

# Mining Shape Expressions with ShapeIt<sup>\*</sup>

Ezio Bartocci<sup>1</sup>, Jyotirmoy Deshmukh<sup>2</sup>, Cristinel Mateis<sup>3</sup>,  
Eleonora Nesterini<sup>1,3</sup>, Dejan Ničković<sup>3</sup>, and Xin Qin<sup>2</sup>

<sup>1</sup> TU Wien, Austria

<sup>2</sup> University of Southern California, USA

<sup>3</sup> AIT Austrian Institute of Technology, Austria

**Abstract.** We present SHAPEIT, a tool for mining specifications of cyber-physical systems (CPS) from their real-valued behaviors. The learned specifications are in the form of *linear shape expressions*, a declarative formal specification language suitable to express behavioral properties over real-valued signals. A linear shape expression is a regular expression composed of parameterized lines as atomic symbols with symbolic constraints on the line parameters. We present here the architecture of our tool along with the different steps of the specification mining algorithm. We also describe the usage of the tool demonstrating its applicability on several case studies from different application domains.

## 1 Introduction

Specification mining [1–3] is the process of inferring likely system properties from observing its execution and the behavior of its environment. This is an emerging research field that supports the engineering of cyber-physical systems (CPS) where computational units are tightly embedded with physical entities such as sensors and actuators controlling a physical process. CPS often operate (autonomously) in sophisticated and unpredictable environments.

In this context, mined properties can be used to complete existing incomplete or outdated specifications, to understand essential properties of black-box components (e.g., machine learning components) and to automate difficult tasks such as fault-localization [4, 5], failure explanation [6] and falsification analysis [7]. The symbolic and declarative nature of formal specification languages provide an high-level and abstract framework that facilitates generalisation. Furthermore, mined specifications are re-usable, data-efficient, compositional and closer to human understanding.

In this paper, we present SHAPEIT, a tool for automatic mining formal specifications from positive examples of time-series data encoding system behaviors or a discrete-time trace of the value of a particular system variable. SHAPEIT

---

<sup>\*</sup> This project has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 956123 and it is partially funded by the TU Wien-funded Doctoral College for SecInt: Secure and Intelligent Human-Centric Digital Technologies.

uses *Linear Shape Expressions* (LSEs) [8], a recent introduced declarative formalism suitable to express expected behaviors over noisy real-valued signals. A linear shape expression is a regular expression composed of parameterized lines as atomic symbols with symbolic constraints on the line parameters.

Given a set of time-series and a maximum error threshold, SHAPEIT implements the specification mining procedure [9] consisting of three steps: (1) **segmentation** of time-series into an optimal piecewise-linear approximation, (2) **abstraction** and **clustering** of linear segments into a finite set of symbols, where each symbol represent a set of similar lines, and (3) **learning** of linear shape expressions from the sequences of symbols generated in the previous step.

In the rest of the paper, we present the specification language and the architecture of the tool. We also show the usage of our tool, demonstrating the applicability to several different examples of time-series taken from the literature. The code of our tool is publicly available at: <https://www.doi.org/10.5281/zenodo.5569447>.

## 2 Shape Expressions

*Linear shape expressions* (LSE) [8] are regular expressions defined over parameterized *linear atomic shapes*, where a linear atomic shape is uniquely determined by three parameters: slope  $a$ , (relative) offset  $b$  and duration  $d$ . LSEs can have additional constraints over these parameters. We use the following syntax to define the fragment of LSEs supported by SHAPEIT.

$$\begin{aligned} \text{shape} &:= \text{line}(a, b, d) \mid \text{shape}_1 + \text{shape}_2 \mid \text{shape}_1 \cdot \text{shape}_2 \mid (\text{shape})^* \\ \text{cst} &:= x \text{ in } [\mathbf{c1}, \mathbf{c2}] \mid \text{cst}_1 \text{ and } \text{cst}_2 \\ \text{SE} &:= \text{shape} : \text{cst} \end{aligned}$$

where  $\mathbf{c1}$  and  $\mathbf{c2}$  are rational constants such that  $\mathbf{c1} \leq \mathbf{c2}$ .

A LSE SE consists of two main components, a regular expression **shape** that captures the qualitative aspect of the specification, and a constraint **cst** imposed on the LSE parameters. Shape expressions are evaluated against finite signals – sequences of (time, value) pairs. The semantics of LSE is defined in terms of a *noisy match* relation. We say that a signal is a  $\nu$ -noisy match of a linear atomic shape, if there exists an ideal line segment with some slope  $a$ , relative offset  $b$  and duration  $d$  such that (1)  $a$ ,  $b$  and  $d$  satisfy the constraint **cst**, and (2) the mean square error (MSE) between the signal segment and the ideal line segment is smaller than or equal to  $\nu$ . This definition is inductively lifted to arbitrary LSEs. In essence, a signal is a  $\nu$ -noisy match of an arbitrary LSE if there exists a sequence of linear atomic shapes with instantiated parameters such that: (1) the sequence is consistent with the qualitative (regular expression) part of the LSE, (2) the instantiated parameters satisfy the LSE constraint, and (3) the signal can be split into the sequence with the same number of segments, such that each signal segment is a  $\nu$ -noisy match of its corresponding atomic shape. The formal syntax and semantics of shape expressions are presented in [9].

### 3 ShapeIt Architecture, Methods and Implementation

The architecture of SHAPEIt is depicted in Figure 1. The tool consists of five components: (1) segmentation, (2) abstraction, (3) clustering, (4) automata learning and (5) translation from automata to regular expressions. SHAPEIt is implemented in Python 3 with the use of external Python and Java libraries.

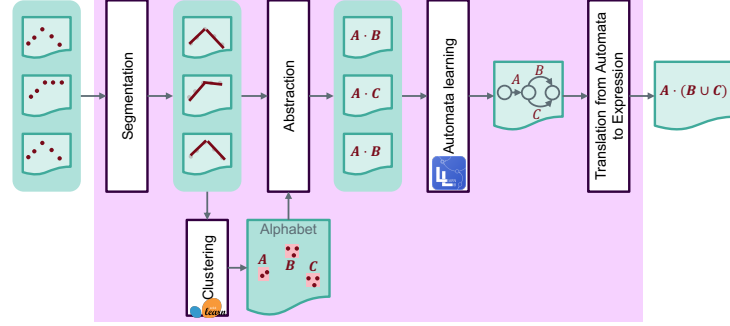


Fig. 1: Overview of SHAPEIt workflow.

*Segmentation* module implements the piecewise-linear approximation algorithm with quadratic complexity from [9] that given a time series and a mean square error (MSE) threshold computes the minimal sequence of segments such that for each segment of data, its linear regression MSE is below the threshold. The input of this module is a set of time-series and the output is a set of line segment sequences, where each line segment is characterized by slope, relative offset and duration parameters.

*Abstraction and clustering* module takes as input the set of line segments (computed by the segmentation module) and uses the k-Means clustering implementation from the `scikit-learn` library<sup>4</sup> to group lines with similar parameters. The user specifies a threshold on the derivative of the Within-Cluster-Sum-of-Squares (WCSS) error measure to determine the optimal number of clusters. The tool defines a finite alphabet in which each letter is associated to a different cluster. Each letter is also assigned the minimal bounding cube that contains all the points in its corresponding cluster. Each line segment is mapped to a letter in the alphabet, resulting in a set of finite words.

*Automata learning* module applies the Regular Positive and Negative Inference (RPNI) algorithm for passive learning from positive examples, implemented in the Java `learnlib` library<sup>5</sup>, to infer a deterministic finite automaton (DFA) from a set of finite words. The integration of the Java library in our Python implementation is done using the `JPy` library.<sup>6</sup>

<sup>4</sup> <https://scikit-learn.org/stable/>

<sup>5</sup> <https://learnlib.de/>

<sup>6</sup> <https://jpy.readthedocs.io/en/latest/>

*DFA to shape expressions* module implements the algorithm for translating DFAs to regular expressions using the state elimination method. The NetworkX library<sup>7</sup> is used to represent and manipulate DFAs during the translation.

## 4 Evaluation

In the following, we evaluate the applicability of SHAPEIT<sup>8</sup> to find temporal patterns over different time-series datasets stored in the UCR Time Series Classification Archive [10]. Our experiments run on a Notebook Dell Latitude 5320, Intel Quad-Core i7-1185G7 (3,00 GHz/Turbo 4,80 GHz), RAM 32 GB. SHAPEIT software components run on Python version 3.8.8 and on Java version 16.0.2. For all the experiments we set to 10 the threshold on the derivative of the WCSS error discussed in Section 3.

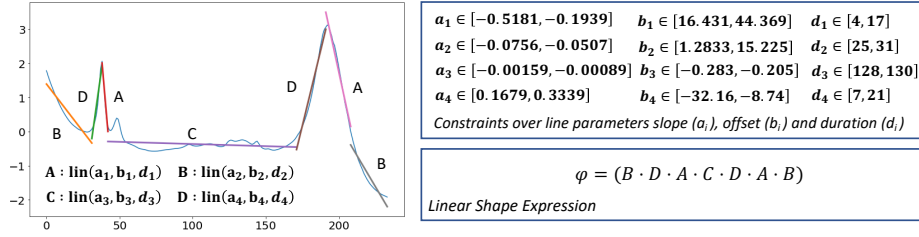


Fig. 2: (Left) An example of piece-wise linear approximation of a trace in *Wine* dataset with  $\varepsilon_{\max} = 0.05$  (Right) Generated Linear Shape Expression.

*Wine dataset* This dataset [10] consists of 111 traces, representing the food spectrograph of two kinds of wine. We consider only one class of wine data, containing 57 traces of length 234 samples (Fig. 2 shows one example).

By setting the maximum error threshold  $\varepsilon_{\max}$  to 0.05 (a little insight into how the learned specification varies depending on this value can be found in Table 2), SHAPEIT obtains an alphabet of four letters, each one describing a set of segments characterized by the values of slope, relative offset and duration reported in Fig. 2.

# traces	$ w $	$t_s(s)$	$t_c(s)$	$t_l(s)$	$t_{\text{total}}(s)$
1	10	$2.205 \cdot 10^{-4}$	$1.100 \cdot 10^{-6}$	$3.499 \cdot 10^{-4}$	$5.724 \cdot 10^{-4}$
1	100	$7.173 \cdot 10^{-2}$	$4.968 \cdot 10^{-3}$	$3.727 \cdot 10^{-4}$	$7.707 \cdot 10^{-2}$
1	234	$4.227 \cdot 10^{-1}$	$5.195 \cdot 10^{-3}$	$4.319 \cdot 10^{-4}$	$4.283 \cdot 10^{-1}$
10	10	$1.993 \cdot 10^{-3}$	$4.932 \cdot 10^{-3}$	$4.175 \cdot 10^{-4}$	$7.281 \cdot 10^{-3}$
10	100	$7.232 \cdot 10^{-1}$	$5.114 \cdot 10^{-3}$	$7.976 \cdot 10^{-4}$	$7.183 \cdot 10^{-1}$
10	234	4.353	$1.176 \cdot 10^{-2}$	$1.537 \cdot 10^{-3}$	4.366
57	10	$1.21 \cdot 10^{-2}$	$7.594 \cdot 10^{-3}$	$6.122 \cdot 10^{-4}$	$2.039 \cdot 10^{-2}$
57	100	4.110	$2.954 \cdot 10^{-2}$	$2.188 \cdot 10^{-2}$	4.161
57	234	$2.934 \cdot 10$	$2.983 \cdot 10^{-2}$	$4.201 \cdot 10^{-3}$	$2.937 \cdot 10$

Table 1: Computational cost of SHAPEIT.

<sup>7</sup> <https://networkx.org/>

<sup>8</sup> commit in the repository used: d92341998d66615cf6a9c4f3bcc419df4cd988b6

The concatenation of letters  $D$  and  $A$  represents the peaks that appear in the shape (see Figure 2), in which  $D$  describes the rising part (with positive slope) and  $A$  the decreasing one (with negative slope). Letter  $C$  represents the approximately constant part of the trace that separates the two peaks, while  $B$  describes the two extremes (they are both decreasing segments but less steep than the ones that come after the peaks' maxima).

In this particular application the values of slopes would be able to distinguish the different letters by their own: the intervals of slopes are indeed disjoint. The same happens for the relative offset but not for the duration.

SHAPEIt generates an LSE specification (see Fig. 2) that captures the two main peaks of the trace, but it is not able to recognize the little one that comes immediately after the first peak. The maximum error threshold  $\varepsilon_{\max}$  should be reduced if one is interested in detecting also this little curve.

In Table 1, we report the time (expressed in seconds) required by the tool to complete the three different phases: segmentation ( $t_s$ ), clustering ( $t_c$ ) and automata learning ( $t_l$ ). In the last column,  $t_{\text{total}}$  summarizes the total time needed. Varying the number of traces and their lengths, we can observe that almost always the segmentation represents the most expensive part of the computation, while the clustering and the automata learning can be both considered negligible in terms of computation time. The only exceptions are the two cases in which the total number of traces is 1 or 10 with traces long only 10: the values of  $t_s$ ,  $t_c$  and  $t_l$  are comparable since the segmentation is very fast due to the low number of samples to approximate.

In Table 2, we compare the specifications learned varying the maximum error threshold  $\varepsilon_{\max}$  from 0.05 to 0.5. The number of clusters does not decrease monotonically when increasing the maximum error allowed in the segmentation, while the specifications become shorter and therefore have less explanatory power.

$\varepsilon_{\max}$	$\varphi$	# clusters
0.05	$B \cdot D \cdot A \cdot C \cdot D \cdot A \cdot B$	4
0.1	$F \cdot E \cdot G \cdot I \cdot H \cdot F$	5
0.5	$K \cdot (L + M)$	3

Table 2: Sensitivity w.r.t.  $\varepsilon_{\max}$

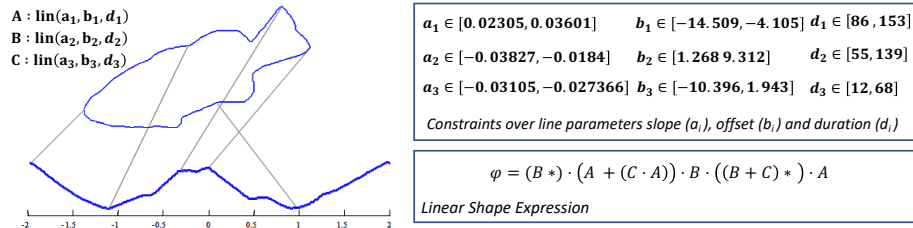


Fig. 3: (Left) From fish contour to time series. (Right) Generated Linear Shape Expression.

*Fish Data Set* This data set [10] is composed by 350 time series representing the shape of seven different species of fishes (chinook salmon, winter coho, brown

trout, Bonneville cutthroat, Colorado River cutthroat trout, Yellowstone cutthroat and mountain whitefish). Starting from 50 images for each class, Lee *et al.* in [11] generated the data set leveraging a novel technique that transforms the contour of the fish into a time series using a turn-angle function illustrated in Fig. 3 on the left. Setting to 0.05 the maximum threshold error, with SHAPEIT we are able to learn a specification (Fig. 3 on the right) from 26 shapes of the same species of fish, each one containing 463 samples.

The concatenation of letters  $B$  and  $A$  represents the predominant shape in the traces: the triangular repeating behavior where, in particular,  $A$  describes the rising part and  $B$  the descending one (see Figure 3). Letter  $C$  is instead used to symbolize the noisy parts, both with positive and negative slope, that eventually separates these longer segments. The choice operator  $(+)$  represents the possibility to have multiple symbols or expressions in different time series. Finally, the Kleene star  $(*)$  is used to indicate that a symbol or an expression can appear zero or more times.

The learned specification provides insights about the relevant shapes in the time series data, displaying them in an human understandable language and therefore offering interpretability to the user. In this example, referring to the fish image in Fig. 3 on the left, we can associate the concatenation of letters  $B$  and  $A$  in the specification to the upper contour of the fish silhouette that is starting from the head and is ending with the tail. Since the same concatenation is then repeated in the specification, we can infer that the contour of the lower part of the fish is not significantly different from the upper one. Finally, letter  $C$  can be interpreted as the presence of a big fin that interrupts the predominant lines described by letters  $A$  and  $B$ .

## 5 Conclusion and Future Work

In this paper, we presented SHAPEIT, a tool for mining specifications that describe the behaviors of CPS. The tool requires a set of real-valued signals generated by the system under study as input and it returns as output the specification that better summarize the properties of the traces in the form of linear shape expression. SHAPEIT is structured in three phases: segmentation (approximating the traces with segments), abstraction and clustering (grouping lines with similar parameters) and automata learning (learning a DFA from words). Two additional values are needed as inputs to regulate the first two processes: a threshold expressing the maximum error allowed by the approximation and a threshold for the WCSS error to find an optimal number of clusters. We demonstrated the applicability of our tool over two different case studies (*Wine* and *Fish*) but other datasets are present in our repository. These data can be used as well to do experiments and gain confidence with SHAPEIT.

As possible future works, we are interested in exploring and learning more general Shape Expressions (not necessarily linear ones), probably gaining explanatory power at the cost of an increasing computation time. We will also study how to automatize the tuning of the two thresholds required by the tool

for the segmentation and the clustering phases. In this paper, the segmentation tool finds automatically the optimal number of segments to be used for the approximation, given a maximum error allowed. However it has already been developed to work in the other way round: receiving the number of required segments as input and then finding the approximation that provides the minimum error. It will be therefore interesting to exploit this feature to embed some domain knowledge (in the form of number of segments) in the specification mining process. A step forward will be adding the possibility to set constraints to the parameters of the lines. Finally, an other direction of work could be trying to generalize the tool in order to make it able to handle online processes instead of only offline ones.

## References

1. L. Nenzi, S. Silveti, E. Bartocci, and L. Bortolussi, “A robust genetic algorithm for learning temporal specifications from data,” in *Proc. of QEST 2018*, vol. 11024 of *LNCS*, pp. 323–338, Springer, 2018.
2. E. Bartocci, L. Bortolussi, and G. Sanguinetti, “Data-driven statistical learning of temporal logic properties,” in *Proc. of FORMATS*, pp. 23–37, 2014.
3. F. Wang, Z. Cao, L. Tan, and H. Zong, “Survey on learning-based formal methods: Taxonomy, applications and possible future directions,” *IEEE Access*, vol. 8, pp. 108561–108578, 2020.
4. E. Bartocci, T. Ferrère, N. Manjunath, and D. Nickovic, “Localizing faults in Simulink/Stateflow models with STL,” in *HSCC*, pp. 197–206, ACM, 2018.
5. X. Jin, A. Donzé, J. V. Deshmukh, and S. A. Seshia, “Mining requirements from closed-loop control models,” *IEEE TCAD*, vol. 34, no. 11, pp. 1704–1717, 2015.
6. E. Bartocci, N. Manjunath, L. Mariani, C. Mateis, and D. Nickovic, “Automatic failure explanation in CPS models,” in *SEFM*, vol. 11724 of *LNCS*, pp. 69–86, 2019.
7. E. Bartocci, J. V. Deshmukh, A. Donzé, G. E. Fainekos, O. Maler, D. Nickovic, and S. Sankaranarayanan, “Specification-based monitoring of cyber-physical systems: A survey on theory, tools and applications,” in *Lectures on Runtime Verification - Introductory and Advanced Topics*, pp. 135–175, Springer, 2018.
8. D. Nickovic, X. Qin, T. Ferrère, C. Mateis, and J. V. Deshmukh, “Shape expressions for specifying and extracting signal features,” in *Proc. of RV*, vol. 11757 of *LNCS*, pp. 292–309, 2019.
9. E. Bartocci, J. Deshmukh, F. Gigler, C. Mateis, D. Nickovic, and X. Qin, “Mining shape expressions from positive examples,” *IEEE Trans. Comput. Aided Des. Integr. Circuits Syst.*, vol. 39, no. 11, pp. 3809–3820, 2020.
10. Y. Chen, E. Keogh, B. Hu, N. Begum, A. Bagnall, A. Mueen, and G. Batista, “The UCR time series classification archive,” July 2015. [www.cs.ucr.edu/~eamonn/time\\_series\\_data/](http://www.cs.ucr.edu/~eamonn/time_series_data/).
11. L. Dah-Jye, J. Archibald, R. Schoenberger, A. Dennis, and D. Shiozawa, *Contour Matching for Fish Species Recognition and Migration Monitoring*, vol. 122, pp. 183–207. 2008.