

SAMOS - A Framework for Model Analytics and Management

Önder Babur^{a,b}, Loek Cleophas^{b,c}, Mark van den Brand^b

^a*Wageningen University & Research
Wageningen, The Netherlands*

^b*Eindhoven University of Technology
Eindhoven, The Netherlands*

^c*Stellenbosch University
Matieland, Republic of South Africa*

Abstract

The increased popularity and adoption of model-* engineering paradigms, such as model-driven and model-based engineering, leads to an increase in the number of models, metamodels, model transformations and other related artifacts. This calls for automated techniques to analyze large collections of those artifacts to manage model-* ecosystems. SAMOS is a framework to address this challenge: it treats model-* artifacts as data, and applies various techniques—ranging from information retrieval to machine learning—to analyze those artifacts in a holistic, scalable and efficient way. Such analyses can help to understand and manage those ecosystems.

Keywords: model-driven engineering, model analytics, software ecosystems, information retrieval, machine learning

1. Introduction

Model-* engineering approaches such as model-driven engineering and model-based engineering promote the use of models, metamodels and model transformations as first-class citizens to tackle the complexity of building and maintaining software-intensive systems. As such approaches have gained popularity and are applied to larger problems, however, the complexity, size

Email addresses: `Onder.Babur@wur.nl` (Önder Babur), `L.G.W.A.Cleophas@tue.nl` (Loek Cleophas), `M.G.J.v.d.Brand@tue.nl` (Mark van den Brand)

7 and variety of the involved artifacts (notably models) increase. This obser-
8 vation holds for open source as well as industrial ecosystems [1, 2].

9 In this paper we present SAMOS (Statistical Analysis of MOdelS), a
10 framework to address the scalability issue in building and managing model-*
11 ecosystems. It treats collections of models as data and offers various large-
12 scale analysis techniques. It has so far been applied for different modelling
13 languages such as metamodels, feature models, statecharts, domain-specific
14 models, and business process models; and for scenarios such as domain anal-
15 ysis, repository management, clustering, outlier detection, language usage
16 analysis, and clone detection [3, 4, 5, 6, 7, 8, 9, 10, 11]. In this initial public
17 version (1.0), it includes the core framework, feature extraction for Ecore
18 metamodels, as well as domain clustering and clone detection functionalities.
19 We have followed the artifact sharing guidelines [12] to maintain high quality
20 in our work.

21 The rest of this paper is organized as follows. Section 2 summarizes some
22 background concepts. Section 3 outlines the SAMOS architecture and func-
23 tionalities. We illustrate SAMOS for domain clustering and clone detection
24 for Ecore metamodels (Section 4). Section 5 concludes with future work.

25 2. Background

26 Information Retrieval (IR) tackles the challenge of efficiently indexing, an-
27 alyzing and searching various forms of content such as text documents [13].
28 IR consists of a typical workflow. First, documents are collected and indexed,
29 e.g. using a Vector Space Model (VSM) with (1) a vector representation
30 of vocabulary occurrence in a document, named *term frequency*, (2) *zones*
31 (e.g. 'author', 'title'), (3) weighting schemes such as inverse document fre-
32 quency (idf), and zone weights, (4) NLP techniques for handling synonyms.

33 With the VSM we effectively transform each document into an n -dimen-
34 sional vector, yielding an $m \times n$ matrix for m documents. Document similarity
35 can be defined as the distance (e.g. Euclidean or cosine) between vectors in
36 the VSM. We can use these distances for identifying similar groups of docu-
37 ments —i.e. clustering as an unsupervised machine learning (ML) technique.

38 Finally, n -grams [14] are used in computational linguistics to build prob-
39 abilistic models of natural language text, e.g. for estimating the next word
40 given a sequence of words, or comparing text collections based on their
41 n -gram profiles. In essence, n -grams represent a linear encoding of struc-
42 tural context. Alternative (and increasingly more complex) fragmentations

of structural context include sub-trees and sub-graphs. We will refer to n-grams and sub-trees while discussing SAMOS in the next section.

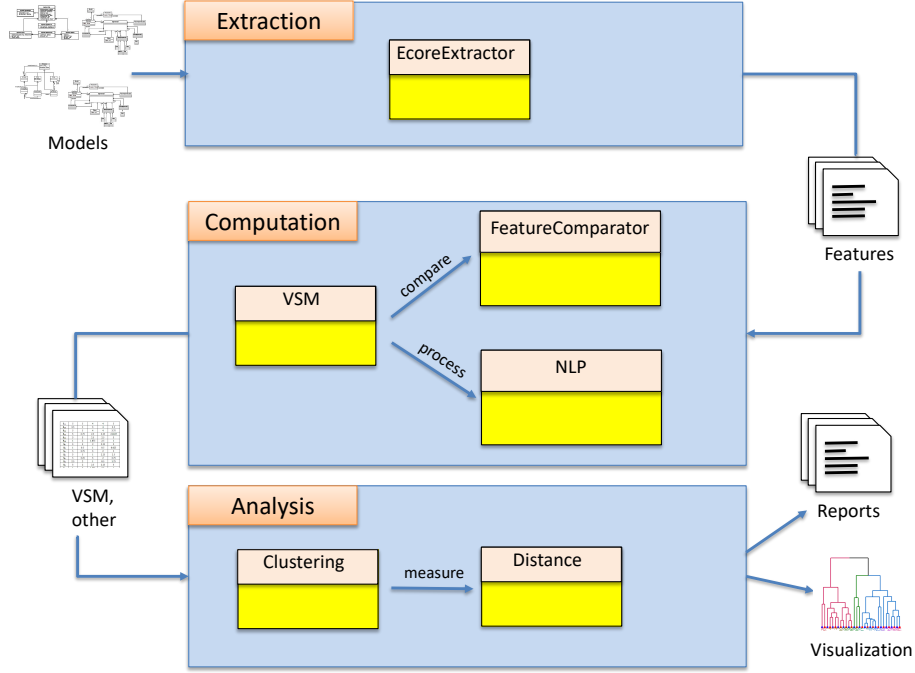


Figure 1: Overview of SAMOS Architecture.

3. SAMOS Framework

SAMOS applies the IR and ML-inspired approaches as summarized above to models. In this section, we outline the overall architecture, implementation details and the current functionalities offered by SAMOS.

3.1. Architecture

The architecture, as depicted in Figure 1, involves three main modules: extraction, computation and analysis. A typical workflow for analytics starts with a metamodel-driven extraction (i.e. in module *Extraction*) of features from a set of input models. Features can be, for instance, singleton names of model elements (very similar to the vocabulary of documents) or larger fragments of the underlying graph structure such as n-grams or sub-trees. In our context, an example n-gram for an EMF metamodel would be for $n = 2$

an EClass containing an EAttribute [5]. The feature files are then fed to the *Computation* module. SAMOS computes a term-frequency based VSM, using feature comparison schemes (for instance maximum similar subsequence for n-grams), weighting schemes (for instance EClass weighted higher than EAttribute) and NLP such as tokenization, filtering, stemming and synonym checking. The results include a vector space model and other metadata files, which are then moved to the *Analysis* module. The fact that the models are effectively transformed to their corresponding numerical vector representations allows the application of various analytics and machine learning techniques. Applying various distance measures suitable to the problem at hand, SAMOS applies different clustering algorithms and can output automatically derived cluster labels or diagrams for visualization and manual inspection and exploration. For more elaboration on the underlying techniques and a detailed example-driven demonstration of the workflow, please refer to our previous work [3, 8].

3.2. Implementation

SAMOS is implemented in Java for feature extraction and VSM generation, and in R for the analytics and machine learning parts. SAMOS public version 1.0 as presented in this paper offers the following major functionalities:

- *Modelling languages:* SAMOS supports Ecore metamodels for feature extraction and subsequent steps of the analytics workflow.
- *Analyses:* SAMOS comes preconfigured to run in predefined settings for two types of analyses: domain clustering and clone detection.
- *Feature extraction:* Users can configure SAMOS to extract unigrams, bigrams and sub-trees with depth 1. Orthogonal to this configuration, they can choose to have basic vertex information (e.g. just the name of the model element) or the full details (i.e. all the attributes).
- *Fragmentation:* Users can specify the extraction and analysis scope: whole model or fragments (e.g. EPackages or EClasses in metamodels).
- *NLP components:* Users can configure NLP settings such as tokenization, lemmatization, WordNet synonym checking and threshold values.
- *Matching schemes:* Users can have type-based weights (e.g. EClass having more weight than EParameter) and idf (see Section 2).
- *Feature comparison:* Users can have different algorithms for feature comparison, e.g. maximum similar subsequence for n-grams or tree edit distance for subtrees [8].

- 94 • *Distance measures:* Different distance measures are used in the R
95 scripts for different goals, e.g. cosine for domain clustering and masked
96 Bray-Curtis for clone detection.
- 97 • *Clustering methods:* Similarly, different clustering methods are used
98 for different goals, e.g. hierarchical clustering for domain clustering and
99 density-based clustering for clone detection.

100 3.3. Extension

101 Beyond what the wide range of available configuration options as outlined
102 in Section 3.2, SAMOS is highly flexible and extendable, due to its simple and
103 modular architecture. The three main modules of extraction, computation
104 and analysis, as depicted in Figure 1, are independent of each other and the
105 communication of data among them is achieved via a file-based pipeline. This
106 approach so far has allowed us to plug-and-play different components, e.g. in-
107 tegrating an extractor for a different modelling language without changing
108 the rest of the architecture.

109 Over the course of the years, parts of the workflow have been extended
110 for various modelling languages and scenarios. Examples for the extraction
111 component would include extending the *IExtractor* abstract class for UML2
112 models (ongoing work), or adding a completely new feature extractor SXFM
113 feature models using its own Java-based parser [7], while reusing the compu-
114 tation and analysis components mostly as-is. Similarly, one can simply add
115 a new statistical analysis component at the end of the workflow, as we have
116 done in previous work by applying topic modelling on top of the vector space
117 model we obtained from the computation module [9]. Additional extensions
118 we have experimented with include attaching machine learning components
119 on top of the computation output (i.e. the vector space model with its nu-
120 merical vector representation), or even deep learning components directly on
121 top of the extracted features. We plan to gradually integrate these extensions
122 into future release versions of SAMOS.

123 4. Illustrative Use Cases

124 While the reference work for SAMOS already includes plenty of demon-
125 stration and evaluation of SAMOS, we provide two introductory scenarios
126 in this version. The sample runs contain a prerequisite step of crawling 19
127 Ecore metamodels from ATL Zoo [15], and running the following scenarios.

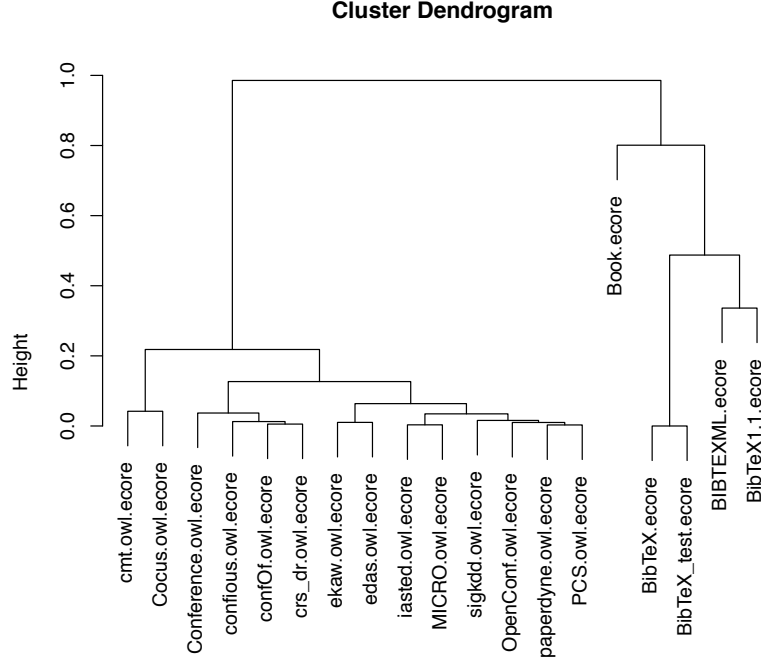


Figure 2: Plot of the dendrogram as output of domain clustering.

128 *Domain clustering scenario.* Users can run SAMOS with the domain clus-
 129 tering setting to see how conceptually related the model domains are. The
 130 predefined settings are: unigrams with just model element name informa-
 131 tion, whole model scoping, full components of NLP, type and idf weighting,
 132 simple n-gram comparison, cosine distance and hierarchical clustering. Clus-
 133 tering yields two types of output: one csv file with the list of cluster labels,
 134 and a pdf file visualizing the whole dataset in the form of a dendrogram.
 135 The dendrogram depicts how similar metamodels are to each other or other
 136 groups. The height values of the connection points of metamodels or groups
 137 map to intra-model or intra-group distances. The interpretation of the den-
 138 drogram in Figure 2 is that most of the models on the left main branch of
 139 the tree are conceptually very similar; after manual investigation, they ap-
 140 pear to be all conference management metamodels. The rightmost branch
 141 in turn has a cluster of bibliography management metamodels, although one
 142 item (Book.ecore) seems to be less similar to the rest of the cluster. For
 143 further reference on domain clustering on models, readers are referred to our
 144 previous work using SAMOS [3, 7, 9] and other related work [16].

145 *Clone detection scenario.* Users can run SAMOS with the clone detection
 146 setting to see whether there are highly similar model fragments (*clones*)
 147 within or among metamodels. The predefined settings are: n-trees (depth
 148 1) with just full model element attributes, EClass scoping, full components
 149 of NLP, type weighting only, Hungarian distance for subtree comparison,
 150 masked Bray-Curtis distance and density-based clustering. Clone detection
 151 yields three csv files corresponding to clusters of different clone types: Type
 152 A for identical fragments except layout and cosmetic differences, Type B
 153 for slightly different fragments (distance threshold 10%), and Type C for
 154 substantially different fragments (distance threshold 30%), with the threshold
 155 values chosen with respect to our previous work [8]. Table 1 depicts an
 156 excerpt for Type B clone fragments resulting from the sample run. Each
 157 row corresponds to a clone cluster with a cluster label (first column), list of
 158 model fragments separated by whitespace (second column), and the average
 159 size of the cluster in terms of total number of model elements. The model
 160 fragments are encoded in the format `MODEL_NAME$FRAGMENT_NAME`, and the
 161 clusters are sorted with respect to their average sizes.

cluster label	models	average cluster size
4	MICRO.owl.ecore\$Person sigkdd.owl.ecore\$Person OpenConf.owl.ecore\$Person ...	50.5
15	PCS.owl.ecore\$Paper OpenConf.owl.ecore\$Paper sigkdd.owl.ecore\$Paper ...	46.5
14	MICRO.owl.ecore\$Conference PCS.owl.ecore\$Conference edas.owl.ecore\$Conference ...	33.25
...

Table 1: Clone detection results for Type B clones in the sample dataset.

162 SAMOS can also crawl our labelled dataset on Zenodo [17]. Users can ap-
 163 ply domain clustering and clone detection on this significantly larger dataset.
 164 While there is no other related work for metamodel clone detection, readers
 165 can refer to literature for clone detection on other (non-meta-)model types,
 166 for instance, for Simulink models [18] and UML models [19].

167 5. Conclusions and Future Work

168 In this paper, we presented the SAMOS framework for model analytics
169 and management, in its public version 1.0. We briefly outlined the moti-
170 vation for developing this framework, its workflow and main functionalities.
171 We also provided two illustrative scenarios for domain clustering and clone
172 detection. As already mentioned in the introduction, SAMOS has been in
173 development for many years with a lot of (partly experimental, not all pub-
174 lished) extensions for different languages, settings and analyses. We plan to
175 gradually integrate those extensions into SAMOS, including but not limited
176 to:

- 177 • Major performance optimizations, including iterative clone detection [8],
- 178 • Support for other languages: feature models [7], UML [6], BPMN [10],
179 Simulink,
- 180 • Other analyses and ML algorithms such as topic modelling [9],
- 181 • High performance computing using Apache Spark [20],
- 182 • Mining software repositories and Neo4j bridges [11].

183 References

- 184 [1] Ö. Babur, L. Cleophas, M. van den Brand, B. Tekinerdogan, M. Aksit,
185 Models, more models, and then a lot more, in: M. Seidl, S. Zschaler
186 (Eds.), *Software Technologies: Applications and Foundations*, Springer
187 International Publishing, Cham, 2018, pp. 129–135.
- 188 [2] B. Tekinerdogan, Ö. Babur, L. Cleophas, M. van den Brand, M. Aksit,
189 Introduction to model management and analytics for large scale systems,
190 in: B. Tekinerdogan et al. (Ed.), *Model Management and Analytics for
191 Large Scale Systems*, Elsevier, 2019, pp. 3–11.
- 192 [3] Ö. Babur, L. Cleophas, M. van den Brand, Hierarchical clustering of
193 metamodels for comparative analysis and visualization, in: *Proc. of the
194 12th European Conf. on Modelling Foundations and Applications*, 2016,
195 pp. 2–18.
- 196 [4] Ö. Babur, Statistical analysis of large sets of models, in: *Proc. of the
197 31st IEEE/ACM Int. Conf. on Automated Software Engineering*, ACM,
198 2016, pp. 888–891.

- 199 [5] Ö. Babur, L. Cleophas, Using n-grams for the automated clustering of
200 structural models, in: 43rd Int. Conf. on Current Trends in Theory and
201 Practice of Computer Science, 2017, pp. 510–524.
- 202 [6] D. Wille, Ö. Babur, L. Cleophas, C. Seidl, M. van den Brand, I. Schaefer,
203 Improving custom-tailored variability mining using outlier and cluster
204 detection, *Science of Computer Programming* 163 (2018) 62–84.
- 205 [7] Ö. Babur, L. Cleophas, M. van den Brand, Model analytics for feature
206 models: case studies for S.P.L.O.T. repository, in: *Proc. of MODELS*
207 *2018 Workshops, co-located with ACM/IEEE 21st Int. Conf. on Model*
208 *Driven Engineering Languages and Systems*, 2018, pp. 787–792.
- 209 [8] Ö. Babur, L. Cleophas, M. van den Brand, Metamodel clone detection
210 with SAMOS, *Journal of Computer Languages* 51 (2019) 57 – 74.
- 211 [9] Ö. Babur, A. Suresh, W. Alberts, L. Cleophas, R. Schiffelers, M. van den
212 Brand, Model analytics for industrial MDE ecosystems, in: *Model Man-*
213 *agement and Analytics for Large Scale Systems*, Elsevier, 2020, pp. 273–
214 316.
- 215 [10] M. S. Nikoo, Ö. Babur, M. van den Brand, Business process model clone
216 detection, Accepted for publication, *PeerJ Computer Science* (2022) .
- 217 [11] Ö. Babur, E. Constantinou, A. Serebrenik, Language usage analysis for
218 EMF metamodels on GitHub, Submitted to a journal (2022) .
- 219 [12] C. D. N. Damasceno, D. Strüder, Quality guidelines for research arti-
220 facts in model-driven engineering, in: *2021 ACM/IEEE 24th Interna-*
221 *tional Conference on Model Driven Engineering Languages and Systems*
222 *(MODELS)*, IEEE, 2021, pp. 285–296.
- 223 [13] C. D. Manning, P. Raghavan, H. Schütze, et al., *Introduction to infor-*
224 *mation retrieval*, Vol. 1, Cambridge University Press, 2008.
- 225 [14] C. D. Manning, H. Schütze, *Foundations of statistical natural language*
226 *processing*, Vol. 999, MIT Press, 1999.
- 227 [15] Atlanmod Zoo, <https://github.com/atlanmod/atlantic-zoo>, ac-
228 cessed: 2022-05-11.

- 229 [16] F. Basciani, J. Di Rocco, D. Di Ruscio, L. Iovino, A. Pierantonio, Auto-
 230 mated clustering of metamodel repositories, in: Int. Conf. on Advanced
 231 Information Systems Engineering, Springer, 2016, pp. 342–358.
- 232 [17] Ö. Babur, A labeled Ecore metamodel dataset for domain clustering
 233 (Mar. 2019). doi:10.5281/zenodo.2585456.
 234 URL <https://doi.org/10.5281/zenodo.2585456>
- 235 [18] M. H. Alalfi, J. R. Cordy, T. R. Dean, M. Stephan, A. Stevenson, Models
 236 are code too: Near-miss clone detection for simulink models, in: 28th
 237 Int. Conf. on Software Maintenance, IEEE, 2012, pp. 295–304.
- 238 [19] H. Störrle, Towards clone detection in UML domain models, Software
 239 & Systems Modeling 12 (2) (2013) 307–329.
- 240 [20] Ö. Babur, L. Cleophas, M. van den Brand, Towards distributed model
 241 analytics with Apache Spark, in: Proc. of the 6th Int. Conf. on Model-
 242 Driven Engineering and Software Development, 2018, pp. 767–772.

243 **Required Metadata**

244 **Current code version**

Nr.	Code metadata description	
C1	Current code version	1.0
C2	Permanent link to code/repository used of this code version	For example: https://github.com/onderbabur/samos
C3	Legal Code License	MIT license
C4	Code versioning system used	git
C5	Software code languages, tools, and services used	Java SE JDK 1.8, Eclipse Modeling-Tools 2021-03, Eclipse plugin "m2e" version 1.18, R version 3.6.1, R package "rJava" version 3.6.2 and custom package "vegan"
C6	Compilation requirements, operating environments & dependencies	MacOS, Linux, Windows
C7	If available Link to developer documentation/manual	https://onderbabur.github.io/samos/
C8	Support email for questions	Onder.Babur@wur.nl

Table 2: Code metadata (mandatory)