
A VISUAL ARTIFACT DETECTOR FOR ITS.APE

AN EXTENSION TOOL FOR ITS.APE TO
PROVIDE VISION FOR THE SUITE

LAB REPORT

by

FELIX ROSSMANN

2471236

submitted to

RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT BONN

INSTITUT FÜR INFORMATIK IV

ARBEITSGRUPPE FÜR IT-SICHERHEIT

in degree course

COMPUTER SCIENCE (M.Sc.)

First Supervisor: Prof. Dr. Michael Meier
University of Bonn

Second Supervisor: Dr. Matthias Frank
University of Bonn

Sponsor: Arnold Sykosch, M.Sc.
University of Bonn

Bonn, 2nd of April 2019

ABSTRACT

An additional tool for the ITS.APE suite is presented to detect whether visual cues of deployed artifacts are visible on a Windows user's desktop. The task is done by image recognition in a screenshot using OpenCV. The recognition process consists of feature detection, feature matching and outlier rejection with the goal to find a transformation matrix between the features of the observed screenshot and provided reference images. The implemented tool is then evaluated with respect to correctness of the matches and resource consumption, proving the successful implementation with about 90% successful matches using reasonable time and resources for the targeted execution environment.

CONTENTS

1	INTRODUCTION	1
2	RELATED WORK	2
3	METHODS	3
3.1	Software Design	3
3.1.1	Requirement analysis	3
3.1.2	Implementation decisions	4
3.2	Technical backgrounds	5
3.2.1	Feature detection with ORB	6
3.2.2	Descriptor matching	7
3.2.3	Outlier rejection	7
4	RESULTS	9
4.1	Evaluation setup	9
4.2	Evaluation results	11
4.2.1	Detection rate	11
4.2.2	Execution resources	13
5	CONCLUSION	16
6	BIBLIOGRAPHY	17
	APPENDICES	19
	LIST OF FIGURES	21
	LIST OF TABLES	22

1 INTRODUCTION

The *IT-Security Awareness Penetration Testing Environment* (ITS.APE) [Syk15] is a tool to analyse the awareness of a computer's user regarding IT-security. It measures the user's response to certain deployed *artifacts* which differ from the user's familiar environment. Most of them are similar to real-world attack scenarios or share certain principles with them.

There are multiple aspects to the user's response, one of them being the reaction speed. For measuring the reaction speed one must determine the point in time from which on it is possible for the user to interact with the artifact. As ITS.APE is a tool for the Microsoft Windows operating systems (which focus on a graphical user interface, a *GUI*), a user is able to interact with an artifact as soon as he or she can visually perceive it on the screen, or in Windows terms: on the *desktop*.

The process of finding this point in time can not be done reliably by observing the currently running processes or even meta information about it. Only the visual information of the current desktop is useful in this case: An extra feature or tool was needed to perform an image detection algorithm on the current screen and decide whether the visual cues of the currently deployed artifact are present on the user's desktop.

The most important measure for this tool would be the *sensitivity*, describing the probability of successful detection of artifacts when they are actually present (*true positives*). This term will be called *detection rate* in the rest of this work and should ideally be at 100%. The counter-measure for this is the *specificity* which describes the rate of *true negatives*. A second, equally important measure is the resource consumption for execution, especially the execution time. As the detection rate is an obvious measure for such a tool, the execution time and resources are almost equally important to not disrupt the user's experience or distort the reaction analysis by ITS.APE. The limitation on the resources is given by the environment that ITS.APE is mainly used in: low-tier office hardware.

The work at hand describes the design considerations and implementation details of such a tool, called *Visual Artifact Detector* (VAD), and presents the results of an evaluation given the previously proposed measures. For this chapter 2 outlines some related publications and the differences to the presented approach. In chapter 3 the software design and technical details of the VAD are explained. The evaluation setup for the given measures and the results can then be found in chapter 4, followed by a conclusion in chapter 5.

As ITS.APE is designed for computers running *Windows 7* in the 32-bit version, the VAD was built for this operating system specifically but can run on newer Windows versions as well. Written in C#, it uses Microsoft's *.NET framework* in version 4.7.2 [NET] and the *OpenCV-wrapper Emgu CV* in version 3.4.3 [EMGU].

2 RELATED WORK

For evaluation of UI patterns Bakaev et al. [Bak+18] use image recognition functions to analyse a website's screenshot and produce a machine-readable representation of it. A grayscale version of the screenshot is used for rectangle detection, text sections are identified and recognized, before detecting special UI element types using a trained feature extractor. The results are then used to analyse composite structures and output a textual representation of the website. While the task seems similar to the one presented in this report, the proposed analysis has different goals and does not suit the requirements for the VAD: The full classification of the screenshot before identifying possible candidates for artifacts would imply a high demand on resources and is not a necessary step. Since those candidates would then be compared to reference images of artifacts, one can skip the classification steps and use an image feature detector on the screenshots directly.

Aradhya et al. [AMH05] present a way to classify image-based e-mails as spam using image analysis. Their approach features a way to analyse complex images without resorting to OCR: First, text regions are extracted by performing connectivity analysis on thresholded intensities of the grayscale image. With these regions identified, they categorize the image based on certain features which are then used to train *Support Vector Machines* (SVMs) to detect different mail classes. The key advantage of this way to classify images is that the spam images have distinct features to distinguish them from other images in mails. On one hand the ITS.APE artifact database is diverse and training SVMs for every artifact type would be demanding. On the other hand it is not easy to distinguish the more advanced ITS.APE artifacts from other applications by a rigidly defined feature set, because most of the artifacts ITS.APE employs try to imitate a usual application.

The goal of automated quality assurance lead Mozgovoy et al. [MP18] to the usage of image recognition on GUIs. They identified a problem for conventional GUI-testing suites for games with non-standard GUI and other elements. Their solution decides against exact bitmap matching, as the elements could be transformed, and instead match them approximately against a template. Since their approach is unable to find scaled matches, they first scale the reference to the observed image and then match them in total. Here, the main difference to the task at hand lies in the references to match against: In their case the matched output should have limited variation compared to the references so that they are able to match the whole screenshot onto the scaled reference. This is not applicable for the use cases in ITS.APE, as the artifacts are usually present in only a part of the whole screenshot. Since the rest of the screenshot is subject to high variation building up and deploying an adequate reference image collection for whole screens is highly infeasible.

Since none of the related work matches the goals for the Visual Artifact Detector, a new tool had to be implemented.

3 METHODS

This chapter begins with a description of the general software design, including a brief requirement analysis and the resulting decisions for the implementation. The second part introduces the technical backgrounds and algorithms used, and especially explains the choice of the image recognition algorithms and their parameters.

3.1 SOFTWARE DESIGN

The client of ITS.APE that gathers the user's reaction to the deployed artifacts runs on a PC as a Windows service. Instead of implementing the Visual Artifact Detector as an additional feature of the ITS.APE service it is implemented as an external tool. This is to keep a modular approach in which software parts can easily be exchanged, and secondly to be able to give the individual processes an independent priority for execution, thus guaranteeing enough resources for the execution of the ITS.APE service.

The main reason the VAD was designed to run on the same machines as the ITS.APE service instead of e.g. a server infrastructure with better hardware or other external hardware was the user's privacy: As ITS.APE aims to collect the data anonymously or pseudonymously, transactions of highly sensitive data (such as screenshots of the current desktop) would contradict this goal.

The information on artifact types to detect are provided by a *recipes repository* on the same computer containing reference images of the artifacts to be deployed. Of those the VAD only has to look for visual cues of a single type per run, as the ITS.APE service will call it with information on the currently deployed artifact type. Additionally, the service will supply the path to a screenshot to detect the type in.

3.1.1 REQUIREMENT ANALYSIS

To implement the VAD, the following main requirements were identified:

1. Robust and high detection rate for visual cues of artifact in screenshots of a Windows desktop.
2. Low runtime of at most few seconds on low-tier office hardware.

These requirements reflect the measures given for this task in the first chapter. The 1. requirement mainly has influence on the choice of image recognition algorithms, while the 2. requirement contradicts the use of a very resource-heavy algorithm (which would provide better results, gener-

ally speaking) and has influence on the parameters of the used algorithm. This lead to a trade-off between the two requirement at some points which are mentioned in section 3.2.

Further analysis of the nature of ITS.APE and refinement of the main requirements lead to the following secondary requirements:

3. No usage of parallel computing on the graphics card (e.g. with the CUDA Toolkit[CUDA]), CPU only.
4. Windows 7 console program (32-bit).

The 3. requirement is a specification of the 2., as low-tier office computers usually don't include a CUDA-enabled graphics card or a non-integrated graphics card at all. This requirement results in better versatility of the tool, but comes at cost of a higher runtime.

Given the deployment area for ITS.APE in the present and near future the 4. requirement is due to the environment ITS.APE is mainly run in. This is also important for the nature of the visual cues of the artifacts which are mostly those of standard Windows GUI objects.

3.1.2 IMPLEMENTATION DECISIONS

It was initially decided to use an existing library for image recognition rather than implementing the necessary algorithms anew. The choice for an image recognition library fell to OpenCV [OCV], as it is licensed as *Open Source* under the BSD license and arguably the most maintained and popular of such libraries with almost 20 million downloads [CVDL19]. Other libraries for similar tasks exist, but they either focus more on machine learning like *Google's Tensorflow* [GTF], are web-based APIs (mostly non-free) like *Google's Vision API* [GVAPI] or *Amazon Rekognition* [ARK], or are not well maintained, such as *VLFeat* [VLF].

Because the runtime performance is so important for this tool it was not suitable to implement the program using an interpreted programming language. The usage of OpenCV limits the available compiled languages to C, C++ or Java, and using the wrapper Emgu CV also allows to use C#.

The advantage of C# is the possibility to use Microsoft's *.NET Framework* [NET]. It allows to efficiently implement programs for the target platforms while the latest *redistributable packages* of .NET needed for executing .NET programs are available per default if the respective Windows updates are installed. Furthermore, the speed differences between the execution of .NET's intermediate language in their *Just-in-Time* compiler are nowadays insignificant to programs written in e.g. native C++. Therefore the VAD is written in C# using .NET 4.7.2 and OpenCV 3.4.3 via the Emgu CV wrapper.

During first evaluation runs it was noticed that the used post-processing function by Emgu CV lead to erroneous results. Therefore it was decided to provide an own implementation of the RANSAC algorithm as a replacement (see subsection 3.2.3 for details).

In addition to the tool itself a Windows installer setup was composed to allow easy deployment of the VAD after compilation.

3.2 TECHNICAL BACKGROUNDS

The VAD expects a screenshot (observed image) and an artifact type to be supplied for each run. It tries to match features of the observed image to each of the artifact types reference images (model images). For the recognition process itself only one model image at a time is tried to be recognized in the observed image. Accordingly, since the given artifact type may have multiple reference images for several characteristics, the recognition process might be performed as many times as there are reference images, but exits as soon as a match between any reference and the screenshot was found. This design is especially important for the performance, as the order of the artifact's reference images are directly affecting the runtime (see chapter 4).

The goal of the recognition process is to find a valid transformation matrix to the model's features for the largest possible subset of the observed image's features. If such a matrix can not be found for a subset of certain size, it is decided that no match exists. There are three major steps involved in the image recognition process: *Feature detection and description*, *feature matching* and *outlier rejection*, while the last step includes finding a transformation matrix.

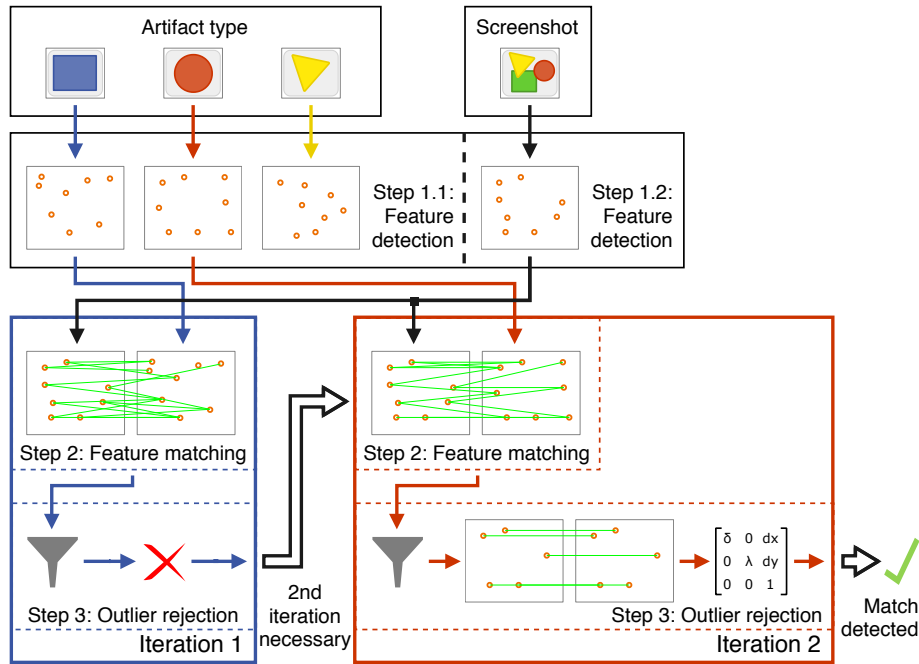


FIGURE 1: Exemplary recognition process diagram.

An example of how the process of image recognition could be executed is shown in figure 1: Given an artifact type (here containing three reference images) and a screenshot the first step yields sets of features extracted from all of the model images (step 1.1) and from the observed image (step 1.2). The features are represented as feature descriptors which are then matched with the first model's descriptors in step 2. This results in a set of candidates for features found in both images. These are then post-processed during outlier rejection (step 3) to find a transformation matrix from one feature subset to the other while matching the most candidates. If there is such a matrix fitting to a certain percentage of the candidates (see subsection 3.2.3) the VAD is finished and returns the successful detection of the artifact (return value 0). Otherwise, as shown in the example, the next reference image will proceed through steps 2 and 3 to find a transformation matrix. If there is no

more reference left and no transformation matrix was found the tool will return that no match was found (return value 1).

The retrieval of the transformation matrix can take advantage of the special properties of the deployed artifact's being mostly standard Windows GUI elements. There are only few possible transformations for those GUI objects: First, they can be moved within the desktop resulting in uniformly translated versions of the artifact's reference images. In this case, the transformation matrix is easy to find and this transformation is the most common one. Secondly the GUI objects can be scaled in a sometimes complex manner, e.g. by shifting sub-parts of the object to different locations and thus changing the aspect ratio between feature points. The combination of both transformation is the third possibility. It is hard to impossible to find a single transformation matrix for the last two of the mentioned transformations, but sometimes sub-parts of the image can be recognized more easily. For this reason the VAD in the current implementation mainly focusses on detecting the first kind of transformation reliably and leaves the recognition of the more complex transformations for future work.

The following sections will describe the choice of algorithms for each step and their configuration details.

3.2.1 FEATURE DETECTION WITH ORB

The set of available algorithms for feature detection were provided by the implementations in Emgu CV, respectively OpenCV. From those, as shown in secondary literature such as by Tareen et al. [TS18], ORB configured to find a limited amount of features provides the best combination of resource usage and result quality compared to the other three algorithms.

The ORB (*Oriented FAST and Rotated BRIEF*) keypoint detector and descriptor is a “very fast binary descriptor based on BRIEF [...], which is rotation invariant and resistant to noise” [Rub+11, p. 1]. It is specifically designed for real-time systems and low computing power, thus being faster than most of the other currently available methods. As the name states, it combines an enhanced version of the *FAST* [RDo6] keypoint detector with *rotation-aware BRIEF* (*rBRIEF*) [Cal+10] descriptors. [Rub+11; TS18; RDo6; Cal+10]

The extended FAST detector employed in ORB is used to detect corners on a pyramid scheme for scale variants [KMo8] of the image. At instantiation of Emgu CV's implementation of ORB one can adjust the scale factor and number of levels for those pyramids. For the VAD the preconfigured scale factor 1.2 with 8 levels was used. Since the VAD currently responds poorly to scaled versions of the artifact's reference images these settings could be adjusted in future versions (see chapter 4).

The feature detection process in the Emgu CV implementation of ORB is done until a configured amount of features are found (if possible), the preconfigured amount in Emgu CV being 500 features. This amount was doubled for the VAD to extract up to 1000 features per image, as a higher value allows for better differentiation between the reference images. An even higher amount of features to extract affects the execution time negatively while not resulting in better matches for the provided reference screenshots.

There are more configuration parameters available for ORB, but the preconfigured values were found to be sufficient for this task.

3.2.2 DESCRIPTOR MATCHING

The descriptor matching step tries to find one or more descriptor of the model image's keypoints (called the *training set*) for each of the observed image's keypoints (called the *query descriptors*). For feature matching, there are only two possibilities implemented in Emgu CV: Brute-force matching and FLANN (*Fast Library for Approximate Nearest Neighbor*) [ML09] based matching. There are different matching algorithms implemented for both of them such as the *k*-nearest-neighbor algorithm (*k*-NN).

FLANN based *k*-NN-matching performs best for large datasets with high dimensionality, but is less likely to find all possible candidates for matches [ML09]. The trade-off between the focus on speed and accuracy mentioned in section 3.1 lead to the choice of FLANN based *k*-NN-matching for better performance. For the VAD *k* was set to 2, so that for each query descriptor, the 2 nearest-neighbor descriptors from the training set are returned – if possible. The choice of *k* allows to apply several post-processing filters (see subsection 3.2.3) which reduce the total time needed for recognition.

3.2.3 OUTLIER REJECTION

The best *k* matches for each query descriptor are filtered in several steps before trying to find a transformation matrix for the rest using the *Random Sample Consensus* (RANSAC) [FB81] algorithm. The match candidates that actually match the correct features are called *inliers*, while the ones that do not are called *outliers*. Therefore the goal of outlier rejection is to separate the outliers from the inliers. The filtering steps try to exclude outliers as early as possible from the given query descriptor set so that the outlier rejection process can be stopped as soon as there are not enough match candidates left to make a match probable.

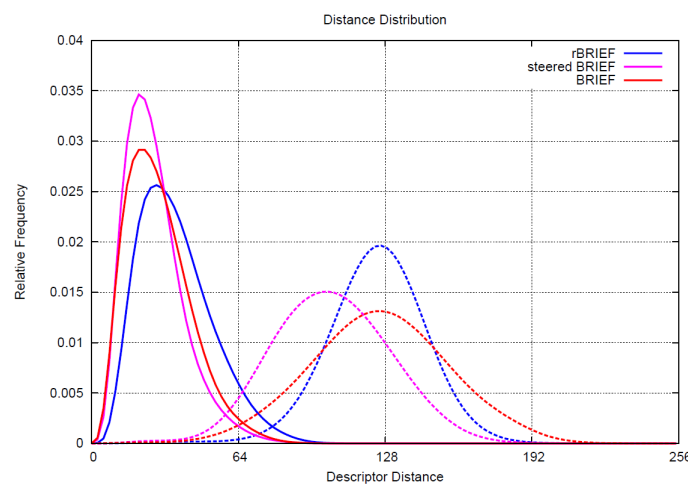


FIGURE 2: Analysis of the relation between descriptor distance and frequency of the match being an inlier or outlier, dotted for outliers and solid for inliers [Rub+11, p. 4].

Since the 2-NN-matching can yield less than two matches for some query descriptors only those descriptors that have two matches are taken into account for the next steps. The first filter applied in the VAD is the distance threshold filter proposed in the original ORB publication [Rub+11]. Their analysis links the probability of a match candidate being an inlier or outlier to its distance from the target. They showed that above a certain distance a match candidate is unlikely to be an inlier. For the VAD a threshold of 80 was chosen, as it is the approximate turning point of the probability when using rBRIEF descriptors (see figure 2). The second filter is an implementation of *Lowe's ratio test* [Lowe04] with the recommended distance ratio of 0.7 between the first and the second match's distance. After these two filters two functions provided by OpenCV/Emgu CV are used to eliminate duplicates and to rule out candidates that do not have the size or orientation of the majority of match candidates.

Between and after those filtering steps the VAD determines the current count of remaining match candidates and checks if there are still enough for a successful match. They are compared to a fixed threshold parameter $min_{matches}$ of 25 divided by the ratio of the observed and model images' areas. The parameter $min_{matches}$ has been found empirically to be the best trade-off between a suitably high detection rate given the currently implemented artifact types in ITS.APE and still skipping most of the matching process early if a match is improbable. The ratio of the image's areas allow model images to be found that have a small image size comparison to the observed image. This is a necessary factor, as the density of the feature descriptors per image area becomes more sparse if the image size becomes larger since the total count of features is set to an absolute limit of 1000.

At the last step of image recognition, if there are enough candidates left, an own implementation of the RANSAC algorithm is performed: Two match candidates are randomly sampled and are used to calculate translation and scaling factors for a hypothetical transformation matrix. This hypothesis is then tested on all match candidates allowing a transformation error of an 11×11 pixel patch due to the inaccuracy of floating point arithmetic for C#. If a certain percentage of candidates support it, given by the confidence parameter, a match was found.

The iteration count of RANSAC was set to 1000 iterations. This is derived from the square of $min_{matches}$ adjusted upward to the next power of ten to allow for an error margin, respectively a high area ratio factor. In the best case scenario this iteration count allows the minimum of needed match candidates to be combined. The confidence factor was set to 85% via empirical analysis, meaning that only 15% of the remaining candidates may be outliers of the hypothetical matrix.

The calculation of the translation and scaling factors in this part is influenced by the observation of possible transformations for GUI elements. In the current implementation a high detection rate despite arbitrary translation of the Windows GUI elements is the main focus while detecting other, more complex transformations is a secondary concern.

4 RESULTS

There are two main focus points for the evaluation, reflecting the measures given in chapter 1: First, an analysis of the detection rate is presented, followed by an evaluation of the used execution resources. For measuring the detection rate the VAD was executed to detect a representative subset of more than half of the artifact types currently provided by the ITS.APE recipes repository. They had to be detected in numerous screenshots so that the measured detection rate can be generalized for the real application domain.

The results of this first setup show large differences in the runtime for some artifact types and even for some test runs for the same artifact type but with a different screenshot. Since the artifact types in the ITS.APE recipes repository and the supplied screenshots are very heterogeneous it is hard to prove hypothetical influences on the runtime using only this setup. Even running the test in this setup several times would not decrease the complexity of confounding variables.

This lead to an additional setup for the second focus point. To have more control over the various possible factors eight specially crafted artifact types were created and used for those test runs, only varying in the order and amount of their reference images. Furthermore, only two different screenshots were queried for each of those artifact types. With this setting it is possible to observe the resources needed in dependency to the artifact types' structure while minimizing the influence of confounding variables. Using the combination of the results for both setups it is possible to extract more generally valid statements about the execution resources needed by the VAD. In section 4.1 the evaluation setups are both explained in detail. The results of the two evaluation parts can be found in section 4.2.

The 2. requirement resulted in the implementation of a *cache* for already processed data of artifact images from previous runs. This leads to drastic runtime improvements in consecutive runs of the same artifact type, saving up to 40% of the total time, but is opposed by another factor: As soon as the size of the cache increases, the duration of loading and decoding increases as well. Even for a small count of artifact types of three or four types the loading time becomes larger than the time saved by caching. Therefore, the cache was discarded early in the evaluation phase and not used during this evaluation. Its concept could still be continued in future work (see chapter 5).

4.1 EVALUATION SETUP

The VAD was evaluated in an environment that should closely resemble the real-world application domain of ITS.APE. For the evaluation Windows 7 was setup on a virtual machine using Oracle VM VirtualBox[VB]. This allowed an easy setup and quick configuration of the available hardware.

All official Windows 7 updates were installed on this VM up to those available at the 25th of March 2019. No additional software was installed other than VirtualBox's *Guest Additions* which are necessary to share data between the host and the virtual machine.

The virtual machine was run with 2GB of RAM available and one core of an Intel Core i5-6600K CPU (3,50 GHz) set to a 40% execution cap, resulting in a single virtual CPU at a speed of 1.4 GHz. These hardware settings were expected to approximate low-tier office hardware sufficiently enough to allow conclusions for the real-world application domain of ITS.APE.

After installing the VAD in the virtual machine using its Windows installer a Batch-script ran the executable multiple times using data supplied in shared folders by the host system. The supplied data was adjusted such that it fitted the respective evaluation focus. After the measurements for a single setup were finished the state of the virtual machine was restored to a snapshot taken before installing the VAD to have a clean testing environment for each test run.

For measuring the detection rate in the first test setup 25 artifact types were selected from ITS.APE's original artifact library of 48 types in total. They have been selected at random and were assumed to be an representative subset. These artifact types have then been numbered A 1 to A 25 according to table 1 in the appendix. Each of their reference images has been merged with a collection of five different Windows desktop screenshots to create a large set of testing screenshots. This merging is done programmatically and places the reference image at a random position on the desktop screenshot such that at least 90% of its area is still inside the images boundaries. Additionally, screenshots of objects that do not belong to any artifact type, but represent common desktop applications have been combined with the desktop screenshots in the same fashion. The resulting screenshots should not be detected as a match for any of the 25 artifact types. In total there 365 generated screenshots for the evaluation of the detection rate and thus leads to a total of 9125 executions of the VAD in this setup.

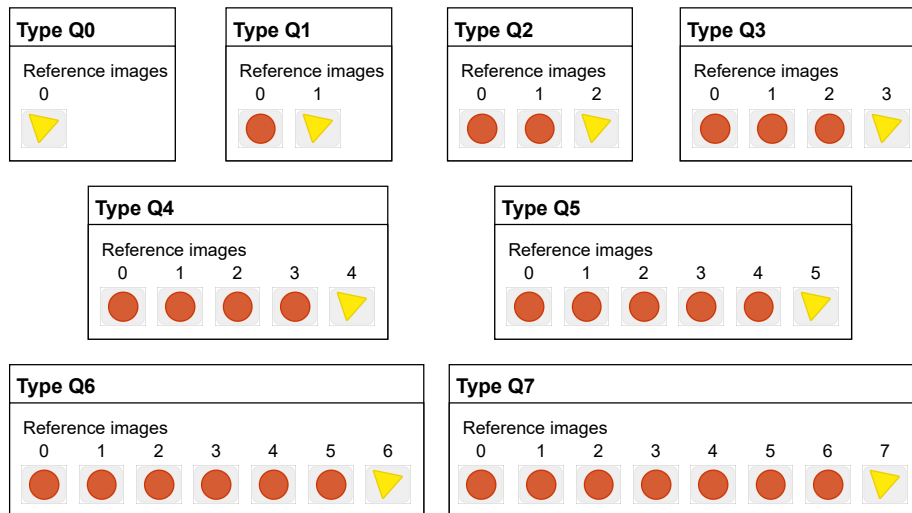


FIGURE 3: Structure of the artifact types Q0 to Q7.

For the second setup eight new artifact types were constructed on the basis of artifact type A 1. Those new types were named Q0 to Q7 and their structure is depicted in figure 3 having reference image indices starting at 0. They have one to eight reference images in increasing order which can

either be a reference image r_a (yellow triangle) or a reference image r_b (red circle) both taken from the previously used type A 1. As shown in figure 3 r_a is only present once in the collection, being always the last image of the set, while the other references are all duplicates of r_b . In addition, only two screenshots were prepared based on the same, empty desktop screenshot. In the first screenshot s_{ap} the artifact type r_a is present, in the second screenshot s_{aa} it is absent. This is done to eliminate a possible influence of screenshot selection on the runtime. Furthermore, the reference image r_b is chosen such that only the last step of outlier rejection is skipped when queried for s_{ap} . This last step is the execution of RANSAC to recover a transformation matrix and with this setup it is possible to measure its influence on the total runtime, as it is assumed to be one of the major factors.

Each of the two screenshots s_{aa} and s_{ap} got queried 1000 times for each artifact type Q0 to Q7, resulting in 16000 executions of the VAD in total for this setup. The high count of repetitions was used to even out randomly occurring influences on the execution time so that the sample is as representative for the execution time as possible. Those random influences could be anything from resource occupation by different (system) processes to temporarily slow reading speed on the shared folders due to virtualisation issues, but since there were not found to be systematic they could not be assessed further in the extent of this evaluation.

4.2 EVALUATION RESULTS

The evaluation results presented in this section prove the successful implementation of the VAD regarding the detection rate and resource usage aimed for, but hint to several aspects to improve in future work. Especially when it comes to the duration of the execution time there is room for improvement shown by the results.

Given the described first setup the detection rate was found to be at about 88% in total, together with a specificity of more than 99%. The details and implications of those results can be found in subsection 4.2.1. While the execution time for artifact types with few or conveniently ordered reference images has a median of less than 500ms per run, the amount and order of reference images have a large influence on the execution time. The measurements of the second setup show that for each additional image that has to be loaded and processed the VAD needs an extra 150ms to 250ms to compute a result. This can lead to execution times deemed too long for application in the ITS.APE environment and thus needs improvement. In subsection 4.2.2 the results for the execution resources can be found, together with further conclusions on this topic.

4.2.1 DETECTION RATE

For the evaluation of the detection rate in the first setup the return value of the VAD has been compared to the correct result. In figure 4 the outcomes grouped by the 25 artifact types can be found, where the y-axis denotes the categories *true positive* (TP), *true negative* (TN), *false positive* (FP) and *false negative* (FN). Each cell shows the relation between the actual results (numerator) and the correct results (denominator). In the last column of figure 4(b) the summed up results over all artifacts is given. The cells are color coded in a spectrum from red at 0% over yellow at 50% to green at 100% for TP and TN and vice-versa for FP and FN.

		A 1	A 2	A 3	A 4	A 5	A 6	A 7	A 8	A 9	A 10	A 11	A 12	A 13
sum measured sum goal	TP	$\frac{10}{10}$	$\frac{9}{10}$	$\frac{15}{15}$	$\frac{7}{10}$	$\frac{9}{10}$	$\frac{4}{5}$	$\frac{20}{20}$	$\frac{8}{15}$	$\frac{4}{5}$	$\frac{11}{15}$	$\frac{5}{5}$	$\frac{4}{5}$	$\frac{4}{5}$
	TN	$\frac{355}{355}$	$\frac{355}{355}$	$\frac{345}{350}$	$\frac{355}{355}$	$\frac{355}{355}$	$\frac{360}{360}$	$\frac{345}{345}$	$\frac{350}{350}$	$\frac{360}{360}$	$\frac{345}{350}$	$\frac{360}{360}$	$\frac{360}{360}$	$\frac{360}{360}$
	FP	$\frac{0}{355}$	$\frac{0}{355}$	$\frac{5}{350}$	$\frac{0}{355}$	$\frac{0}{355}$	$\frac{0}{360}$	$\frac{0}{345}$	$\frac{0}{350}$	$\frac{0}{360}$	$\frac{5}{350}$	$\frac{0}{360}$	$\frac{0}{360}$	$\frac{0}{360}$
	FN	$\frac{0}{10}$	$\frac{1}{10}$	$\frac{0}{15}$	$\frac{3}{10}$	$\frac{1}{10}$	$\frac{1}{5}$	$\frac{0}{20}$	$\frac{7}{15}$	$\frac{1}{5}$	$\frac{4}{15}$	$\frac{0}{5}$	$\frac{1}{5}$	$\frac{1}{5}$
		Artifact												

(a) Results for artifacts A1 to A12.

		A 14	A 15	A 16	A 17	A 18	A 19	A 20	A 21	A 22	A 23	A 24	A 25	Total
sum measured sum goal	TP	$\frac{5}{5}$	$\frac{5}{5}$	$\frac{20}{20}$	$\frac{11}{20}$	$\frac{5}{5}$	$\frac{18}{20}$	$\frac{10}{10}$	$\frac{40}{40}$	$\frac{10}{10}$	$\frac{5}{5}$	$\frac{20}{20}$	$\frac{14}{15}$	$\frac{273}{305}$
	TN	$\frac{360}{360}$	$\frac{360}{360}$	$\frac{345}{345}$	$\frac{345}{345}$	$\frac{360}{360}$	$\frac{345}{345}$	$\frac{355}{355}$	$\frac{325}{325}$	$\frac{355}{355}$	$\frac{360}{360}$	$\frac{345}{345}$	$\frac{350}{350}$	$\frac{8810}{8820}$
	FP	$\frac{0}{360}$	$\frac{0}{360}$	$\frac{0}{345}$	$\frac{0}{345}$	$\frac{0}{360}$	$\frac{0}{345}$	$\frac{0}{355}$	$\frac{0}{325}$	$\frac{0}{355}$	$\frac{0}{360}$	$\frac{0}{345}$	$\frac{0}{350}$	$\frac{10}{8820}$
	FN	$\frac{0}{5}$	$\frac{0}{5}$	$\frac{0}{20}$	$\frac{9}{20}$	$\frac{0}{5}$	$\frac{2}{20}$	$\frac{0}{10}$	$\frac{0}{40}$	$\frac{0}{10}$	$\frac{0}{5}$	$\frac{0}{20}$	$\frac{1}{15}$	$\frac{32}{305}$
		Artifact												

(b) Results for artifacts A13 to A24 and their totals.

FIGURE 4: Detection rate for all runs by artifacts A 1 to A 25 and their totals.

The results show that the sensitivity of the result is about 88% for the summed up total, while the specificity is at more than 99%. Although the sensitivity is above 50% for all tests some artifact types show a significantly lower sensitivity. This is especially prominent for A 4, A 8, A 10 and A 17, for which the VAD has a detection rate between 53% and 70%. This is most probably due to these types' reference images, which have very few distinct features. Anomalies in the specificity are only found for A 4 and A 10, with each of them having five false positive matches. Cross-checking the screenshot collection showed that a highly similar reference picture was used for both artifact types, only differing in a small text region and resulting in five evaluation screenshots each. Those were then identified as a match by the opposite artifact type. In conclusion, both error sources can be depleted by choosing more distinct and feature-rich reference images for the artifact types.

Since the k-NN-matching with FLANN is not necessarily deterministic the test was run a second time under the same conditions. It showed the exact same detection rates for all executions, proving stable matching results for the VAD. [MLog]

For measuring the detection rate of scaled versions of the artifact types a second collection of screenshots was generated containing not only translated, but additionally stretched versions of the reference images. This resulted in a sensitivity below 50%, showing that the algorithm does not reliably detect such scaling. But since for real-world use-cases the artifacts would not usually be uniformly stretched as they are in this collection this result is not substantial for the conclusion of an overall success.

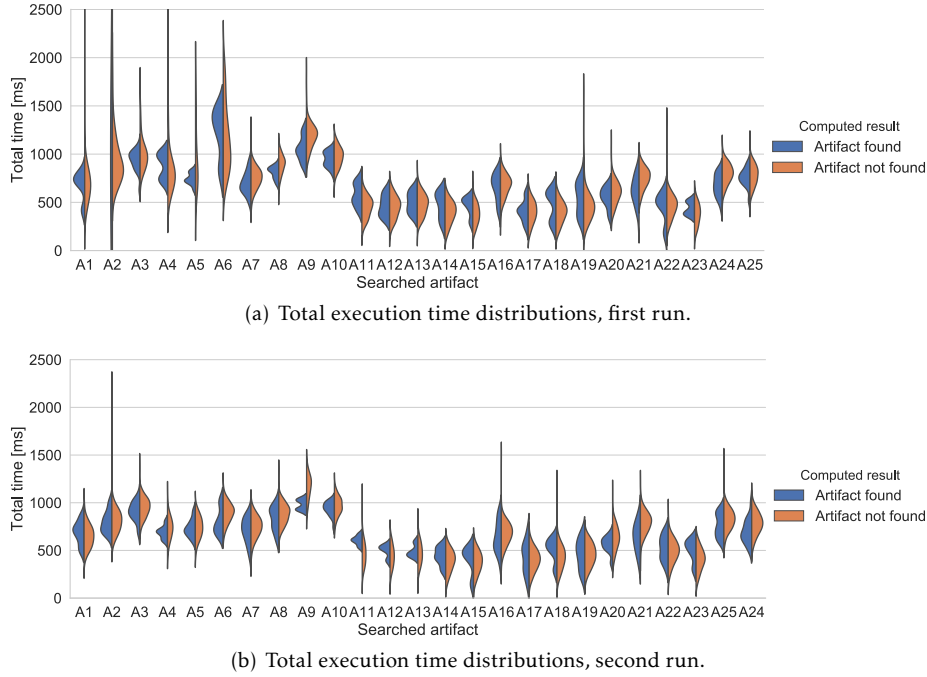


FIGURE 5: Total execution time distributions of the image recognition process for artifacts A 1 to A 25.

4.2.2 EXECUTION RESOURCES

The results of measuring the execution speed in the first evaluation setup are depicted in figure 5 using a violin plot with a kernel density estimator of 0.5. The distribution of the total execution times is plotted in the direction of the y-axis for each artifact type on the x-axis. Since the first setup was run twice there are two plots in this figure, showing similar medians for each artifact type across both measurements.

The violins for each type are split by the computed result to show that the outcome of the computation has only little influence on the execution time on a larger scale. This is not self-explanatory, as the outlier rejection process is designed such that it should skip a substantial part of the process early if a match is improbable. If the median execution times for both cases have such little difference there are multiple possible explanations: For example the rules for skipping may be not strict enough so that they skip too few runs, or the skipped part is not actually that substantial for execution time, or other factors counterbalance the saved execution time.

Both plots show that the median execution time lies within 250ms and 1500ms while for most of the types it is fairly below 1000ms. There also are some divergences for few measurements with randomly occurring execution time peaks of up to about 3000ms and a single run taking 7870ms.

To further assess possible influences on the runtime, the results from the run of the second evaluation setup is shown in figure 6 using the same representation as for figure 5.

The first plot in figure 6(a) shows the execution time needed for loading all of the type's reference images and extracting their features. The increase of the duration dependent on the image count can be assumed to be linear given the implementation. This is an expected behaviour and matches the results of the second plot in figure 6(b) which is showing the duration of the same process for

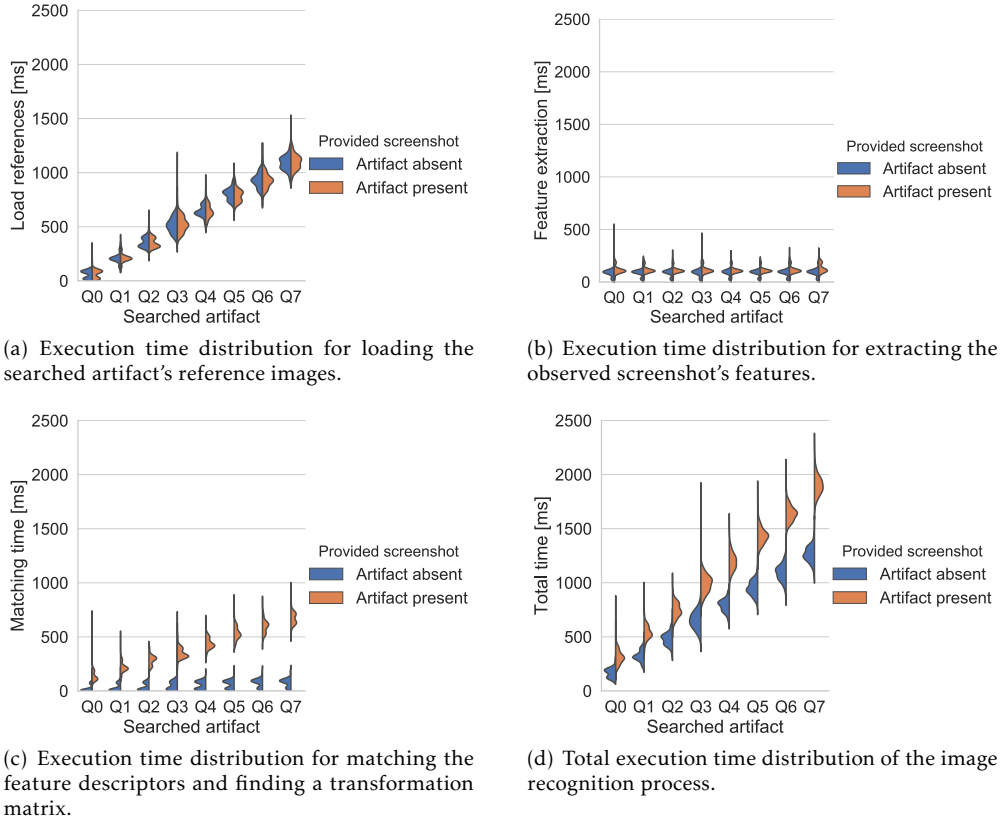


FIGURE 6: Execution time distributions of the image recognition process and parts of it for artifacts Q0 to Q7.

the observed screenshot. Since for all runs there is only one observed screenshot to extract features from, the time needed for this task is constant within a certain margin of error. Both plots also show the independence of these steps from the actual absence or presence of an artifact in the screenshot.

Plot three in figure 6(c) shows the duration of the last two step of the recognition process (from here: *matching time*), namely feature matching and outlier rejection. One can see how the presence of the artifact in the image drastically increases the matching time. With only eight of such artifact types it is not possible to distinguish the type of increase depending on the reference count, but a linear increase is probable. The difference in execution time for s_{ap} and s_{aa} show how the threshold checking of the match count prevents the execution of further filtering steps and thus saves execution time. It also shows a large influence by RANSAC on the execution time.

In figure 6(d) the sum of the previous three values is plotted. It shows that the image recognition process of the VAD takes at least 100ms (s_{aa}) to 250ms (s_{ap}) and can take up to 2000ms or more, if there are multiple reference images in an inconvenient order. The plot also shows that there is a significant difference for s_{ap} and s_{aa} in execution time which is only due to the difference in matching time.

There are major outliers for the measured times of each step, being up to double to triple the time of the median. The reason for this is likely to be a delay in access to the file storage if it

occurs while loading and extracting the references or screenshot. For all steps it is also possible, that an increased CPU-usage by some other program decreased the resources available for the VAD, although such an execution was tried to be suppressed in the virtual machine.

Since the tests are run in a virtual machine, the built-in tools for VirtualBox have been used for measuring the resource consumptions at random points during the execution. The RAM usage of one of VAD's instances was between 30MB and 45MB consistently, while the CPU usage was at almost 100% for all runs.

5 CONCLUSION

The evaluation showed the importance of the quality of reference images for the detection rate, although for most types in the recipes repository the reference images were *good enough* to yield a detection rate of more than about 90%. This detection rate is most probably going to suffice for ITS.APE if it can be confirmed in the real-world and therefore fulfils the 1. requirement of a robust and high detection rate. The chosen parameters of the VAD can still be adjusted based on new results to achieve a higher detection rate. An important aspect for successful employment for the tool is the low rate of false positives being less than 1%, which means that there are a minimum of false alarms when using the VAD.

The measured execution time can certainly be decreased in future work. The measurements of the first evaluation setup prove that on average the VAD fulfils the 2. requirement of a low runtime, but also shows some rare, but extreme divergences randomly occurring. Possible solutions to prevent them and to generally decrease the runtime could include a refined version of the implemented caching mechanism. The new version could e.g. only cache a single artifact type at once, which is then lazy-loaded only if needed for the current execution. The choice of parameters for skipping filtering steps during outlier rejection can be fine-tuned as well to achieve a better performance. Especially the implemented RANSAC algorithm could probably be optimized for execution speed.

In conclusion, an overall successful implementation of the VAD has been accomplished in regards to the measures of the detection rate and resource consumption (see chapter 1) and given the test setup. A real-world evaluation is nevertheless highly advisable to e.g. show the influence of other running processes on the execution time.

6 BIBLIOGRAPHY

- [AMHo5] H. B. Aradhye, G. K. Myers, and J. A. Herson. “Image analysis for efficient categorization of image-based spam e-mail”. In: *Eighth International Conference on Document Analysis and Recognition (ICDAR’05)*. Aug. 2005, 914–918 Vol. 2. doi: [10.1109/ICDAR.2005.135](https://doi.org/10.1109/ICDAR.2005.135).
- [ARK] Amazon.com, Inc. *Amazon Rekognition*. Website. <https://aws.amazon.com/de/rekognition/>. Mar. 2019.
- [Bak+18] Maxim Bakaev et al. “HCI Vision for Automated Analysis and Mining of Web User Interfaces”. In: *ICWE*. Vol. 10845. Lecture Notes in Computer Science. Springer, 2018, pp. 136–144.
- [Cal+10] Michael Calonder et al. “BRIEF: Binary Robust Independent Elementary Features”. In: *Computer Vision – ECCV 2010*. Ed. by Kostas Daniilidis, Petros Maragos, and Nikos Paragios. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 778–792. ISBN: 978-3-642-15561-1.
- [CUDA] NVIDIA Corporation. *CUDA Toolkit*. Website. <https://developer.nvidia.com/cuda-toolkit>. Mar. 2019.
- [CVDL19] Slashdot Media. *OpenCV Download Statistics: All Files*. Website. <https://sourceforge.net/projects/opencvlibrary/files/stats/timeline>. Mar. 2019.
- [EMGU] Canming Huang. *Emgu CV: OpenCV in .NET (C#, VB, C++ and more)*. Website. http://www.emgu.com/wiki/index.php/Emgu_CV. Dec. 2018.
- [FB81] Martin A. Fischler and Robert C. Bolles. “Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography”. In: *Commun. ACM* 24.6 (June 1981), pp. 381–395. ISSN: 0001-0782.
- [GTF] Google Brain Team, Google LLC. *TensorFlow*. Website. <https://www.tensorflow.org/>. Mar. 2019.
- [GVAPI] Google LLC. *Vision API*. Website. <https://cloud.google.com/vision/>. Mar. 2019.
- [KMo8] Georg Klein and David Murray. “Improving the Agility of Keyframe-Based SLAM”. In: *Computer Vision – ECCV 2008*. Ed. by David Forsyth, Philip Torr, and Andrew Zisserman. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 802–815. ISBN: 978-3-540-88688-4.
- [Low04] David G. Lowe. “Distinctive Image Features from Scale-Invariant Keypoints”. In: *International Journal of Computer Vision* 60.2 (Nov. 2004), pp. 91–110. ISSN: 1573-1405.

- [ML09] Marius Muja and David G. Lowe. “Fast approximate nearest neighbors with automatic algorithm configuration”. In: *In VISAPP International Conference on Computer Vision Theory and Applications*. 2009, pp. 331–340.
- [MP18] Maxim Mozgovoy and Evgeny Pyshkin. “Unity Application Testing Automation with Appium and Image Recognition”. In: *Tools and Methods of Program Analysis*. Ed. by Vladimir Itsykson, Andre Scedrov, and Victor Zakharov. Cham: Springer International Publishing, 2018, pp. 139–150. ISBN: 978-3-319-71734-0.
- [NET] Preeti Krishna. *Announcing the .NET Framework 4.7.2*. Website. <https://devblogs.microsoft.com/dotnet/announcing-the-net-framework-4-7-2/>. Apr. 2018.
- [OCV] The OpenCV team. *OpenCV library*. Website. <https://opencv.org/>. Mar. 2019.
- [RDo6] Edward Rosten and Tom Drummond. “Machine Learning for High-Speed Corner Detection”. In: *Computer Vision – ECCV 2006*. Ed. by Ale Leonardis, Horst Bischof, and Axel Pinz. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, pp. 430–443. ISBN: 978-3-540-33833-8.
- [Rub+11] Ethan Rublee et al. “ORB: An Efficient Alternative to SIFT or SURF”. In: *Proceedings of the 2011 International Conference on Computer Vision*. ICCV ’11. Washington, DC, USA: IEEE Computer Society, 2011, pp. 2564–2571. ISBN: 978-1-4577-1101-5. DOI: [10.1109/ICCV.2011.6126544](https://doi.org/10.1109/ICCV.2011.6126544). URL: <http://dx.doi.org/10.1109/ICCV.2011.6126544>.
- [Syk15] Arnold Sykosch. *IT-Security Awareness Penetration Testing – ITS.APT*. Website. <https://itsec.cs.uni-bonn.de/itsapt>. July 2015.
- [TS18] S. A. K. Tareen and Z. Saleem. “A comparative analysis of SIFT, SURF, KAZE, AKAZE, ORB, and BRISK”. In: *2018 International Conference on Computing, Mathematics and Engineering Technologies (iCoMET)*. Mar. 2018, pp. 1–10. DOI: [10.1109/ICOMET.2018.8346440](https://doi.org/10.1109/ICOMET.2018.8346440).
- [VB] Oracle. *Oracle VM VirtualBox*. Website. <https://www.virtualbox.org/>. Mar. 2019.
- [VLF] The VLFeat Authors. *VLFeat - Home*. Website. <http://gs.statcounter.com/screen-resolution-stats>. Mar. 2019.

Appendices

TABLE 1: Mapping of the abbreviated artifact names to their full names

Abbreviation	Full name
A 1	o1_Jibberish-Mittel
A 2	o2_Link_Only-Schwer
A 3	o3_Targo_Bank-Schwer
A 4	o4_Targeted_IT_Abteilung-Einfach
A 5	o4_Targeted_IT_Abteilung-Schwer
A 6	o5_Newsletter-Schwer
A 7	o6_Versichertenkarte-Schwer
A 8	o7_Paypal-Mittel
A 9	o8_Weiterbildung-Schwer
A 10	12_IT_Ticket-Mittel
A 11	browser_advertisement
A 12	browser_defacing
A 13	email-bounce-exchange-de
A 14	exe_anti_virus_extended
A 15	exe_anti_virus_simple
A 16	exe_file_scanner
A 17	exe_login_window
A 18	exe_self_remove
A 19	exe_updater_generic
A 20	exe_updater_human
A 21	exe_updater_java
A 22	exe_updater_remote
A 23	exe_updater_simple
A 24	ms_word_macro
A 25	ms_word_protected_view

LIST OF FIGURES

1	Exemplary recognition process diagram.	5
2	Analysis of the relation between descriptor distance and frequency of the match being an inlier or outlier, dotted for outliers and solid for inliers [Rub+11, p. 4]. . .	7
3	Structure of the artifact types Q ₀ to Q ₇	10
4	Detection rate for all runs by artifacts A ₁ to A ₂₅ and their totals.	12
5	Total execution time distributions of the image recognition process for artifacts A ₁ to A ₂₅	13
6	Execution time distributions of the image recognition process and parts of it for artifacts Q ₀ to Q ₇	14

LIST OF TABLES

1	Mapping of the abbreviated artifact names to their full names	20
---	-------------------------------------------------------------------------	----

STATEMENT OF AUTHORSHIP

I hereby confirm that the work presented in this lab report has been performed and interpreted solely by myself except where explicitly identified to the contrary. I declare that I have used no other sources and aids other than those indicated. This work has not been submitted elsewhere in any other form for the fulfilment of any other degree or qualification.

Bonn, 2nd of April 2019

Felix Rossmann