

Seminar paper on Probabilistic Near-Duplicate Detection Using Simhash

Arne Beer, MN 6489196, University of Hamburg

08.07.2019

1 Introduction

In the age of the modern Internet, many services depend in large parts on crawlers and proper document detection and duplication elimination, near duplicate document detection becomes a necessity. Real time detection of which website has already been visited and whether a website is new or just has been edited are important tasks during crawling [3].

Standard hashing functions are often inefficient and operate in $O(n^2)$ space requirements for RAM and computing time [2]. At the same time, the size of available documents grow steadily. Google's website index alone has multiple hundreds of billions of web pages and over 100 Petabyte in size, according to their information website [1].

This paper attends to the paper *Seminar paper on Probabilistic Near-Duplicate Detection Using Simhash* [5]. The main challenge tackled by this paper is to find all matching pairs of fingerprints withing a certain Hamming distance h . At this time, the fasted implementation for this procedure has been *Block-Permuted Hamming Search* (BPHS), which requires RAM space at least four times the size of the whole dataset. The authors of [5] aim to design a new algorithm that allows significantly faster online and batch document comparison and furthermore reduces RAM requirements, in exchange for a small percentage of recall loss.

In the first chapters, the basics for understanding this topic will be explained. Afterwards the proposed algorithm will be looked at and the authors' findings will be discussed.

2 Conventional Hashing

Hashing is a technique, which is used to map data of an arbitrary size to a fingerprint with some fixed size. This procedure could be seen as a function $f(i) \rightarrow j$, which produces a value j from any value i , where $j \in H$ and H is the set of values of the fixed length s with $s \in \mathbb{N}$. Well-known hash functions are, for instance, *MD5* or *SHA256*. These hashing functions are commonly used to check whether two files are absolutely identical or, for instance, to verify that a file has not been corrupted during transport. This is possible, since these hashing functions are designed to flip half of the output hash bits on average, if an input bit changes [2]. Without this property it would be easier to change the input without the hash signature being modified. This would allow malicious third parties to, for instance, change code in a binary, without users being able to detect the change with the help of this hash and would require a full byte level comparison between the original and the copied file to verify its integrity.

If, on the other hand, one's goal is to find near duplicates, which are identical for the most part, but sometimes only differ by a few bits or bytes, this hashing approach immediately becomes useless, due to this property. Due to the need for a hashing algorithm, that creates a fingerprint based on the features and structure of the input data, *Simhash* has been created.

3 Simhash

Simhash is a procedure used to create a fingerprint of a any kind of data. This fingerprint can then be used to, for instance, inspect two files for similarity.

The process for creating such a fingerprint can completely differ depending on the features in the hashed data one is interested in. In case one wants to find similar binary files, it would be reasonable to split the data into equal chunks and use these chunks as features. For websites or documents, looking at the composition and structure of text could be a viable approach to select features for hashing. The original data can then be seen as a high dimensional vector of features.

The size of available features can vary significantly and is completely in the hand of the designer for a *Simhash* implementation. It's important to note, that there exists no clear guideline on which features of a data set are interesting and which features can be discarded. The performance of a *Simhash* implementation thereby also depends on the chosen features and the respective properties of the dataset. Such features can be for instance binary chunks, file extension [4], individual words, tags or URLs [5].

After determining in which way feature are extracted from the original data, each

feature is

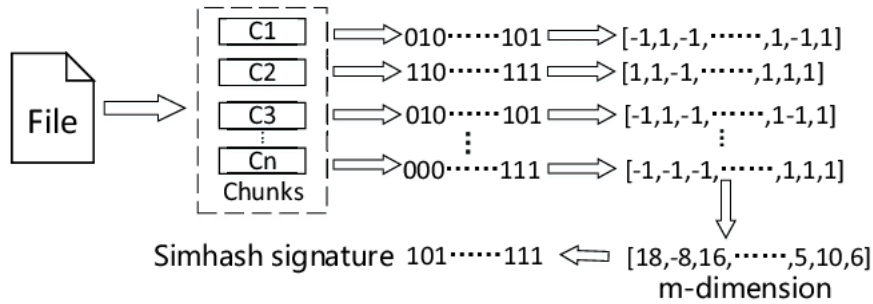


Figure 1: Process of calculating a *Simhash* fingerprint. A binary file split into chunks, which are then hashed. Combined all hashes results in the desired fingerprint. [6]

[4]

3.1

3.2 Achievements of Session Juggler

4 Impact in the scientific community

5 Relevance as of 2018

References

- [1] Google. *How Search organizes information*. 2019. URL: <https://www.google.com/search/howsearchworks/crawling-indexing/> (visited on 07/02/2019).
- [2] A.G. Konheim. *Hashing in Computer Science: Fifty Years of Slicing and Dicing*. July 2010. Chap. 8. DOI: 10.1002/9780470630617.
- [3] Hsin-Tsang Lee et al. "IRLbot: Scaling to 6 Billion Pages and Beyond". In: *ACM Transactions on the Web (TWEB)* 3 (Jan. 2008), p. 8. DOI: 10.1145/1541822.1541823.
- [4] Caitlin Sadowski and Greg Levin. "SimHash: Hash-based Similarity Detection". In: Dec. 2007. DOI: 10.1.1.473.7179.

- [5] Sadhan Sood and Dmitri Loguinov. “Probabilistic Near-Duplicate Detection Using Simhash”. In: Oct. 2011, pp. 1117–1126. DOI: 10 . 1145 / 2063576 . 2063737.
- [6] Yongtao Zhou et al. “EPAS: A Sampling Based Similarity Identification Algorithm for the Cloud”. In: *IEEE Transactions on Cloud Computing* PP (Feb. 2016), pp. 1–1. DOI: 10 . 1109 / TCC . 2016 . 2527646.