

How to Write a Great Guided Research

And why should I do it?

Dr. Roman Haas

With material from Dr. Elmar Juergens

In close cooperation with the Academic Advisors at TUM Computer Science

2011 – 2017



2017 – now



Research collaboration
with Prof. Pretschner

thesisguide.org

- Slides
- Video
- Detailed Essays
- FAQ



THESIS GUIDE

[START HERE](#)

[PREFACE](#)

[CONTENT](#)

[CONTRIBUTE](#)

[ABOUT ME](#)



Agenda

1. Motivation
2. Preparation
3. Doing the work

Guided Research



```
graph TD; A[Guided Research] --> B[Guidance]; A --> C[Your own (small) research project];
```

- Guidance
 - Supervisor has research experience, helps you on your way
 - Examiner must be from TUM Informatics or affiliated with the CIT
- Your own (small) research project
 - Related Work
 - Implementation?
 - Proof?
 - Evaluation?
- Document and present your work
- Insights into real scientific work

Guided Research

- Voluntary
- 6 months, 10 ECTS
- Effort comparable to a more labor-intensive lab course
- Approx. 40 students/semester

Master's Thesis

- Mandatory
- 6 months, 30 ECTS
- Full-Time
- Approx. 100 students/semester

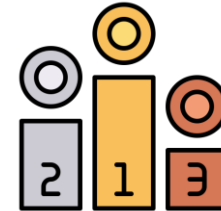
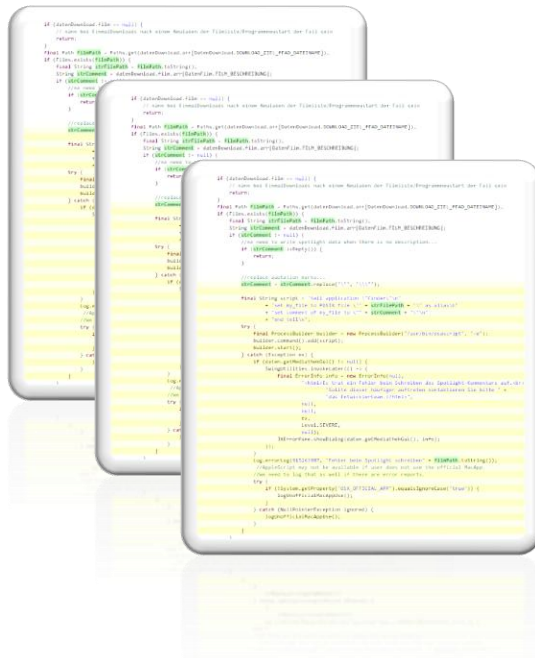
Less Formal than a Thesis

- Written document is „just“ a scientific report on your results (8-12 pages in English) which you need to send to your supervisor/examinor only
- You have to present your work
 - At the chair
 - Or at a „scientific event“

There are some formalia, though...

- You have to be enrolled in a Master's program (Informatics, Data Engineering & Analytics, Information Systems, Games Engineering, Robotics)
- Registration can be done anytime [online](#)
- Submission no later than the first lecture week of the next semester (6 months duration)
- Cannot be extended
- No transfer of credits, you need an internal examiner (with whom you may work together abroad)

Learning to Rank Extract Method Refactoring Suggestions for Long Methods



Result

which is described in more detail by Aravamdan and Nekrashev [28](#), and measures the goodness of the ranking list (obtained by the application of the scoring function). Mistakes in the top-most ranks have a bigger impact on the scoring function value. This is useful and important to us because we will not suggest all possible refactoring candidates, but only the highest-ranked ones. Given a long method, m , with refactoring candidates, C ; suppose the m is the ranking list on C , and p_0 the set of manually determined grades, then the DCG at position k is defined as $DCG(k) = \sum_{i=1, i \in C} G(i) / D(p_0(i))$, where $G(i)$ is an exponential gain function, $D(i)$ is a position discount function, and $\pi(i)$ is the position of refactoring candidate, $c_{i,j}$, in p_0 . We set $G(i) = 2^{i-1} - 1$ and $D(i) = \frac{1}{\log(1 + i)}$. To normalize the DCG, and to make it comparable with measures of other long methods, we divide this DCG by the DCG that a perfect ranking would have obtained. Therefore, the NDCG of a candidate ranking will always be in $[0, 1]$, where the NDCG of 1 can only be obtained by perfect rankings. In our evaluation, we consider the NDCG value of the last position so that all ranks are taken into account. See Jiang [46](#) for further details.

1.3 Approach

We discuss our approach to improve the scoring function in order to find the best suggestions for extract method refactoring.

1.3.1 Extract Method Refactoring Candidates

In our previous work [46](#), we presented an approach to derive extract method refactoring suggestions automatically for long methods. The main steps are: generating valid extract method refactoring candidates, ranking the candidates, and pruning the candidate list.

In the following, a *refactoring candidate* is a sequence of statements that can be extracted from a method into a new method. The *remainder* is the method that contains all the statements from the original method after applying the refactoring, plus the call of the extracted method. The suggested refactorings will help to improve the readability of the code and reduce its complexity, because these are main reasons for developers to initiate code refactoring [48](#).

We derived refactoring candidates from the control and data flow graph of a method using the Continuous Quality Assessment Toolchain (ConQAT) open source software. We filtered out all invalid candidates, that is those that violate preconditions that need to be fulfilled for extract method refactoring (for details, see [46](#)). The second step of our approach was to rank the valid

¹ www.conqat.org

US-3, whereas for SVM-rank it is 0.0790. Therefore, the scoring function based by ListMLE performed better than the scoring function based by SVM-rank.

Table 1.2: Coefficients of Variation for Learned Coefficients

Method	CV
ListMLE	0.0567
SVM-rank	0.0513

RQ2: How stable are the learned scoring functions?

Table [1.2](#) shows the average, minimum and maximum coefficients of variation (CV) for the learned coefficients for ListMLE and for SVM-rank. Small CV indicates that the results from the single run of the 10-folds fold procedure did not vary a lot, whereas big CVs indicate big differences between the learned coefficients. As the CVs of the single features from ListMLE are much smaller than those of SVM-rank, the ranking results from ListMLE are much more stable compared with SVM-rank. SVM-rank scores with a big variance between the single iterations of the validation process; that is, despite the heavy overlapping of the training sets, the learned coefficients vary a lot and can hardly be generalized.

RQ3: Can the scoring function be simplified?

Figure [1.4](#) shows a plot of the averaged NDCG measure for all 12 runs. Remember that we actually had three length measures, and we considered the absolute and the relative values for all of them. As the reduction of the number of statements led to a higher NDCG for ListMLE (which outperformed SVM-rank with respect to NDCG), we chose to use it as our length measure. In practice, that seems sensible since, while LoC also count empty and commented lines, the number of statements only counts real code.

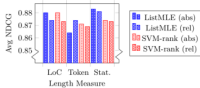


Fig. 1.4: Averaged NDCG When Considering Only One Length Measure

by filtering out very similar candidates, in order to obtain essentially different suggestions.

In the present paper, we focus on the ranking of candidates, and especially on the scoring function that defines that ranking.

1.3.2 Scoring Function

We aimed for an optimized scoring function that is capable of ranking extract method refactoring candidates, so that top-most ranked candidates are most likely to be chosen by developers for an extract method refactoring. The scoring function is a linear function that calculates the dot product of a coefficient vector, c , and a feature value vector, v , for each candidate. Candidates are arranged in decreasing order of their score.

In this paper, we use a basis of 20 features for the scoring function. In the following, we give a short overview about the features. There are three categories of feature: complexity-related features, parameters, and structural information.

We illustrate the feature values with reference to two example refactoring candidates (C_1 and C_2) that were chosen from the example method given in Figure [1.1](#). The gray areas show the nesting areas, which is defined below. The white numbers specify the nesting depth of the corresponding statement.



Fig. 1.1: Example Method with Nesting Areas of Statements And Example Candidates

Complexity-related features

We mainly focused on reducing complexity and increasing readability. For complexity indicators, we used length, nesting and data flow information. For

on the ranking performance and removed it in the next iteration. A scoring function that only considered the number of input parameters and length and nesting areas reduction still had an average NDCG of 0.885.

RQ1: How does the learned scoring function compare with our manually determined one?

The scoring function that we presented in [46](#) achieved a NDCG of 0.891, which is better than the best scoring function learned in this evaluation.

1.4 Discussion

Our results show that, in the initial run of the learning to rank tools, features indicating a reduction of complexity are much more relevant for the ranking. This is because that the manually determined scoring function was based on the 10-folds fold procedure did not vary a lot, whereas big CVs indicate big differences between the learned coefficients. As the CVs of the single features from ListMLE are much smaller than those of SVM-rank, the ranking results from ListMLE are much more stable compared with SVM-rank. SVM-rank scores with a big variance between the single iterations of the validation process; that is, despite the heavy overlapping of the training sets, the learned coefficients vary a lot and can hardly be generalized.

The results for RQ3 show that it is possible to achieve a great simplification without losing readability in the ranking performance. The biggest influences on the ranking performance were the reduction of the number of statements, the reduction of nesting areas (both are complexity indicators), and the number of input parameters.

Manual improvement As already mentioned, the learned scoring functions did not outperform the manually determined scoring function from our previous work. Obviously, the learning tools were not able to find optimal coefficients for the features. To improve the scoring function from our previous work, we did manual experiments that were influenced by the results of ListMLE and SVM-rank, and evaluated the results using the whole learning data set.

We were able to find several scoring functions that had only a handful of features and a better ranking performance than our scoring function from previous work (column 'Previous' in Table [1.3](#)). In addition to the three most important features that we obtained in the answer to RQ3 (features #3, #7, #10), we also took the count features (#14-17) into consideration. The main differences between the previous scoring function and the manually improved one from this paper are the length reduction features, the omission of nesting depth, and the number of output parameters.

By taking the results of ListMLE and SVM-rank into consideration, we were able to find a coefficient vector such that the scoring function achieved a NDCG of 0.894 (see Table [1.3](#)). That means that we were able to find a better scoring function when we combined the findings of our previous work with the learned coefficients from this paper.

Learning to Rank Extract Method Refactoring Suggestions for Long Methods

Roman Haas¹ and Benjamin Hummel²

¹ Technical University of Munich, Leitenbergstr. 8, Garching, Germany

roman.haas@tum.de
² QSE GmbH, Leitenbergstr. 8, Garching, Germany
hummel@qse.com

Summary. Extract method refactoring is a common way to shorten long methods in software development. It improves code readability, reduces complexity, and is one of the most frequently used refactorings. Nevertheless, refactorings are often refrains from applying it because identifying an appropriate set of statements that can be extracted into a new method is error-prone and time-consuming.

In a previous work, we presented a method that could be used to automatically derive extract method refactoring suggestions for long Java methods, that generated useful suggestions for developers. The approach was based on a scoring function that ranks all valid refactoring possibilities (that is, all *candidates*) to identify suitable candidates for an extract method refactoring that could be suggested to developers. Even though the evaluation has shown that the suggestions are useful for developers, there is a lack of understanding of the scoring function. In this paper, we present research on the scoring function and its features. We analyze the ranking capability. In addition, we evaluate the ranking capability of the suggested scoring function, and derive a better and less complex one using learning to rank techniques.

Key words: Learning to Rank, Refactoring Suggestion, Extract Method Refactoring, Long Method

1.1 Introduction

A long method is a bad smell in software systems [49](#), and makes code harder to read, understand and test. A straight-forward way of shortening long methods is to extract parts of them into a new method. This procedure is called 'extract method refactoring', and is the most often used refactoring practice [48](#).

The process of extracting a method can be partially automated by using modern development environments, such as Eclipse IDE or IntelliJ IDEA, that can put a set of extractable statements into a new method. However, developers still need to find this set of statements by themselves, which takes

attention of the method length (with respect to the longest method after the refactoring). We considered length based on the number of lines of code (LoC), on the number of tokens, and on the number of statements – all of them as both absolute values and relative to the original method length.

We consider highly nested methods as more complex than moderately nested ones, and use two features to represent the reduction of nesting: reduction of nesting depth and reduction of nesting area. The nesting area of a method with statements S_1 to S_n which having a nesting depth of d_i , is defined to be $\sum_{i=1}^n d_i \cdot |S_i|$. The nesting area of a method is the sum of the nesting area of all statements. The nesting area of a method is the sum of the nesting area of all statements. The nesting area of a method is the sum of the nesting area of all statements.

Dataflow information can also indicate complexity. We have features representing the number variables that are read, written or read and written.

Parameters

We considered the number of input and output parameters as an indicator of data coupling between the original and the extracted methods, which we want to keep low using our suggestions. The more parameters that are needed for a set of statements to be extracted from a method, the more the statements will depend on the rest of the original method.

Structural information

Finally, we have some features that represent structural aspects of the code. A design principle for code is that methods may only use one value [50](#). Methods that follow this principle are easier to understand. As developers often put blank lines or comments between blocks of code that process something else, we use features representing the existence and the number of blank or commented lines at their beginning, or at their end. Additionally, for first statement of the candidate, we check to see whether the type of the preceding is the same, and for the last statement, we check to see whether the type of the following statement is the same. Our last feature considers a structural complexity indicator – the number of branching statements in the candidate.

1.3.1 Training and Test Data Generation

To be able to learn a scoring function, we need training and test data. We derived this data by manually ranking approximately 1,000 extract method refactoring suggestions. To obtain this learning data, we selected 13 Java open source systems from various domains, and of different sizes. We consider a method to be 'long' if it has more than 40 LoC. From each project we randomly selected 15 long methods. For each method, we manually selected valid refactoring candidates, where the number of candidates depended on the method length.

	Previous	Previous Improved
Train	1000	1000
Test	1000	1000
Validation	1000	1000
Learning	1000	1000
Test	1000	1000
Validation	1000	1000
Learning	1000	1000

1.5 Threats to Validity

Learning from data sources that are either too similar or too small means that there is a chance that no generalization of the results is possible. To have enough data to enable us to learn a scoring function that can rank extract method refactoring candidates, we used 13 Java open source projects from various domains and from each project we randomly selected 15 long methods.

We manually reviewed the long methods, and filtered out those that were not appropriate for the extract method refactoring. The ranking was manually done five to nine valid refactoring suggestions, depending on the method length. We ensured that our learning data did not contain any code clones to avoid learning from redundant code.

The manual ranking was performed by a single individual, which is a threat to validity since there is no commonly agreed way on how to shorten a long method, and therefore to single ranking criterion exists. The ranking was done very carefully, with the aim of reducing the complexity and increasing the readability and understandability of the code as much as possible; so, the scoring function should provide a ranking such that we can make further refactoring suggestions with the same aim.

We relied on two learning to rank tools, which represents another threat to validity. The learned scoring functions heavily depend on the tool. As the learned scoring functions vary, it is necessary to have an independent way of evaluating the ranking performance of the learned scoring functions. We used the widely used measure NDCG to evaluate the scoring functions, and applied a 10-fold cross validation procedure to obtain a meaningful evaluation of the ranking performance of the learned scoring functions.

A threat to external validity is the fact that we derived our learning data from 13 open source Java systems. Therefore, results are not necessarily generalizable.

1.6 Related Work

In our previous work [46](#), we presented an automatic approach to derive extract method refactoring suggestions for long methods. We obtained valid

non-exhaustive datasets and some statements that cannot be extracted but supported by the programming language [49](#).

The refactoring process can be improved by suggesting to developers which statements could be extracted into a new method. The literature presents several approaches that can be used to extract method refactorings. In a previous work, we suggested a method that could be used to automatically find good extract method refactoring candidates for long Java methods [46](#). Our first prototype, which was derived from manual experiments on several open source systems, implemented a scoring function to rank refactoring candidates. The result of our evaluation has shown that this first prototype finds suggestions that are followed by experienced developers. The results of our first prototype have been implemented in an industrial software quality analysis tool.

Problem statement. The scoring function is an essential part of our approach to derive extract method refactoring suggestions for long methods. It is decisive for the quality of our suggestions, and also important for the complexity of the implementation of the refactoring suggestion. However, it is currently unclear how good the scoring function actually performs in ranking refactoring suggestions and how much complexity will be needed to obtain useful suggestions. Therefore, in order to enhance our work, we need a deeper understanding of the scoring function.

Contribution. We do further research on the scoring function of our approach to derive extract method refactoring suggestions for long Java methods. We use learning to rank techniques in order to learn which features of the scoring function are useful for ranking refactoring suggestions, and to keep the scoring function as simple as possible. In addition, we evaluate the ranking performance of our previous scoring function, and compare it with the new scoring function that we learned. For the machine learning setting, we use 177 training and testing data sets that we obtained from 13 well-known open source systems by manually ranking five to nine randomly selected valid refactoring candidates.

In this paper, we show how we derived better extract method refactoring suggestions than in our previous work using learning to rank tools.

1.2 Fundamentals

We use learning to rank techniques to obtain a scoring function that is able to rank extract method refactoring candidates, and we normalized discounted cumulative gain (NDCG) metrics to evaluate the ranking performance. In this section, we explain the techniques, tools and metrics that we use in this paper.

into the code. Therefore, in the pruning step of our approach, we usually enter our candidates that need more than three input parameters, thus avoiding the 'long parameter list' mentioned by Fowler [51](#). To avoid learning that too many input parameters are bad, we considered only candidates that had less than four input parameters.

We ranked the selected candidates manually with respect to complexity reduction and readability improvement. The higher the ranking we gave a candidate, the better the suggestion was for us.

Selected to the randomly selected methods were not suitable for an extract method refactoring. That was most commonly the case when the code would not benefit from the extract method, but from other refactorings. In addition, in some methods, we could not derive a scoring function, because the methods were only very weak candidates. That is why we did not use 18 of the 195 randomly selected long methods to learn our scoring function.

1.4 Evaluation

In this section, we present and evaluate the results from the learning procedure.

1.4.1 Research Questions

RQ1: What are the results of the learning tools? In order to get a scoring function that is capable of ranking the extract method refactoring candidates, we decided to use two learning to rank tools that implement different approaches, and that had performed well in previous studies.

RQ2: How stable are the learned scoring functions? To be able to derive implications for a real-world scoring function, the coefficients of the learned scoring function should not vary a lot during the 10-fold cross evaluation procedure.

RQ3: Can the scoring function be simplified? For practical reasons, it is useful to have a scoring function with a limited number of features. Additionally, reducing the search space may increase the performance of the learning to rank tools – resulting in better scoring functions.

RQ4: How does the learned scoring function compare with our manually determined one? In our previous work, we derived a scoring function from manual experiments. Now we use our learning data set to evaluate the ranking performance of the previously defined scoring function, and to compare it with the learned one.

¹ On http://ja.tum.de/~haas/12r_src_data.zip we provide our rankings and the corresponding code base from which we generated the training data.

All valid refactoring candidates were ranked by a manually-determined scoring function that aims to reduce code complexity and increase readability. In the present work, we have put the scoring function on more solid ground by learning a scoring function from many manually ranked, and manually ranked refactoring suggestions.

In the literature, there are several approaches that learn to suggest the most beneficial refactorings – usually for code clones. Wang and Gollrey [52](#) propose an automated approach to recommend clones for refactoring by training a decision-tree based classifier, C4.5. They use 15 features for decision-tree model training, where four consider the cloning relationship, four the content of the clone, and seven relate to the code of the clone. In the present paper, we have used a similar approach, but with a different aim: instead of clones, we have learned on long methods.

Mondal et al. [53](#) rank clones for refactoring through mining association rules. Their idea is that clones which are often changed together, and thus have a similar function, are worthy candidates for refactoring. Their prototype tool, MARC, identifies clones that are often changed together in a similar way, and suggests association rules among clones. A second step of their evaluation on thirteen software systems is that clones that are highly ranked by MARC are important refactoring possibilities. We used learning to rank techniques to find a scoring function that is capable of ranking extract method refactoring candidates from long methods.

1.7 Conclusion and Future Work

In this paper, we have presented an approach to derive a scoring function that is able to rank extract method refactoring suggestions by applying learning to rank tools. The scoring function can be used to automatically rank extract method refactoring candidates, and thus present a set of best refactoring suggestions to developers. The resulting scoring function needs less parameters than previous scoring functions but has a better ranking performance.

In the future, we would like to suggest sets of refactorings, especially those that remove clones from the code.

We would also like to find out whether the scoring function provides good suggestions for object-oriented programming languages other than Java and whether other features need to be considered in that case.

Acknowledgments

Thanks to the anonymous reviewers for their helpful feedback. This work was partially funded by the German Federal Ministry of Education and Research (BMBF), grant "Q-Effekt, 01IS1503A3". The responsibility for this article lies with the authors.

to a ranking task [46](#).

There are several learning to rank approaches, where the pairwise and the listwise approach usually perform better than common pointwise regression approaches [45](#). The pairwise approach learns by comparing two training objects and their given ranks ('ground truth'), whereas in the case the listwise approach learns from the list of all given rankings of refactoring suggestions for a long method. Liu et al. [45](#) pointed out that the pairwise and the listwise approaches usually perform better than the pointwise approach. Therefore, we do not rely on a pointwise approach but use pairwise and listwise learning to rank tools.

Qin et al. [45](#) constructed a benchmark collection for research on several learning to rank tools on the Learning To Rank (LETOR) data set. Their results support the hypothesis that pointwise approaches perform badly compared with pairwise and listwise approaches. In addition, listwise approaches often perform better than pairwise. However, SVM-rank, a pairwise learning to rank tool by Tsochantzidis et al. [48](#), performs quite well and the first experiments on our data set showed that SVM-rank may lead to us interesting results. We set the parameter κ to 0.5 and the parameter σ to 5,000 as a trade-off between time consumption and learning performance.

Beside SVM-rank, we used a listwise learning to rank tool, ListMLE by Xia et al. [44](#). In their evaluation, they showed that ListMLE performs better than ListNet by Coyot et al. [43](#), which was also considered to be good by Qin et al. Liu et al. [45](#) improved the learning capability of ListMLE, but did not provide literature or source code, so we were unable to use the improved version.

ListMLE needs to be assigned a tolerance rate and a learning rate. In a series of experiments we performed, we found that the optimal ranking performance on our data set was with a tolerance rate of 0.001 and a learning rate of 1E-15.

1.2.2 Training and Testing

The learning process consisted of two steps: training and testing. We applied cross-validation [46](#) with 10 sets, that is, we split our learning data into 10 sets of (nearly) equal size. We performed 10 iterations using these sets, where nine of the sets were considered to be training data and one set was used as test data.

Test data is used to evaluate the ranking performance of the learned scoring function by comparing the grade of a refactoring candidate determined by the learned scoring function with its grade given by the learning data. We use NDCG metric to compare different scoring functions and their performances.

To answer RQ1 and RQ2, we used the learning to rank tools SVM-rank and ListMLE to perform a 10-fold cross validation on our training and test data set of 177 long methods, and a total of 1,185 refactoring suggestions. We illustrate the stability of the single coefficients by using box plots that show the coefficients are distributed over the ten iterations of the 10-fold cross validation.

To answer RQ3, we simplified the learned scoring function by omitting features, where the selection criterion for the omitted features was the stability of the ranking capability of the scoring function. Our initial feature set contained six different measures of length. For the sake of simplicity, we would like to have only one measure of length in our scoring function. To find out which measure best fits in with our training set, we re-ran the validation procedure (again using ListMLE and SVM-rank), but this time with only one length measurement, using each of the length measurements one at a time. We continued with the feature set reduction until only one feature was left.

1.4.3 Results

The following paragraphs answer the research questions.

RQ1: What are the results of the learning tools?

Figures [1.2](#) and [1.3](#) show the results of the 10-fold cross validation for ListMLE and for SVM-rank, respectively. For each single feature, i , there is a box plot of the corresponding coefficient, c_i .

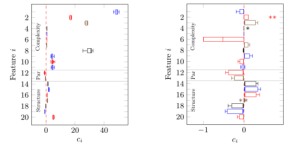


Fig. 1.2: Learning Result From ListMLE With All Features

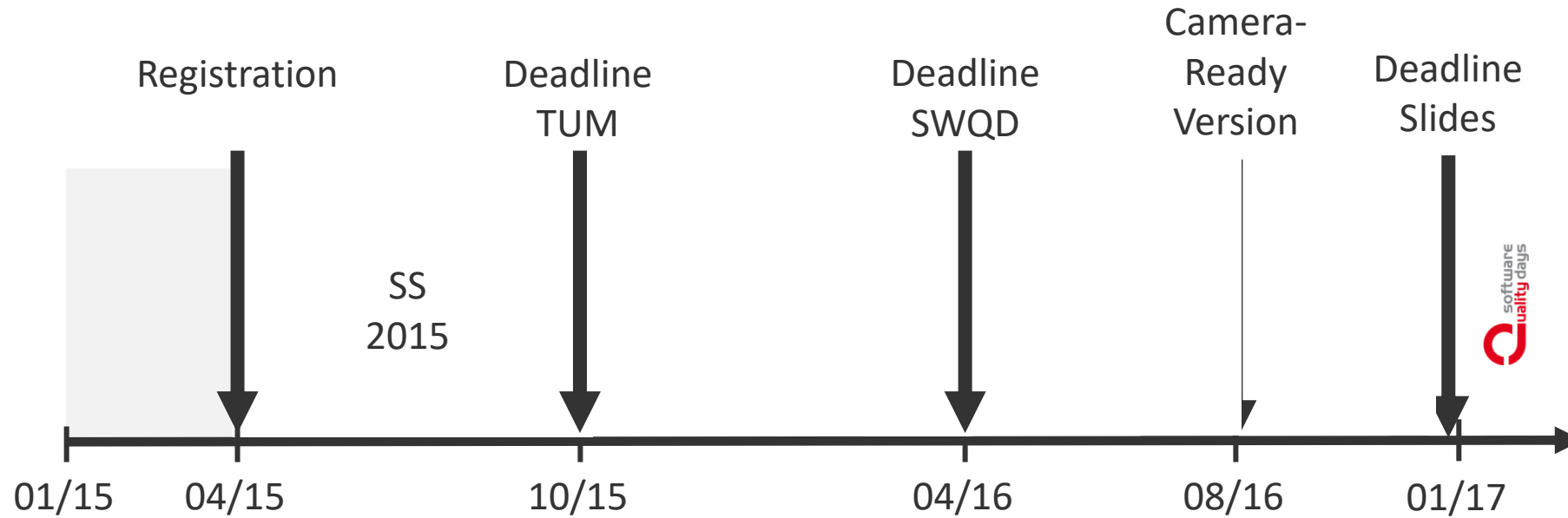


1. Z. Cao, T. Qin, T.-Y. Liu, M.-F. Tian, and H. Li. Learning to rank: From pairwise approach to listwise approach. In *24th ICML*, 2007.
2. M. Fowler. *Refactoring: Improving the Design of Existing Code*. Addison-Wesley object technology series. Addison-Wesley, Reading, PA, 1999.
3. R. Haas and B. Hummel. Deriving extract method refactoring suggestions for long methods. In *38QSO*, 2016.
4. L. Jiang. A short introduction to learning to rank. *IEEE Transactions on Information and Systems*, 98(10):1854–1862, 2011.
5. K. Järvelin and J. Kekkonen. IR evaluation methods for retrieving highly relevant information. In *2nd SIGIR*, 2002.
6. M. Kim, T. Zimmermann, and N. Nagappan. A field study of refactoring challenges and benefits. In *2008 International Symposium on the FSE*, 2012.
7. Y. Luo, Y. Zhao, S. Xia, and X. Cheng. Position paper: Position paper: A sequential learning process for ranking. In *2008 Conference on F4I*, 2014.
8. T.-Y. Liu. Learning to rank for information retrieval. *Foundations and Trends in Information Retrieval*, 3(3):225–331, 2009.
9. B. C. Martin. *Clean Code: A Handbook of Agile Software Craftsmanship*. Robert Martin series. Prentice Hall, Upper Saddle River, NJ, 2009.
10. M. Mondal, C. K. Roy, and K. Schneider. Automatic ranking of clones for refactoring through mining association rules. In *CSMR-ICRSE*, 2014.
11. E. Murphy-Hill and A. P. Black. Why don't people use refactoring tools? In *1st WET*, 2007.
12. E. Murphy-Hill and A. P. Black. Breaking the barriers to successful refactoring: Observations and tools for extract method. In *2006 ICSE*, 2006.
13. W. F. Opdyke. *Refactoring: Object-Oriented Fundamentals*. PhD thesis, University of Illinois at Urbana-Champaign, 1992.
14. A. Osti. *A Manual and Automatic Approach for Recommending Software Refactoring*. PhD thesis, Université de Montréal, 2015.
15. T. Qin, T.-Y. Liu, J. Xu, and H. Li. Letor: A benchmark collection for research on learning to rank information retrieval. *Information Retrieval*, 18(4):386–374, 2010.
16. C. Sammut, editor. *Encyclopedia of machine learning*. Springer, New York, 2011.
17. N. Tassatilis and A. Tsochantzidis. Ranking refactoring suggestions based on historical stability. In *25th ECSMR*, 2011.
18. I. Tsochantzidis, T. Joachims, T. Hofmann, and Y. Altun. Large margin methods for structured and semistructured output variables. *Journal of Machine Learning Research*, 6:1453–1484, 2005.
19. W. Wang and M. W. Gollrey. Recommending clones for refactoring using design, context, and history. In *ICSM*, 2014.
20. D. Whiting, U. F. Kaba, and S. Rowson. An empirical evaluation of refactoring. *e-Information*, 1(1):27–42, 2007.
21. F. Xia, T.-Y. Liu, J. Wang, and W. Zhang. In the Listwise approach to learning to rank: Theory and Theory. In *25th ICML*, 2008.



	Track A	Track B	Track C	Scientific-Track	Solution Provider Forum I	Solution Provider Forum II
08:00	Registration					
09:00	Opening session / Begrüßung & Konferenzzeröffnung					
09:20	Keynote: Managing for Happiness					
10:20	Coffee break & networking in the exhibition area / Kaffeepause & Networking im Ausstellungsbereich					
10:50						
10:55	Software Engineering Complexity and Challenges of Software Engineering (Result Planning Limited, Kolbotn, Norwegen)	Testen ist Unfug!...aber ist frühe QS in Form von statischer Analyse wirklich so einfach - Über die soziologischen Aspekte (Zöhlke Engineering (Austria) GmbH, Vienna, AT)	Continuous Integration für Mobile Apps (Zöhlke Engineering (Austria) GmbH, Vienna, AT)	Improve your software models with search-based techniques (TU Wien, Wien, AT) Englisch, Fortgeschrittene	Ranorex in the Agile World (Ranorex GmbH, Graz, Österreich) Deutsch, Einsteiger	Testumgebungen auf einen Klick - zeitgemäßes Testumgebungsmanagement als Herausforderung und Lösung (ANECON Software Design und Beratung GmbH, Wien, AT) Deutsch, Fortgeschrittene
11:20	Traceability in a Fine Grained Software Configuration Management System (Vector Informatik GmbH, Stuttgart, DE) Englisch, Fortgeschrittene	Projektbericht „Optimierte Testautomatisierung bei Vienna Insurance Group“ (BIAC - Business Insurance Application Consulting GmbH, Wien) (Tricentis GmbH, Wien, AT) Deutsch, Einsteiger				Agiles Requirements Management – eine effiziente Umsetzung mit agosense.fidella (agosense GmbH, Kornwestheim, DE) Deutsch, Einsteiger
11:50						
12:20						
12:50						
13:20						
13:50						
14:25						
14:35	Kontinuierliche Architekturanalyse (Software Quality Lab, Linz) Deutsch, Einsteiger	Strukturierte Tests bei defizitärer Dokumentation - Wie man zwei Fliegen mit einer Klappe schlägt (SRC Security Research & Consulting GmbH, Bonn, DE) Deutsch, Fortgeschrittene	Continuous Delivery - Feel your Quality - Every Day (Automic Software GmbH (CA Technologies), Wien, AT) (Automic Software GmbH, Wien, AT) Englisch, Fortgeschrittene	A portfolio of internal quality metrics for software architects (University of Gothenburg, Gothenburg, SE) Englisch, Fortgeschrittene	Scrum in Embedded Systems (Software Quality Lab GmbH, Linz, AT) (ENGEL AUSTRIA GmbH, Schwertberg, AT) Deutsch, Fortgeschrittene	Zertifizierung Quality Engineer für das Internet der Dinge (ISQI GmbH, Potsdam, DE) Englisch, Experte
14:55				Validating converted Java Code via Symbolic Execution (ZT Prentner-IT, Wien, AT) Englisch, Fortgeschrittene		

Chronological Overview

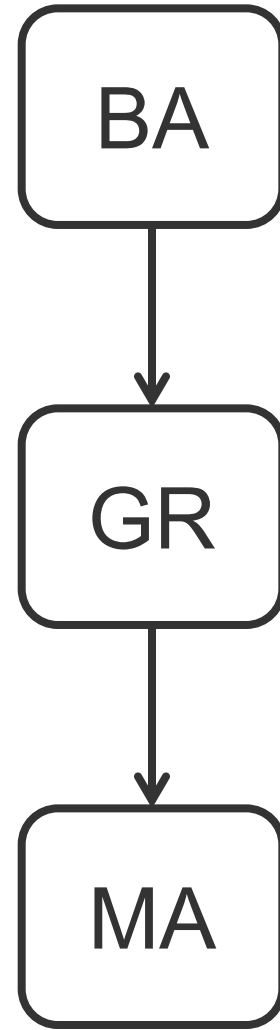


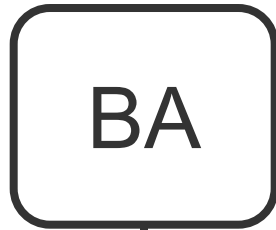
What is Different to Other Study Projects?

- More Freedom
 - Topic
 - Own research
 - You define schedule and pace
- Requires high level of self-organization
- Better opportunities for personal growth

Personal Conclusion

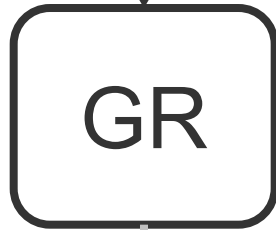
- My GR was on my „mental Stack“ during my entire studies in the Master's program
- GR got me out of my comfort zone
- Learned a lot on research methodologies and practical application of machine learning techniques
- Working on my research topic was fun for me
- I would do it again 😊





Timo Pawelka

Automatische Erkennung der Sprache von Quelltext-Kommentaren
Bachelor's Thesis, not published



Timo Pawelka, Elmar Juergens:

Is This Code Written in English? A Study of the Natural Language of Comments and Identifiers in Practice.
Proceedings of the 31st International Conference on Software Maintenance and Evolution (ICSME'15), 2015.



BA



GR

Raphael Nömmner, Roman Haas

Test Suite Minimization

Guided Research, to be published in Conference Proceedings of SWQD '20



MA

Raphael Nömmner

Design and Evaluation of Regression Test Suite Minimization Techniques

Master's Thesis

Funding

Costs 1k€ – 5k€

- Travel and accomodation costs
- Conference fee

Funding sources (often mixed)

- Travel Subsidies
- Chairs
- DAAD scholarships
- e.g., CQSE

Decision processes take long, so organize this early!

Agenda

1. Motivation
- 2. Preparation**
3. Doing the work

Get the Most out of your GR?!

- GR provides the opportunity to publish scientific work at a scientific venue.
- Nevertheless, formally, you do not need to publish anything
- My recommendation: aim for a scientific publication

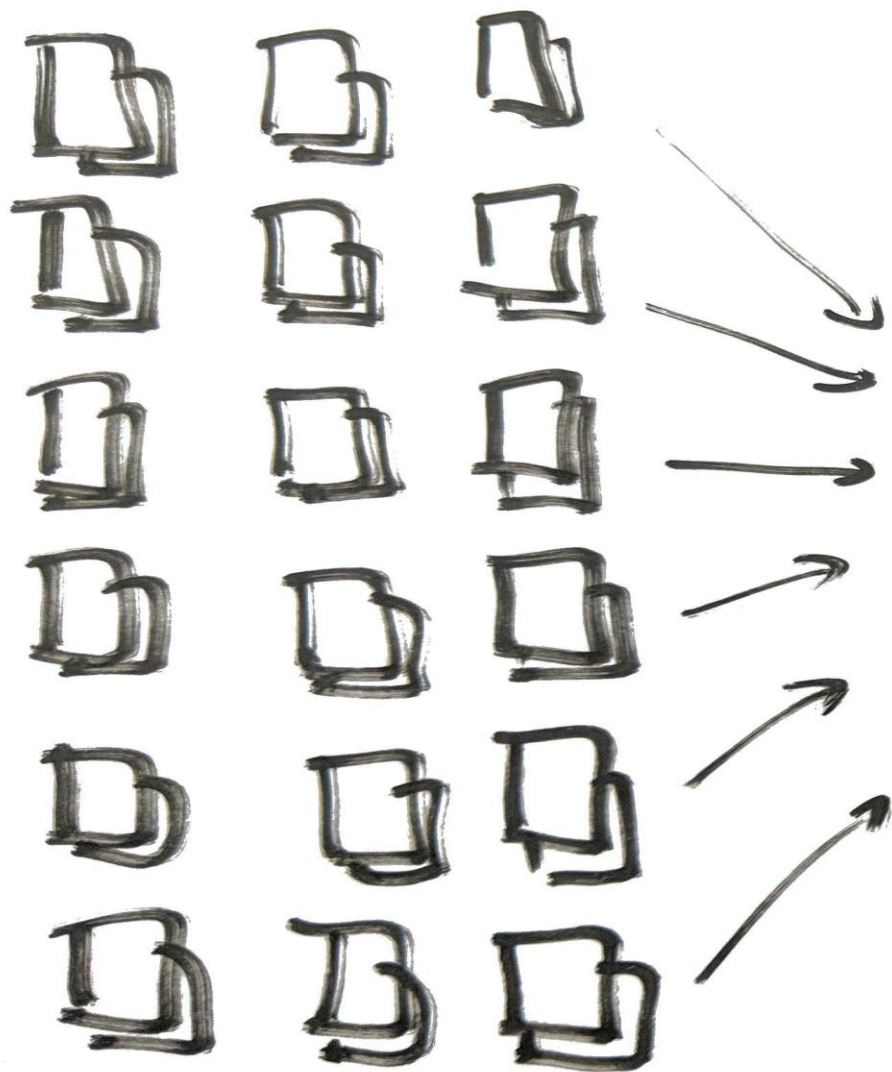


	Track A	Track B	Track C	Scientific-Track	Solution Provider Forum I	Solution Provider Forum II
08:00	Registration					
09:00	Opening session / Begrüßung & Konferenzzeröffnung					
09:20	Keynote: Managing for Happiness					
10:20	Coffee break & networking in the exhibition area / Kaffeepause & Networking im Ausstellungsbereich					
10:50						
10:55	Software Engineering Complexity and Challenges of Software Engineering (Result Planning Limited, Kolbotn, Norwegen)	Testen ist Unfug!...aber ist frühe QS in Form von statischer Analyse wirklich so einfach - Über die soziologischen Aspekte (Zöhlke Engineering (Austria) GmbH, Vienna, AT)	Continuous Integration für Mobile Apps (Zöhlke Engineering (Austria) GmbH, Vienna, AT)	Improve your software models with search-based techniques (TU Wien, Wien, AT) Englisch, Fortgeschrittene	Ranorex in the Agile World (Ranorex GmbH, Graz, Österreich) Deutsch, Einsteiger	Testumgebungen auf einen Klick - zeitgemäßes Testumgebungsmanagement als Herausforderung und Lösung (ANECON Software Design und Beratung GmbH, Wien, AT) Deutsch, Fortgeschrittene
11:20	Traceability in a Fine Grained Software Configuration Management System (Vector Informatik GmbH, Stuttgart, DE) Englisch, Fortgeschrittene	Projektbericht „Optimierte Testautomatisierung bei Vienna Insurance Group“ (BIAC - Business Insurance Application Consulting GmbH, Wien) (Tricentis GmbH, Wien, AT) Deutsch, Einsteiger				Agiles Requirements Management – eine effiziente Umsetzung mit agosense.fidella (agosense GmbH, Kornwestheim, DE) Deutsch, Einsteiger
11:50						
12:20						
12:50						
13:20						
13:50						
14:25						
14:35	Kontinuierliche Architekturanalyse (Software Quality Lab, Linz) Deutsch, Einsteiger	Strukturierte Tests bei defizitärer Dokumentation - Wie man zwei Fliegen mit einer Klappe schlägt (SRC Security Research & Consulting GmbH, Bonn, DE) Deutsch, Fortgeschrittene	Continuous Delivery - Feel your Quality - Every Day (Automic Software GmbH (CA Technologies), Wien, AT) (Automic Software GmbH, Wien, AT) Englisch, Fortgeschrittene	A portfolio of internal quality metrics for software architects (University of Gothenburg, Gothenburg, SE) Englisch, Fortgeschrittene	Scrum in Embedded Systems (Software Quality Lab GmbH, Linz, AT) (ENGEL AUSTRIA GmbH, Schwertberg, AT) Deutsch, Fortgeschrittene	Zertifizierung Quality Engineer für das Internet der Dinge (ISQI GmbH, Potsdam, DE) Englisch, Experte
14:55				Validating converted Java Code via Symbolic Execution (ZT Prentner-IT, Wien, AT) Englisch, Fortgeschrittene		

Submissions

Selection Procedure

Agenda



Pecking Order



Conference
10%-25%

Acronym	Full Name	Date
CHASE	11th International Workshop on Cooperative and Human Aspects of Software Engineering	27-May
CSI-SE	5th International Workshop on Crowd Sourcing in Software Engineering	27-May
MET	International Workshop on Metamorphic Testing	27-May

Workshop
40%-60%

RAISE	SoHeal	MISE	GE	SQUADE	SE4COG	SER&IP	SE4Science
SEAD	WETSEB	SEHS	RoSE	AST	FairWare	SESoS	RET
SEsCPS	GREENS	CESI	SEFAIAS	SBST	RCoSE	GI	SEEM

Aim: Submission to workshops



Author



Organizer



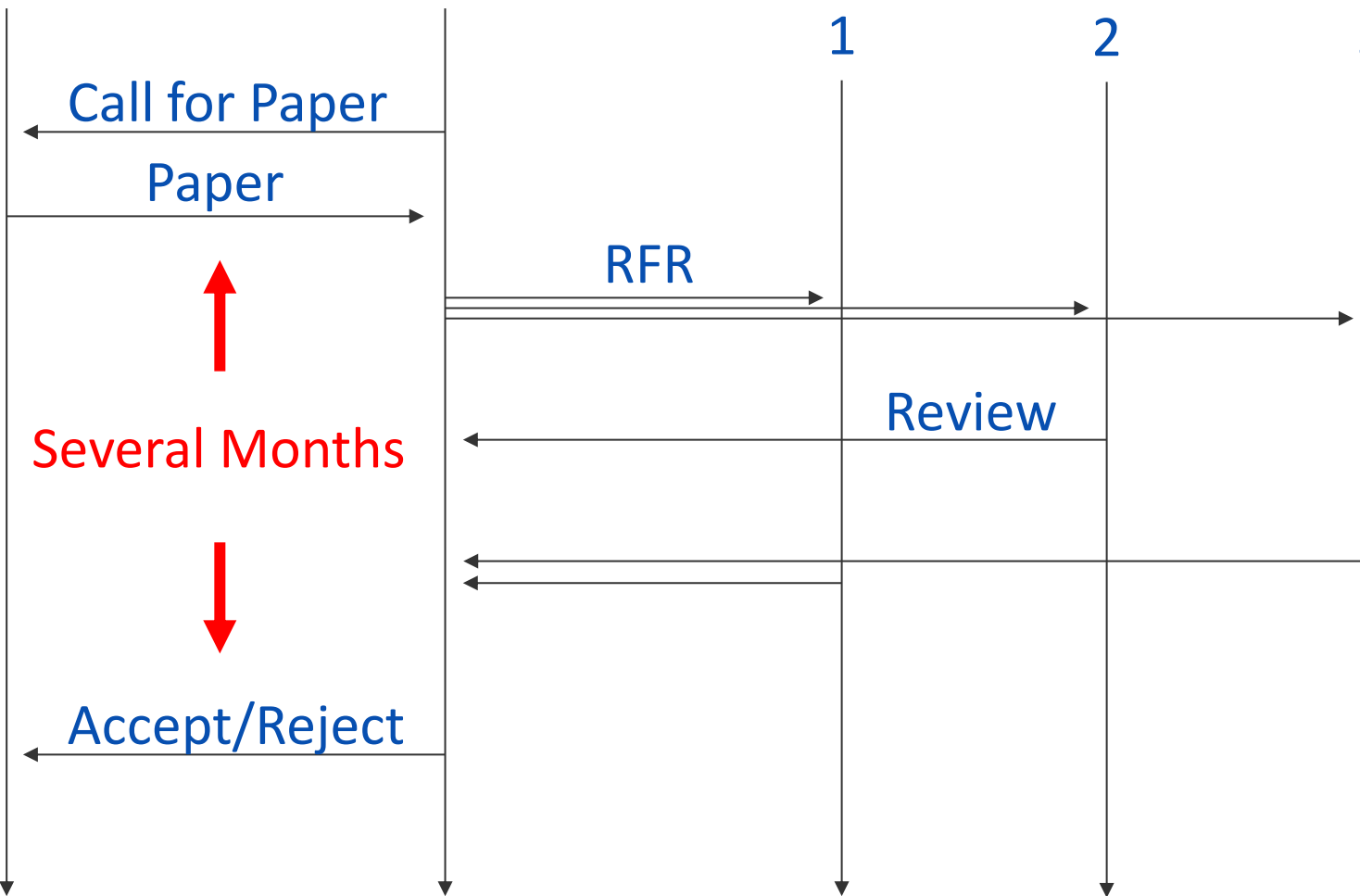
Reviewer
1



Reviewer
2



Reviewer
3



Call for Papers

12th International Workshop on Software Clones (IWSC 2018)

Co-located with the [25th IEEE International Conference on Software Analysis, Evolution, and Reengineering \(SANER 2018\)](#)

March 20, 2018, Campobasso, Italy

Software clones are often a result of copying and pasting as an act of ad-hoc reuse by programmers, and can occur at many levels, from simple statement sequences to blocks, modules, classes, packages, models, requirements or architectures. They are a common phenomenon in software development today.

IWSC series of events has provided a forum for researchers to discuss the state-of-the-art in software clones.

IWSC aims to bring researchers and practitioners together to discuss the state-of-the-art in software clones.

In particular, we expect the in-depth discussions and presentations about IWSC 2018 are here on this page.

TOPICS OF INTEREST:

Topics of interest include but not limited to:

- Use cases for clones and clones detection
- Experiences with clones analysis
- Types and nature of clones
- Causes and effects of clones
- Techniques and algorithms
- Clone and clone pattern visualization
- Tools and systems for clones detection
- Applications of clones detection
- System architecture and clones
- Effect of clones to system evolution
- Clone analysis in families of software
- Measures of code similarity
- Economic and trade-off models
- Evaluation and benchmarking
- Licensing and plagiarism issues
- Clone-aware software design
- Refactoring through clones
- Higher-level clones in models
- Clone evolution and variability
- Role of clones in software evolution

PAPERS SOUGHT:

Each paper will be reviewed by at least three members of the program committee following a full double-blind process. Authors must adhere to SANER's double blind guidelines - <http://saner.unimol.it/restrack>. The following types of papers are sought:

- Full papers (7 pages maximum)
- Position papers (2 pages maximum)
- Tool demonstration papers (4 pages maximum)

SUBMISSION:

Papers must conform to the [IEEE proceedings paper format guidelines](#). If the paper is accepted, at least one author must attend the workshop and present the paper. Accepted papers will be published in the [IEEE Xplore Digital Library](#) along with the SANER proceedings.

All submissions must be in PDF and must be submitted online by the deadline via the IWSC 2018 EasyChair conference management system.

[Submit your papers here >>> EasyChair<<<](#)

IMPORTANT DATES:

- Abstract submission deadline: January 19, 2018 AoE
- Paper submission deadline: January 26, 2018 AoE
- Notifications: February 16, 2018
- Camera Ready deadline: ** February 22, 2018 **
- Workshop day: March 20 2018

GENERAL CHAIR:

TBD

PROGRAM CO-CHAIRS:

- [Ying \(Jenny\) Zou](#) (ying.zou@queensu.ca), Queen's University, Canada
- [Matthew Stephan](#) (stephamd@miamioh.edu), Miami University, USA

STEERING COMMITTEE:

- [James R. Cordy](#), Queen's University, Canada
- [Katsuro Inoue](#), Osaka University, Japan
- [Rainer Koschke](#), University of Bremen, Germany

Call for Papers

12th International Workshop on Software Clones (IWSC 2018)

Co-located with the [25th IEEE International Conference on Software Analysis, Evolution, and Reengineering \(SANER 2018\)](#)

March 20, 2018, Campobasso, Italy

Software clones are often a result of copying and pasting as an act of ad-hoc reuse by programmers, and can occur at many levels, from simple statement sequences to blocks, models, requirements or architectures today.

IWSC series of events has provided

IWSC aims to bring researchers

particular, we expect the in-depth

about IWSC 2018 are here on this

TOPICS OF INTEREST:

Topics of interest include but not

- Use cases for clones and clones
- Experiences with clones and clones
- Types and nature of clones
- Causes and effects of clones
- Techniques and algorithms
- Clone and clone pattern vis
- Tools and systems for detecting
- Applications of clone detection
- System architecture and clones
- Effect of clones to system
- Clone analysis in families of
- Measures of code similarity
- Economic and trade-off models
- Evaluation and benchmarking
- Licensing and plagiarism issues
- Clone-aware software design
- Refactoring through clones
- Higher-level clones in models
- Clone evolution and variability
- Role of clones in software

PAPERS SOUGHT:

Each paper will be reviewed by at least three members of the program committee following a full double-blind process. Authors must adhere to SANER's double blind guidelines - <http://saner.unimol.it/restrack>. The following types of papers are sought:

- Full papers (7 pages maximum)
- Position papers (2 pages maximum)
- Tool demonstration papers (4 pages maximum)

SUBMISSION:

Papers must conform to the [IEEE proceedings paper format guidelines](#). If the paper is accepted, at least one author must attend the workshop and present the paper. Accepted papers will be published in the [IEEE Xplore Digital Library](#) along with the SANER proceedings.

All submissions must be in PDF and must be submitted online by the deadline via the IWSC 2018 EasyChair conference management system.

[Submit your papers here >>> EasyChair<<<](#)

IMPORTANT DATES:

- Abstract submission deadline: January 19, 2018 AoE
- Paper submission deadline: January 26, 2018 AoE
- Notifications: February 16, 2018
- Camera Ready deadline: ** February 22, 2018 **
- Workshop day: March 20 2018

GENERAL CHAIR:

TBD

PROGRAM CO-CHAIRS:

- [Ying \(Jenny\) Zou](#) (ying.zou@queensu.ca), Queen's University, Canada
- [Matthew Stephan](#) (stephamd@miamioh.edu), Miami University, USA

STEERING COMMITTEE:

- [James R. Cordy](#), Queen's University, Canada
- [Katsuro Inoue](#), Osaka University, Japan
- [Rainer Koschke](#), University of Bremen, Germany

Call for Papers

12th International Workshop on Software Clones
Co-located with the 25th IEEE International Conference on Software Maintenance
March 20, 2018, Campobasso, Italy

Software clones are often a result of evolution of statement sequences to blocks, models, requirements or architectures today.

IWSC series of events has provided a forum for researchers to present their work.

IWSC aims to bring researchers and practitioners together.

In particular, we expect the in-depth discussions and presentations.

For more information about IWSC 2018 are here on the website.

TOPICS OF INTEREST:

Topics of interest include but not limited to:

- Use cases for clones and clones
- Experiences with clones and clones
- Types and nature of clones
- Causes and effects of clones
- Techniques and algorithms
- Clone and clone pattern visualization
- Tools and systems for clone detection
- Applications of clone detection
- System architecture and clones
- Effect of clones to system
- Clone analysis in families of clones
- Measures of code similarity
- Economic and trade-off models
- Evaluation and benchmarking
- Licensing and plagiarism issues
- Clone-aware software development
- Refactoring through clones
- Higher-level clones in models
- Clone evolution and variability
- Role of clones in software evolution

PAPERS SOUGHT:

Each paper will be reviewed by at least two reviewers using double blind guidelines - <http://seer.cs.cmu.edu/>

- Full papers (7 pages maximum)
- Position papers (2 pages maximum)
- Demonstration papers (4 pages maximum)

Program Committee

Name	Institution	Country
Toshihiro Kamiya	Shimane University	Japan
Daqing Hou	Clarkson University	USA
Tien Nguyen	University of Texas at Dallas	USA
Nils Göde	CQSE GmbH	Germany
Jens Krinke	University College London	UK
Otavio Lemos	ICT-UNIFESP	Brazil
Manishankar Mondal	University of Saskatchewan	Canada
Ravindra Naik	Tata Consultancy Services	India
Robert Tairas	Vanderbilt University	USA
Minhaz Zibran	University of New Orleans	USA
Eunjong Choi	Nara Institute of Science and Technology	Japan
Michael Godfrey	University of Waterloo	Canada
Yoshiki Higo	Osaka University	Japan
Foutse Khomh	Ecole Polytechnique de Montréal	Canada
Nicholas A. Kraft	ABB Corporate Research	USA
Chanchal Roy	University of Saskatchewan	Canada
Hitesh Sajjani	Microsoft	USA
Suresh Thummalapenta	Microsoft	USA
Xiyou Wang	University of Texas at San Antonio	USA
Norihiro Yoshida	Nagoya University	Japan

attend the workshop and

management system.

What If I have no Topic in Mind?

- Ask potential supervisors for ideas
 - Supervisor from Bachelor's Thesis
 - Lectures
 - Seminars
 - Lab courses
- As a supervisor, I do **not** expect
 - Students to come up with thesis topics
 - Students to apply only for documented topics
- If you have a rough idea, discuss it with potential supervisors

CQSE



Development
Operations

Services

Audits
Quality Control

Research

Software Quality
e.g., Coding, Testing



Requirements for a GR topic

- Is there a clear problem statement?
- Can different solutions be evaluated objectively?

Why?

- Decision making while you work on it
- Easier to convince supervisor
- Easier to convince program chair

Even more important for a GR than BA/MA

More info: www.thesisguide.org

Wann ist ein Thema Schrott?

Wenn sich nicht klar beurteilen lässt, ob eine Lösung besser ist, als eine andere.

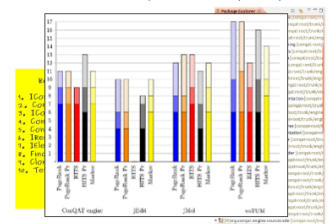
Wichtigste Faktoren:

- Gibt es ein klares Problem Statement?
- Kann ich Alternative Lösungen objektiv bewerten?

Warum?

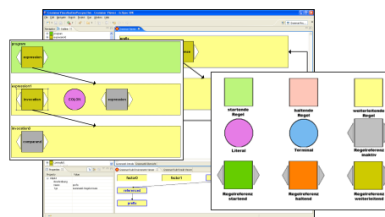
- Entscheidungsfindung während Bearbeitung
- Einfacher, Betreuer zu überzeugen
- Betreuer kann Professor einfacher überzeugen

Using Network Analysis for Recommendation of Central Software Classes (Daniela Steidl, 2012)



Grafiken aus Foliensatz von Daniela

Unterstützung von Sprachentwicklung durch Visualisierung



Grafiken aus Foliensatz von Ludwig

Themen-Antipatterns

- Search my Literature
- Implementation only
- Choose my Tool

- Wenig objektive Bewertungskriterien
- Kein eigenes Feedback während Arbeit
- Veröffentlichung sehr schwierig

What Makes a Good Guided Research Supervisor



- Needs to have publishing experience
- Has already successfully published (ideally on the same workshop if you aim for a publication)
- Sources: scholar.google.com, DBLP, personal webpage



Roman Haas 

CQSE GmbH
Bestätigte E-Mail-Adresse bei cqse.eu





<input type="checkbox"/> TITEL  	ZITIERT VON	JAHR
<input type="checkbox"/> How can manual testing processes be optimized? developer survey, optimization guidelines, and case studies R Haas, D Elsner, E Juergens, A Pretschner, S Apel Proceedings of the 29th ACM Joint Meeting on European Software Engineering ...	37	2021
<input type="checkbox"/> Is static analysis able to identify unnecessary source code? R Haas, R Niedermayr, T Roehm, S Apel ACM Transactions on Software Engineering and Methodology (TOSEM) 29 (1), 1-23	29	2020
<input type="checkbox"/> An evaluation of test suite minimization techniques R Noemmer, R Haas International Conference on Software Quality, 51-66	25	2019
<input type="checkbox"/> Deriving extract method refactoring suggestions for long methods R Haas, B Hummel International Conference on Software Quality, 144-155	25	2015
<input type="checkbox"/> Teamscale: tackle technical debt and control the quality of your software R Haas, R Niedermayr, E Juergens 2019 IEEE/ACM International Conference on Technical Debt (TechDebt), 55-56	16	2019
<input type="checkbox"/> Learning to rank extract method refactoring suggestions for long methods R Haas, B Hummel Software Quality. Complexity and Challenges of Software Engineering in ...	6	2017

Agenda

1. Motivation
2. Preparation
3. **Doing the work**

View as an Supervisor







































-  Regular meeting
-  Meeting on demand





























































































ICSE 2021

“ ICSE 2021 received 615 submissions. Of these, 13 were desk rejected for double-blind or formatting violations. The remaining 602 papers went through a thorough review process, with at least three reviewers, one meta-reviewer, and an area chair per paper. Following an online discussion, the program committee decided to accept 138 papers, including 30 conditional ones. We will announce the acceptance rate after finalizing all conditional decisions.”

Program Board ICSE Technical Papers

 Tevfik Bultan Track Chair University of California, Santa Barbara	 Jon Whittle Track Chair Monash University Australia	 Sven Apel University of Passau Germany	 Andrew Begel Microsoft Research United States	 Antonia Bertolino CNR-ISTI Italy	 Eric Bodden Heinz Nixdorf Institut Paderborn University and Fraunhofer IEM	 Yuriy Brun University of Massachusetts Amherst United States	 Margaret Burnett Oregon State University United States
 Jordi Cabot ICREA - UOC	 Cristian Cadar Imperial College London United Kingdom	 Marsha Chechik University of Toronto	 Jane Cleland-Huang University of Notre Dame United States	 Daniela Damian University of Victoria Canada	 Laura Dillon Michigan State University	 Bernd Fischer Stellenbosch University South Africa	 Alessandro Garcia PUC-Rio Brazil
 Dimitra Giannakopoulou NASA Ames Research Center	 Sung-hun Kim Hong Kong University of Science and Technology Hong Kong	 Andrew J. Ko University of Washington	 Claire Le Goues Carnegie Mellon University United States	 David Lo Singapore Management University Singapore	 Darko Marinov University of Illinois at Urbana- Champaign United States	 Mira Mezini TU Darmstadt, Germany Germany	 Richard Paige McMaster University Canada
 Corina S Pasareanu Carnegie Mellon University Silicon Valley, NASA Ames Research Center	 Lori Pollock University of Delaware, USA United States	 Michael Pradel TU Darmstadt Germany	 Abhik Roychoudhury National University of Singapore Singapore	 Julia Rubin University of British Columbia Canada	 Koushik Sen University of California, Berkeley United States	 Eleni Stroulia University of Alberta Canada	 Lin Tan Purdue University United States
 Frank Tip Northeastern University United States	 Chao Wang University of Southern California United States	 Dongmei Zhang Microsoft Research, China China	 Andrea Zisman The Open University United Kingdom				

Program Committee ICSE Technical Papers

 Jonathan Aldrich Carnegie Mellon University United States	 Aldeida Aleti Monash University Australia	 Dalal Alrajeh Imperial College London United Kingdom	 Samik Basu Iowa State University	 Benoit Baudry KTH Royal Institute of Technology, Sweden Sweden	 Gabriele Bavota Università della Svizzera Italiana (USI) Switzerland	 Nelly Bencomo Aston University United Kingdom	 Ayse Bener Ryerson University Canada
 Domenico Bianculli University of Luxembourg Luxembourg	 Christian Bird Microsoft Research United States	 Kelly Blincoe University of Auckland New Zealand	 Barbara Buhnova Masaryk University	 Marcel Bohme Monash University Australia	 Maria Christakis MPI-SWS Germany	 Siobhán Clarke Trinity College Dublin, Ireland Ireland	 James Clause University of Delaware United States
 Rob DeLine Microsoft Research	 Danny Dig School of EECS at Oregon State University	 Yvonne Dittrich IT University of Copenhagen, Denmark	 Hakan Erdogmus Carnegie Mellon University United States	 Robert Feldt Chalmers University of Technology	 Maria Angela Ferraro University of London United Kingdom	 Antonio Filieri Imperial College London United Kingdom	 Thomas Fritz University of Zurich, University of British Columbia Switzerland
 Diego Garbervetsky University of Buenos Aires, Argentina Argentina	 Jaco Geldenhuys University of Stellenbosch, South Africa South Africa	 Milos Gligoric University of Texas at Austin United States	 Michael W. Godfrey University of Waterloo, Canada Canada	 Alex Groce Northern Arizona University United States	 John Grundy Monash University Australia	 Lars Grunske Humboldt-Universität zu Berlin Germany	 Paul Grünbacher JKU Linz, Austria Austria
 Arie Gurfinkel University of Waterloo Canada	 William G.J. Halfond University of Southern California United States	 Tracy Hall Brunel University United Kingdom	 Sylvain Hallé Université du Québec à Chicoutimi, Canada Canada	 Dan Hao Peking University China	 Mark Harman Facebook and University College London United Kingdom	 Rachel Harrison University of Oxford United Kingdom	 Emily Hill Drew University United States
 Rashina Hoda The University of Auckland New Zealand	 Reid Holmes University of British Columbia Canada	 Jennifer Horkoff Chalmers and University of Gothenburg Sweden	 Jeff Huang Texas A&M University	 James Jones University of California, Irvine United States	 Sarfraz Khurshid University of Texas at Austin United States	 Moonzoo Kim KAIST Korea, South	 Dimitris Kolovos University of York United Kingdom
 Patricia Lago Vrije Universiteit Amsterdam Netherlands	 Wei Le Iowa State University	 Emmanuel Letier University College London United Kingdom	 Grace Lewis Carnegie Mellon Software Engineering Institute United States	 Antónia Lopes University of Lisbon Portugal	 Sam Malek University of California, Irvine United States	 Shahar Maoz Tel Aviv University Israel	 Wes Masri American University of Beirut Lebanon
 Na Meng Virginia Tech United States	 Marija Mikic Google United States	 Raffaella Mirandola Politecnico di Milano Italy	 Henry Muccini University of L'Aquila, Italy Italy	 Sarah Nadi University of Alberta Canada	 Nachiappan Nagappan Microsoft Research United States	 Shiva Nejati SRI Centre/University of Luxembourg Luxembourg	 Tien Nguyen University of Texas at Dallas United States
 Liliana Pasquale University College Dublin & Lero, Ireland Ireland	 Marco Pistoia IBM Research United States	 Adam Porter University of Maryland United States	 Paul Ralph University of Auckland New Zealand	 Cindy Rubio-Gonzalez University of California, Davis United States	 Caitlin Sadowski	 Anita Sarma Oregon State University United States	 Federica Sarro University College London, UK United Kingdom
 Ina Schaefer Technische Universität Braunschweig Germany	 Carolyn Seaman University of Maryland Baltimore County United States	 Alexander Serebrenik Eindhoven University of Technology Netherlands	 Jocelyn Simmonds University of Chile Chile	 Kathryn Stolee North Carolina State University United States	 Zhendong Su ETH Zurich Switzerland	 Gabriele Taentzer Universität Marburg Germany	 Christoph Treude The University of Adelaide Australia
 Burak Turhan Brunel University	 Daniel Varro McGill University / Budapest University of Technology and Economics Canada	 Bogdan Vasilescu Carnegie Mellon University United States	 Helene Waeselynck LAAS-CNRS France	 Westley Weimer University of Michigan United States	 Jim Whitehead University of California, Santa Cruz United States	 Xin Xia Monash University Australia	 Tao Xie University of Illinois at Urbana- Champaign
 Yingfei Xiong Peking University China	 Tuba Yavuz University of Florida United States	 Cemal Yilmaz Sabanci University Turkey	 Xiangyu Zhang Purdue University	 Minghui Zhou Peking University China	 Ying Zou Queen's University, Kingston, Ontario Canada	 Marcelo d'Amorim Federal University of Pernambuco Brazil	 Cleudson de Souza Vale Institute of Technology and Federal University of Pará Belém, Brazil
 Arie van Deursen Delft University of Technology Netherlands	 Andre van der Hoek University of California, Irvine United States						

Write for the Reviewer

- Make problem statement and contribution very clear
- Use established outline (e.g., see [thesisguide](#))
- Make text easily readable. This is hard and exhausting work. But you can learn it, this is no issue of talent.

My Personal Best Practices

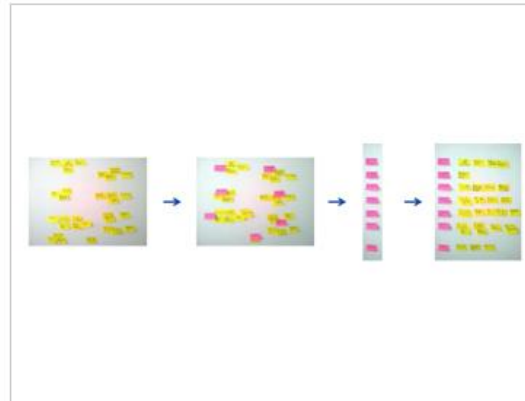
- Block writing time
- Begin with outline
- Separate writing from improving
- Write complete paragraphs before improving them
- Let text „cool down“ and proof-read it later again
- There is not the one silver-bullet way of writing

English Writing Center

- Free one-to-one consulting with native English speakers
 - GR, Thesis, Homework, CV etc.
 - Text needs not to be ready

<https://www.sprachenzentrum.tum.de/sz/sprachen/englisch/english-writing-center/>

Prepare Presentation

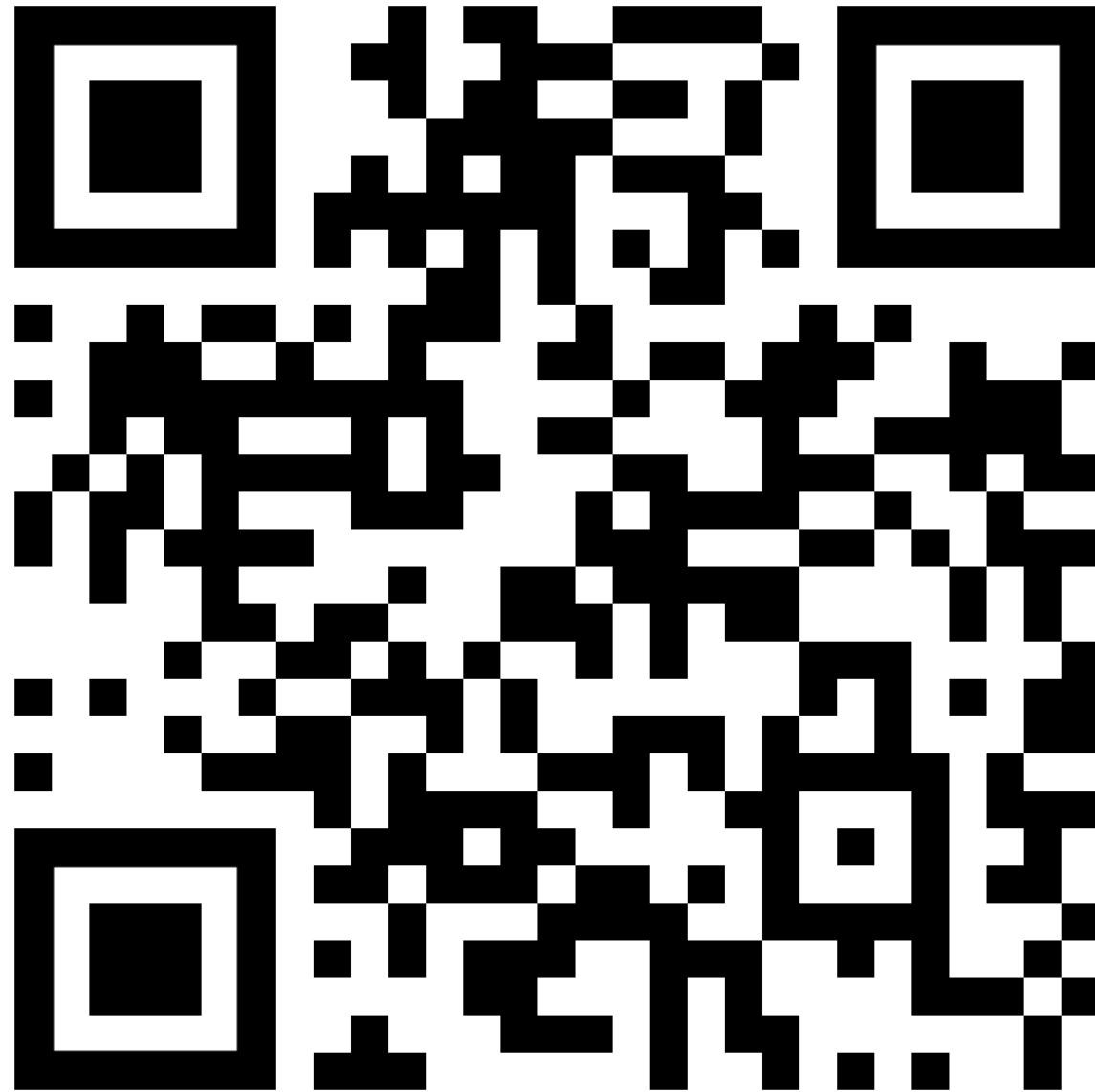


<https://thesisguide.org/2015/03/04/how-to-draft-your-presentation/>

Presentation Differences to BA/MA

- Rehearsal talk with supervisor
- Practice it in English
- Formulate starting sentences and learn them by heart
- Backup slides for questions (e.g., more details)

Do you want to do your own research and get to know the research community? Then a guided research is the best you can do!



<https://cqse.eu/feedback-tum-talk>

Thank you!

If you are interested in a guided research in the field of software analysis and testing, please let me know:

haas@cqse.eu

More Info:

www.thesisguide.org

Feedback:



<http://cqse.eu/feedback-tum-talk>