *behavior* of a new program [19, 53, 58]. In the IFPR context, students already know the problem and (roughly) what the program is supposed to do. Reading others' code therefore has different goals: notably, to understand the structure that someone else brought to the problem and to contrast that with one's own. This is a less burdensome task than asking if the program *matches* an existing specification.

Studies on program comprehension comparing experts and novices have found that experts engage in more metacognitive behavior [17], so the metacognitive context of review may put students in the right frame of mind to understand programs in the first place. Other work on program comprehension suggests that the process has a lot to do with understanding the high-level plan of a program [58]. Since in the in-flow context students have at least constructed *a* plan of their own for the same or a similar problem, they may at least be able to determine if the solution they are viewing matches their plan, or is doing something different.

# 8.7 Increasing Socialization in Programming-Oriented

IFPR has potential to foster a collaborative and social atmosphere in programming assignments. We discussed some of the motivations for a more social CS course when discussing IFPR for non-majors (section 6.5). There are other approaches to meeting these goals that IFPR can learn from.

Garvin-Doxas and Barker emphasize the importance of the classroom climate, emphasizing that courses that reward "hero" programmers and individual accomplishment give rise to a defensive atmosphere that can be counterproductive to learning for students with less prior ability [22]. In later work, Barker and Garvin-Doxas describe the outcome of running an IT course more like a fine arts course than a traditional engineering course [4]. This included projects that were more meaningful, public critique of results, and routine collaboration. This approach created a classroom culture where learning is a social and community process, rather than isolated, and the result was a greater retention of female students than the traditional engineering teaching approach.

The technique from Barker et al. most relevant for in-flow peer review, though not completely the same, is the approach to knowledge sharing during lab work. Students actively solicit help from any student, for example, by yelling questions out, resulting in a fluid exchange of ideas and techniques. Even though the projects considered in these labs were run in an open, collaborative setting, cheating was avoided by using individualized assignments.

Another working group discussed design decisions in computermediated collaborative (CMC) educational settings [73]. That report emphasizes goals that CMC can help reach, including encouraging peer review, having teamwork experiences, developing selfconfidence, and improving communication skills. A course using IFPR that performs review through an online tool is certainly an instance of a CMC setting, and addresses many of the same goals.

#### 8.8 Existing Uses of In-flow Peer Review

Others have used strategies for peer review that fall under the umbrella of in-flow peer review, even though they did not go by that name.

In the implementation of a multi-stage compiler, students in Søndergaard's course review one another's work between stages [60]. The evaluation in that work was only in the form of surveys after the assignment, but shows generally positive attitudes from students indicating that they felt the review had helped.

Expertiza [48] (discussed in section 5.4.1) is used for large, multi-stage collaborative projects. This includes assessment of the reviews themselves as an explicit motivator for giving good feedback. It is notable that in Expertiza, students often review other students' components of a larger whole, which can be a task that the reviewer didn't complete him or herself. In several of our case studies (1, 3, 4, 13, and 14), students review an instance of the *same* work that they just did themselves.

CaptainTeach [47] supports in-flow peer review for programming assignments. The Web-based tool supports test-first, datastructure-first, and one-function-at-a-time stagings of programming problems. It uses asynchronous reviewing, where students see 2-3 reviews from the most recent students to submit, combined with random known-good and known-bad solutions provided by the staff. It supports a fixed set of open-ended review prompts combined with Likert scales for each of tests, implementation, and data structures. Students are also allowed to give (optional) feedback on reviews they received. Politz et al. report that students engaged with the process, submitting stages early enough to get reviews (more than 24 hours before the deadline), and receiving review feedback promptly (a few hours) [47]. In another analysis of data on reviews of test suites in CaptainTeach, Politz et al. report that students were more likely to add missing tests for a feature after review if they reviewed or were reviewed by a student who had tested for that

Informa's "Solve and Evaluate" approach integrates a simple form of peer review into a software-based class room response system [28]. During a lecture the instructor poses a problem, and each student solves it by creating a solution in Informa. Informa is not limited to multiple-choice problems; it also allows a variety of problem types, including free text (e.g., code snippets), or drawings (e.g., diagrams of the structure or state of a program). Students submit their solution as soon as they are done, and they immediately are assigned a solution of a peer for evaluation. They evaluate a solution simply by scoring it as correct or incorrect. In Informa, a key reason for including an evaluation phase is to keep the faster students engaged while the slower students are still solving the problem. While the results of peer review are not shown to the authors of the submissions, they are used by the instructor to estimate the level of understanding of the class, and to focus the class discussion that follows the evaluation phase. A lecture using Informa often consists of multiple stages, and often the problems in subsequent stages build on each other (e.g., the first problem asks students to draw a control-flow graph of a program with conditionals, and the second problem asks for a control-flow graph including loops). In such a scenario, a lecture with Informa is an instance of in-flow peer review.

#### 8.9 Actionable Peer Review

Some other uses of peer review on large projects are related to inflow peer review because they allow students to improve their work in response to review. These uses don't necessarily stage assignments into reviewable pieces, instead performing review on entire intermediate artifacts. For example, Clark has students exercise the functionality of one another's projects, and lets groups improve their work based on the feedback their classmates give them [10]. Similarly, Wang, et al. [68], Zeller [74], Papadopoulos et al. [45], and the Aropä system [25] use assignment structures that allow students to update revisions of entire submissions that were reviewed by peers. Other studies have students write test cases (or manually test) one another's work as part of a review [50, 57]. These tests are most often on entire systems, rather than on pieces of a project that build up along with reviews. Students do, however, have the chance to improve their projects in response to their peers' feedback.

## 9. Conclusion

IFPR is a highly-configurable mechanism for making peer review more actionable. It leverages the fact that many problems, both programming and otherwise, can be split into several steps that occur at key moments for triggering reflection, and uses those moments as vehicles for peer feedback. It encourages reflective and critical thinking in both reviewers and reviewees, and prepares students for professional activities in judging others' work and incorporating feedback into their own.

IFPR has a lot in common with existing peer review approaches, and in collaborative and participatory pedagogic styles in general. All of those benefits, from enhancing a sense of community to improving communication skills, are also goals of IFPR. The main new idea is to engage students in the collaborative process by better integrating feedback into the flow of assignments.

This report outlines a large space for designing in-flow peer review assignments. We encourage practitioners to consider many of these factors, but not to be intimidated by them, or to be concerned that there are too many challenges to tackle in adopting IFPR. The key task is to pick good moments for reflection in the middle of assignments, and use those moments to get the most out of peer feedback (which we already know has many benefits).

Not surprisingly, our discussions raised several questions for future research or in-class experimentation. In addition to obvious questions about which configurations of IFPR are most useful in various contexts, there are questions about overall logistics of IFPR. One of our reviewers posed two good examples:

- How extensively should one incorporate IFPR into a course in order to maximize the benefits?
- If the answer to the previous question includes multiple assignments, how can this be made to work in a single term (quarter or semester) without over-stretching both the students and resources?

We hope to see future projects and papers explore these and other questions.

#### 10. Acknowledgments

Many students and course staff participated in our experiments with peer review prior to the working group meeting. We appreciate their feedback, humor, and patience; their experiences influenced many working-group discussions. Our reviewers provided extremely useful and detailed feedback. We regret that the revision window was too narrow for us to incorporate more of their recommendations.

## **Bibliography**

- [1] Christopher Alvin, Sumit Gulwani, Rupak Majumdar, and Supratik Mukhopadhyay. Synthesis of geometry proof problems. In *Proc. AAAI Conference on Artificial Intelligence*, 2014
- [2] Robert K. Atkinson, Sharon J. Derry, Alexander Renkl, and Donald Wortham. Learning from Examples: Instructional Principles from the Worked Examples Research. *Review of Educational Research* 70(2), pp. 181–214, 2000.
- [3] Alberto Bacchelli and Christian Bird. Expectations, Outcomes, and Challenges of Modern Code Review. In Proc. International Conference on Software Engineering, 2013.
- [4] Lecia J. Barker, Kathy Garvin-Doxas, and Eric Roberts. What Can Computer Science Learn from a Fine Arts Approach to Teaching? In *Proc. Special Interest Group on Computer Science Education*, pp. 421–425, 2005.
- [5] B. S. Bloom and D. R. Krathwohl. Taxonomy of Educational Objectives: The Classification of Educational Goals. Handbook I: Cognitive Domain. Longmans, 1956.
- [6] David Boud, Ruth Cohen, and Jane Sampson. Peer learning and assessment. Assessment & Evaluation in Higher Education 24(4), pp. 413–426, 1999.
- [7] Kwangsu Cho and Charles MacArthur. Student revision with peer and expert reviewing. *Learning and Instruction* 20(4), pp. 328–338, 2010.
- [8] Kwangsu Cho and Christian D. Sun. Scaffolded writing and rewriting in the discipline: A web-based reciprocal peer review system. *Computers and Education* 48(3), pp. 409–426, 2005
- [9] Jenny C. C. Chung and Susanna M. K. Chow. Promoting student learning through a student-centred problem-based learning subject curriculum. *Innovations in Education and Teach*ing International 41(2), pp. 157–168, 2004.
- [10] Nicole Clark. Peer testing in software engineering projects. In *Proc. Australasian Computing Education Conference*, 2004.
- [11] Tony Clear. Thinking Issues: Managing Mid-project Progress Reviews: A Model for Formative Group Assessment in Capstone Projects. *ACM Inroads* 1(1), pp. 14–15, 2010.
- [12] Jason Cohen, Steven Teleki, and Eric Brown. Best Kept Secrets of Peer Code Review. SmartBear Software, 2013.
- [13] The College Board. AP Computer Science Principles, Draft Curriculum Framework. 2014.
- [14] Catherine H. Crouch and Eric Mazur. Peer instruction: Ten years of experience and results. *American Journal of Physics* 69(9), pp. 970–977, 2001.
- [15] Elizabeth A. Davis and Maria C. Linn. Scaffolding students' knowledge integration: prompts for reflection in KI. *Interna*tional Journal of Science Education 22(8), 2000.
- [16] Paul Denny, John Hamer, Andrew Luxton-Reilly, and Helen Purchase. PeerWise: Students Sharing Their Multiple Choice Questions. In Proc. Fourth International Workshop on Computing Education Research, 2008.
- [17] Anneli Eteläpelto. Metacognition and the Expertise of Computer Program Comprehension. Scandinavian Journal of Educational Research 37, pp. 243–254, 1993.
- [18] M. E. Fagan. Design and code inspections to reduce errors in program development. *IBM Systems Journal*, pp. 182–211, 1976.

- [19] Vikki Fix, Susan Wiedenbeck, and Jean Scholtz. Mental Representations of Programs by Novices and Experts. In *Proc. INTERACT*, 1993.
- [20] John R. Frederiksen and Barbara Y. White. Cognitive Facilitation: A Method for Promoting Reflective Collaboration. In Proc. Computer Support for Collaborative Learning, 1997.
- [21] Ursula Fuller, Joyce Currie Little, Bob Keim, Charles Riedesel, Diana Fitch, and Su White. Perspectives on Developing and Assessing Professional Values in Computing. *SIGCSE Bull.* 41(4), pp. 174–194, 2010.
- [22] Kathy Garvin-Doxas and Lecia J. Barker. Communication in Computer Science Classrooms: Understanding Defensive Climates As a Means of Creating Supportive Behaviors. J. Educ. Resour. Comput. 4(1), 2004.
- [23] Rich Gerber and Paolo Gai. The START V2 ConferenceManager. 2014.
- [24] John Hamer, Quintin Cutts, Jana Jackova, Andrew Luxton-Reilly, Robert McCartney, Helen Purchase, Charles Riedesel, Mara Saeli, Kate Sanders, and Judithe Sheard. Contributing Student Pedagogy. SIGCSE Bulletin 40(4), pp. 194–212, 2008
- [25] John Hamer, Catherine Kell, and Fiona Spence. Peer Assessment using Aropä. In Proc. Australasian Computing Education Conference, 2007.
- [26] John Hamer, Kenneth T. K. Ma, and Hugh H. F. Kwong. A method of automatic grade calibration in peer assessment. In Proc. Australasian Conference on Computing Education, 2005
- [27] John Hamer, Helen Purchase, Andrew Luxton-Reilly, and Paul Denny. A comparison of peer and tutor feedback. In Proc. Assessment & Evaluation in Higher Education, 2014.
- [28] Matthias Hauswirth and Andrea Adamoli. Teaching Java Programming with the Informa Clicker System. Science of Computer Programming, 2011.
- [29] Geert Hofstede, Gert Jan Hofstede, and Michael Minkov. Cultures and Organizations: Software of the Mind. McGraw-Hill, 2005.
- [30] Simon Hooper and Michael J Hannafin. The Effects of Group Composition on Achievement, Interaction, and Learning Efficiency During Computer-Based Cooperative Instruction. *Journal of Educational Computing Research* 4, pp. 413–424, 1988.
- [31] Christopher D. Hundhausen, Anukrati Agrawal, and Pawan Agarwal. Talking About Code: Integrating Pedagogical Code Reviews into Early Computing Courses. *Transactions on Computing Education* 13(3), 2013.
- [32] Eddie Kohler. Hot Crap! In Proc. Workshop on Organizing Workshops, Conferences, and Symposia in Computer Systems, WOWCS at NSDI'08, 2008.
- [33] Alfie Kohn. Punished By Rewards. Houghton Mifflin Company, 1999.
- [34] Shriram Krishnamurthi. The CONTINUE Server (or, How I Administered PADL 2002 and 2003). In Proc. International Symposium on Practical Aspects of Declarative Languages, 2003.
- [35] Chinmay Kulkarni, Steven P. Dow, and Scott R. Klemmer. Early and Repeated Exposure to Examples Improves Creative Work. In *Proc. Cognitive Science*, 2012.

- [36] Chinmay Kulkarni, Koh Pang Wei, Huy Le, Daniel Chia, Kathryn Papadopoulos, Justin Cheng, Daphne Koller, and Scott R. Klemmer. Peer and Self Assessment in Massive Online Classes. ACM Transactions on Computer-Human Interaction, 2013.
- [37] Cynthia Baily Lee, Saturnino Garcia, and Leo Porter. Can Peer Instruction Be Effective in Upper-division Computer Science Courses? *Transactions on Computing Education* 13(3), pp. 1–22, 2013.
- [38] Ngar-Fun Liu and David Carless. Peer feedback: the learning element of peer assessment. *Teaching in Higher Education* 11, pp. 279–290, 2006.
- [39] Andrew Luxton-Reilly. A Systematic Review of Tools that Support Peer Assessment. Computer Science Education 19(4), pp. 209–232, 2009.
- [40] Eric Mazur. Peer Instruction: A User's Manual. 1996.
- [41] Amanda Miller and Judy Kay. A Mentor Program in CS1. In Proc. ACM SIGCSE Conference on Innovation and Technology in Computer Science Education, 2002.
- [42] Mark Montague, Jeremiah Konkle, and Luke Montague. Linklings. 2014.
- [43] M. M. Nelson and C. D. Schunn. The nature of feedback: How different types of peer feedback affect writing performance. *Instructional Science* 27(4), pp. 375–401, 2009.
- [44] Annemarie Sullivan Palinscar and Ann L. Brown. Reciprocal Teaching of Comprehension-Fostering and Comprehension-Monitoring Activities. *Cognition and Instruction*, pp. 117– 175, 1984.
- [45] Pantelis M. Papadopoulos, Thomas D. Lagkas, and Stavros N. Demetriadis. How to Improve the Peer Review Method: Free-selection vs Assigned-pair Protocol Evaluated in a Computer Networking Course. *Computing & Education* 59(2), pp. 182–195, 2012.
- [46] Joe Gibbs Politz, Shriram Krishnamurthi, and Kathi Fisler. In-flow Peer Review of Tests in Test-First Programming. In Proc. Innovation and Technology in Computer Science Education, 2014.
- [47] Joe Gibbs Politz, Daniel Patterson, Shriram Krishnamurthi, and Kathi Fisler. CaptainTeach: Multi-Stage, In-Flow Peer Review for Programming Assignments. In *Proc. Innovation* and Technology in Computer Science Education, 2014.
- [48] Lakshmi Ramachandran and Edward F. Gehringer. Reusable learning objects through peer review: The Expertiza approach. In *Proc. Innovate: Journal of Online Education*, 2007.
- [49] Lakshmi Ramachandran and Edward F. Gehringer. Automated Assessment of Review Quality Using Latent Semantic Analysis. In Proc. IEEE International Conference on Advanced Learning Technologies, 2011.
- [50] K. Reily, P. L. Finnerty, and L. Terveen. Two peers are better than one: Aggregating peer reviews for computing assignments is surprisingly accurate. In *Proc. ACM International Conference on Supporting Group Work*, 2009.
- [51] Jochen Rick and Mark Guzdial. Situating CoWeb: A scholarship of application. *International Journal of Computer-Supported Collaborative Learning* 1(1), pp. 89–115, 2006.
- [52] Dieter Rombach, Marcus Ciolkowski, Ross Jeffery, Oliver Laitenberger, Frank McGarry, and Forrest Shull. Impact of Research on Practice in the Field of Inspections, Reviews and Walkthroughs: Learning from Successful Industrial Uses. SIGSOFT Softw. Eng. Notes 33(6), pp. 26–35, 2008.

- [53] Carsten Schulte, Tony Clear, Ahmad Taherkhani, Teresa Busjahn, and James H. Paterson. An Introduction to Program Comprehension for Computer Science Educators. In Proc. Proceedings of the 2010 ITiCSE Working Group Reports, 2010.
- [54] Beth Simon, Sarah Esper, Leo Porter, and Quintin Cutts. Student Experience in a Student-centered Peer Instruction Classroom. In Proc. ACM Conference on International Computing Education Research, 2013.
- [55] Gerald K. Sims. Student peer review in the classroom: a teaching and grading tool. *Journal of Agronomic Education* 18(2), pp. 105–108, 1989.
- [56] R. Singh, S. Gulwani, and S. Rajamani. Automatically generating algebra problems. In *Proc. AAAI Conference on Artificial Intelligence*, 2012.
- [57] Joanna Smith, Joe Tessler, Elliot Kramer, and Calvin Lin. Using Peer Review to Teach Software Testing. In *Proc. International Computing Education Research Conference*, 2012.
- [58] Elliot Soloway and Kate Ehrlich. Empirical Studies of Programming Knowledge. *IEEE Transactions of Software Engineering* 10(5), pp. 595–609, 1984.
- [59] Karen Swan, Jia Shen, and Starr Roxanne Hiltz. Assessment and Collaboration in Online Learning. *Journal of Asyn*chronous Learning, 2006.
- [60] Harald Søndergaard. Learning from and with Peers: The Different Roles of Student Peer Reviewing. In Proc. ACM SIGCSE Conference on Innovation and Technology in Computer Science Education, 2009.
- [61] Keith Topping. Peer Assessment Between Students in Colleges and Universities. *Review of Educational Research* 68(3), pp. 249–276, 1998.
- [62] Nancy M. Trautmann. Designing Peer Review for Pedagogical Success: What Can We Learn from Professional Science? *Journal of College Science Teaching* 38(4), pp. 14–19, 2009.
- [63] Richard van de Stadt. CyberChair: A Web-Based Groupware Application to Facilitate the Paper Reviewing Process. 2014.
- [64] Susan van Rooyen, Nick Black, and Fiona Godlee. Development of the review quality instrument (RQI) for assessing peer reviews of manuscripts. *Journal of Clinical Epidemiology* 52(7), pp. 625–629, 1999.
- [65] Andrei Voronkov. EasyChair Conference Management System. 2014.
- [66] Lawrence G. Votta Jr. Does every inspection need a meeting? In *Proc. Foundations of Software Engineering*, 1993.
- [67] Lev Vygotsky. Interaction between learning and development. *Mind and Society*, pp. 79–91, 1978.
- [68] Yanqing Wang, Hang Li, Yanan Sun, Jiang Yu, and Jie Yu. Learning outcomes of programming language courses based on peer code review model. In *Proc. International Conference* on Computer Science & Education, 2011.
- [69] Gerald M. Weinberg. The Psychology of Computer Programming. Van Nostrand Reinhold, 1971.
- [70] Patrick Wessa. How Reproducible Research Leads to Non-Rote Learning within Socially Constructivist Statistics Education. *Electronic Journal of e-Learning* 7(2), pp. 173–182, 2009.
- [71] Patrick Wessa and Antoon De Rycker. Reviewing peer reviews-A rule-based approach. In *Proc. 5th International Conference on E-Learning (ICEL)*, 2010.

- [72] Barbara Y. White, John R. Frederiksen, T. Frederiksen, E. Eslinger, and A. Collins. Inquiry Island: Affordances of a Multi-Agent Environment for Scientific Inquiry and Reflective Learning. In Proc. International Conference of the Learning Sciences (ICLS), 2002.
- [73] Ursula Wolz, Jacob Palme, Penny Anderson, Zhi Chen, James Dunne, Göran Karlsson, Atika Laribi, Sirkku Männikkö, Robert Spielvogel, and Henry Walker. Computer-mediated Communication in Collaborative Educational Settings (Report of the ITiCSE '97 Working Group on CMC in Collaborative Educational Settings). In Proc. The Supplemental Proceedings of the Conference on Integrating Technology into Computer Science Education: Working Group Reports and Supplemental Proceedings, 1997.
- [74] Andreas Zeller. Making Students Read and Review Code. In *Proc. ACM SIGCSE Conference on Innovation and Technology in Computer Science Education*, 2000.