

Software Trajectory Analysis: An empirically based method for automated software process discovery

Pavel Senin
IRISA, campus de Beaulieu
35042 Rennes Cedex, France
pavel.senin@inria.fr

ABSTRACT

A process defines a set of routines which allow to organize, manage and improve activities to reach a goal. With an expert intuition and apriori knowledge software processes were modeled for the long time resulting in Waterfall, Spiral and other development models. Later, with a wide use of SCM systems and public availability of primitive software process artifact trails, formal methods such as Petri Nets, State Machines and others were applied for recurrent processes discovery and control. Recent advances in metrication effort, increased use of continuous integration and extensive documentation of the performed process make information-rich fine-grained software process artifacts trails available for analysis. Potentially this fine-grained data may shed more light on the software process. In this work I propose to investigate an automated technique for the discovery and characterization of recurrent behaviors in software development - “programming habits” either on an individual or a team level.

1. INTRODUCTION AND MOTIVATION

A *software process* is a set of activities performed in order to design, develop and maintain software systems. Examples of such activities include design methods; requirements collection and creation of UML diagrams; testing and performance analysis. The intent behind a software process is to structure and coordinate human activities to achieve the goal - deliver a software system successfully; while the objective of software process research is to design a process consisting of a set of recurrent behaviors which allows to reproduce such a success consistently and improve the process performance.

Much work has been done in software process research resulting in a number of industrial standards for process models (CMM, ISO, PSP etc. [3]) which are widely accepted by many government and industrial institutions. Nevertheless,

software development remains error-prone and more than a half of all software development projects ending up failing or being very poorly executed. Some of them are abandoned due to running over budget, some are delivered with such low quality or so late that they are useless, and some, when delivered, are never used because they do not fulfill requirements [10]. The cost of this lost effort is enormous and may in part be due to our incomplete understanding of software process.

There is a long history of software process improvement through proposing specific patterns of software development. For example, the Waterfall Model process proposes a sequential pattern in which developers first create a Requirements document, then create a Design, then create an Implementation, and finally develop a Tests. On the contrary, the Test Driven Development process proposes an iterative pattern in which the developer must first write a test case, then write the code to implement that test case, then refactor the system for maximum clarity and minimal code duplication. Probably the main problem with this traditional top-down approach to process development is that it requires the developer or manager to notice a recurrent pattern of behavior in the first place [3].

As an alternative to that, in my research, I am applying knowledge discovery and data mining techniques to the domain of software engineering in order to evaluate their ability to automatically notice interesting recurrent patterns of behavior from collected software process artifacts. While I am not proposing to be able to infer a complete and correct software process model, my system will provide its users with a formal description of recurrent behaviors in their software development. As a simple example, consider a development team in which committing code to a repository triggers a build of the system. Sometimes the build passes, and sometimes the build fails. To improve the productivity of the team, it would be useful to be aware of any recurrent behaviors of the developers. My system might generate one recurrent pattern consisting of a) implementing code b) running unit tests, c) committing code and d) a passed build: $i \rightarrow u \rightarrow c \rightarrow s$, and another recurrent pattern consisting of a) implementing code, b) committing code, and c) a failed build: $i \rightarrow c \rightarrow f$. The automated generation of these recurrent patterns can provide actionable knowledge to developers; in this case, the insight that running test cases prior to committing code reduces the frequency of build failures.

The rest of the paper is organized as follows: next I will briefly introduce the most relevant prior approaches for the

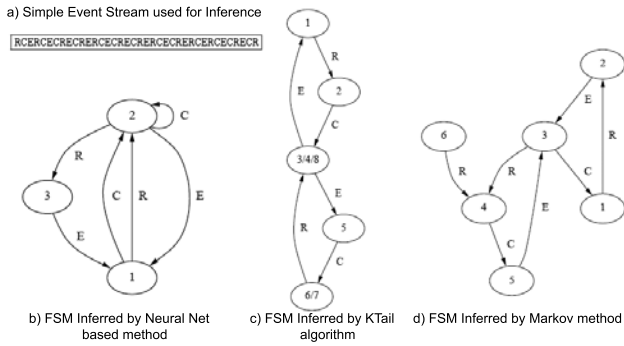


Figure 1: Process discovery through the grammar inference: panel a) a sample event stream (simple process involving three types of events: Edit, Review, and Checkin); and FNA results obtained by applying three methods of process discovery from Cook & Wolf [6].

problem. Further, I will present research objectives, approach and data analysis methods. Pilot study results, future directions for exploratory research and empirical evaluation will conclude the paper.

2. RELEVANT PRIOR WORK

Although process mining in the business domain is a well-established field with much software developed up to date (ERP, WFM and other systems), “Business Process Intelligence” tools usually do not perform process discovery and typically offer relatively simple analyzes that depend upon a correct a-priori process model [25] [23] [1]. This fact restricts direct application of business domain process mining techniques to software engineering, where processes are usually performed concurrently by many agents, are more complex and typically have a higher level of noise. Taking this fact in account, I will review only the approaches to the process mining for which applicability to software process was expressed.

Perhaps, the research most relevant to my own was done by Cook & Wolf in [6]. The authors developed a “*process discovery*” techniques intended to discover process models from event streams. The authors did not attempt to generate a complete model, but rather to generate sub-models that express the most frequent patterns in the event stream. They designed a framework which collects process data from history logs, and generates a set of recurring patterns of behavior characterizing observed process. In this work they extended two methods of *grammar inference* from previous work: purely statistical (neural network based *RNet*) and purely algorithmic (*KTail*) as well as developing their own Markovian method (*Markov*).

The first method extended by the authors, the neural-network based grammar inference, *RNet* algorithm, defines a recurrent neural network architecture which is trained by the sequences of events. After training, this neural net is able to characterize a current system state by looking on past behavior. The authors extract the FSM from the trained neural network by presenting different strings to it and extracting the hidden neurons activity through observations. Due to the nature of Neural Net, closely related activation patterns are clustered into the same state; therefore, by not-

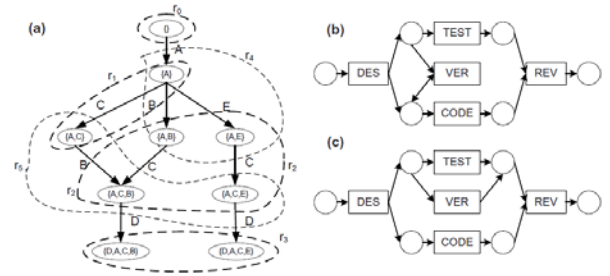


Figure 2: Illustration of the “Generation and Synthesis Approach” from [24]: a) Transition System with regions shown; b),c) Petri Nets synthesized from the Transition System.

ing the current pattern, the input token, and the next activation pattern, transitions are recorded and compiled into the inferred FSM.

The second method investigated, is a purely algorithmic *KTail* method, which was taken from the work of Biermann & Feldman [2]. The idea is that a current state is defined by what future behaviors can occur from it. The *future* is defined as the set of next k tokens. By looking at a window of successor events, the *KTail* algorithm can build the equivalence classes that compose the process model. The authors extensively modified the original *KTail* algorithm improving the folding in the mined model making it more robust to noise.

The Markov based method developed by the authors is based on both algorithmic and statistical approaches. It takes to account past and future system behavior in order to guess the current system state. Assuming that a finite number of states can define the process, and that the probability of the next state is based only on the current state (Markov property), the authors built a n^{th} -order Markov model using the first and second order probabilities. Once built, the transition probability table corresponding to the Markov model is converted into FSM which is in turn reduced based on the user-specified cut-off threshold for probabilities.

The authors implemented all three of these algorithms in a software tool called DAGAMA as a plugin for larger software system called Balboa [4]. By performing benchmarking, Cook & Wolf found that the Markov algorithm was superior to the two others. *RNet* was found to be the worst of the three algorithms. The software tool was applied to a real-world process data and demonstrated an abstraction of the actual process executions and ability to capture important properties of the process behavior. The major drawback of the approach, as stated by the authors, lies in the inability of the FSMs to model concurrency of processes which limits its applicability to the software development process. Later, Cook et al. in [5] addressed this limitation by using Petri-nets and Moore-type FSM.

Another set of findings relevant to my research approach was developed by Rubin et al. [22] and van der Aalst et al. [23] and is called *incremental workflow mining*. The authors not only designed sophisticated algorithms but built a software system using a business process mining framework called ProM by van Dongen et al. [24] which synthesizes a Petri Net corresponding to the observed process. The system was tested on SCM logs and while the process artifacts

retrieved from the SCM system are rather high-level, the approach discussed is very promising for the modeling of software processes from the low-level product and process data.

The algorithm input is an event chain constructed through the “abstraction on the log level”, which aggregates basic events into single high-level entities, is treated with the *Generate* part of the “*Generate and Synthesis*” [23] algorithm in order to generate a *Transition System* which represents an ordered series of events. This algorithm looks at the history (prefix) and the future (suffix) sequences of events related to the current one in order to discover transitions. When applied to the abstracted log information, the algorithm generates a rather large Transition System graph where edges connect to abstracted events. This transition system is then successively simplified by using various reduction strategies. At the last step of the incremental workflow mining approach, Transition Systems are used to *Synthesize* labeled Petri nets (where different transition can refer to the same event) with the help of “*regions theory*” [7]. As with the Transition System generation, the authors investigate many different strategies of Petri nets synthesis, showing significant variability in the results achieved. (see Figure 2). The significant contribution of this research is in the generality of the method. It was shown that by tuning the “Generate” and “Synthesize” phases it is possible to tailor the algorithm to a wide variety of processes. In particular, as mentioned before, Rubin et al. successfully applied this framework to the SCM logs and audit trails analysis.

In addition to discussed work latest trends in software process research emphasize mining of software process artifacts and behaviors [25] [12] [19].

It worth noting that while reviewed work demonstrate general approaches to modeling of concurrent processes, to the best of my knowledge their application to the real world data was very limited. This may be partially due to the high level of noise and computational limitations which made inferred models too complex and impractical to use; or maybe due to the lack of the means of fine-grained software process artifacts collection which made models too abstract. In my approach I am planning to address both issues by leveraging the ability of the Hackstat system [13] to collect such a fine grained data and by the application of advanced symbolic and temporal data mining techniques resistant to noise.

3. RESEARCH OBJECTIVES

As shown by previous research, it is possible to infer and successively formalize sequential software process by observing its artifacts, and in particular, recurrent behavioral patterns. The problem of finding such patterns is the cornerstone of my research. Solving this will extend previous research with new knowledge that will support improvements in our understanding of software process.

The main research objectives of my work is a design, development and evaluation of a previously unexplored approach to discovering of recurrent behaviors in software process through the data mining of low-level process and product artifacts.

To reach this ultimate goal I am planning to perform a number of steps. At first I am exploring the applicability of proposed pattern-mining techniques to the various levels of software process artifacts complexity. This exploratory study will result in a set of well-defined data mining work-

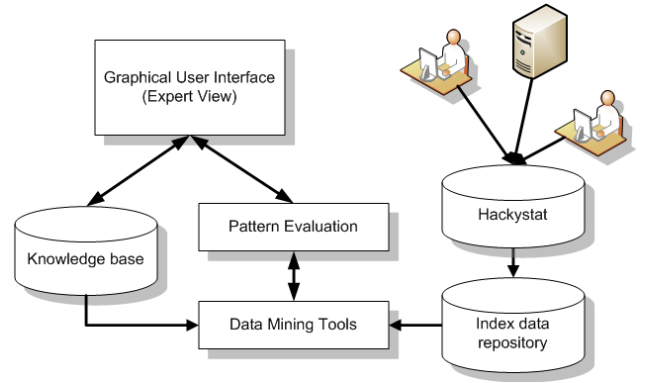


Figure 3: The high-level system overview. Software engineering process and product data are collected and aggregated by Hackstat and then used to generate temporal symbolic indexes. Data mining tools constrained by software engineering domain knowledge are then used for unsupervised patterns discovery. The GUI provides an interface to the discovered patterns and aids in investigation of a discovered phenomena.

flows designed for particular kinds of software process data and research goals. In turn, workflow implementations will be used in the an expert software package aiding recurrent behavior discovery. The software and results of its application will be empirically evaluated at the final phase of my research.

4. RESEARCH APPROACH

My approach to this problem rests on the application of data-mining techniques to symbolic time-point and time-interval series constructed directly from the software process artifacts such as SCM logs, software audit trails or telemetry streams provided by Hackstat or by a similar in-process software development monitoring system.

To investigate the requirements for a software tool that aids in the discovery of recurrent behavioral patterns in software process, I am designing and developing the “Software Trajectory” framework. A high-level overview of the framework is shown in Figure 3 and resembles the flow of the “Knowledge Discovery in Database” process discussed by Han et al. in [11]. As shown, the data collected by Hackstat is transformed into a symbolic format and then indexed for further use in data-mining. The tools, designed for data-mining, have a specific restrictions placed on the search space by domain and context knowledge in an attempt to limit the amount of reported patterns to useful ones. I am planning to design a GUI in a way that will allow easy access and modification of these restrictions.

While I plan to use Hackstat as a primary data collection system in my future experimental setup, for the current exploratory study I am working with various ways of data collection, extraction and abstraction. The ability to abstract into symbolic representation various software artifact streams such as SCM logs and audit trails is crucial for the exploratory phase enabling me to experiment with much broader spectrum of existing software process data collections.

5. DATA AND ANALYSIS METHODS

As mentioned before, this research based on the mining of recurrent behavioral patterns from a symbolic representation of software processes as temporal sequences of events performed concurrently by agents - individual developers or automated tools. Various types and complexity of data are available for such software process analysis. For example, simplest data extracted from SCM logs allow to track software change and perform a basic software process reconstruction - however a very little information can be recalled about individual development activities and behaviors from such models. When development audit trails added to analyzes - enriched data allow to shed more light on the performed process, while adding considerable amount of noise to the models. By using logs of an automated build system it is possible to collect a broader spectrum of artifacts and software metrics such as Cyclomatic Complexity or Coverage, moreover centralized build system creates a set of singular points which are tying together in time processes performed concurrently by developers. Finally, collecting temporal data about atomic development events by the means of Hackstat and similar systems - such as background compilation and test results, IDE buffer transfers along with quantification of the actual effort enables characterization of the dynamic behavior of a software process in great detail. However, as shown by previous experience - the amount of noise will be unbearable high for conventional methods of process inference.

Having this limitation in mind, approach I am taking in my research differs from previous attempts by its computational simplicity and inspired by recent advances in two research fields: in temporal patterns mining from time-series and in motif finding in biological sequence analysis.

Current state of the art approach in temporal data mining called Symbolic Aggregate approXimation and was proposed by Lin et al. in [18]. This method extends the PAA-based approach, inheriting algorithmic simplicity and low computational complexity, while providing satisfiable sensitivity and selectivity in range-query processing. Moreover, as mentioned by the authors, the use of a symbolic representation opens the door to the existing wealth of data-structures and string-manipulation algorithms in computer science such as hashing, regular expression pattern matching, suffix trees etc. As showed by the authors, SAX outperforms all previously known methods - DFT, DWT, DTW and similar - by using SAX it is possible to index and find recurrent and "surprise" patterns in vast amounts of data in almost linear time and space [15] [17] [14]. Looking at the SAX-performed data reduction one might notice that due to the PAA step of the algorithm high level of noise is reduced dramatically - leaving only major time-series trends for "symbolization". In my opinion, this protocol can be adopted for the complexity reduction of the software process artifacts streams. For example, if symbolic approximation of performed process would be $i \rightarrow i \rightarrow i \rightarrow u \rightarrow i \rightarrow i \rightarrow i \rightarrow u \rightarrow c \rightarrow s$ indicating long cycle of development with intermediate unit test at some point (position 4) it is possible to dismiss this unit test event and collapse in time all development events into the single one - $i \rightarrow u \rightarrow c \rightarrow s$ reducing the pattern complexity. Another idea borrowed from the SAX and Bioinformatics sequence alignment is in the distance function implementation - which will score sequences different by a single symbol (single substitution) as equal if the both

symbols are approximately equal to each other in the context of actions. For example, one can say that running static code analysis tools such as CheckStyle, PMD or running Unit test is equal in the meaning - checking code quality and one common symbol would be assigned to all such events.

Multiple alignments of protein sequences are important in many biological applications including phylogenetic tree estimation, secondary structure prediction and critical residue identification. Recent advances in MSA algorithms [8] include a new method for design and evaluation of objective functions for improving alignments by its profiling which is done by statistically inferred weighting schemas. Within the multiple alignment of biological sequences algorithms usually introducing gaps and substitution within possibly similar sequences to find the best possible score for an alignment, clever weighting of such discrepancies was shown to improve alignments. This strategy can be also adopted within the search of software development motifs - for example by ignoring (scoring low) presence or lack of code analysis events within a development cycle - it will immediately boost the evidence for both types of patterns - making them much more noticeable while preserving the overall contextual correctness.

6. PRELIMINARY RESULTS

During my work on the pilot version of Software Trajectory framework, I began a set of small exploratory experiments in order to aid in the architectural design and algorithms implementation. In addition, these experiments helped me to outline the boundaries of applicability of my approach to certain problems in software engineering. I call these experiments *Pilot study*.

One of the experiments was performed to evaluate the ability of SAX approximations and indexing to capture a temporal specificity of development time telemetry streams through the discovery of recurrent temporal patterns. Knowing about the frequently misleading results of a time-series clustering [16], I did not expect to capture many interesting facts, nevertheless the results were encouraging. The streams under analysis were composed by quantifying development events within a day: for example for the BuildSuccess stream a single value 5 means that five builds succeeded within that day. The data used in this study was collected from eight students and represents Hackstat metrics collected during sixty days of a classroom project. The two clustering experiments were conducted using the distance between vectors of motif frequencies (note that temporal ordering not accounted in these analysis) extracted by indexing of telemetry streams. The time span of a motif in this experiment was set to one week.

At first, by clustering of development time telemetry streams collected from individual developers I was able to group developers with similar behavioral patterns within clusters, which indicates the feasibility of the classification by effort approach. Figure 5 depicts results of this analysis.

Clustering of product-related telemetry streams however was not successful. I was able to group telemetry streams, but while these groups look intuitively meaningful - for example clustering together filemetrics and development time - the close examination of the stream features suggests that this grouping happened due to the similar temporal behavior on the short stretches rather than positive correlation. This result, while proving the correctness of approach, in-

icates its limitation, pointing that instead of using just motif frequencies, some temporal ordering should be taken into account. Figure 4 displays results of this analysis.

Another result was also achieved by the direct application of indexing of symbolic approximation of the data streams. For this experiment I was using data collected from my own concurrent development of Trajectory software package and JMotif library. Within a symbolic approximation paradigm the change of symbols can be interpreted as growth or drop in the underlying data stream: for example consecutive *ab* symbols indicate that growth happened and *ba* correspond to a drop. The idea was to find if there were any drop or growth coincidences within two streams. For this purpose I defined a sliding window consisting of three days and searched such coinciding within a window events. Events were found within a development time stream, churn, commit and others. Detailed analysis of coinciding growth events revealed that changes within JMotif data manipulation routines almost always followed by changes within Trajectory data analyzes routines which heavily relying on the JMotif API. While within my software project portfolio such a discovery is a quite obvious fact one can see the immediate application of such a method to large portfolio or software system - by performing such analyzes it will be possible to build a repertoire of frequent consecutive changes and also estimate an effort required for performing such changes extrapolating previous experience.

7. EXPLORATORY STUDY DIRECTIONS

Previous section describes two finished experiments which provide an inside on the applicability of SAX approximation of Telemetry streams to discovery of recurrent behaviors in software process. Work on both experiments resulted in a database schema which supports symbolic indexing of telemetry streams. Such a storage solution provides an immediate and structured access to the symbolic streams reflecting software process allowing to perform time-range queries with use of Regular Expressions (in SQL language) and extract sequences of interest which for example include successful Unit-Tests, Commits, or Cyclomatic Complexity above an arbitrary threshold value. Such a database solution is somewhat similar to the CVSanaly tool [21] providing an abstraction layer over the raw, mostly unstructured data. One of the direction I am working on right now is the complimentary use of both databases under unifying API in order to evaluate a capacity of such approach in the recurrent behavior discovery for the MSR challenge [9].

Motif finding in protein sequences is a well-established research field which deals with vast amount of data as well as with high level of noise. This noise problem is inevitable in biology and reflects a natural way of evolution in life through mutations. Recent advances in biological data mining resulted in a wealth of algorithms and software tools. According to [20] there are two categories of approaches for motif finding in biological sequences: while the first is based on the similar to reviewed concepts such as HMM, the second category “uses patterns with ‘mismatch representation’ and define a signal to be a consensus pattern allowing up to a certain number of mismatches to occur in each instance of the pattern”. The second approach looks very promising as an extension of exact motifs search explained before. By allowing mismatches the noise from non-frequent events will not decrease the evidence value for the similar patterns allow-

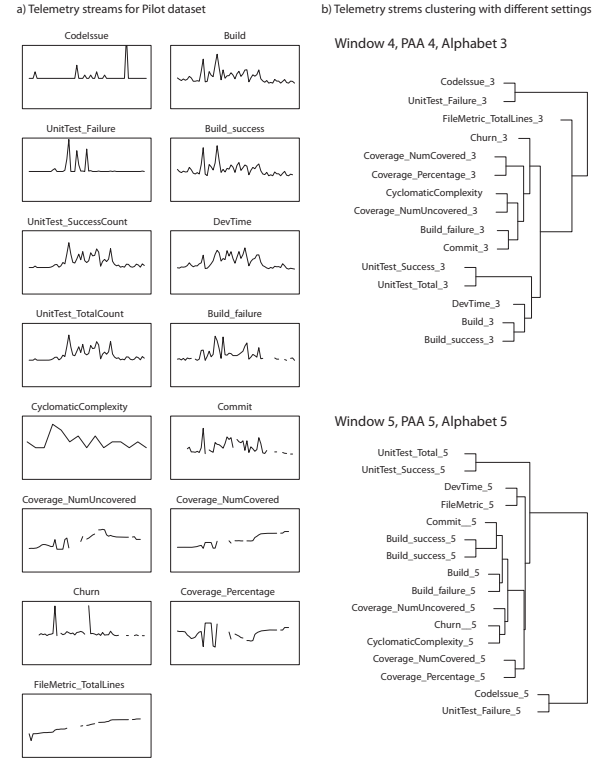


Figure 4: Clustering of telemetry streams for classroom pilot dataset using symbolic approximation and vectors of motif frequencies. While it seems to be meaningful to find correlation between *UnitTest_Failure* and *CodeIssue* streams unit test, this grouping happened due to the similarity of behavior pattern - short, high amplitude bursts; but note, there is no correlation of features in time .

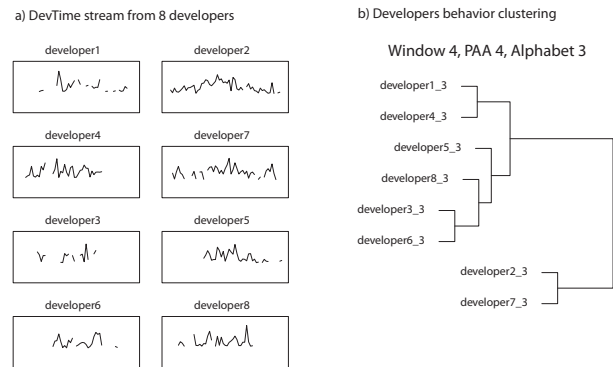


Figure 5: Clustering of developers behavior using symbolic approximation and vectors of motif frequencies. This analysis captured similar development behavior among developers. Developers #2 and #7 were consistent (no bursts observed) in both, coding and measuring effort during whole time interval, while all others can be characterized with bursty, inconsistent effort.

ing to identify key events of the observed process (as similar to the key amino-acids residues within protein sequences). I am planning to perform an experiment consisting of abstraction of software process artifacts into the symbolic representation mimicking protein encoding and application of motif-finding packages after evaluating their ability to recall frequent patterns from such a data.

Concurrent software process occurs within a team working on a common goal. Within such a process developers interact at many levels, but in the context of this work the focus put on the role of SCM and build automation (such as continuous integration) in their ability to serve as a place which persists among all levels and creates means of interaction for all agents. I plan to perform analyzes on telemetry streams coming from individual agents using SCM and automation events as the time-points breaking such individual continuous streams onto short intervals representing fragment sets of concurrent process. I expect to find difference among such a fragment sets reflecting different phases of development as well as differences among sets from individual agents reflecting team-assigned roles.

8. FUTURE EMPIRICAL STUDY DESIGN

It is proposed to conduct two case studies: *Public data case study*, and *Classroom case study* in order to empirically evaluate the capabilities and performance of Software Trajectory framework. These studies differ in the granularity of data used, and in the approaches for evaluation.

My intent behind these empirical studies is to assess the ability of Software Trajectory framework to recognize well known recurrent behavioral patterns and software processes (for example Test Driven Development), as well as its ability to discover new ones. In addition, these studies will support a classification and extension of the current Hackystat sensor family in order to improve Software Trajectory's performance. It is quite possible that some of the currently collected sensor data will be excluded from the Software Trajectory datasets, while some new ones will be designed and developed in order to capture important features from the studied software development data streams.

The proposed public data case study is based on the use of publicly available Software Configuration Management (SCM) audit trails of the big, ongoing software projects such as Eclipse, GNOME etc. Mining of SCM repositories is a well-developed area of research with much work published [9]. SCM repositories contain coarse software product artifacts which are usually mined with a purpose of discovering of various characteristics of software evolution and software process. I am using a mixed-method approach in this study. In the first phase of this study, I plan to perform SCM audit trail data mining following published work and using Software Trajectory as a tool in order to discover confirmed patterns in software process artifacts, and thus quantitatively evaluate Software Trajectory's performance when compared to existing tools. In the second phase, I will develop my own pre-processing and taxonomy mapping of software process artifacts into temporal symbolic series. By using this data and Software Trajectory framework, I plan to develop a new approach for SCM audit trail mining and possibly discover new evolutionary behaviors within software process.

The classroom case study is based on a more comprehensive data set. This data will be collected by Hackystat from Continuous Integration and from individual developers and

will contain fine-grained information about performed software process which may or may not reflect a development practice restriction placed by the lecturer. The approach I am taking in this study is very similar to the public data case study. I will develop my own taxonomy for mapping of software process artifacts into symbolic temporal data and will apply Software Trajectory analyzes to this data in order to discover recurrent behaviors. In turn, these discovered knowledge will be evaluated through interviewing for usefulness and meaningfulness.

9. RESEARCH RISKS

While I strongly believe that my research will result in the discovery of novel software processes as well as Software Trajectory package will contribute to the research community, it may be the case that I underestimate the overall complexity of the problem and the approach taken will not be able to solve it. In this section I will assess the risks associated with my research plan as well as explain my thoughts about these issues.

First of all, one of the big risks for my research lies in the question if I will be able to find any of the recurrent behaviors while mining software process artifacts. The failure to find patterns maybe due to the nature of the software development as a human activity. It is possible that there are no recurrent behaviors within such a process, in other words, it may be that software development is a craft - a practical art, which in turn makes every programmer an artist with a very own, unique behavior. While such a result will be "terminal" for my research, the impact it makes is worth getting my research done. Another possibility which yields a similar result - no recurrent behaviors - may be due to the inefficiency of process monitoring and data collection means. If this will be a case - my research will outline directions for further development of Hackystat and similar systems by identifying what kind of the activities need to be captured to support future development in the software process research. Finally, it may be possible that patterns exist as well as they are represented within the data, but methods chosen are not sensitive or not efficient enough to find them. By detecting this case through the classroom experiment setup and interviewing I will investigate and outline limitations of implemented methods.

Secondly, I may discover some novel behavioral patterns through mining, but will not be able to confirm them experimentally. This situation maybe due to the number of reasons - for example, interviewee may not notice and confirm such behaviors or will hesitate to confirm "not good ones" intentionally. If this happens, it is hard to see at this point which way to proceed further, but one of the vital options will be to set up additional targeted experiments for evaluation of findings.

And finally, it maybe that I will find patterns and will confirm them empirically, but they will be obvious and only trivial ones - for example building a system before committing changes or doing an update before running a build. This result is very similar to the case when no patterns are identified and will be resolved in the similar fashion.

10. CONCLUSIONS

This paper presented an approach to discovering of novel recurrent behaviors in software process. Summarizing the

experience collected within the exploratory study I can only see that properties of this approach and its current implementation in the Software Trajectory framework appear to be very promising. Application of current advances in temporal symbolic data mining and Bioinformatics allows to overcome many computational limitations in existing approaches for mining of software process artifacts; while Hackystat provides the ability to capture fine-grain software product and process metrics providing a richness of data, which, potentially, might reveal new insights.

11. REFERENCES

- [1] R. Agrawal, D. Gunopulos, and F. Leymann. Mining process models from workflow logs. pages 467–483. 1998.
- [2] A. W. Biermann and J. Feldman. On the synthesis of finite-state machines from samples of their behavior. *IEEE Trans. Computers*, pages 592–597, 1972.
- [3] R. Conradi. SPI frameworks: TQM, CMM, SPICE, ISO 9001, QIP experiences and trends - norwegian SPIQ project.
- [4] J. E. Cook. *Process discovery and validation through event-data analysis*. PhD thesis, Boulder, CO, USA, 1996.
- [5] J. E. Cook, Z. Du, C. Liu, A. L. Wolf, and Er. Discovering models of behavior for concurrent workflows, 2004.
- [6] J. E. Cook and A. L. Wolf. Discovering models of software processes from event-based data. *ACM Trans. Softw. Eng. Methodol.*, 7(3):215–249, July 1998.
- [7] J. Cortadella, M. Kishinevsky, L. Lavagno, and A. Yakovlev. Deriving petri nets from finite transition systems. *Computers, IEEE Transactions on*, 47(8):859–882, 1998.
- [8] R. Edgar. Muscle: a multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics*, 5(1):113+, August 2004.
- [9] H. Gall, M. Lanza, and T. Zimmermann. 4th international workshop on mining software repositories (MSR 2007). In *Software Engineering - Companion, 2007. ICSE 2007 Companion. 29th International Conference on*, pages 107–108, 2007.
- [10] W. W. Gibbs. Software’s chronic crisis. *Scientific American*, September 1994.
- [11] J. Han and M. Kamber. *Data Mining, Second Edition: Concepts and Techniques (The Morgan Kaufmann Series in Data Management Systems) (The Morgan Kaufmann Series in Data Management Systems)*. Morgan Kaufmann, January 2006.
- [12] C. Jensen and W. Scacchi. Guiding the discovery of open source software processes with a reference model. pages 265–270. 2007.
- [13] P. M. Johnson, S. Zhang, and P. Senin. Experiences with hackystat as a service-oriented architecture. Technical Report CSDL-09-07, Department of Information and Computer Sciences, University of Hawaii, Honolulu, Hawaii 96822, Los Angeles, California, February 2009.
- [14] E. Keogh, K. Chakrabarti, M. Pazzani, and S. Mehrotra. Dimensionality reduction for fast similarity search in large time series databases. *Knowledge and Information Systems*, 3(3):263–286, 2001.
- [15] E. Keogh, J. Lin, and A. Fu. Hot sax: efficiently finding the most unusual time series subsequence. In *Data Mining, Fifth IEEE International Conference on*, pages 8 pp.+, 2005.
- [16] E. Keogh, J. Lin, and W. Truppel. Clustering of time series subsequences is meaningless: Implications for previous and future research. In *ICDM ’03: Proceedings of the Third IEEE International Conference on Data Mining*, Washington, DC, USA, 2003. IEEE Computer Society.
- [17] E. Keogh, S. Lonardi, and B. . Chiu. Finding surprising patterns in a time series database in linear time and space. In *KDD ’02: Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 550–556, New York, NY, USA, 2002. ACM.
- [18] J. Lin, E. Keogh, L. Wei, and S. Lonardi. Experiencing SAX: a novel symbolic representation of time series. *Data Mining and Knowledge Discovery*, 15(2):107–144, October 2007.
- [19] D. Lo, H. Cheng, J. Han, S. C. Khoo, and C. Sun. Classification of software behaviors for failure detection: a discriminative pattern mining approach. In *KDD ’09: Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 557–566, New York, NY, USA, 2009. ACM.
- [20] C. K. Reddy, Y.-C. C. Weng, and H.-D. D. Chiang. Refining motifs by improving information content scores using neighborhood profile search. *Algorithms for molecular biology : AMB*, 1(1):23+, November 2006.
- [21] G. Robles, S. Koch, J. M. González-Barahona, and J. Carlos. Remote analysis and measurement of libre software systems by means of the cvsanaly tool. In *In Proceedings of the 2nd ICSE Workshop on Remote Analysis and Measurement of Software Systems (RAMSS)*, pages 51–55, 2004.
- [22] V. Rubin, C. Günther, W. van der Aalst, E. Kindler, B. van Dongen, and W. Schäfer. Process mining framework for software processes. pages 169–181. 2007.
- [23] W. van der Aalst, V. Rubin, H. Verbeek, B. van Dongen, E. Kindler, and C. Günther. Process mining: a two-step approach to balance between underfitting and overfitting. *Software and Systems Modeling*, 2009.
- [24] B. F. van Dongen, A. K. A. de Medeiros, H. M. W. Verbeek, A. J. M. M. Weijters, and W. M. P. van der Aalst. The ProM framework: A new era in process mining tool support. pages 444–454. 2005.
- [25] W. Vanderaalst, B. Vandongen, J. Herbst, L. Maruster, G. Schimm, and A. Weijters. Workflow mining: A survey of issues and approaches. *Data & Knowledge Engineering*, 47(2):237–267, November 2003.