

**NIST Special Publication 800-168**

---

**Approximate Matching:  
Definition and Terminology**

---

Frank Breitingger  
Barbara Guttman  
Michael McCarrin  
Vassil Roussev  
Douglas White

<http://dx.doi.org/10.6028/NIST.SP.800-168>

**NIST**  
**National Institute of  
Standards and Technology**  
U.S. Department of Commerce



**NIST Special Publication 800-168**

# **Approximate Matching: Definition and Terminology**

Frank Breitingger  
*University of New Haven  
West Haven, CT*

Douglas White  
Barbara Guttman  
*Software and Systems Division  
Information Technology Laboratory*

Michael McCarrin  
*Naval Postgraduate School  
Monterey, CA*

Vassil Roussev  
*University of New Orleans  
New Orleans, LA*

<http://dx.doi.org/10.6028/NIST.SP.800-168>

May 2014



U.S. Department of Commerce  
*Penny Pritzker, Secretary*

National Institute of Standards and Technology  
*Patrick D. Gallagher, Under Secretary of Commerce for Standards and Technology and Director*

## Authority

This publication has been developed by National Institute of Standards and Technology (NIST) to further its statutory responsibilities under the Federal Information Security Management Act (FISMA), Public Law (P.L.) 107-347. NIST is responsible for developing information security standards and guidelines, including minimum requirements for Federal information systems, but such standards and guidelines shall not apply to national security systems without the express approval of appropriate Federal officials exercising policy authority over such systems. This guideline is consistent with the requirements of the Office of Management and Budget (OMB) Circular A-130, Section 8b(3), *Securing Agency Information Systems*, as analyzed in Circular A-130, Appendix IV: *Analysis of Key Sections*. Supplemental information is provided in Circular A-130, Appendix III, *Security of Federal Automated Information Resources*.

Nothing in this publication should be taken to contradict the standards and guidelines made mandatory and binding on Federal agencies by the Secretary of Commerce under statutory authority. Nor should these guidelines be interpreted as altering or superseding the existing authorities of the Secretary of Commerce, Director of the OMB, or any other Federal official. This publication may be used by nongovernmental organizations on a voluntary basis and is not subject to copyright in the United States. Attribution would, however, be appreciated by NIST.

National Institute of Standards and Technology Special Publication 800-168  
Natl. Inst. Stand. Technol. Spec. Publ. 800-168, 21 pages (May 2014)  
<http://dx.doi.org/10.6028/NIST.SP.800-168>  
CODEN: NSPUE2

Certain commercial entities, equipment, or materials may be identified in this document in order to describe an experimental procedure or concept adequately. Such identification is not intended to imply recommendation or endorsement by NIST, nor is it intended to imply that the entities, materials, or equipment are necessarily the best available for the purpose.

There may be references in this publication to other publications currently under development by NIST in accordance with its assigned statutory responsibilities. The information in this publication, including concepts and methodologies, may be used by Federal agencies even before the completion of such companion publications. Thus, until each publication is completed, current requirements, guidelines, and procedures, where they exist, remain operative. For planning and transition purposes, Federal agencies may wish to closely follow the development of these new publications by NIST.

Organizations are encouraged to review all draft publications during public comment periods and provide feedback to NIST. All NIST Computer Security Division publications, other than the ones noted above, are available at <http://csrc.nist.gov/publications>.

### Comments on this publication may be submitted to:

National Institute of Standards and Technology  
Attn: Software and Systems, Information Technology Laboratory  
100 Bureau Drive (Mail Stop 8970) Gaithersburg, MD 20899-8970  
Email: [barbara.guttman@nist.gov](mailto:barbara.guttman@nist.gov)

## **Reports on Computer Systems Technology**

The Information Technology Laboratory (ITL) at the National Institute of Standards and Technology (NIST) promotes the U.S. economy and public welfare by providing technical leadership for the Nation's measurement and standards infrastructure. ITL develops tests, test methods, reference data, proof of concept implementations, and technical analyses to advance the development and productive use of information technology. ITL's responsibilities include the development of management, administrative, technical, and physical standards and guidelines for the cost-effective security and privacy of other than national security-related information in Federal information systems. The Special Publication 800-series reports on ITL's research, guidelines, and outreach efforts in information system security, and its collaborative activities with industry, government, and academic organizations.

### **Abstract**

This document provides a definition of and terminology for approximate matching. Approximate matching is a promising technology designed to identify similarities between two digital artifacts. It is used to find objects that resemble each other or to find objects that are contained in another object. This can be very useful for filtering data for security monitoring, digital forensics, or other applications. The purpose of this document is to provide a definition and terminology to describe approximate matching in order to promote discussion, research, tool development and tool acquisition.

### **Keywords**

approximate matching; bitwise matching; syntactic matching; semantic matching; fuzzy hashing

## **Acknowledgements**

This document is based on the work of the NIST Approximate Matching Working Group, which consists of Frank Breiting, John Delaroderie, Simson Garfinkel, Barbara Guttman, John Kelsey, Jesse Kornblum, Mary Laamanen, Michael McCarrin, Vassil Roussev, Clay Shields, John Tebbutt, Douglas White, and Joel Young.

# Table of Contents

TABLE OF CONTENTS .....	v
INTRODUCTION .....	1
1.1 PURPOSE AND SCOPE .....	1
1.2 AUDIENCE .....	1
2. DEFINITION AND TERMINOLOGY .....	2
2.1 USE CASES .....	2
2.2 TERMINOLOGY .....	3
2.3 ESSENTIAL REQUIREMENTS .....	4
2.4 RELIABILITY OF RESULTS .....	5
3. TESTING BYTEWISE APPROXIMATE MATCHING ALGORITHMS .....	7
3.1 MOTIVATION .....	7
3.2 TEST DATA .....	7
3.3 PROPERTIES OF SPECIAL INTEREST .....	8
REFERENCES .....	10
APPENDIX 1: REFERENCE TEST METHODOLOGY FOR BYTEWISE APPROXIMATE MATCHING .....	11
<b>TABLE 1.</b> CUMULATIVE EMPIRICAL FILE SIZE DISTRIBUTION IN THE <b>govdoc</b> -CORPUS .....	11
APPENDIX 2: APPROXIMATE LONGEST COMMON SUBSTRING .....	13
<b>TABLE 2.</b> EMPIRICAL PROBABILITY DISTRIBUTION FUNCTION ( <i>PDF</i> ) AND CUMULATIVE DISTRIBUTION FUNCTION ( <i>CDF</i> ) FOR $D_R$ .....	14

# Introduction

## 1.1 Purpose and Scope

Approximate matching is a promising technology designed to identify similarities between two digital artifacts. It is used to find objects that resemble each other or to find objects that are contained in another object. This can be very useful for filtering data for security monitoring, digital forensics, or other applications. The purpose of this document is to provide a definition and terminology to describe approximate matching in order to promote discussion, research, tool development and tool acquisition. Approximate matching has the potential to provide valuable filtering in a world inundated with information, but the technique is not widely used. The goal of the document is to provide infrastructure to support advances in approximate matching and its use in forensics and security.

## 1.2 Audience

The intended audience of this document is security digital forensics programmers and other technical professionals with a need to determine, build, or use technology to identify similarity.



## 2. Definition and terminology

Approximate matching is a generic term describing any technique designed to identify similarities between two digital artifacts. In this context, an *artifact* (or an *object*) is defined as an arbitrary byte sequence, such as a file, which has some meaningful interpretation.

Different approximate matching methods may operate at different levels of abstraction. At the lowest level, generic techniques may detect the presence of common byte sequences (substrings) without any attempt to interpret the artifacts. At higher levels, approximate matching can incorporate more abstract analysis. In general lower level methods are expected to be faster and more generic in their applicability, whereas higher level methods are typically more targeted and require more processing.

One common approach in security and forensic analysis is to find identical objects using cryptographic hashing. Approximate matching can be viewed as a generalization of that idea in that, instead of providing a yes/no {0, 1} answer to a comparison, it provides a range of outcomes, [0, 1], with the result interpreted as a measure of similarity.

### 2.1 Use cases

Broadly, there are two types of similarity queries that are of interest – *resemblance* and *containment* [1]. In the case of resemblance, we compare two similarly sized objects and interpret the result as a measure of the commonality between them; for example, two successive versions of a piece of code are likely to resemble each other substantially. When the compared objects differ in size significantly, such as a file and a whole-disk image, the test for commonality is interpreted as a *containment* query because it addresses the question of whether the large object contains the smaller one.

An approximate matching algorithm should address at least one of the following problems (divided according to whether the query type is (R)esemblance or (C)ontainment) [2, 3]:

*Object similarity detection (R)*: identify related artifacts, e.g., different versions of a document.

*Cross correlation (R)*: identify artifacts that share a common object, e.g., a Microsoft Word document and a PDF document containing the same image, or other embedded object.

*Embedded object detection (C)*: identify a given object inside an artifact, e.g., an image within a compound document or an executable inside a memory capture.

*Fragment detection (C)*: identify the presence of traces/fragments of a known artifact, e.g., identify the presence of a file in a network stream based on individual packets.

In most analytical scenarios, approximate matching is used to *filter* data in, or out, based on a known reference set. In security monitoring applications, approximate matching could potentially be used to *blacklist* known bad artifacts, and (by extension) anything closely resembling them. However, approximate matching is not nearly as useful when it comes to *whitelisting* artifacts as malicious content can often be quite similar to benign content; e.g., a backdoored `ssh` server would approximately match regular one.

## 2.2 Terminology

Although the common language definition of ‘similarity’ is sufficient to give an intuitive sense of the term, the multitude of ways in which two artifacts can be said to be similar poses a challenge when attempting to describe the purpose and behavior of approximate matching algorithms. For example, two strings ‘ababa’ and ‘cdcdc’ might be considered similar in that they both have five characters ranging over two alternating values, or they might be treated as dissimilar because they have no common characters. To resolve this ambiguity, approximate matching algorithms define similarity in terms of *features* that represent characteristics of the artifacts pertinent to the algorithm’s method of comparison.

*Features.* The basic elements through which artifacts are compared. Comparison of two *features* always yields a binary {0, 1} outcome indicating a match or non-match; because features are defined as the most basic comparison unit that the algorithm considers, partial matches are not permitted. Generally, a *feature* can be any value derived from an artifact. Each approximate matching algorithm must define the structure of its features and the method by which they are derived. For example, an algorithm might define a feature as a (byte, offset) pair produced by reading the value of a byte and storing it along with the offset at which it was read.

*Feature set.* The set of all features associated with a single artifact is its feature set. Each algorithm must include a criteria by which candidate features are selected for inclusion in this set. For example, an algorithm might select all the (byte, offset) pairs produced by reading every 16th byte in the artifact.

*Similarity.* The similarity of two artifacts, as measured by a particular approximate matching algorithm, is defined as an increasing monotonic function of the number of matching features contained in their respective feature sets.

Based on the level of abstraction of the similarity analysis performed, approximate matching methods can be placed in one of three main categories [4]:

*Bytewise* matching relies only on the sequences of bytes that make up a digital object, without reference to any structures within the data stream, or to any meaning the byte stream may have when appropriately interpreted. Such methods have the widest applicability as they can be applied to any piece of data; however, they also carry the implicit assumption that artifacts that humans perceive as similar have similar byte-level encodings. This assumption is not universally valid. Analyst expertise is necessary to evaluate the significance of a byte-level match.

*Syntactic* matching uses internal structures present in digital objects. For example, the structure of a TCP network packet is defined by an international standard and matching tools can make use of this structure during network packet analysis to match the source, destination or content of the packet. Syntax-sensitive similarity measurements are specific to a particular class of objects that share an encoding but require no interpretation of the content to produce meaningful results.

*Semantic* matching uses contextual attributes of the digital object to interpret the artifact in a manner that more closely corresponds with human perceptual categories. For example, perceptual hashes allow the matching of visually similar images and are unconcerned with the low-level details of how the images are

persistently stored. Semantic methods tend to provide the most specific results but also tend to be the most computationally expensive ones.

In current literature, researchers use a number of terms to refer to various approximate matching methods: *fuzzy hashing* and *similarity hashing* denote bitwise approximate matching; *perceptual hashing* and *robust hashing* denote semantic approximate matching. There is no widely-used pre-existing terminology for syntactic approximate matching as it is mostly viewed as pre-processing (to separate the features) before hashing, or applying a bitwise approximate matching algorithms. For example, network flows are usually reconstructed before any processing is done on them.

Bitwise approximate matching algorithms work in two phases. In the first, a *similarity digest* representation (also referred to as a *signature* or *fingerprint*) is generated from the original data. In the second phase, digests are compared to produce a *similarity* score. More precisely:

*Similarity digest.* A similarity digest is a (compressed) representation of the original data object's feature set that is suitable for comparison with other similarity digests created by the same algorithm. In most cases, the digest is much smaller than the original artifact and the original object is not recoverable from the digest.

Every bitwise approximate matching technique requires at least two core functions:

*Feature extraction function:* identifies and extracts features/attributes from each digital artifact. The mechanism by which features are picked and interpreted depends on the approximate matching algorithm. The representation of this collection is the *similarity digest* of the object.

*Similarity function:* compares two similarity digests and outputs a score. The recommended approach is to assign a score  $s$  in the  $0 \leq s \leq 1$  range, where 0 indicates no similarity and 1 indicates high similarity. This score represents a normalized estimate of the number of matching features in the feature sets corresponding to the artifacts from which the similarity digests were created.

*Normalization strategy:* The similarity function can follow one of two normalization strategies, depending on whether the algorithm describes resemblance or containment. For resemblance queries, the number of matching features will be weighed against the total number of features in both objects. In the case of containment queries, the algorithm may disregard unmatched features in the larger of the objects' two-feature sets.

Because features and feature sets can be arbitrarily complex and, furthermore, deal with byte-level structures to which meaning is not clearly assigned, the interpretation of the similarity score can prove challenging. To address this problem, some approximate matching algorithms make use of an empirically determined *threshold* value to attempt to correlate bitwise similarity scores with higher-level properties of interest. In such cases, the similarity score can be treated as a *confidence score*, where results above the threshold value are considered likely to exhibit common human-recognizable traits.

## 2.3 Essential requirements

Like traditional hash functions, there are several defining characteristics that approximate matching functions should exhibit. Each algorithm should define how it incorporates each of

these properties and how it satisfies the reporting requirements for those properties, where appropriate.

*Similarity preservation:* Similarity digests must be constructed such that the outcome of a comparison between any two digests is uniquely determined by the similarity of the artifacts from which they were produced. That is, if  $A'$  is a similarity digest created from artifact  $A$  and  $B'$  is a similarity digest created from artifact  $B$ , the results of comparing  $A'$  and  $B'$  should be uniquely determined by the similarity of  $A$  and  $B$ .

*Self-evaluation:* The similarity measure should be accompanied by a measure of the accuracy of the matching technique under the circumstances in which it is used, e.g., a margin of error or confidence level. The description of the output score should also state whether a score of 1 indicates an exact match.

*Compression:* A compact similarity digest is desired as it normally allows a faster comparison and requires less storage space. In the best case, it will have a fixed length like the output of traditional hash functions. If the efficiency and reliability of the results remains unchanged, then a shorter similarity digest is preferable.

*Ease of computation:* First, the algorithm description should include the results of testing the runtime efficiency of the feature extraction function and of the similarity function. The former might be expressed relative to a standard hashing algorithm, such as SHA-1.

Second, the algorithm description should state the theoretical complexity for a similarity digest comparison in big O notation. For instance, common lookup complexities for comparing a single digest against a database with  $n$  entries, are:

- $O(1)$  for fixed-length digests stored in hash tables (e.g. dictionaries)
- $O(\log_2 n)$  for fixed-length digests stored in binary trees or a sorted list
- $O(n)$  for fixed-length digests stored in an unsorted list, or other scenarios in which no indexing or sorting is possible

## 2.4 Reliability of results

The reliability of the results for a given approximate matching technique depends on three factors. Each algorithm should define how it incorporates these factors and how it satisfies their reporting requirements.

*Sensitivity & robustness:* The algorithm should provide some measure of its robustness. A technique's robustness will define the operating conditions in which it can function effectively, also called its *performance envelope*. For example, robustness addresses the minimum and maximum sizes of objects that an algorithm can reliably distinguish between.

*Precision & recall:* The algorithm should include a description of the methods used to determine its reliability and to select the test data. Specifically, it should indicate whether test data is culled from existing collections or developed solely to support testing. Test results may include precision, recall, F-score, and other relevant measures of effectiveness.

*Security of results:* The algorithm should indicate whether it includes security properties designed to prevent attacks. Such attacks include manipulation of the matching technique or input data such that a data object appears dissimilar to another object to which it is similar or similar to another object with which it has little in common.

### 3. Testing bitwise approximate matching algorithms

The section gives a general overview of the motives and methods behind testing approximate matching algorithms. Defining a universal testing criteria is difficult; the effectiveness of a given approximate matching algorithm depends largely on the particular task against which it is evaluated. As such, rather than attempt to enumerate a comprehensive evaluation strategy, we focus on key considerations that arise during test design and evaluation. We also describe several properties of approximate matching algorithms that may be useful for characterizing their behavior.

#### 3.1 Motivation

There are at least three major motivations for testing bitwise approximate matching algorithms. The first and most basic is to understand the classes of objects that the algorithm identifies as similar. For example, an algorithm may define its feature set such that it measures similarity as a function of the number of certain common byte sequences shared between two digital artifacts. Because this comparison occurs at a very low level, its semantic interpretation is not obvious. Testing across various tasks and types of data is necessary to determine whether such an algorithm is useful for identifying, say, html documents that contain similar content, or pdfs with the same embedded font.

Second, as an algorithm is developed, a standard suite of tests, such as those proposed by Breitingger, Stivaktakis, and Baier [3], allows new versions to be compared against previous versions. Testing against a fixed set of tasks helps to highlight any changes in the algorithm's similarity score as a function of its input. Improvements in speed or space efficiency can also be measured.

Finally, testing allows approximate matching algorithms to be compared to each other. The use of a standard suite of tasks can be beneficial in this testing scenario also. Additional measures must be taken, however, to ensure that tests produce comparable results. Because the algorithms define their similarity score independently, direct comparison of scores is rarely meaningful. One solution for this problem is to use a threshold or some other criteria to translate scores from their  $[0,1]$  range back to a binary  $\{0,1\}$  output which can be evaluated against known ground truth to produce confusion matrices (Roussev's comparison of sddash and ssdeep offers an example of this approach [2]). Alternately, the tester might define other high-level criteria for the expected behavior of an algorithm in a specific use case of interest, and evaluate the algorithms against this expectation.

#### 3.2 Test Data

Approximate matching algorithms can be tested against either controlled (synthetic) data [7] or real data (i.e., data originally created for purposes other than testing). The former typically consists of randomly generated blocks or files; its main advantage is that ground truth is constructed and, therefore, precisely known. This allows tests to be run automatically and the results to be interpreted with standard statistical measures.

The drawback to using randomly generated data is that it is not necessarily representative. For example, random data tends to be unique, high-entropy and internally varied. These properties may be appropriate for use cases involving compressed or encrypted data, but could prove

misleading for other file types (such as text documents). A potential workaround is to develop more sophisticated randomization that preserves the characteristic properties of the target file type.

The obvious advantage of using real data is that the results can be directly be related to practical applications. However, the challenges of defining a representative sample, establishing the ground truth, and running experiments at scale (without a human in the loop) are non-trivial. See Appendix A1 Example test methodologies for both controlled and real data.

### 3.3 Properties of Special Interest

The following properties may be especially useful for testing and comparing approximate matching algorithms. This list is not meant to be exhaustive. Note that the properties enumerated here are adapted from those proposed by Breitingner, Stivaktakis, and Baier [3]. See their work on the FRASH testing framework for an example of a complete description (including full implementation details and suggested test values) of a testing framework that incorporates these properties.

#### 3.3.1 Efficiency

The efficiency of an algorithm is measured in terms of the space and time it requires; the latter property is further subdivided according to time needed for generating similarity digests and time need for comparing them.

*Space efficiency.* Traditional hash functions return a fixed length fingerprint; in contrast, the length of similarity digests is sometimes variable and proportional to the input length. If the digest is of variable length, space efficiency measures the ratio between input and the digest and returns a percentage value. That is,

$$\text{space efficiency} = \frac{\text{digest length}}{\text{input length}} \quad (1)$$

*Generation efficiency.* Generation efficiency measures the throughput of an algorithm when producing a similarity digest from raw input. To enable useful comparisons across different architectures, it is recommended that the throughput of a standard algorithm implementation, such as SHA-1 in *openssh*, be included as a reference point.

*Comparison efficiency.* The comparison efficiency measures the rate at which similarity digests can be compared. It is useful to have both a formal analysis, which provides the theoretical complexity of the comparison (in O- notation) and an empirical evaluation based on a reference data set.

A related property of note is the ability of the algorithm to leverage parallel computational resources; these may include conventional multi-core CPU architectures, as well as massively parallel ones, such as GPUs. To that end, tests should include scalability analysis, which shows speedup as a function of available hardware concurrency.

#### 3.3.2 Sensitivity

*Sensitivity* is a measure of the ability of an approximate matching algorithm to find correlations among objects based on fine-grained commonality—the smaller the

features being correlated, the more sensitive the algorithm is. Clearly, there is a threshold below which the sensitivity will be too high and all objects will appear similar; it is up to the algorithm designer to identify that threshold and incorporate it into the implementation.

Two examples of tests for measuring an algorithm's sensitivity are fragment detection and single-common-block correlation.

- *Fragment detection.* Fragment detection quantifies the length of the shortest sample from a data object, for which the similarity tool reliably correlates the sample and the whole object. Common uses include correlating a disk block, or network packet to file.
- *Single-common-block correlation.* The single-common-block correlation test is designed to characterize the behavior of an algorithm in the case where two files share a single common object. That is, given two files  $f_1$  and  $f_2$  that share a common object  $O$  (but are otherwise dissimilar), what is the smallest  $O$  for which the similarity tool reliably correlates the two targets?

### 3.3.3 Robustness

*Robustness* is a measure of the ability of an approximate matching algorithm to find correlation among related objects in the presence of noise and routine transformations. Common transformations include internal fragmentation (e.g., during network transmission) and misalignment (adding content during artifact editing).

Alignment robustness and white noise resistance are examples of tests that measure an algorithm's robustness.

- *Alignment robustness.* Alignment robustness is an attempt to quantify the sensitivity of an algorithm to alignment shifts. Specifically, the test analyzes the impact of inserting sequences of bytes at the beginning of one of the compared artifacts.
- *White noise resistance.* This test measures the amount of (uniformly) random noise that can be added to an object before the approximate matching algorithm becomes unable to correlate the original and the modified version. For example, for `ssdeep` [5] it was shown that a few changes distributed over the input are sufficient to prevent a match [6].



## References

1. A. Z. Broder. "On the resemblance and containment of documents," in *In Compression and Complexity of Sequences (SEQUENCES'97)*. IEEE Computer Society, 1997, pp. 21–29.
2. V. Roussev. "An evaluation of forensic similarity hashes," *Digital Forensic Research Workshop*, vol. 8, pp. 34–41, 2011.
3. F. Breitinger, G. Stivaktakis, and H. Baier. "FRASH: A framework to test algorithms of similarity hashing," in *13th Digital Forensics Research Conference (DFRWS'13)*, Monterey, August 2013.
4. F. Breitinger, H. Liu, C. Winter, H. Baier, A. Rybalchenko, and M. Steinebach. "Towards a process model for hash functions in digital forensics," *5th International Conference on Digital Forensics & Cyber Crime*, September 2013.
5. J. Kornblum. "Identifying almost identical files using context triggered piecewise hashing," *Digital Forensic Research Workshop (DFRWS)*, vol. 3S, pp. 91–97, 2006.
6. H. Baier and F. Breitinger. "Security aspects of piecewise hashing in computer forensics," *IT Security Incident Management & IT Forensics (IMF)*, pp. 21–36, May 2011.
7. F. Breitinger, G. Stivaktakis, and V. Roussev. "Evaluating Detection Error Trade-offs for Bitwise Approximate Matching Algorithms," *5th ICST Conference on Digital Forensics & Cyber Crime (ICDF2C)*, September 2013.
8. L. C. Noll. (2001) Fowler / Noll / Vo (FNV) Hash. [Online]. Available: <http://www.isthe.com/chongo/tech/comp/fnv/index.html>

## Appendix 1: Reference test methodology for bitwise approximate matching

This appendix presents detailed example methodologies for testing bitwise approximate matching algorithms with controlled and real data.

**Testing with controlled data.** The main purpose of controlled data experiments is to know exactly the ground truth by carefully constructing the test cases. In this case, the goal is to build artifacts—files—that have known common substrings. The most practical way to accomplish this is to use (pseudo-)random data.

The first step is to determine appropriate sizes for the constructed files. This will vary depending on the intended application, but if the size distribution of the almost 1,000,000 files in the `govdoc`-corpus can be treated as representative, it may be sufficient to evaluate algorithms against six reference file sizes—1, 4, 16, 64, 256 and 1024 KiB. As shown in Table 1, nearly 91% of all files in this set are smaller than 1 MiB.

**Table 1.** Cumulative empirical file size distribution in the `govdoc`-corpus.

File size range (KiB)	$\leq 4$	$\leq 16$	$\leq 64$	$\leq 256$	$\leq 1024$
Cumulative probability (%)	5.4	20.71	52.54	75.82	90.60

*Test methodology.* The proposed approach consists of four basic steps: build a set of unique files, create mutations of them, run approximate matching comparisons between the original and modified versions (for all algorithms), and summarize the results with appropriate statistics. In regards to producing the file modifications, we provide four example *modification methods*, below. Each corresponds to one of the four test examples listed in sections 3.3.2 and 3.3.3. In each,  $f_1$  refers to the original randomly generated file,  $f_2$  refers to a copy of that file that is transformed according to the method, and  $X$  refers to a set of object sizes to be tested with that method.

- *Fragment detection:*  $f_2$  is a fragment of  $f_1$ , selected randomly for each run of the test, where the size of  $f_2$  ranges the values  $X$ , which are selected such that each element of  $X$  is less than the size of  $f_1$ . For reference, the default settings of FRASH [3] recommend  $X = \{.01s_1, .02s_1, .03s_1, .04s_1, .05s_1, .1s_1, .15s_1, .2s_1, .3s_1, .5s_1\}$ , where  $s_1$  is the size of  $f_1$ .
- *Single-common-block correlation:*  $f_1$  and  $f_2$  have equal size and share a common byte string (block), the size of which ranges across the values in  $X$ . The position of the common block and its content are chosen randomly for each file/run combination. For reference, the default settings of FRASH [3] recommend  $X = \{.01s_1, .02s_1, .03s_1, .04s_1, .05s_1, .1s_1, .15s_1, .2s_1, .3s_1, .5s_1\}$ , where  $s_1$  is the size of  $f_1$ .
- *Alignment robustness:*  $f_2$  is prefixed with a random byte string of size  $s$  where  $s$  ranges across the values in  $X$ . The content of the prefix is randomized across runs. For reference, the default settings of FRASH [3] recommend  $X = \{.01s_1, .02s_1, .03s_1, .04s_1, .05s_1, .1s_1, .2s_1\}$ , where  $s_1$  is the size of  $f_1$ .
- *Random-noise resistance:*  $f_2$  is edited with some number,  $n$ , of single-byte edit operations (either additions, deletions or replacements) applied to it at random, where  $n$  ranges across the values in  $X$ . For reference, the default settings of FRASH [3] recommend  $X = \{.005s_1,$

$.01s_1, .015s_1, .02s_1, .025s_1\}$ , where  $s_1$  is the size of  $f_1$  in bytes and the values of  $X$  are rounded to the nearest whole number.

For every choice of file size and modification method, each test has two additional parameters: *file count* and *number of runs*. The former specifies the number of files in the test set; the latter specifies the number of independent test runs to be executed (where each run creates its own new test set).

**Testing with real data.** As already mentioned, two of the main challenges in testing with real data are the choice of representative samples, and the establishment of ground truth. In respect to the former the choice depends, to some degree, on the expected characteristics of the target data. For example, for general evaluation of artifacts found on the Internet, the **govdocs**-corpus is a good starting point.

The focus of this section is to provide an approach for establishing ground truth using automated means. The proposed approach is to use the longest common substring (LCS) as the reference metric and to characterize the behavior of bitwise approximate matching algorithms with respect to this metric.

Using a string comparison algorithm as a reference is a natural choice given that the algorithms treat the data objects as plain strings with no attempt to parse or interpret them. LCS should be considered a first-order approximation; two objects may have a lot more in common than what the LCS result suggests, so further refinements are to be expected at a later stage.

Given an unordered pair of files  $(f_1, f_2)$ , define the absolute ( $L_a$ ) and relative ( $L_r$ ) results as follows:

$$L_a = \text{alcs}_a(f_1, f_2), \text{ where } 0 \leq L_a \leq \min(|f_1|, |f_2|). \quad (2)$$

$$L_r = \frac{\lfloor \frac{100 \times L_a}{\min(|f_1|, |f_2|)} \rfloor}{100}, \text{ where } 0 \leq L_r \leq 1. \quad (3)$$

where  $|f|$  denotes the file size in bytes.

Broadly, any two strings sharing a substring are related; however, we suggest a more practical lower bound on the minimum amount of commonality to declare two files related. Specifically, we require that the absolute size  $L_a$  is at least 100 (bytes) and that the relative result  $L_r$  exceeds 0.5% of the size of the smaller file. More formally, the true positive function  $TP_{lcs}(f_1, f_2)$  is defined as

$$TP_{lcs}(f_1, f_2) \equiv L_a \geq 100 \wedge L_r \geq .01 \quad (4)$$

(Note: result of  $L_r$  is rounded and thus 0.5 is equal to 1.)

Clearly, the true negative function  $TN_{lcs}(f_1, f_2) = \neg TP_{lcs}(f_1, f_2)$ .

*Approximate ground truth* LCS is a well-studied problem and has known solutions of quadratic time complexity— $O(mn)$ , where  $m$  and  $n$  are the string lengths. Given that files can be large, the exact solution quickly becomes too burdensome to be practical. Therefore, we suggest an approximation of the longest common substring which, by design, provides a lower bound on LCS; details are given in Appendix 2.

## Appendix 2: Approximate longest common substring

The basic idea of the approximate longest common substring metric (aLCS) is not to compare files byte-by-byte but rather block-by-block. To identify the blocks, we apply the rolling hash from **ssdeep**. Our settings aim at having blocks of  $\approx 80$  bytes. Instead of comparing blocks byte-wise, each one is hashed and compared using the 64-bit FNV-1a hash [8]. Besides the hash value, we also store the entropy and length for each block in a final linear list called *alcs-digest*; a reference implementation is publicly available.<sup>1</sup>

Let  $L_a$  denote the absolute longest common substring of two *alcs-digests*. Comparing two *alcs-digests* is equal to comparing two linear lists. If the hash of an item on list *A* has the same value as the hash of an item on list *B*, we are convinced that  $L_a$  is greater than or equal to the length of the blocks corresponding to the hashes. If two consecutive items on list *A* have the same hash values as two consecutive items on list *B*, we sum up the length of both blocks to receive  $L_a$ . Of course, the usage of hash functions implies the possibility of false positives. Nevertheless, this is an easy and fast method to get a good estimation of the longest common substring.

*Implementation details.* The tool is implemented in C and proceeds in three steps: reading, hashing and comparing, which are declared in the main function. Because it is a command line tool, it can be executed by `./aLCS <dir>`.

First, all files in *dir* are read. Out of the file names, we create “hash-tasks” which are added to a thread pool. A hash-task contains the path to a file and denotes “hash file”. Depending on the number of threads, these tasks are processed. Once all *alcs-digests* are created, we perform an all-against-all comparison. Therefore, we create compare-tasks (compare *file<sub>1</sub>* against *file<sub>2</sub>*) which are again added to the thread pool. The output is printed to the standard output.

The reference implementation has three main settings configurable in `header/config.h`. **MIN\_LCS** is the minimum  $L_a$  length which is printed to `stdio` and is by default 0 (all comparisons are printed). The **THREAD\_POOL\_QUEUE\_SIZE** is the length of the queue and should be  $\frac{\text{fileamount} \times (\text{fileamount} - 1)}{2}$ . **NUMTHREADS** is the number of threads which should be equal to the number of cores.

*Verification of ground truth.* To verify the correctness of our approximate longest common substring, we compared the results against LCS for a subset of *t5*. In order to do this, we implemented a parallelized LCS tool written in C++.<sup>2</sup> The output is a summary file structured similarly to our aLCS output: **file1 | file2 | LCS**. A small, ruby script is used to compare the LCS- summary and aLCS-summary.

Our subset consists of 201 randomly selected files. We compare these files using aLCS as well as LCS and finally compare both summaries. All  $\frac{(200) \times (201)}{2} = 20,100$  comparisons yield *alcs* scores in the correct range, *i.e.*,  $0 \leq alcs \leq lcs$ .

---

<sup>1</sup> <https://www.dasec.h-da.de/staff/breitinger-frank/#downloads> (last accessed 2013-05-09).

<sup>2</sup> <https://www.dasec.h-da.de/staff/breitinger-frank/#downloads> (last accessed 2013-05-09).

We also consider the distribution of the differences between the LCS and aLCS scores. Specifically, we define  $d_r$  for files  $f_1$  and  $f_2$  as follows:

$$d_r = \left\lceil \frac{[lcs(f_1, f_2) - aLCS(f_1, f_2)]}{\min(|f_1|, |f_2|)} \right\rceil, d_r \in 0, 1, \dots, 100.$$

In other words, we consider the score difference relative to the size of the smaller of the two files, and build the empirical distribution in Table 2. As we can see, upwards of 95% of the observed differences do not exceed 3% of the size of the smaller files – we consider this a reasonable starting point for our purposes (further research may refine this). If anything, this should give tools a slight boost as the available commonality would be underestimated.

**Table 2.** Empirical probability distribution function (*pdf*) and cumulative distribution function (*cdf*) for  $d_r$ .

$X$	0	1	2	3	4	5	10	15	20
$P_r\{d_r = X\}$	0.8869	0.0449	0.0155	0.0040	0.0047	0.0116	0.0062	0.0001	0.0000
$P_r\{d_r \leq X\}$	0.8869	0.9318	0.9473	0.9513	0.9561	0.9677	0.9834	0.9992	0.9999