

Accepted Manuscript

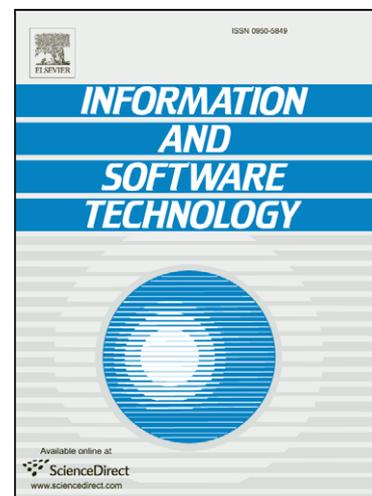
The Effectiveness of Pair Programming: A Meta-Analysis

Jo E. Hannay, Tore Dybå, Erik Arisholm, Dag I.K. Sjøberg

PII: S0950-5849(09)00012-3
DOI: [10.1016/j.infsof.2009.02.001](https://doi.org/10.1016/j.infsof.2009.02.001)
Reference: INFSOFT 4907

To appear in: *Information and Software Technology*

Received Date: 20 October 2008
Revised Date: 9 February 2009
Accepted Date: 9 February 2009



Please cite this article as: J.E. Hannay, T. Dybå, E. Arisholm, D.I.K. Sjøberg, The Effectiveness of Pair Programming: A Meta-Analysis, *Information and Software Technology* (2009), doi: [10.1016/j.infsof.2009.02.001](https://doi.org/10.1016/j.infsof.2009.02.001)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

The Effectiveness of Pair Programming: A Meta-Analysis

Jo E. Hannay^{a,b,*} Tore Dybå^{a,c} Erik Arisholm^{a,b} Dag I.K. Sjøberg^{a,b}

^aSimula Research Laboratory, Department of Software Engineering, Pb. 134, NO-1325 Lysaker, Norway

^bUniversity of Oslo, Department of Informatics, Pb. 1080 Blindern, 0316 Oslo, Norway

^cSINTEF Information and Communication Technology, NO-7465 Trondheim, Norway

Abstract

Several experiments on the effects of pair versus solo programming have been reported in the literature. We present a meta-analysis of these studies. The analysis shows a small significant positive overall effect of pair programming on quality, a medium significant positive overall effect on duration, and a medium significant negative overall effect on effort. However, between-study variance is significant, and there are signs of publication bias among published studies on pair programming. A more detailed examination of the evidence suggests that pair programming is faster than solo programming when programming task complexity is low and yields code solutions of higher quality when task complexity is high. The higher quality for complex tasks comes at a price of considerably greater effort, while the reduced completion time for the simpler tasks comes at a price of noticeably lower quality. We conclude that greater attention should be given to moderating factors on the effects of pair programming.

Key words: Pair programming, Evidence-based software engineering, Systematic review, Meta-analysis, Fixed effects, Random effects

1. Introduction

Much of the current interest in pair programming is likely due to the popularity of extreme programming (XP), of which pair programming is one of the core practices [2]. Common-sense, but scientifically unwarranted claims as to both the benefits and the adversities of pair programming abound. Advocates of pair programming claim that it has many benefits over individual programming when applied to new code development or when used to maintain and enhance existing code. Stated benefits include higher-quality code in about half the time as individuals, happier programmers, improved teamwork, improved knowledge transfer, and enhanced learning [53]. There are also expectations with respect to the benefits and drawbacks of various kinds of pairing, e.g., that “expert-expert pairings seem to be especially accelerated” (*ibid.*, p. 102) and that “novice-novice is the most problematic” (*ibid.*, p. 120). Stephens and Rosenberg [49] claim that pair programming is clearly the least beneficial of XP practices and also that novice-novice pairing is obviously an undesirable combination, because it is akin to “the blind leading the blind”.

With the considerable interest in pair programming by the software industry and academia, it is important to determine scientifically, whether, and if so when, pairing two programmers is beneficial in terms of important cost drivers such as duration, effort, and quality. The finding across decades of small group research is that groups usually fall short of reasonable expectations to improved performance [48,35,33]. An important question, therefore, is whether the claims regarding pair programming can be substantiated by empirical evidence and how pair programming relates to such group research. Several empirical studies have now been conducted that set out to examine the effects of pair programming in a systematic manner. This article provides a systematic review in terms of a meta-analysis of all published experiments on the effectiveness of pair programming and subsequently offers recommendations for evidence-based practice [18]. This review extends the intermediate analysis reported in [15], which briefly summarized pair programming experiments published up and until 2006. In the present full-scale analysis, we take into account between-study variance, subgroup differences and publication bias. We also take into account studies published in 2007.

Section 2 presents the methods for systematic review and statistical meta-analysis. Section 3 presents the results

* Corresponding author. E-mail: johannay@simula.no

from the analysis, and Section 4 discusses implications for theory and practice. Section 5 concludes.

2. Method

Informed by the general procedures for performing systematic reviews [34] and the established methods of meta-analysis [39,45], we undertook the meta-analysis in distinct stages: identification of inclusion and exclusion criteria, search for relevant studies, identification and selection of studies, data extraction, and statistical analysis and synthesis, see details below.

The meta-analysis focused on combining quantitative effects on three central outcome constructs that were investigated in the included studies. We did not assess the quality of the included studies in terms of, e.g., the appropriateness of the chosen effect size measures [30], the appropriateness of randomization procedures [31], subject/task selection and validity issues [46], statistical power [17], the use of theory [25], the approach to realism [24], etc. Future meta-analyses might incorporate study quality, but at present, it is not clear how to aggregate this multifaceted concept into a single measure to be used in a meta-analysis.

2.1. Inclusion and exclusion criteria

Studies were eligible for inclusion in the meta-analysis if they presented quantitative data on the effectiveness of pair programming in which a comparison was made between pairs and individuals, possibly in a team context. The subjects could be either students or professional software developers. Included studies had to report one of the primary outcomes *Quality*, *Duration*, or *Effort*. We did not want to put any restrictions on the operationalization of these outcome constructs. Furthermore, the studies had to be reported in English.

2.2. Data sources and search strategy

The search strategy included electronic databases and hand searches of conference proceedings. We searched in the following electronic databases: ACM Digital Library, Compendex, IEEE Xplore, and ISI Web of Science. We did not perform separate searches in the SE-specific databases Kluwer Online, ScienceDirect, SpringerLink, and Wiley Inter Science Journal Finder, because previous experience with systematic search strategies has shown that articles retrieved from these databases are also returned by either ISI Web of Science or Compendex [16]. In addition to the electronic databases, we hand-searched all volumes of the following thematic conference proceedings: XP, XP/Agile Universe, Agile, and Agile Development Conference. We used the basic search string “pair programming” OR “collaborative programming” to conduct the searches.

2.3. Study identification and selection

The identification and selection process consisted of three stages. At Stage 1, the second author applied the search string to the titles, abstracts, and keywords of the articles in the included electronic databases and conference proceedings. All retrieved articles were published, or accepted for publication, before or in 2007. Excluded from the search were editorials, prefaces, article summaries, interviews, news items, correspondence, discussions, comments, reader’s letters, and summaries of tutorials, workshops, panels, and poster sessions. This search strategy resulted in a total of 236 unique citations.

At Stage 2, the first and second authors went through the titles and abstracts of all the studies resulting from stage 1 for relevance to the review. If it was unclear from the title, abstract, and keywords whether a study conformed to our inclusion criteria, it was included for a detailed review. This screening process resulted in 57 citations that were passed on to the next stage.

At Stage 3, the full text of all 57 citations from Stage 2 were retrieved and reviewed in detail by the first and second authors. This resulted in 23 included articles according to the inclusion criteria. Five of these did not report enough information to compute standardized effect sizes and were excluded. Thus, 18 studies (all experiments) met the inclusion criteria and were included in the review (see Appendix).

2.4. Data extraction and checking

We collected data from the 18 articles, including type of treatment, type of system, type of tasks, duration of the study, number of groups, group assignment, type of subjects and their experience with pair programming, number of pairs, number of individuals, outcome variable, means, standard deviations, counts, percentages, and p -values. Every article included in the review was read in detail and the data was extracted and cross-checked by the first, second and third authors. Discrepancies were resolved by discussion among all four authors.

2.5. Statistical analysis

We used Comprehensive Meta-Analysis *v2* to calculate effect size estimates for all the tests in the 18 articles.¹ In order to be comparable across studies, effect sizes must be standardized. In this meta-analysis, we used Hedges’ g as the standardized measure of effect size. Like Cohen’s d and Glass’ Δ , Hedges’ g is simply the difference between the outcome means of the treatment groups, but standardized with respect to the pooled standard deviation, s_p , and corrected for small sample bias [37]:

¹ Comprehensive Meta-Analysis is a trademark of Biostat Inc.

$$\text{Hedges' } g = \frac{\bar{x}_1 - \bar{x}_2}{s_p} \quad (1)$$

The pooled standard deviation is based on the standard deviations in both groups, s_1 , s_2 :

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{(n_1 - 1) + (n_2 - 1)}} \quad (2)$$

Hedges' g , Cohen's d , and Glass' Δ have the same properties in large samples (i.e., they are equivalent in the limit $(n_1 + n_2) \rightarrow \infty$), but Hedges' g has the best properties for small samples when multiplied by a correction factor that adjusts for small sample bias [26]:

$$\text{Correction factor for Hedges' } g = 1 - \frac{3}{4(N - 2) - 1}, \quad (3)$$

where N = total sample size.

An effect size of 0.5 thus indicates that the mean of the treatment group is half a standard deviation larger than the mean of the control group. Effect sizes larger than 1.00 can be assumed to be *large*, effect sizes of 0.38–1.00 *medium*, and effect sizes of 0–0.37 *small* [30].

2.6. Synthesis

We conducted separate meta-analyses for the three outcome constructs *Quality*, *Duration*, and *Effort*. Some studies applied several tests on the same outcome construct. In these cases, we used the mean of the effect sizes over these tests to give only one effect size per outcome per study. Comprehensive Meta-Analysis calculates these mean effect sizes as an option.

Because we expected considerable heterogeneity, we decided to calculate the resulting meta-analytic effect sizes both under the assumption of the random-effects model and under the assumption of the fixed-effects model. These models lead to different significance test and confidence intervals for meta-analytic results [29].

The fixed-effects model assumes an unknown and fixed population effect size that is estimated by the studies in the meta-analysis. All the studies in the meta-analysis are seen as drawn from the same population, and variances in effect sizes between individual studies are viewed as due to subject variability [27].

The random-effects model, on the other hand, assumes an unknown and stochastic population effect-size distribution. That is, the true effect size of pair programming varies around a mean μ . This caters for the view that the effect of pair programming depends on situational variables and other factors (both known and unknown) that are not taken into consideration in the analysis. Variance in effect sizes are then seen as due to subject variability, and also to inter-study variability, since each study is seen as approximating a different part of the true effect size distribution [27,4].

Both models relate to an unknown population parameter. The approaches are hence referred to as *unconditional* [37]. The choice of which model to use is made prior to the analysis based on theory, past empirical findings, and on insights as to what the included studies describe.

However, one may also use the analysis techniques associated with the two models merely to characterize the studies relatively to each other without any reference to a population effect size. Which model to use in this case, is determined from the observed data based on heterogeneity measures, which are calculated under the assumption of a fixed-effects model. If heterogeneity is non-significant, a fixed-effects model is an appropriate characterization of data. Otherwise, a random-effects model best characterizes the data. The results from such *conditional* approaches should, however, not be confounded with statements regarding population parameters.

We conducted our analysis from both an unconditional and a conditional perspective. For the unconditional perspective, we chose the random-effects model. Hedges and Vevea state that "In the case of random-effects models, for example, some individual effect-size parameters may be negative even though μ is positive. That corresponds to the substantive idea that some realizations of the treatment may actually be harmful even if the average effect is beneficial" [27]. Results in [1] suggest that the effects of pair programming may be positive or negative dependent on other factors (e.g., expertise and task complexity). Also, the outcome constructs *Quality*, *Duration*, and *Effort* are not yet well-defined in software engineering. These constructs are operationalized in very diverse ways (Section 3.1), and for the time being, it is reasonable to view these diverse operationalizations as different aspects of a construct (a so-called *formative* measurement model [3,6,12,40]). Under these circumstances, and until the constructs are better understood, it is reasonable to relate to a random-effects model.

In the conditional approach, we tested whether there were genuine differences underlying the results of the studies (heterogeneity), or whether the variation in the findings was compatible with chance alone.

In the following, we give an overview of the technical details of the meta-analysis. For further elaborations, see e.g., [27,4,39]. Let k be the number of studies in the meta-analysis. Let T_i be the standardized effect size estimate (in our case, Hedges' g) of study i . In the fixed-effects model, the estimate \bar{T}_\bullet of the assumed fixed population effect size, and the estimate's variance v_\bullet , are

$$\bar{T}_\bullet = \frac{\sum_{i=1}^k w_i T_i}{\sum_{i=1}^k w_i} \quad v_\bullet = \frac{1}{\sum_{i=1}^k w_i} \quad (4)$$

where $w_i = 1/v_i$ is the weight assigned to study i , i.e., the reciprocal of the variance v_i for study i . Thus, \bar{T}_\bullet is a weighted mean over the effect sizes of the individual studies, where studies with less variance are given greater weight. In the random-effects model, the weights

Table 1
Characteristics of the Included Studies.

Study	Subjects	N_{Tot}	N_{Pair}	N_{Ind}	Study setting
Arisholm <i>et al.</i> (2007) ¹	Professionals	295	98	99	10 sessions with individuals over 3 months and 17 sessions with pairs over 5 months (each of 1 day duration, with different subjects). Modified 2 systems of about 200-300 Java LOC each.
*Baheti <i>et al.</i> (2002) ²	Students	134	16	9	Teams had 5 weeks to complete a curricular OO programming project. Distinct projects per team.
Canfora <i>et al.</i> (2005) ³	Students	24	12	24	2 applications each with 2 tasks (run1 and run2).
Canfora <i>et al.</i> (2007) ⁴	Professionals	18	5(4)	8(10)	Study session and 2 runs (totalling 390 minutes) involving 4 maintenance tasks (grouped in 2 assignments) to modify design documents (use case and class diagrams).
Domino <i>et al.</i> (2007)	Professionals Students	88	28	32	Run as several sessions during a period of two months. Pseudocode on "Create-Design" tasks Test-driven development.
*Heiberg <i>et al.</i> (2003) ⁵	Students	84(66)	23(16)	19(17)	4 sessions over 4 weeks involving 2 programming tasks to implement a component for a larger "gamer" system.
Madeyski (2006) ⁶	Students	188	28	31(35)	8 laboratory sessions involving 1 initial programming task in a finance accounting system (27 user stories).
Madeyski (2007)	Students	98	35	28	Java course project of developing a 27 user story accounting system over 8 laboratory sessions of 90 minutes each. Test-driven development.
Müller (2005) ⁷	Students	38	19	23	2 runs of 1 programming session each on 2 initial programming tasks (Polynomial and Shuffle-Puzzle) producing about 150 LOC.
Müller (2006) ⁸	Students	18(16)	4(5)	6	1 session involving initial design + programming tasks on an elevator system.
Nawrocki & Wojciechowski (2001) ⁹	Students	15	5	5	4 lab sessions over a winter semester, as part of a University course. Wrote four C/C++ programs ranging from 150-400 LOC.
Nosek (1998) ¹⁰	Professionals	15	5	5	45 minutes to solve 1 programming task (database consistency check script).
*Phongpaibul & Boehm (2006) ¹¹	Students	95	7	7	12 weeks to complete 4 phases of development + inspection.
*Phongpaibul & Boehm (2007)	Students	36	5	4	Part of a team project to extend a system. 13 weeks to complete 4 phases of development + inspection.
Rostaher & Hericko (2002) ¹²	Professionals	16	6	4	6 small user stories filling 1 day.
*Vanhanen & Lassenius (2005) ¹³	Students	20	4	8	9-week student project in which each subject spent a total of 100 hours (400 hours per team). A total of 1500-4000 LOC was written.
Williams <i>et al.</i> (2000) ¹⁴	Students	41	14	13	6-week course where the students had to deliver 4 programming assignments.
Xu & Rajlich (2006) ¹⁵	Students	12	4	4	2 sessions with pairs and 1 session with individuals. 1 initial programming task producing around 200-300 LOC.

are based on between-study variance in addition to within-study variance v_i . Specifically, the estimate \bar{T}_\bullet^* of the mean μ of the assumed population effect size distribution, and the estimate's variance v_\bullet^* , are

$$\bar{T}_\bullet^* = \frac{\sum_{i=1}^k w_i^* T_i}{\sum_{i=1}^k w_i^*} \quad v_\bullet^* = \frac{1}{\sum_{i=1}^k w_i^*} \quad (5)$$

where $w_i^* = 1/v_i^*$, for $v_i^* = v_i + \tau^2$. Here, τ^2 is the additional between-study variance:

$$\tau^2 = \begin{cases} \frac{Q - df}{C} & \text{if } Q > df \\ 0 & \text{if } Q \leq df \end{cases} \quad (6)$$

where the degrees of freedom $df = k - 1$, and Q represents the total variance:

$$Q = \sum_{i=1}^k w_i (T_i - \bar{T}_\bullet)^2 \quad (7)$$

In Equation (6), C is simply a scaling factor that ensures that τ^2 has the same denomination as within-study variance, i.e., $C = \sum w_i - \sum w_i^2 / \sum w_i$.

In fact, Q is a statistic that indicates heterogeneity, and one that we used for this purpose. A significant Q rejects the null hypothesis of homogeneity and indicates that the variability among the effect sizes is greater than what is likely to have resulted from subject-level variability alone [27]. We also calculated the I^2 -statistic, which indicates heterogeneity in percentages:

$$I^2 = 100\%(Q - df)/Q \quad (8)$$

A value of 0 indicates no observed heterogeneity, 25% low, 50% moderate, and 75% high heterogeneity [28].

2.7. Assumption of Parametric Normality

Hedges' g is a parametric effect size measure that is based, firstly, on the assumption that one wishes to relate to a population distribution, and, secondly, that each sub-population (solo, pair) is normally distributed on the response variable. This assumption also underlies the meta-analytic procedures and estimates that we used. It is well-known that *Duration* and *Effort* measures have Gamma distributions, and the various *Quality* measures in this study may or may not have normal distributions. Thus, one may question whether other standardized effect-size measures should have been used for the meta-analysis. Non-parametric effect-size measures would merely allow one to characterize data and not allow one to relate to a population. On the other hand, specifying other distributions (e.g., Gamma) would demand that we had access to the individual studies' raw data. Since means are meaningful measures for all three outcome constructs, we therefore decided to use Hedges' g as a best compromise for conducting the meta-analysis. Note also that even if the assumed population distribution is not evident in a (small) sample, this reason alone should not lead one to abandon the model (unless the sole aim is to characterize data). Small samples will of course, lead to less confident parameter estimates but these confidence estimates will be calculated correctly according the assumed population distribution.

3. Results

We first present characteristics of the 18 studies that were included in the review and meta-analysis. We then give the meta-analytic results in terms of overall effects, subgroup effects and publication bias.

3.1. Description of studies

Characteristics of the 18 studies included in the meta-analysis are summarized in Table 1 in alphabetical order, while full citations are provided in the Appendix.

Of the 18 studies, 11 were from Europe and seven from North America. The number of subjects in the studies varied from 12 to 295. Thirteen of the studies used students as subjects, while four used professionals. One used both professionals and students. Five studies made the comparison within a team context, that is, teams of pairs versus teams of individuals (marked with an asterisk in the first column). The studies often administered several tests and the number of data points may have varied across tests (numbers in parentheses in Table 1).

All studies used programming tasks as the basis for comparison. In addition, Madeyski (2006) and Madeyski (2007) included test-driven development; Müller (2005), Phongpaibul and Boehm (2006), and Phongpaibul and

Boehm (2007) included inspections; and Müller (2006) included design tasks.

Quality was typically reported as the number of test cases passed or number of correct solutions of programming tasks, but student grades, delivered functionality, and metrics for code complexity were also used as measures of *Quality*. *Duration* was reported mainly in two modes: as the total time taken to complete all tasks considered (all solutions), or as the total time taken to complete the tasks that had been assessed as having passed a certain quality standard (checked solutions). For comparisons between pairs and solo programmers, pair *Effort* was reported as twice the duration of each individual in the pair. For team-based comparisons, e.g., teams of individuals versus teams of pairs, *Effort* was reported as the total effort spent by the respective groups.

Thus, the studies included in the meta-analysis do not all apply the same measures or have similar context variables. Rather, they investigate the effectiveness of pair programming with respect to different aspects of the constructs *Quality*, *Duration*, and *Effort*. As such, the studies may be seen as differentiated replications [36], and any measure from a particular study in the meta-analysis is but one indicator of, perhaps, many indicators regarding one of these constructs.

3.2. Effects of pair programming

Table 2 summarizes the meta-analysis. The g column shows the meta-analytic estimates \bar{T}_\bullet and \bar{T}_\bullet^* , in terms of Hedges' g , of the population effect size parameter in the fixed-effects model and random-effects model, respectively, along with 95% confidence intervals and p -values. Also given are heterogeneity measures calculated under the assumption of a fixed-effects model. The *Overall effects* analysis will be described in detail in this section, and the *Subgroup* analysis will be given in detail in Section 3.3.

Overall, 14 studies used *Quality* as an outcome construct, 11 used *Duration*, and 11 used *Effort*. Figure 1 shows Forest plots² of the standardized effects for each of the three outcome constructs. The studies are sorted according to the relative weight that a study's effect size receives in the meta-analysis. Relative weights are normalized versions of the weights w_i or w_i^* used in calculating the meta-analytic estimates \bar{T}_\bullet or \bar{T}_\bullet^* . The rightmost columns in Figure 1 show these weights according to the fixed-effects and random-effects models. Estimates from larger studies will usually be more precise than the estimates from smaller studies; hence, larger studies will generally be given greater weight.

The squares indicate the effect size estimate for each study. The size of each square is proportional to the relative weight of the study according to the random-effects

² The plots were generated with PSTricks postscript macros in LaTeX within an Excel spreadsheet using data produced by Comprehensive Meta-Analysis.

Table 2
Summary of Meta-analysis.

Analysis	<i>k</i>	Model	Effect Size			Heterogeneity			
			<i>g</i>	95% CI	<i>p</i>	<i>Q</i>	<i>I</i> ²	τ^2	<i>p</i>
<i>Overall effects</i>									
<i>Quality</i>	14	fixed	.23	.09 .37	.00	35.97	63.86	.14	.00
		random	.33	.07 .60	.01				
<i>Duration</i>	11	fixed	.40	.21 .59	.00	33.57	70.21	.28	.00
		random	.54	.13 .94	.01				
<i>Effort</i>	11	fixed	-.73	-.94 -.51	.00	66.36	84.93	.87	.00
		random	-.52	-1.18 .13	.12				
<i>Subgroups</i>									
<i>Students</i>									
<i>Quality</i>	11	fixed	.22	.06 .38	.01	32.97	69.66	.18	.00
		random	.32	-.01 .65	.06				
<i>Duration</i>	7	fixed	.58	.33 .84	.00	12.28	51.13	.13	.06
		random	.63	.24 1.02	.00				
<i>Effort</i>	8	fixed	-.59	-.88 -.30	.00	48.85	85.67	1.17	.00
		random	.04	-.82 .91	.92				
<i>Professionals</i>									
<i>Quality</i>	3	fixed	.26	-.05 .57	.10	2.97	32.56	.07	.23
		random	.37	-.10 .85	.12				
<i>Duration</i>	4	fixed	.16	-.13 .46	.28	16.87	82.22	.83	.00
		random	.50	-.55 1.54	.35				
<i>Effort</i>	3	fixed	-.90	-1.22 -.58	.00	15.48	87.08	1.54	.00
		random	-1.99	-3.56 -.41	.01				
<i>Teams</i>									
<i>Quality</i>	3	fixed	.19	-.22 .60	.36	.23	.00	.00	.89
		random	.19	-.22 .60	.36				
<i>Duration</i>	2	fixed	.34	-.13 .81	.16	1.46	31.55	.06	.23
		random	.31	-.27 .89	.30				
<i>Effort</i>	2	fixed	.74	-.13 1.61	.09	11.25	91.11	4.14	.00
		random	.99	-1.96 3.94	.51				
<i>No teams</i>									
<i>Quality</i>	11	fixed	.24	.09 .39	.00	35.70	71.99	.18	.00
		random	.38	.05 .70	.02				
<i>Duration</i>	9	fixed	.41	.20 .63	.00	32.02	75.02	.37	.00
		random	.63	.13 1.13	.01				
<i>Effort</i>	9	fixed	-.82	-1.05 -.60	.00	43.36	81.55	.61	.00
		random	-.85	-1.48 -.23	.01				

model. The relative weights of a random-effects model are generally more uniform than those of fixed-effects models, due to the incorporation of between-study variance into all the studies' weights. The horizontal lines indicate the 95% confidence intervals for each study's effect size estimate according to the random-effects model.

The diamonds at the bottom give the meta-analytic effect size estimate according to the fixed-effects and the random-effects model, i.e., \bar{T}_\bullet and \bar{T}_\bullet^* , respectively. The diamonds' centers and widths indicate the estimates and their 95% confidence interval, respectively.

Figure 2 shows *one-study-removed* analyses for each of the three outcome constructs. The plots show the meta-analytic effect size estimate when each study is removed from the meta-analysis. The resulting deviation from the full analysis indicates the *sensitivity* of the full analysis with respect to each study, that is, how much difference a given study makes to the meta-analysis.

3.2.1. Effects of pair programming on Quality

Fourteen studies compared the effects of pair programming on *Quality* in a total of 38 tests. These studies used a total of 1160 subjects, although in some studies, not all subjects were used in these particular tests. The subjects were distributed to study units (pairs, solo, teams of pairs, teams of solo) in various ways (Section 3.1). The meta-analytic effect size estimate is .23 in the fixed effects model and .33 in the random effects model.

Both the fixed-effects model and the random-effects model suggest that there is a small³ positive overall effect of pair programming on *Quality* compared with solo programming. Only one study showed a negative effect of pair programming on *Quality* (Domino *et al.* 2007). All the other studies showed mostly small to medium positive effects. The three studies by Domino *et al.* (2007), Arisholm *et al.* (2007), and Madeyski (2006) contribute more than 50% of the total weight in the meta-analysis for *Quality*. The one-study-removed analysis shows that the meta-analysis is most sensitive to the inclusion/exclusion of Williams *et al.* (2000). Heterogeneity is significant at a medium level ($Q = 35.97$; $p < .01$; $I^2 = 63.86\%$).

3.2.2. Effects of pair programming on Duration

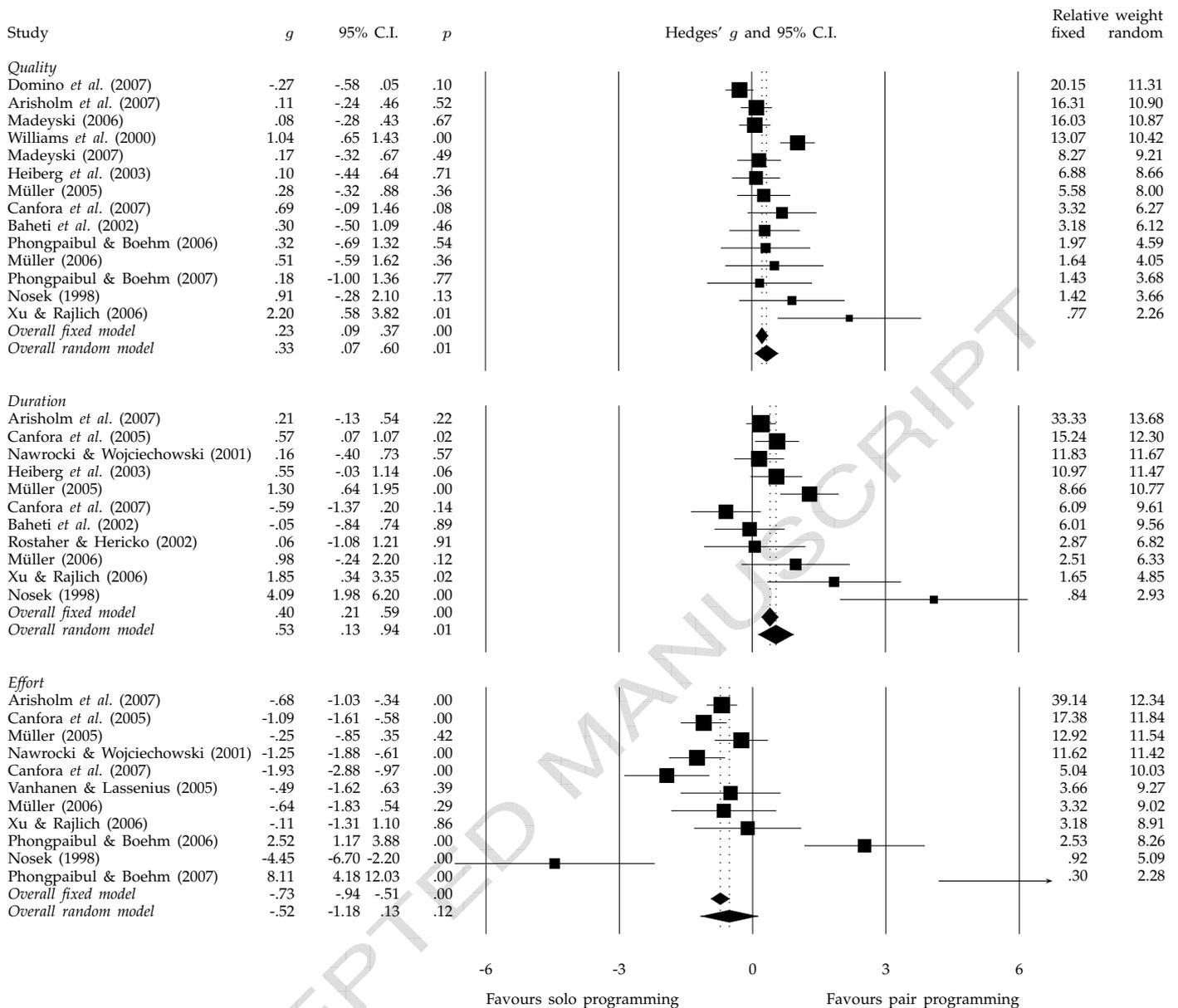
Eleven studies reported effects on *Duration* in a total of 21 tests. These studies used a total of 669 subjects. Both the fixed-effects model and the random-effects model suggest that there is a medium positive overall effect of pair programming on *Duration*.

Compared with *Quality*, the studies on *Duration* show a more mixed picture; two of the 11 studies show a negative effect, while the remaining nine show a positive effect. In addition, the smaller studies show medium to large contradictory effects. The three studies by Arisholm *et al.* (2007), Canfora *et al.* (2005), and Nawrocki and Wojciechowski (2001) contribute more than 60% of the total weight in the meta-analysis of *Duration*. The one-study-removed analysis shows that the meta-analysis is most sensitive to the inclusion/exclusion of Nosek (1998). Heterogeneity is significant at a medium level ($Q = 33.57$; $p < .01$; $I^2 = 70.21\%$).

3.2.3. Effects of pair programming on Effort

Eleven studies reported effects on *Effort* in a total of 18 tests. These studies used a total of 586 subjects. Both the fixed-effects model and the random-effects model suggest that there is a medium negative overall effect of pair programming on *Effort* compared with solo programming. All the included studies show a negative effect on effort, apart from the two studies by Phongpaibul and Boehm (2006,2007). However, the results of those studies are not directly comparable, because the researchers compared pair programming teams with teams of individu-

³ The effect size is small compared with the effect sizes reported in other software engineering experiments [30].

Fig. 1. Forest Plots for Meta-Analysis of *Correctness*, *Duration*, and *Effort*.

als who also performed inspections. The three studies by Arisholm *et al.* (2007), Canfora *et al.* (2005), and Müller (2005) contribute almost 70% of the total weight in the meta-analysis of *Effort*. The one-study-removed analysis shows that the meta-analysis is most sensitive to the inclusion/exclusion of either of Phongpaibul and Boehm (2006,2007). Heterogeneity is significant and high ($Q = 66.36$; $p < .01$; $I^2 = 84.93\%$).

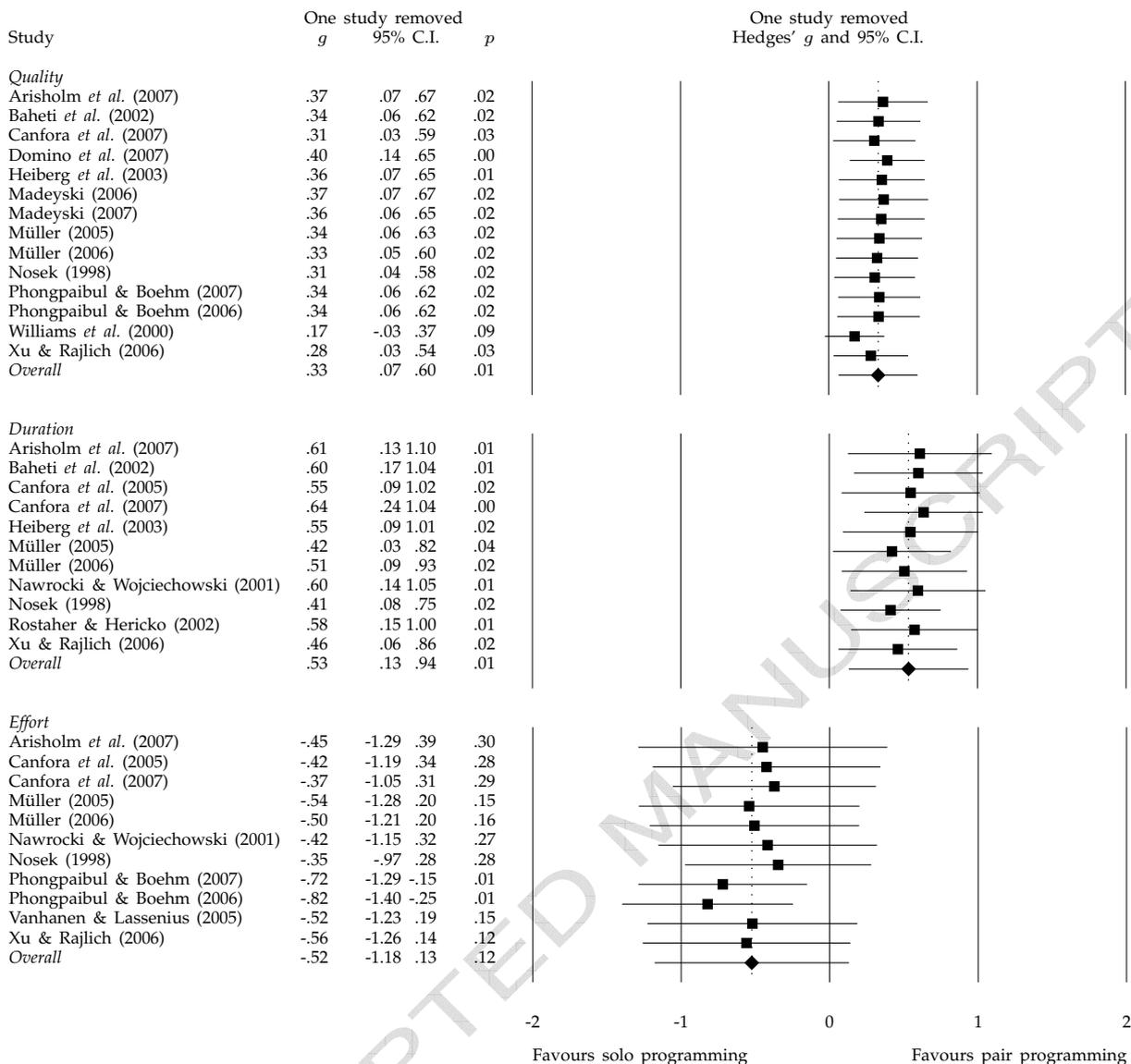
3.3. Subgroup analyses

Because of medium to high heterogeneity, we decided to conduct subgroup analyses as a step to identify possible immediate sources of heterogeneity. Two subgroup types stand out due to surface dissimilarities: the type of developers (students or professionals) and the type of comparison (isolated pairs vs. individuals, or teams of

pairs vs. teams of individuals). The results of these analyses are summarized in Table 2.

The most dramatic result of the subgroup analysis is the reversal of effect on *Effort* for the *Teams* subgroup (from $-.52$ in the overall analysis to $.99$ (non-significant) in the subgroup analysis, for the random-effects model). This is due to the two studies Phongpaibul and Boehm (2006,2007), which we remarked upon above. Without these two studies (*No teams* subgroup), the effect size increases to from $-.52$ to $-.85$ in the random-effects model.

Apart from this reversal, effect sizes in the subgroups remain in the same order of magnitude as in the overall analysis, except for the following cases: For *Quality*, the *No teams* subgroup reaches a medium level in the random model's estimate ($.38$). For *Duration*, the *Professionals* subgroup declines to a small level in the fixed model's estimate ($.16$). For *Effort*, the *Students* subgroup declines

Fig. 2. Forest Plots for One-Study-Removed Meta-Analysis of *Correctness*, *Duration*, and *Effort*.

to small (.04) (but non-significant, $p = .92$), and the *Professionals* subgroup increases to a large effect size (-1.99) in their respective random-effects models.

Heterogeneity levels remain medium to high. *Professionals* and *No teams* increase to high heterogeneity for *Duration* (82.22) and *Teams* decreases to small (0.00) (but non-significant, $p = .89$) for *Quality*.

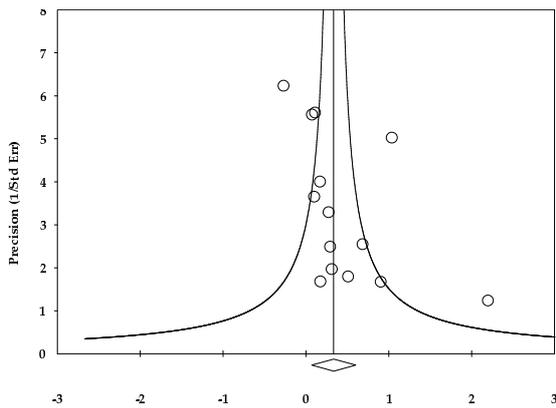
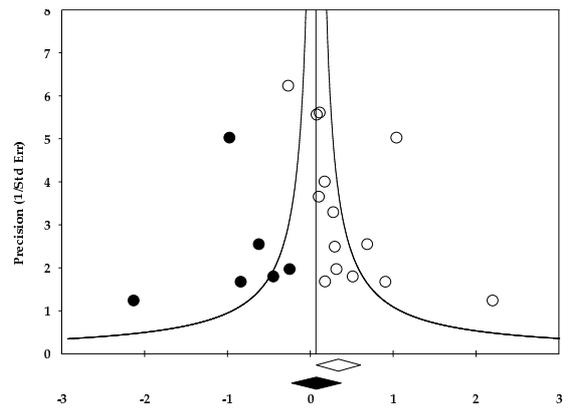
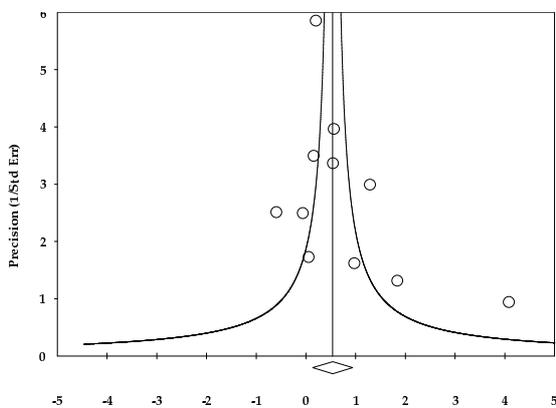
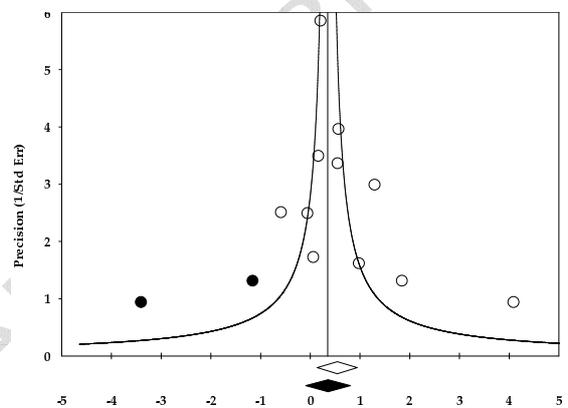
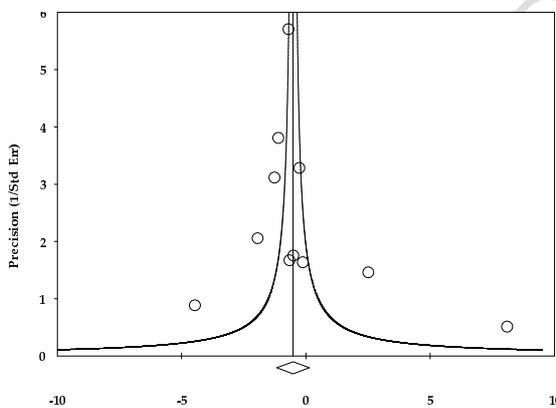
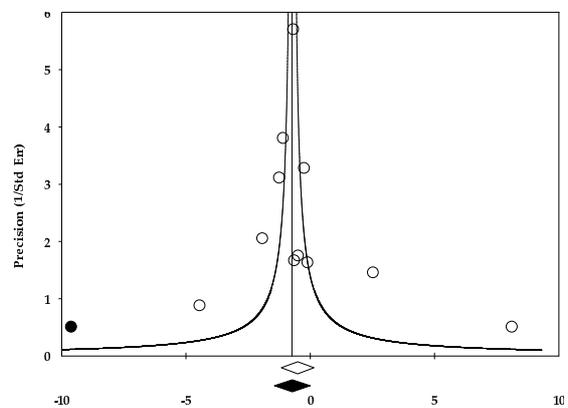
Significance levels for both effect size estimates and heterogeneity decrease dramatically in several instances in the subgroups, probably due to small k . This is particularly the case in the *Professionals* and *Teams* subgroups. Note however, that the effect on *Effort* actually turns significant in the *Professionals* subgroup. Note that the *Professionals* subgroup excludes Phongpaibul and Boehm (2006,2007) since they have student subjects.

3.4. Publication Bias

Publication bias captures the idea that studies that report relatively large treatment effects are the ones that were most likely to be initiated in the first place, and/or submitted and accepted for publication. The effect size estimated from a biased collection of studies will tend to overestimate the true effect. We used Comprehensive Meta-Analysis to assess the likely extent of publication bias, and its potential impact on the conclusions.

Figs. 3–5 show funnel plots for *Quality*, *Duration* and *Effort*, respectively. These graphs plot each study according to a measure of study size called *precision* (which is simply $1/\text{StdError}$) on the vertical axis, and according to effect size (Hedges' g) on the horizontal axis.

Hence, large studies tend to appear toward the top of the graph, and smaller studies tend toward the bottom of the graph. Since there is more inter-study variance in ef-

Fig. 3. Funnel Plot *Quality*.Fig. 6. Funnel Plot *Quality* Trim and Fill Fixed-Random.Fig. 4. Funnel Plot *Duration*.Fig. 7. Funnel Plot *Duration* Trim and Fill Fixed-Random.Fig. 5. Funnel Plot *Effort*.Fig. 8. Funnel Plot *Effort* Trim and Fill Fixed-Random.

fect size estimates among the smaller studies, these studies will be dispersed wider, while larger studies will tend to cluster near the meta-analytic estimate of effect size.

In the absence of publication bias, studies should be distributed symmetrically on either side of the meta-analytic estimate of effect size. In the presence of bias, the bottom part of the plot should show a higher concentration of studies on one side of the meta-analytic estimate than the other. This would reflect the fact that smaller studies are more likely to be published if they have larger than average effects, which makes them more likely to meet the criterion for statistical significance.

Various statistical measures can be used to complement the picture given by the funnel plots.

Duval and Tweedie's *trim and fill* method [13,14], takes the basic idea behind the funnel plot one step further and imputes (computes and inserts) assumed missing studies to obtain symmetry if the funnel plot is asymmetric and then recomputes the meta-analytic estimate of effect size. The method initially trims the asymmetric studies from the biased side to locate the unbiased effect (in an iterative procedure), and then fills the plot by re-inserting the trimmed studies on the biased side as well as their imputed counterparts on the opposite side of the meta-

analytic estimate of effect size. Figs. 6–8 show the plots of Figs. 3–5 with imputed studies (black dots) and adjusted meta-analytic estimates of effect size (black diamonds).

These analyses suggest that there is indeed some publication bias in our sample. For example, six studies are imputed to obtain symmetry for *Quality*. This might be taken as an incentive to publish more high-quality studies that, perhaps, exhibit low effects or opposite effects to those that are expected. This might include promoting gray literature to be readily accessible in journals and conferences.

However, it is important to note that publication bias may not be the only cause of funnel plot asymmetry. For instance, asymmetry may be caused by between-study heterogeneity, study quality, or other factors [44].

In the presence of heterogeneity (which we are assuming), there are three meaningful ways to impute the assumed missing studies. (1) A random-effects model is used to trim and fill, and then the adjusted meta-analytic estimate of effect size is calculated from the filled data using a random-effects model (a so-called *random-random* approach) [13,14]. However, in a random-effects model meta-analysis, smaller studies are given added weight in the synthesis. Thus, if publication bias exists, the meta-analytic estimate is likely to be more biased than that obtained from a fixed-effects model meta-analysis [44]. Thus, an alternative is: (2) A fixed-effects model is used to trim and fill, and then the adjusted meta-analytic estimate of effect size is calculated from the filled data using a fixed-effects model (*fixed-fixed*). It is recommended to report both fixed-fixed and random-random analyses [51]. The fixed-fixed approach may be unsatisfying if one wishes to view adjustments relative to the original random-effects model meta-analysis. (3) A third approach is to use a fixed-effects model to trim and fill the meta-analysis, but a random-effects model to calculate the meta-analytic effect size estimate from the filled data (*fixed-random*). The fixed-effects trim and fill process is less likely to be influenced by any publication bias than the random-effects model approach, but the resulting estimate from the random effects model is likely to be more conservative than if a fixed effects model is used [44]. The plots in Figs. 6–8 were generated using the fixed-random approach.

Table 3 presents the results from all perspectives. For example, for *Correctness* a fixed-effects model approach to trim and fill yields 6 imputed studies and adjusts the estimate in the fixed-effects model from .23 to .03 (fixed-fixed), while it adjusts the estimate in the random-effects model from .33 to .07 (fixed-random).

The trim and fill procedure searches for missing studies on one side at a time. For *Correctness* and *Duration* missing studies were found on the left (Table 3), i.e., in the opposite direction of the estimated effect, which is consistent with the assumption of publication bias. For *Effort*, the fixed approach found a missing study on the left, while the random approach found missing studies

Table 3
Trim and Fill Analysis

	Imputed	fixed			random			Q
		g	95% C.I.		g	95% C.I.		
<i>Correctness</i>								
Observed		.23	.09	.37	.33	.07	.60	35.97
fixed	6	.03	-.09	.16	.07	-.23	.37	81.42 left
random	0	.23	.09	.37	.33	.07	.60	35.97 left
<i>Duration</i>								
Observed		.40	.21	.59	.53	.13	.94	33.57
fixed	2	.34	.15	.54	.35	-.10	.80	49.77 left
random	1	.37	.18	.56	.44	-.01	.88	45.01 left
<i>Effort</i>								
Observed		-.73	-.94	-.51	-.52	-1.18	.13	66.36
fixed	1	-.75	-.97	-.54	-.74	-1.46	-.01	86.01 left
random	2	-.58	-.78	-.37	-.07	-.79	.65	104.46 right

on the right. The missing study on the left is not consistent with the assumption of publication bias and is due to the overly large effort estimate of one of the studies in the opposite direction of the estimate.

Sterne and Egger [50] note that trim and fill plots merely detect a relationship between sample size and effect size, not a causal mechanism between the two. The effect size may be larger in small studies due to publication bias. However, it is also possible that small studies estimate a different part of the the population effect size distribution than do large studies, e.g., because smaller studies use different operationalizations of the outcome constructs, or have different situational variables than large studies. Because one does not know for certain whether funnel plot asymmetry is really caused by publication bias, it is recommended to use the trim and fill method mainly as a form of sensitivity analysis [44].

3.5. Moderating effects of task complexity and expertise

Due to the interdisciplinary and complex nature of industrial software engineering, it is usually not reasonable to test an hypothesis that considers only one independent (predictor) variable and one dependent (criterion) variable. Hypotheses related to software engineering should typically include additional variables and test more complex relationships in order to provide a more accurate description of reality. Indeed, the relatively small overall effects and large between-study variance (heterogeneity) indicate that one or more moderator variables might play a significant role.

Only two of the studies in the review, Vanhanen and Lassenius (2005) and Arisholm *et al.* (2007), tested explicitly for moderator effects, while only one of the other studies discussed the potential influence of such effects; Williams *et al.* (2000) suggested that the relative improvement by pairs after the first programming task in their experiment was due to “pair jelling”.

The literature on group dynamics (e.g., [5,21]) suggests that the extent to which group performance exceeds that of individuals, and the mechanisms by which such gains in performance may be achieved, depend upon the com-

position of the group and the characteristics of the tasks. Vanhanen and Lassenius (2005) found that task complexity did not affect the differences in effort between solo and pair programming; the Pearson correlation between task complexity and effort difference was as low as $r = -0.02$.

On the other hand, Arisholm *et al.* (2007) found moderating effects of both task complexity and expertise. The results are shown in Fig. 9. Overall, the results showed that the pairs had an 8% decrease in duration ($g = 0.21$) with a corresponding 84% increase in effort ($g = -0.68$) and a 7% increase in correctness ($g = 0.11$) (Fig. 9a). However, the main effects of pair programming were masked by the moderating effect of system complexity, in that simpler designs had shorter duration, while more complex designs had increased correctness (Fig. 9e).

Furthermore, when considering the moderating effect of programmer expertise, junior pairs had a small (5%) increase in duration and thus a large increase in effort (111%), and a 73% increase in correctness (Fig. 9b). Intermediate pairs had a 28% decrease in duration (43% increase in effort) and a negligible (4%) increase in correctness (Fig. 9c). Senior pairs had a 9% decrease in duration (83% increase in effort) and an 8% decrease in correctness (Fig. 9d). Thus, the juniors benefited from pair programming in terms of increased correctness, the intermediates in terms of decreased duration, while there were no overall benefits of pair programming for seniors. When considering the combined moderating effect of system complexity and programmer expertise on pair programming, there appears to be an interaction effect: Among the different treatment combinations, junior pairs assigned to the complex design had a remarkable 149% increase on correctness compared with individuals (Fig. 9f). Furthermore, intermediates and seniors experienced an effect of pair programming on duration on the simpler design, with a 39% (Fig. 9g) and 23% (Fig. 9h) decrease, respectively. However, the cost of this shorter duration was a corresponding decrease in correct solutions by 29% and 13%, respectively.

4. Discussion

The unconditional part of the meta-analysis suggests that the population effect size mean μ for each of the outcome constructs are small to medium.

More interestingly, the conditional part of the meta-analysis showed large and partly contradictory differences in the reported overall effects of pair programming, specifically with respect to *Duration* and *Effort*. Our subgroup analyses do not suggest that these differences are due to differences in samples (e.g., students or professionals). Differences in organization (e.g., teams or no teams) were apparent from the subgroup analysis, but these differences were due to two particular studies that, when removed, did not decrease overall heterogeneity.

At the outset, we anticipated between-study variance

(heterogeneity) due to moderating factors such as expertise and task complexity and to the fact that the outcome constructs of the meta-analysis are operationalized by indicators that are different aspects of the constructs. Such heterogeneity was taken into account by using a random-effects model for the unconditional interpretation.

Nevertheless, contradictory differences still manifested themselves. Thus, it seems clear that moderating factors play an important role and should be investigated further. The study in [1] corroborates this conclusion for expertise and task complexity, but other moderating factors, such as amount of training in pair programming, motivation, team climate, etc., are likely to be relevant as well.

Consequently, the question of whether pair programming is better than solo programming is not precise enough to be meaningful, since the answer to that question in the present context is both “yes” and “no”. On the basis of the evidence from this review, the answer is that “it depends”: It depends on other factors, for example, the expertise of the programmers and on the complexity of the system and tasks to be solved.

In the literature, expertise and task complexity are perhaps the most central situation-independent predictors of performance. (Situation-dependent factors, on the other hand, include more dynamic factors such as motivation, team climate, organizational issues, etc.) Theory predicts that experts perform better on complex tasks than do novices because experts’ level of understanding corresponds to the deep structure [10,8,9,19] of a complex task. That is, experts perceive (objectively) complex tasks as subjectively less complex, i.e., less *variable* and more *analyzable* [43,38]. Conversely, experts perceive (objectively) less complex tasks as more *variable* and less *analyzable* since such tasks do not match the expert’s deep-level understanding of the problem. Novices, on the other hand, do have an understanding that matches the surface structure of a non-complex task, and are expected to do better on non-complex tasks than on complex tasks. Moreover, they may be expected to even outperform experts on non-complex tasks [23].

These effects are evident in [1], in which the levels of correctness for individual juniors, intermediates, and seniors on the non-complex system were 63%, 87%, and 86%, respectively, whereas the levels on the complex system were 34%, 41%, and 81%, respectively. Thus, the performance drop was much higher for juniors than for intermediates and seniors, when moving from the non-complex to the complex system, although juniors did not outperform higher expertise groups on the non-complex system, see also [23].

In our context, one of the most interesting observations is that the pairing up of individuals seems to elevate the junior pairs up to near senior pair performance. Thus, pair collaboration might compensate for juniors’ lack of deep understanding, for example, by inducing an expert-like strategy.

Change tasks rely heavily on program comprehension.

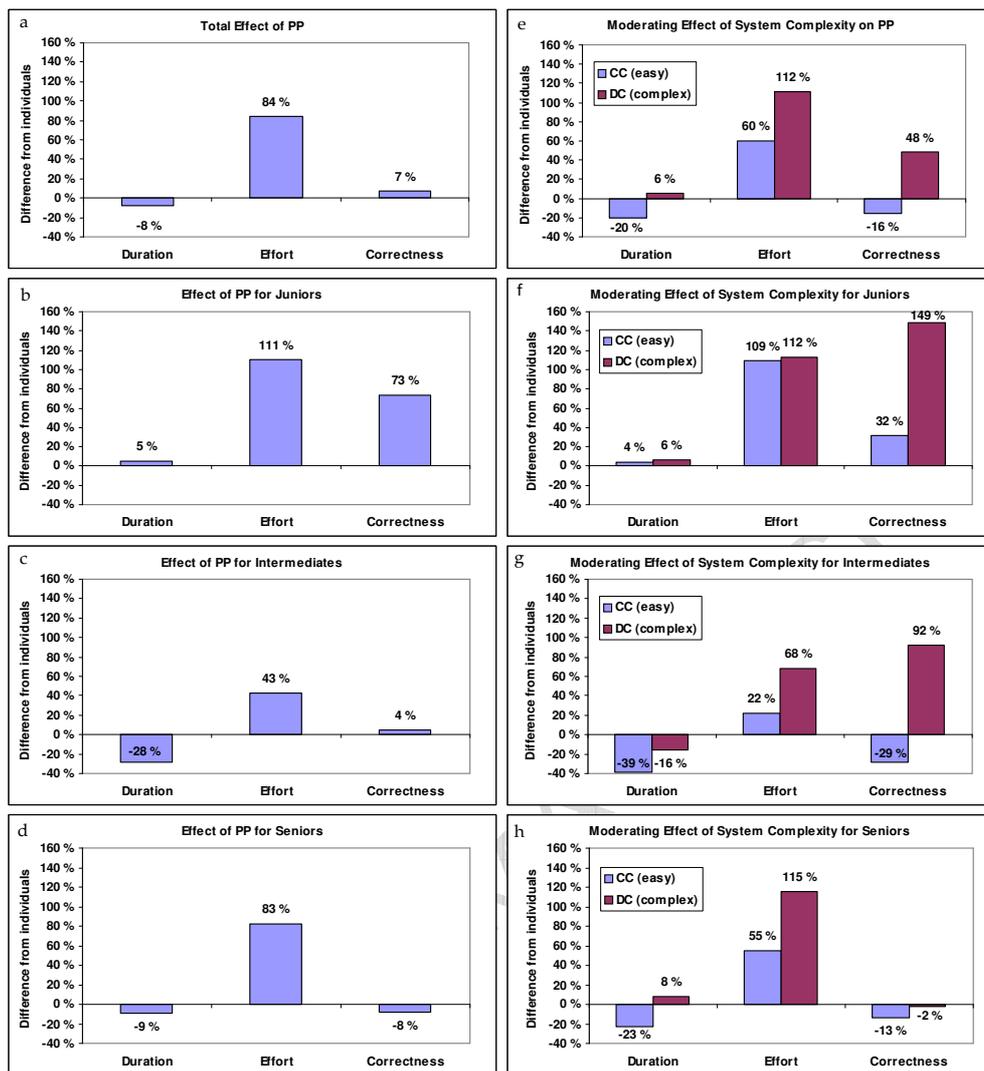


Fig. 9. The moderating effects of programmer expertise (a-d) and system complexity (e-h) on the relation of pair programming on duration, effort, and correctness (Arisholm *et al.* 2007).

To comprehend code, experts use a top-down model [47] that short-cuts the available information by only investigating details as dictated by domain knowledge. This approach to comprehension is more efficient than bottom-up comprehension, which builds understanding from details, and which is the approach found to be used by programmers encountering totally unfamiliar code [41,42].

By forcing junior peers to rationalize their ideas to each other, junior pairs might adopt a top-down strategy to comprehension, rather than getting lost in code on their own. The mere act of thinking aloud whilst solving problems has been shown to increase performance, when the verbalization is intended to reason or explain action e.g., [7] (Type 3 verbalization, in Ericsson and Simon's terminology) [20]. Several studies have concluded that apparent successes of pair programming are not due to the particularities of pair programming (such as the specific roles of driver and navigator), but rather to the sheer amount of verbalization that the pair programming situation ne-

cessitates [11,22].

Group performance not only relies on task complexity but also on the collaborative nature of a task. In fact, the appropriateness of each of the familiar adages "two heads are better than one", "many hands make light work", and "a chain is only as strong as its weakest link" depends on whether a task is additive, compensatory, disjunctive or conjunctive [48]. For example, the chain analogy is appropriate for conjunctive tasks, where all group members must contribute to the solution, but is inappropriate for disjunctive tasks for which it suffices that one group member has the ability to complete the task. It is not obvious what sort of task pair programming is in this respect.

The precise collaborative nature of pair programming also influences what social mechanisms (social loafing, social labouring, social facilitation, social inhibition, social compensation, etc) are applicable. However, these social mechanisms also depend on a host of other factors. In a meta-analysis of social loafing (the phenomenon that

individuals tend to expend less effort when working collectively than when working individually), Karau and Williams [32] identified several conditions in which such loafing is eliminated (e.g., by high group cohesion) and some in which the opposite phenomenon, social loafing [5], could be observed (i.e., greater effort on group tasks). Social loafing seems to occur when complex or highly involving tasks are performed, or when the group is considered important for its members, or if the prevailing values favor collectivism rather than individualism [5].

5. Conclusion

Our meta-analysis suggests that pair programming is not uniformly beneficial or effective, that inter-study variance is high, and that perhaps publication bias is an issue. Hence, if further investigations are to be undertaken on pair programming, then the focus should be on untangling the moderating factors of the effect of pair programming.

However, with respect to the central factors expertise and task complexity, the current state of knowledge suggests that pair programming is beneficial for achieving correctness on highly complex programming tasks. Pair programming may also have a time gain on simpler tasks. By cooperating, programmers may complete tasks and attain goals that would be difficult or impossible if they worked individually. Junior pair programmers, for example, seem able to achieve approximately the same level of correctness in about the same amount of time (duration) as senior individuals. However, the higher quality for complex tasks comes at a price of a considerably higher effort (cost), while the reduced completion time for the simpler tasks comes at a price of a noticeably lower quality. This fact confirms Voas' [52] contention that you cannot expect faster and better and cheaper. These relationships give rise to a few evidence-based guidelines for the use of pair programming for professional software developers. If you do not know the seniority or skill levels of your programmers, but do have a feeling for task complexity, then employ pair programming either when task complexity is low and time is of the essence, or when task complexity is high and correctness is important.

In the future, we intend to investigate deeper into the theoretical and empirical underpinnings of collaboration in pairs, e.g., by studying group dynamics and analyzing pair dialogues to obtain insights into subjects' learning and reasoning processes. Only by understanding what makes pairs work, and what makes pairs less efficient, can steps be taken to provide beneficial conditions for work and to avoid detrimental conditions; or to avoid pairing altogether when beneficial conditions cannot be provided.

Acknowledgements

The authors are grateful to the anonymous referees for insightful remarks.

Appendix

The following articles were included in the meta-analysis:

- Arisholm *et al.* (2007). E. Arisholm, H. Gallis, T. Dybå, and D.I.K. Sjøberg, "Evaluating pair programming with respect to system complexity and programmer expertise," *IEEE Trans. Software Eng.*, vol. 33, pp. 65–86, Feb. 2007.
- Baheti *et al.* (2002). P. Baheti, E. Gehringer, and D. Stotts, "Exploring the efficacy of distributed pair programming," in *Proc. XP/Agile Universe 2002*, ser. Lecture Notes in Computer Science, vol. 2418. Springer Verlag, 2002, pp. 208–220.
- Canfora *et al.* (2007). G. Canfora, A. Cimitlie, F. Garcia, M. Piattini, and C.A. Visaggio, "Evaluating performances of pair designing in industry," *J. Systems and Software*, vol. 80, no. 8, pp. 1317–1327, Aug. 2007.
- Canfora *et al.* (2005). G. Canfora, A. Cimitlie, and C.A. Visaggio, "Empirical study on the productivity of the pair programming," in *Proc. XP 2005*, ser. Lecture Notes in Computer Science, vol. 3556. Springer Verlag, 2005, pp. 92–99.
- Domino *et al.* (2007). M.A. Domino, R. Webb Collins, and A.R. Hevner, "Controlled experimentation on adaptations of pair programming," *Information Technology and Management*, vol. 8, no. 4, pp. 297–312, Dec. 2007.
- Heiberg *et al.* (2003). S. Heiberg, U. Puus, P. Salumaa, and A. Seeba, "Pair-programming effect on developers productivity," in *Proc. XP 2003*, ser. Lecture Notes in Computer Science, vol. 2675. Springer Verlag, 2003, pp. 215–224.
- Madeyski (2006). L. Madeyski, "The impact of pair programming and test-driven development on package dependencies in object-oriented design—an experiment," in *Proc. PROFES 2006*, ser. Lecture Notes in Computer Science, vol. 4034. Springer Verlag, 2006, pp. 278–289.
- Madeyski (2007). L. Madeyski, "On the effects of pair programming on thoroughness and fault-finding effectiveness of unit tests," in *Proc. PROFES 2007*, ser. Lecture Notes in Computer Science, vol. 4589. Springer Verlag, 2007, pp. 207–221.
- Müller (2005). M. Müller, "Two controlled experiments concerning the comparison of pair programming to peer review," *J. Systems and Software*, vol. 78, no. 2, pp. 169–179, 2005.
- Müller (2006). M. Müller, "A preliminary study on the impact of a pair design phase on pair programming and solo programming," *Information and Software Technology*, vol. 48, no. 5, pp. 335–344, May 2006.
- Nawrocki and Wojciechowski (2001). J. Nawrocki and A. Wojciechowski, "Experimental evaluation of pair programming," in *Proc. European Software Control and Metrics Conference (ESCOM'01)*, 2001, pp. 269–276.
- Nosek (1998). J. Nosek, "The case for collaborative programming," *Comm. ACM*, vol. 41, no. 3, pp. 105–108, Mar. 1998.
- Phongpaibul and Boehm (2006). M. Phongpaibul and B. Boehm, "An empirical comparison between pair development and software inspection in Thailand," in *Proc. Int'l Symposium on Empirical Software Engineering (ISESE'06)*, 2006, pp. 85–94.
- Phongpaibul and Boehm (2007). M. Phongpaibul and B. Boehm, "A replicate empirical comparison between pair development and software development with inspection," in *Proc. 1st Int'l Symp. Empirical Software Engineering and Measurement (ESEM'07)*. IEEE Computer Society, 2007, pp. 265–274.
- Rostaher and Herick (2002). M. Rostaher and M. Hericko, "Tracking test first pair programming—an experiment," in *Proc. XP/Agile*

- Universe 2002*, ser. Lecture Notes in Computer Science, vol. 2418. Springer Verlag, 2002, pp. 174–184.
- Vanhanen and Lassenius (2005). J. Vanhanen and C. Lassenius, “Effects of pair programming at the development team level: An experiment,” in *Proc. Int’l Symposium on Empirical Software Engineering (ISESE’05)*, 2005, pp. 336–345.
- Williams *et al.* (2000). L. Williams, R.R. Kessler, W. Cunningham, and R. Jeffries, “Strengthening the case for pair programming,” *IEEE Software*, vol. 17, no. 4, pp. 19–25, 2000.
- Xu and Rajlich (2006). S. Xu and V. Rajlich, “Empirical validation of test-driven pair programming in game development,” in *Proc. Int’l Conference on Computer and Information Science and Int’l Workshop on Component-Based Software Engineering, Software Architecture and Reuse (ICIS-COM SAR’06)*, 2006, pp. 500–505.
- ## References
- [1] E. Arisholm, H. Gallis, T. Dybå, and D.I.K. Sjøberg, “Evaluating pair programming with respect to system complexity and programmer expertise,” *IEEE Trans. Software Eng.*, vol. 33, pp. 65–86, Feb. 2007.
- [2] K. Beck and C. Andres, *Extreme Programming Explained: Embrace Change*, 2nd ed. Addison-Wesley, 2003.
- [3] K. Bollen and R. Lennox, “Conventional wisdom on measurement: A structural equation perspective,” *Psychological Bull.*, vol. 110, no. 2, pp. 305–314, 1991.
- [4] M. Borenstein, L.V. Hedges, and H. Rothstein, “Meta-analysis—fixed effects versus random effects,” *Biostat Inc.*, Tech. Rep., 2007. [Online]. Available: meta-analysis.com/downloads/Meta-analysis_fixed_effect_vs_random_effects.pdf
- [5] R. Brown, *Group Processes: Dynamics within and between Groups*, 2nd ed. Blackwell, 2000.
- [6] C. Burke jarvis, S.B. Mackenzie, and P.M. Podsakoff, “A critical review of construct indicators and measurement model misspecification in marketing and consumer research,” *J. Consumer Research*, vol. 30, pp. 199–218, Sept. 2003.
- [7] M.T.H. Chi, N. de Leeuw, M.H. Chiu, and C. LaVancher, “Eliciting self-explanations improves understanding,” *Cognitive Science*, vol. 18, pp. 439–477, 1994.
- [8] M.T.H. Chi, P.J. Feltovich, and R. Glaser, “Categorization and representation of physics problems by experts and novices,” *Cognitive Science*, vol. 5, pp. 121–152, 1981.
- [9] M.T.H. Chi, R. Glaser, and E. Rees, “Expertise in problem solving,” in *Advances in the Psychology of Human Intelligence*, R. Sternberg, Ed. Lawrence Erlbaum Associates, Inc., 1982, pp. 17–76.
- [10] N. Chomsky, *Syntactic Structures*, 2nd ed. Mouton, 2002.
- [11] J. Chong and T. Hurlbutt, “The social dynamics of pair programming,” in *Proc. 29th Int’l Conf. Software Engineering*, 2007.
- [12] A. Diamantopoulos and J.A. Siguaw, “Formative versus reflective indicators in organizational measure development: A comparison and empirical illustration,” *British J. Management*, vol. 17, pp. 263–282, 2006.
- [13] S. Duval and R. Tweedie, “A nonparametric “trim and fill” method of accounting for publication bias in meta-analysis,” *J. American Statistical Association*, vol. 95, no. 445, pp. 89–98, 2000.
- [14] S. Duval and R. Tweedie, “Trim and fill: A simple funnel-plot-based method of testing and adjusting for publication bias in meta-analysis,” *Biometrics*, vol. 56, no. 2, pp. 455–463, June 2000.
- [15] T. Dybå, E. Arisholm, D.I.K. Sjøberg, J.E. Hannay, and F. Shull, “Are two heads better than one? on the effectiveness of pair programming,” *IEEE Software*, pp. 12–15, November/December 2007.
- [16] T. Dybå, T. Dingsøy, and G.K.S. Hanssen, “Applying systematic reviews to diverse study types: An experience report,” in *Proc. 1st Int’l Symp. Empirical Software Engineering and Measurement (ESEM’07)*. IEEE Computer Society, 2007, pp. 225–234.
- [17] T. Dybå, V.B. Kampenes, and D.I.K. Sjøberg, “A systematic review of statistical power in software engineering experiments,” *J. Information & Software Technology*, vol. 48, 2006.
- [18] T. Dybå, B.A. Kitchenham, and M. Jørgensen, “Evidence-based software engineering for practitioners,” *IEEE Software*, vol. 22, pp. 58–65, 2005.
- [19] K.A. Ericsson and N. Charness, “Expert performance—its structure and acquisition,” *American Psychologist*, vol. 49, pp. 725–747, 1994.
- [20] K.A. Ericsson and H.A. Simon, *Protocol Analysis*, revised ed. The MIT Press, 1993.
- [21] D.R. Forsyth, *Group Dynamics*, 4th ed. Thomson Wadsworth, 2006.
- [22] S. Freudenberg (née Bryant), P. Romero, and B. du Boulay, “‘talking the talk’: Is intermediate-level conversation the key to the pair programming success story?” in *Proc. AGILE 2007*, 2007.
- [23] T. Hærem and D. Rau, “The influence of degree of expertise and objective task complexity on perceived task complexity and performance,” *J. Applied Psychology*, vol. 92, no. 5, pp. 1320–1331, 2007.
- [24] J.E. Hannay and M. Jørgensen, “The role of deliberate artificial design elements in software engineering experiments,” *IEEE Trans. Software Eng.*, vol. 34, pp. 242–259, mar/apr 2008.
- [25] J.E. Hannay, D.I.K. Sjøberg, and T. Dybå, “A systematic review of theory use in software engineering experiments,” *IEEE Trans. Software Eng.*, vol. 33, pp. 87–107, Feb. 2007.
- [26] L.V. Hedges and I. Olkin, *Statistical Methods for Meta-Analysis*. Academic Press, Inc., 1985.
- [27] L.V. Hedges and J.L. Vevea, “Fixed- and random-effects models in meta-analysis,” *Psychological Methods*, vol. 3, no. 4, pp. 486–504, 1998.
- [28] J.P. Higgins, S.G. Thompson, J.J. Deeks, and D.G. Altman, “Measuring inconsistency in meta-analyses,” *British Medical Journal*, vol. 327, no. 4, pp. 557–560, 2003.
- [29] J.E. Hunter and F.L. Schmidt, “Fixed effects vs. random effects meta-analysis models: Implications for cumulative research knowledge,” *International J. Selection and Assessment*, vol. 327, no. 8, pp. 272–292, Dec. 2000.
- [30] V.B. Kampenes, T. Dybå, J.E. Hannay, and D.I.K. Sjøberg, “A systematic review of effect size in software engineering experiments,” *Information and Software Technology*, vol. 49, no. 11–12, pp. 1073–1086, Nov. 2007.
- [31] V.B. Kampenes, T. Dybå, J.E. Hannay, and D.I.K. Sjøberg, “A systematic review of quasi-experiments in software engineering,” *Information and Software Technology*, vol. 49, no. 11–12, pp. 1073–1086, Nov. 2007.
- [32] S.J. Karau and K.D. Williams, “Social loafing: A meta-analytic review and theoretical integration,” *J. Personality and Social Psychology*, vol. 65, no. 4, pp. 681–706, 1993.
- [33] N.L. Kerr and R.S. Tindale, “Group performance and decision making,” *Annual Review of Psychology*, vol. 55, pp. 623–655, 2004.
- [34] B.A. Kitchenham and S. Charters, “Guidelines for performing systematic literature reviews in software engineering,” Keele University, EBSE Technical Report, EBSE-2007-01, Tech. Rep., 2007.
- [35] J.M. Levine and R.L. Moreland, “Progress in small group research,” *Annual Review of Psychology*, vol. 41, pp. 585–634, 1990.
- [36] R.M. Lindsay and A.S.C. Ehrenberg, “The design of replicated studies,” *The American Statistician*, vol. 47, no. 3, pp. 217–228, Aug. 1993.
- [37] M.W. Lipsey and D.B. Wilson, *Practical Meta-Analysis*. Sage, 2001.
- [38] J.G. March and H.A. Simon, *Organizations*, 2nd ed. Wiley-Blackwell, 1993.
- [39] J. Miller, “Applying meta-analytical procedures to software engineering experiments,” *J. Systems and Software*, vol. 54, no. 1–2, pp. 29–39, 2000.
- [40] J.C. Nunnally and I. Bernstein, *Psychometric Theory*, 3rd ed. McGraw-Hill, Inc., 1994.
- [41] N. Pennington, “Comprehension strategies in programming,” in *Proc. Second Workshop Empirical Studies of Programmers*, 1987, pp. 100–113.

- [42] N. Pennington, "Stimulus structures and mental representations in expert comprehension of computer programs," *Cognitive Psychology*, vol. 19, pp. 295–341, 1987.
- [43] C. Perrow, "A framework for the comparative analysis of organizations," *American Sociological Review*, vol. 32, pp. 194–208, 1967.
- [44] J.L. Peters, A.J. Sutton, D.R. Jones, K.R. Abrams, and L. Rushton, "Performance of the trim and fill method in the presence of publication bias and between-study heterogeneity," *Statistics in Medicine*, vol. 26, pp. 4544–4562, 2007.
- [45] R. Rosenthal and M.R. DiMatteo, "Meta-analysis: Recent developments in quantitative methods for literature reviews," *Annual Review of Psychology*, vol. 52, pp. 59–82, 2001.
- [46] D.I.K. Sjøberg, J.E. Hannay, O. Hansen, V.B. Kampenes, A. Karahasanović, N.K. Liborg, and A.C. Rekdal, "A survey of controlled experiments in software engineering," *IEEE Trans. Software Eng.*, vol. 31, no. 9, pp. 733–753, Sept. 2005.
- [47] E. Soloway, B. Adelson, and K. Ehrlich, "Knowledge and processes in the comprehension of computer programs," in *The Nature of Expertise*, M. Chi, R. Glaser, and M. Farr, Eds. Lawrence Erlbaum Assoc., 1988, pp. 129–152.
- [48] I.D. Steiner, *Group Process and Productivity*. New York and London: Academic Press, 1972.
- [49] M. Stephens and D. Rosenberg, *Extreme Programming Refactored: The Case Against XP*. APress, 2003.
- [50] J.A.C. Sterne and M. Egger, "Regression methods to detect publication and other bias in meta-analysis," in *Publication Bias in Meta-analysis: Prevention, Assessment and Adjustments*, H.R. Rothstein, A.J. Sutton, and M. Borenstein, Eds. Wiley, 2005, ch. 6, pp. 99–110.
- [51] A.J. Sutton, "Evidence concerning the consequences of publication and related biases," in *Publication Bias in Meta-analysis: Prevention, Assessment and Adjustments*, H.R. Rothstein, A.J. Sutton, and M. Borenstein, Eds. Wiley, 2005, ch. 10, pp. 175–192.
- [52] J. Voas, "Faster, better, and cheaper," *IEEE Software*, vol. 18, no. 3, pp. 96–99, May/June 2001.
- [53] L. Williams and R.R. Kessler, *Pair Programming Illuminated*. Addison-Wesley, 2002.