

# WikiWho: Precise and Efficient Attribution of Authorship of Revisioned Content

Fabian Flöck  
Institute AIFB  
Karlsruhe Institute of Technology, Germany  
fabian.floeck@kit.edu

Maribel Acosta  
Institute AIFB  
Karlsruhe Institute of Technology, Germany  
maribel.acosta@kit.edu

## ABSTRACT

Revisioned text content is present in numerous collaboration platforms on the Web, most notably Wikis. To track authorship of text tokens in such systems has many potential applications; the identification of main authors for licensing reasons or tracing collaborative writing patterns over time, to name some. In this context, two main challenges arise. First, it is critical for such an authorship tracking system to be precise in its attributions, to be reliable for further processing. Second, it has to run efficiently even on very large datasets, such as Wikipedia. As a solution, we propose a graph-based model to represent revisioned content and an algorithm over this model that tackles both issues effectively. We describe the optimal implementation and design choices when tuning it to a Wiki environment. We further present a gold standard of 240 tokens from English Wikipedia articles annotated with their origin. This gold standard was created manually and confirmed by multiple independent users of a crowdsourcing platform. It is the first gold standard of this kind and quality and our solution achieves an average of 95% precision on this data set. We also perform a first-ever precision evaluation of the state-of-the-art algorithm for the task, exceeding it by over 10% on average. Our approach outperforms the execution time of the state-of-the-art by one order of magnitude, as we demonstrate on a sample of over 240 English Wikipedia articles. We argue that the increased size of an optional materialization of our results by about 10% compared to the baseline is a favorable trade-off, given the large advantage in runtime performance.

## Categories and Subject Descriptors

I.7.1 [Document and Text Processing]: Document and Text Editing—*Version control*; K.4.3 [Computers and Society]: Organizational Impacts—*Computer-supported collaborative work*; K.4.1 [Computers and Society]: Public Policy Issues—*Intellectual property rights*

## Keywords

Wikipedia; authorship; version control; content modeling; community-driven content creation; collaborative writing; online collaboration

## 1. INTRODUCTION

Collaborative authoring of text-based content, such as Wiki pages, shared editable documents or software code is a common sight on the Web today. Most of these systems keep track of the different revisions (i.e., versions) of the content created with every edit (or commit). In this paper, we propose an approach to efficiently and precisely assign the revision of origin – and with it, its author – to particular text tokens (mostly whole words delimited by whitespaces). While line-level tracking of changes is a feature of many code revisioning systems, this level of attribution can prove insufficient in the case of a community that produces large amounts of collaboratively written natural language text. A word-level tracking that is proven to be trustable in its attributions and that can trace reintroductions and moves of chunks of text in an acceptable runtime for end-users can prove very useful, especially on a platform like Wikipedia, as has been previously discussed [6, 7]. While research shows that Wikipedians are motivated by the recognition by their peers that comes with authoring content [8], more practical needs also exist [6]. To reuse a Wikipedia article under the CC-BY-SA license, for example, might require to list the main authors of the article, which are not easily retrievable as there exists not straightforward way in the Mediawiki software to show authors of single pieces of text for a particular revision.<sup>1</sup> Authorship tracking in articles can further raise awareness by editors and readers for editing dynamics, concentration of authorship [7], tracing back certain viewpoints or generally understanding the evolution of an article. Recently, Wikimedia Deutschland e.V. introduced the “Article Monitor”,<sup>2</sup> aiming to assist users with these exact issues and making use of the results of a basic authorship algorithm [7] whose general concept we use as a foundation in this work. The Wikipedia community has come up with a number of intended solutions related to the authorship attribution problem on word level, which highlights the utility of such a solution for Wikipedians.<sup>3</sup>

As outlined in previous work [6], the attribution problem at this fine-grained level and in highly dynamic environments like Wikipedia is not trivial, as we will discuss when introducing our content model in Section 3.1. Frequent reintroductions, content moves and other actions can be hard to monitor. In code revisioning, similar issues can emerge and finer-grained attribution techniques can have similar merits, as small changes of a few characters might have great effects, just as (re)introducing larger code chunks.

Copyright is held by the International World Wide Web Conference Committee (IW3C2). IW3C2 reserves the right to provide a hyperlink to the author’s site if the Material is used in electronic media.  
WWW’14, April 7–11, 2014, Seoul, Korea.  
ACM 978-1-4503-2744-2/14/04.  
<http://dx.doi.org/10.1145/2566486.2568026>.

<sup>1</sup>[https://en.wikipedia.org/wiki/Wikipedia:Reusing\\_Wikipedia\\_content](https://en.wikipedia.org/wiki/Wikipedia:Reusing_Wikipedia_content), CC-BY-SA: <http://creativecommons.org/licenses/by-sa/3.0/>  
<sup>2</sup><http://tools.wmflabs.org/render/stools/>  
<sup>3</sup>[https://en.wikipedia.org/wiki/Wikipedia:Tools#Page\\_histories](https://en.wikipedia.org/wiki/Wikipedia:Tools#Page_histories), cf. also Keppmann et al. [10]

Against this background, the main contributions of this paper are: the model for revisioned content we propose (Section 3), the algorithm we build upon that model (Section 4), the generation of a gold standard for precision testing (Section 5.1.1) and the experimental evaluation of precision, runtime and materialization size in comparison to the state-of-the-art (remainder of Section 5).

Although we use the example of the English Wikipedia as inspiration and testing ground, the proposed model and algorithm can be understood as components of a more generally applicable method for revisioned content. We are convinced that many of the assumptions made for the use case of Wikipedia also hold true not only for other Wikis but also for other revisioned content systems.

## 2. RELATED WORK

In the context of Software Engineering, content attribution has been studied in terms of code ownership. In programming, lines of code are still used for measuring technical quality of source code [4], as well as a basic unit to identify contributors. Therefore, solutions to trace code ownership are designed to operate on a coarse-grained level [13, 14]. Decentralized Source Code Management systems such as Apache Subversion [13] or Git<sup>4</sup> provide a feature to keep track of changes line-by-line. This functionality is denominated `blame` or `praise` depending on the system. When a contributor performs a change on a line of code, the system attributes the whole line to that user. In this way, `blame` allows to identify who last modified each line in a given revision of a file, but the information about the origin of the content is unaccounted for. The `blame` approach is a suitable solution to detect defects in collaborative software development [15] as well as to select expert developers for implementing required changes in programs [12], yet does not provide an appropriate mechanism to trace the *first* author of the content at a more fine-grained level such as single words or special characters, to which we refer as “tokens” in the remainder.

To detect authorship information in Wikipedia article text, several analysis approaches have been employed. *HistoryFlow* by Viegas et al. [17] assigns sentences of a text to the editor who created or changed them. It doesn’t however acknowledge deleted content that was later reconstructed as being written by the original editor. More importantly, by operating on a sentence level, small changes like spelling mistake corrections lead to wrongly recognizing the correcting editor as the author of the whole sentence.

*Wikitrust* generates a visual mark-up of trusted and untrusted passages in any Wikipedia article [3, 1, 2].<sup>5</sup> To track authorship, longest matches for all word sequences of the current revision are searched for in the preceding revision and in previously existing, but now deleted word-chunks. In this way, *Wikitrust* can as well detect reintroduced words and assign the original author – an important feature, as “reverts” to formerly existing revisions are commonplace in Wikipedia. The underlying algorithm is, however, a variation of a greedy algorithm [1], known to look for local optima, which in the case of determining word authorship can lead to grave misinterpretations when word sequences are moved rather than inserted or deleted only [5].

Flöck and Rodchenko [7] introduce an authorship attribution approach for Wikipedia based on a tree model of paragraphs and sentences. The precision of the results according to the evaluation lies at 59.2% compared to 48.4% for *Wikitrust*, values that are rather unsatisfactory for productive usage requiring users to place confidence in the computed attributions. The work presented here builds

on the foundations laid by [7], but formalizes a model based on a  $k$ -partite graph to represent paragraphs, sentences and tokens much more efficiently compared to the tree model of [7]. We also gain over 30% in precision in respect to [7] by refining the conceptualization of authorship and tokenization of text.

The most relevant related work, by de Alfaro and Shavlovsky [6], proposes an algorithm for attributing authorship to tokens in revisioned content based on the concept of processing content as sequences of neighboring tokens and finding “relevant” matches in the revision history given a parameterized rarity function, for the work in question defined as the length of the sequence. This means that a token is uniquely identified solely by its local neighbours. To store the authorship attributions the annotated revision history is remembered by means of a “trie” structure. A trie is a tree where the authorship information is stored in the leaves, while the intermediate nodes are empty. The arcs of this trie are labeled with tokens. All the arcs leaving a given node correspond to neighbors of the token(s) represented in the label of the arc incoming to that node. The algorithm is tested in terms of runtime and the size of a materialization (storage in secondary memory) of the results, but not precision. It takes into account reintroduction of text and can keep track of authorship on a word or smaller token level and therefore conforms exactly to the goals set out for our work.

## 3. MODELLING REVISIONED CONTENT

The following subsections outline our model for representing revisioned content.

### 3.1 A Model Based on Observations of Real-world Revisioned Writing

When observing collaborative writing in systems that rely on long-term, incremental refinement by a large group of users, namely Wikipedia and its sister-projects, certain patterns become salient. In the following we list our conclusions from these observations and from studying related literature (e.g., [9, 11, 17]). These assertions build the conceptual foundation for the content model developed in Section 3.2.

The first assessment is that a considerable part of editing activity after the initial basic writing phase consists of moving, deleting and reintroducing content, while not adding much new subject matter per edit. A notable number of edits consists of reintroductions, which are mostly reverts due to vandalism; another reason for reverts is, e.g., an “edit war” between disagreeing factions. Moves of text sequences are also a regular sight, where a sentence or paragraph gets shifted to another position in the article without alterations. Another sizeable amount of edits is predominantly changing only very small parts of the article per edit, incrementally revising the text. This pattern is occasionally interrupted by a burst of new content; for instance, in case of a larger addition or a fundamental re-write. Still, very often the changes implemented per edit do not span more than one section or paragraph – frequently, they don’t even transgress the boundary of a sentence. These assertions point out that methodically keeping track of reused or relocated content plays an important role when intending to efficiently monitor authorship over large data in such a system.

Regarding the conceptual definition of “authorship”, the larger context of a token plays a crucial role. The paragraph or the section it is embedded in can be as important for the interpretation of its meaning as its immediate token neighbours in a sequence. The same exact string of tokens, even up to the length of a sentence (e.g., a figure of speech), might mean something completely different in one section of a text (e.g., an introduction) than in another segment, where it was potentially introduced for a different

<sup>4</sup><http://git-scm.com/>

<sup>5</sup><http://www.wikitrust.net/>

purpose. It can therefore not necessarily be seen as an exact copy of the same sequence in another position, entailing the attribution of the same author. An example: One editor writes the sequence “A theory is that” in front of a factual statement  $A$  at the top of an article. Later, a different editor adds the same four words ahead of a completely different statement  $B$ . Both authors use the same terms and add the assertion that the subsequent statements are mere theories instead of proven facts. Yet, the declaration that “statement  $B$  is a theory” can only be attributed to the later editor as the first author used the same chain of words in a different context and with a completely different goal (namely to call statement  $A$  a theory). This also applies, e.g., if two authors use an identical literature reference in order to prove different facts. Basically, just by comparing local neighbours of words, as it is done by de Alfaro and Shavlovsky [6], who in practice use four-word sequences, the larger context of the tokens is not taken into account. Extending these neighbour-tracking sequences to sentence or paragraph length, on the other hand, is not constructive, as this would contradict the initial idea of exact word-level author attribution.

Tracking provenance hierarchically, by assigning tokens to a larger enclosing unit (sentences), and linking these to another superordinate element like a paragraph provides a more exact identification of tokens than mere local contextualization. Another key advantage compared to the method of [6] is that not all tokens in the text have to be analyzed if the enclosing unit has already been identified as unchanged or as reintroduced from an earlier revision. This enables a more rapid processing of the data, as changes, like mentioned above, often only affect fractions of the whole article.

### 3.2 A Graph-based Model for Revised Text Content

We propose a model to represent revised content as a  $k$ -partite graph, where the content is partitioned into units of discourse in writing, i.e., paragraphs, sentences and tokens (which can consist of words or single characters as we will explain in Section 4.3).

In order to illustrate the representation of revised content with the proposed model, consider a Wiki page with three revisions  $r_0$ ,  $r_1$  and  $r_2$  as depicted in Figure 1. The first revision,  $r_0$ , contains a single paragraph,  $p_0$ , which is composed of only one sentence,  $s_0$ , with two tokens ( $t_0$  and  $t_1$ ). The labels over the arcs represent the relative position of the nodes. For instance, the token  $t_0$  and  $t_1$  are located in positions 0 and 1 of the sentence  $s_0$ , respectively. The second revision,  $r_1$ , creates two new paragraphs,  $p_1$  and  $p_2$ . The paragraph,  $p_1$ , is written by reusing the sentence  $s_0$  from revision  $r_0$  followed by a new sentence,  $s_1$ . The third revision,  $r_2$ , reuses paragraph  $p_2$  from the previous revision and creates two new paragraphs,  $p_3$  and  $p_4$ . In addition,  $p_3$  contains a new sentence which reuses the token  $t_2$  originally inserted in the previous revision.

*Definition 1. (A graph-based model for revised content).* Given a revised content document, it can be represented as a  $k$ -partite graph, with  $k = 4$ ,  $G = (V, E, \mathbb{N}_0, \phi)$  defined as follows:

- The set of vertices  $V$  in  $G$  is composed of four different subsets  $R, P, S, T$ , i.e.,  $V = R \cup P \cup S \cup T$ . The subset  $R$  represents the revisions of a given document,  $P$  the paragraphs of the document,  $S$  the sentences that compose the paragraphs, and  $T$  the tokens (words, special characters, etc.) in the sentences. The subsets  $R, P, S, T$  are pairwise disjoint.
- The set of arcs  $E$  in  $G$  is partitioned into  $k - 1$  cuts as follows:  $E = \langle R, P \rangle \cup \langle P, S \rangle \cup \langle S, T \rangle$ . The arcs in  $G$  represent the relationship of *containment*, e.g., if  $p \in P$ ,  $s \in S$  and  $(p, s) \in E$  then the paragraph  $p$  contains the sentence  $s$ .

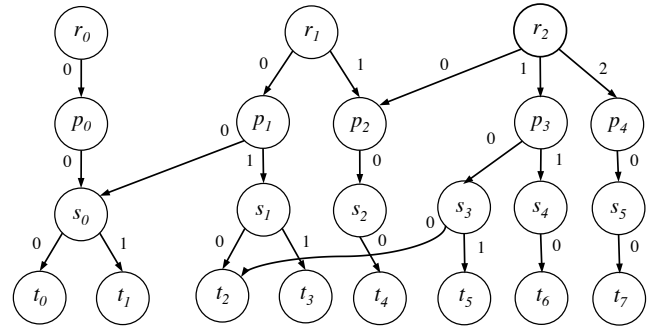


Figure 1: Example of the revised content graph. Revisions are represented by nodes  $r$ , paragraphs by  $p$ , sentences by  $s$ , and tokens by  $t$ . Arcs between nodes correspond to the *containment* relation.

- A labelling mapping  $\phi : E \rightarrow \mathbb{N}_0$  over the arcs in  $G$  represents the *relative position* of a token, sentence or paragraph in a sentence, paragraph or revision, respectively. Each arc is labeled only once, therefore these labels are not updated.

Additionally, it is necessary to keep record of the sequence in which the revisions were generated. Since adding arcs between revision nodes violates the partite graph definition, we represent this information by annotating the revision nodes with an identifier (*label*) such that if revision  $r_i$  is a predecessor of revision  $r_j$ , the condition following condition is met:  $label(r_i) < label(r_j)$ .

### 3.3 Restrictions Over the Model

In the following we present the restrictions to guarantee consistency within the model. We refer to paragraphs, sentences or tokens as ‘content elements’.

The first property refers to the number of content elements within a revision. Particularly, this property allows the definition of an empty revision (with no content elements).

*Property 1.* The number of arcs that leaves a revision vertex (denoted  $deg^+(\cdot)$ ) must be greater than or equal to 0.

$$\forall v \in R (deg^+(v) \geq 0) \quad (1)$$

The second property restricts the existence of empty paragraphs or sentences, i.e., each paragraph or sentence must contain at least one content element.

*Property 2.* The number of arcs that leave a paragraph or sentence vertex must be greater than 0.

$$\forall v, v \in P \vee v \in S (deg^+(v) > 0) \quad (2)$$

The following property establishes that paragraphs, sentences or tokens must be associated to at least one revision, paragraph or sentence, respectively.

*Property 3.* The number of incoming arcs of a paragraph, sentence or token vertex must be greater than 0.

$$\forall v, v \in P \vee v \in S \vee v \in T (deg^-(v) > 0) \quad (3)$$

The last property refers to the labelling of the arcs that leave a given vertex. This property states that each content element (paragraph, sentence or token) can only occupy a single relative position.

*Property 4.* The label of an arc that leaves a vertex in  $R, P$  or  $S$  must be between 0 and the number of arcs that leaves that vertex.

The set of arcs that leave a vertex is denoted as  $d^+(\cdot)$ . Moreover, the labels of the arcs that leave a given vertex must be unique.

$$\forall v, v \in R \vee v \in P \vee v \in S (\exists e \in d^+(v) (0 \leq \phi(e) < \text{deg}^+(v)) \wedge \forall e_1, e_2 \in d^+(v) (e_1 \neq e_2 \rightarrow \phi(e_1) \neq \phi(e_2))) \quad (4)$$

### 3.4 Operations Over the Model

We define four different operations over the model that correspond to the actions that can be performed by editing a document. In the following,  $\text{path}(a, b)$  is defined as the set of paths from vertex  $a$  to vertex  $b$ . The first operation defines the creation of a new (initially empty) revision, which consists of adding a vertex to the set of revision vertices.

*Definition 2. (Creation of a new revision).* Let  $r_{i-1}$  be the last revision in the graph. The operation  $\text{createRevision}(r_i)$  represents the creation of a new revision (denominated *current revision*) and is defined as follows:

$$R := \{r_i\} \cup R \quad (5)$$

After the newly created revision  $r_i$  is added to the set of revisions, the corresponding label of  $r_i$  is assigned as follows:

$$\text{label}(r_i) := |R| - 1 \quad (6)$$

The following operation allows creating a new paragraph, sentence or token in a certain position of a given revision, paragraph or sentence, respectively. This operation consists of adding a vertex to the corresponding vertex partition, and an edge in the corresponding arc cut annotated with the position of the element.

*Definition 3. (Creation of content).* Let  $x$  and  $y$  be content elements such that  $y$  is a new element to be added in  $x$  at position  $\alpha$ . The operation  $\text{createContent}(x, y, \alpha)$ , with  $\alpha = \phi((x, y))$ , is defined as follows:

$$\begin{cases} P := \{y\} \cup P \wedge \langle R, P \rangle := \{((x, y), \alpha)\} \cup \langle R, P \rangle & \text{if } x \in R \\ S := \{y\} \cup S \wedge \langle P, S \rangle := \{((x, y), \alpha)\} \cup \langle P, S \rangle & \text{if } x \in P \\ T := \{y\} \cup T \wedge \langle S, T \rangle := \{((x, y), \alpha)\} \cup \langle S, T \rangle & \text{if } x \in S \end{cases} \quad (7)$$

In addition, the creation of an element  $y$  in a given revision  $r_i$  meets the following condition:

$$\exists \rho (\rho \in \text{path}(r_i, y)) \quad (8)$$

The third operation defines the action of copying, or reintroducing, content from an old revision. This operation consists of creating an arc from the content element of the current revision to the copied element, and labeling the arc with the relative position of the element in the current revision.

*Definition 4. (Copying content from an old revision).* Let  $r_i$  ( $i > 0$ ) be the current revision and  $r_{i-k}$  ( $0 < k \leq i$ ) an old revision. Let  $x$  and  $y$  be content elements such that  $y$  is an element copied from revision  $r_{i-k}$  in the element  $x$  of revision  $r_i$  at position  $\alpha$ . The operation  $\text{copyContent}(x, y, \alpha)$ , with  $\alpha = \phi((x, y))$ , is defined as follows:

$$\begin{cases} \langle R, P \rangle := \{((x, y), \alpha)\} \cup \langle R, P \rangle & \text{if } x \in R, y \in P \\ \langle P, S \rangle := \{((x, y), \alpha)\} \cup \langle P, S \rangle & \text{if } x \in P, y \in S \\ \langle S, T \rangle := \{((x, y), \alpha)\} \cup \langle S, T \rangle & \text{if } x \in S, y \in T \end{cases} \quad (9)$$

In addition, copying an element  $y$  from revision  $r_{i-k}$  to revision  $r_i$  meets the following condition:

$$\exists \rho (\rho \in \text{path}(r_{i-k}, y)) \wedge \exists \rho' (\rho' \in \text{path}(r_i, y)) \quad (10)$$

The last operation is the deletion of content, which models the case when content from the previous revision is removed. This operation requires no alteration on the structures of the model, since elements are never removed from revisioned content.

*Definition 5. (Deletion of content).* Let  $r_i$  ( $i > 0$ ) be the current revision and  $y$  the element from the previous revision to be removed. The deletion of  $y$  in  $r_i$  meets the following condition:

$$\exists \rho (\rho \in \text{path}(r_{i-1}, y)) \wedge \nexists \rho' (\rho' \in \text{path}(r_i, y)) \quad (11)$$

## 4. AUTHORSHIP ALGORITHM

This section describes the implementation of an authorship attribution algorithm based on the presented model.

### 4.1 The Authorship Attribution Problem

The authorship attribution problem consists of identifying for each token the revision in which the token originated. This problem has been previously introduced [6], where each token is annotated with an origin label denoted as  $\Theta$ . In the following we devise a theoretical solution to the authorship attribution problem for a given token, built on top of the proposed graph-based model.

**THEOREM 1.** (A solution to the authorship attribution problem). Let  $G = (V, E, \mathbb{N}_0, \phi)$  be the graph of a given revisioned content, modelled according to Definition 1. The authorship of a token  $t$  can be determined by identifying all the revisions where the token occurs and selecting the revision that was generated first in sequential order, i.e., the revision with the minimum label.

$$\forall t \in T(\Theta(t) := \min \{ \text{label}(r_i) \mid \exists \rho (\rho \in \text{path}(r_i, t)) \wedge r_i \in R \})$$

**PROOF.** We want to demonstrate that if  $t$  ( $t \in T$ ) originated in revision  $r_i$  ( $r_i \in R$ ), then  $\Theta(t) = \text{label}(r_i)$ . By contradiction, let's assume that  $t$  originated in  $r_i$ , but  $\Theta(t) \neq \text{label}(r_i)$ . Therefore, we have two cases:  $\Theta(t) < \text{label}(r_i)$  or  $\Theta(t) > \text{label}(r_i)$ . Furthermore, there exists a revision  $r_j$  ( $r_j \in R$ ) such that  $\Theta(t) = \text{label}(r_j)$ . By hypothesis,  $t$  didn't originate in  $r_j$  but in  $r_i$ . Therefore,  $t$  must have been copied from  $r_i$  to  $r_j$ . According to Definition 4, an element can only be copied from an old revision, thus the case  $\Theta(t) < \text{label}(r_i)$  is discarded. In the other case, we can affirm that  $r_i$  is a predecessor of  $r_j$ , therefore  $\min(\text{label}(r_i), \text{label}(r_j)) = \text{label}(r_i)$ . By the definition of  $\Theta(t)$ , the only possibility for not selecting  $r_i$  as the origin of  $t$  is that there does not exist a path from  $r_i$  to  $t$  (contradiction to Definition 3).  $\square$

### 4.2 Implementation of the Proposed Solution

We have demonstrated that our proposed model provides a straightforward solution to the authorship attribution problem in revisioned content. In the following we devise an algorithm to build this model while generating origin labels of tokens simultaneously. The source code and further information are available online.<sup>6</sup>

Algorithm 1 outlines our proposed solution, WikiWho, an algorithm that constructs a graph according to Definition 1 to represent a document with revisioned content. WikiWho follows a breadth-first search (BFS) strategy to build the graph structures for each revision and assigns the corresponding origin labels to each token.

In order to illustrate the execution of Algorithm 1 consider the revisioned content presented in Figure 2. In this example, the document  $D$  is composed of three revisions: Revision 0, Revision 1 and Revision 2. The algorithm starts processing Revision 0, and creates the corresponding revision node  $r_0$  (Algorithm 1, line 4). Then, the content is split into paragraphs; in our example there is only one

<sup>6</sup><http://people.aifb.kit.edu/ffl/wikiwho/>

### Algorithm 1 WikiWho algorithm

**Input:** A document  $D$  with revisioned content  $r_0, r_1, \dots, r_{n-1}$ .  
**Output:** A graph  $G = (V, E, \mathbb{N}_0, \phi)$  representing the revisioned content from  $D$ .

```

1: create an empty graph  $G = (V, E, \mathbb{N}_0, \phi)$ 
2: create an empty queue  $Q$ 
3: for  $i$  in  $0, 1 \dots n - 1$  do
4:    $G.createRevision(r_i)$ 
5:    $label(r_i) \leftarrow i$ 
6:    $y' \leftarrow tokenize(r_i)$ 
7:   enqueue  $(r_i, y)$  onto  $Q$  for all  $y$  in  $y'$ 
8:    $x_{prev} \leftarrow NULL$ 
9:    $diffed \leftarrow FALSE$ 
10:  while  $Q$  is not empty do
11:     $(x, y) \leftarrow Q.dequeue()$ 
12:    if  $x$  is a sentence  $\wedge !diffed$  then
13:      calculate  $diff$  of unmarked tokens of  $r_{i-1}$  against unmarked tokens
        of  $r_i$  ( $i > 0$ )
14:       $diffed \leftarrow TRUE$ 
15:    end if
16:    if  $x = x_{prev}$  then
17:       $\alpha \leftarrow \alpha + 1$ 
18:    else
19:       $\alpha \leftarrow 0$ 
20:       $x_{prev} \leftarrow x$ 
21:    end if
22:    if  $y \in V \wedge y$  is not marked then {detects reintroduction of content}
23:       $G.copyElement(x, y, \alpha)$ 
24:      mark all the nodes reachable from  $y$ , including  $y$ 
25:    else
26:       $G.createElement(x, y, \alpha)$ 
27:      if  $y$  is a token then
28:         $\Theta(y) \leftarrow label(r_i)$ 
29:      else
30:         $z' \leftarrow tokenize(y)$ 
31:        enqueue  $(y, z)$  onto  $Q$  for all  $z$  in  $z'$ 
32:      end if
33:    end if
34:  end while
35:  unmark all the marked nodes
36: end for
37: return  $G$ 

```

paragraph ( $p_0$ ), which is split into a single sentence ( $s_0$ ). Once the algorithm has tokenized all the sentence nodes, it proceeds to calculate the  $diff$  operation (line 13) between the current text and token nodes from the previous revision.<sup>7</sup> For the first revision, this operation corresponds to  $diff('', 'One house .')$ , i.e., empty content diffed vs. the tokens of  $r_0$ . The  $diff$  output states that all the tokens in revision  $r_0$  are new and the algorithm creates the corresponding token nodes (line 26), and annotates them with  $\Theta = 0$ . The current state of the graph is presented as the leftmost of the three sections of Figure 3.

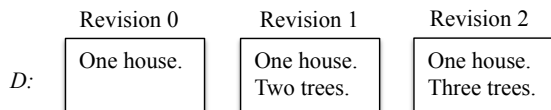


Figure 2: Example of a document  $D$  with revisioned content.  $D$  contains three revisions, each one with a single paragraph.

After processing Revision 0 in Figure 2, in the next iteration the algorithm creates the revision node  $r_1$ . In this revision the paragraph has changed w.r.t. to the previous revision, therefore the node  $p_1$  is created. One of the sentences of  $p_1$  corresponds to  $s_0$  – created in revision  $r_0$  – and the algorithm marks all the nodes

<sup>7</sup>For the actual implementation, the *difflib* library of Python was used: <http://docs.python.org/2/library/difflib.html>

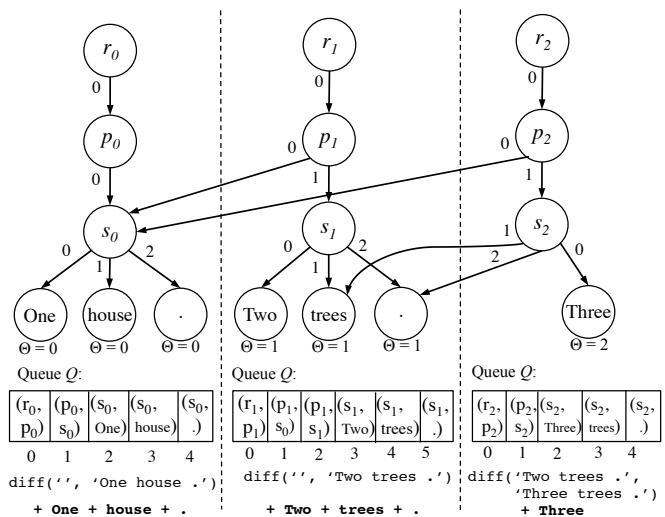


Figure 3: Execution of WikiWho for the example from Figure 2. Sections delimited by dashed lines represent the state of the graph after each revision. At the bottom, the progress of the queue  $Q$  and the output of the  $diff$  for each revision iteration are depicted.

reachable from  $s_0$ , including  $s_0$ . This is a case of reintroduction of content and it is detected by the Algorithm 1 on line 22, since the analysed vertex is already in the graph and has not been used previously in the current revision since it is not marked. The other sentence in  $p_1$  is new, therefore the node  $s_1$  is created. At this step, the  $diff$  is calculated over the sentence 'Two trees .' and the set of unmarked token nodes from  $r_1$  (which is  $\emptyset$ ). The three tokens are identified as new and annotated with  $\Theta = 1$ . The state of the graph at the end of this iteration is illustrated by the combination of the left and the middle section of Figure 3.

In the last iteration of the example, the node  $r_2$  is created. This revision contains a new paragraph  $p_2$ , composed of  $s_0$  and a new sentence  $s_2$ . After processing the sentences, the algorithm calculates  $diff('Two trees .', 'Three trees .')$ . Note that the sentence 'One house .' is not considered by the  $diff$ , since these nodes were marked when  $s_0$  was processed. According to the  $diff$ , only the token 'Three' is new and is annotated with  $\Theta = 2$ . Figure 3 depicts the current state of the graph.

#### 4.2.1 Representation of Content Nodes

As explained earlier, revision nodes are uniquely annotated with a *label* (line 5 of Algorithm 1), usually representing the sequential order in which the revisions were generated. In a Wiki environment, the revision identifier provided may serve as a label in WikiWho. Paragraph and sentence nodes are identified by a hash value. The hash value is the result of applying the MD5 algorithm [16] over the content of the paragraph or sentences, respectively. Token nodes contain the actual text of the revisions, i.e., the single tokens that compose the revisioned content. WikiWho annotates the token nodes with their corresponding origin label ( $\Theta$ ), as shown on line 28 of Algorithm 1. This avoids the calculation of all the paths from revision nodes to a token to retrieve its authorship information.

#### 4.2.2 Implementation of Operations

One crucial operation of WikiWho is determining whether a certain node  $y$  belongs already to the graph (line 22 of Algorithm 1). Depending on the type of the node, this step can be implemented differently. When  $y$  is a paragraph or a sentence, the algorithm only checks the corresponding node partition, i.e., if  $y \in P$  or  $y \in S$ ,

respectively. When  $y$  is a token, the decision whether the token is new or not relies on the output of the `diff` operation (line 13).

On lines 6 and 30, Algorithm 1 performs the tokenization of the text unit that is currently being processed. Tokenization refers to the process of splitting the text into more fine-grained units. We define a token as the smallest unit, be it indivisible. Details regarding the definition of grammatical units are discussed in Section 4.3.

Once the graph is built, accessing the authorship information for each computed revision is straightforward. The origin labels of the content in  $r_i$  can be retrieved by traversing the graph  $G$  with any search approach, e.g., depth-first search (DFS), considering  $r_i$  as the root node and visiting the nodes in the order induced by  $\phi$ .

### 4.3 Design Issue: Tokenization

To achieve authorship attributions that correspond to the “original” author of a given token, it is important to choose the tokenization very carefully in respect to the specific context of the system and its usage. An overly fine-grained tokenization – e.g., on character level – might prove counter-productive if it is not necessary to determine the origin of single characters and make the interpretation of the results too complex for end-users. On the other hand, using only demarcations such as periods to identify sentences can prove too coarse. Whole tokens can in that case spuriously be re-attributed to new authors even when only minor adjustments are applied. As the optimal tokenization can vary immensely between different contexts, we concentrated here on a Wiki environment, especially Wikipedia. We believe, however, that the presented design choices are applicable as well for other Wikis that follow roughly the same patterns of editing and collaborative writing.

When processing the complete source text of articles, as we do with WikiWho, it is not only important to take into account the intricacies of natural language (as it appears on the article front-end) and a corresponding optimal tokenization, but also the function-oriented “Wiki markup” language of the Mediawiki software.<sup>8</sup> Small changes in the markup can entail important changes in the front-end of an article, be it for instance the inclusion of a template via only a few pasted characters or the setting of links. Consider the following example, where a contributor writes the word `Germany` and another contributor in a subsequent revision adds markup content to represent the same word as an internal link (`[[Germany]]`) to the respective article. Using only white spaces as separators would lead to counting the latter string as a new token; using a similarity-metric like a Levenshtein distance might lead to an attribution of the whole string to the former author. What actually happens here is that the word “Germany” was written by one author and the link (specified with “[[” and “]]”) was set by another. Both are elementally important and distinct actions in Wikipedia. Many more of these examples could be given, pertaining to templates, language links, references and numerous others.

Besides white spaces we chose the most commonly used functional characters of the Wiki markup as delimiters, such as “|”, “[”, “]”, “=”, “<”, “>”, to name some. We further split sentences (as defined in our model) at common sentence delimiters such as “.”, “?” etc., and paragraphs at double line breaks, which are used in Wikipedia to begin a new paragraph in the text. All delimiters were also treated as tokens, as they fulfill important functions in the text. We determined all of these demarcations after extensive testing with real article data until we reached a splitting we deemed optimal to achieve the best balance of precision and efficiency.<sup>9</sup>

<sup>8</sup>[http://en.wikipedia.org/wiki/Help:Wiki\\_markup](http://en.wikipedia.org/wiki/Help:Wiki_markup)

<sup>9</sup>The list appears in `Text.py`, part of the algorithm implementation.

## 4.4 Optimization for Wiki Environments: Vandalism Detection

The most expensive operation of the proposed algorithm is the `diff`. We are interested in detecting those vandalism attacks that change large amounts of content from one revision to another, significantly affecting the performance of the `diff` operation. There are different types of vandalism in Wikipedia, such as removing large parts of a page, or modifying a page in a way that adds a lot of vandalistic content.<sup>10</sup> In order to avoid the computation of the `diff` in the previous cases, we implemented two simple vandalism detection techniques that don’t impose a large computational overhead and filter out only the most obstructive cases that would increase runtime notably.

### *Percentage of size change from one revision to another.*

This mechanism is triggered when a large amount of content gets removed at once, by comparing the current content size versus the size of the previous revision. An example of this type of vandalism is *page blanking*, which signifies deleting all the content of a Wiki page.<sup>11</sup> Since the size of the early revisions of a Wiki page can fluctuate notably, this technique is fired only when the revisions have reached a certain size. To not filter out revisions where much content is moved to a different article in good faith, we analyze the edit comment log provided in the article history dumps.

### *Token density.*

This proposed technique aims at detecting vandalism that consists of adding large amounts of disruptive content, often composed of the same text repeated numerous times. For these cases, we propose a measurement called *token density*, defined as follows. Consider the bag of tokens of a revision  $r$  as the result of splitting the revision’s content with a tokenization mechanism. This bag can be formally represented as a multiset  $T_r$ , which consists of a set  $T'_r$  – constructed from removing duplicates in  $T_r$  – and a function  $m : T'_r \rightarrow \mathbb{N}_0$  that denotes the number of times an element of  $T'_r$  occurs in  $T_r$ . We calculate the token density as follows:

$$\text{tokenDensity}(T_r) = \begin{cases} \frac{\sum_{t \in T'_r} m(t)}{|T'_r|} & \text{if } T_r \neq \emptyset \\ 0 & \text{if } T_r = \emptyset \end{cases}$$

A high token density suggests that the content is composed of the same words. From this computation we discard tokens corresponding to Wiki markup elements, since they appear several times in an edit and might be misinterpreted as vandalism. These vandalism filters must be configured with very relaxed thresholds such that no false positives are generated, which was successfully achieved with the values set in the experiments of this work (cf. Section 5.2.2). Note that the objective of implementing these techniques is to avoid the computation of the `diff` operation on large amounts of irrelevant content. We do not aim at applying these techniques as general solutions for the problem of vandalism detection in Wikis.

## 5. EXPERIMENTAL STUDY

We empirically analyzed the performance of the proposed algorithm WikiWho and compared it to the algorithm “A3” by de Alfaro and Shavlovsky [6], which can be regarded as the current bench-

<sup>10</sup>[http://en.wikipedia.org/wiki/Wikipedia:Vandalism\\_types](http://en.wikipedia.org/wiki/Wikipedia:Vandalism_types)

<sup>11</sup>[http://en.wikipedia.org/wiki/Wikipedia:Page\\_blanking](http://en.wikipedia.org/wiki/Wikipedia:Page_blanking)

mark for the given task.<sup>12</sup> In our experiments we report on the execution time of the evaluated algorithms as well as their precision in finding the correct revision of origin for a token. Regarding precision, this is the first evaluation for both algorithms. The datasets, gold standard and further details of the experimental results are available online.<sup>13</sup> For all articles analyzed in the following evaluations, the full history in XML format was retrieved from the English Wikipedia via the MediaWiki `Special:Export` mechanism.<sup>14</sup>

## 5.1 Evaluation of Precision

In the following we explain the three-step procedure to construct and validate the gold standard. We measured the quality of the evaluated algorithms by comparing their authorship attributions with the results of the gold standard.

### 5.1.1 Creating a Gold Standard for Authorship in Revisioned Content

To create the gold standard we selected 40 English Wikipedia articles. Ten articles each were randomly picked from the following four revision-size ranges:<sup>15</sup> articles with i) over 10,000 revisions, ii) 5,000-10,000 revisions, iii) 500-5,000 revisions and iv) 100-500 revisions. The reason for this grouped sampling process was to include a sufficient number of articles that present a challenge to the algorithms when picking the correct revisions of origin, as a higher number of revisions naturally increases the difficulty of the task, as more candidate solutions exist.<sup>16</sup> The latest revision at the point of retrieval of the articles was the “starting revision” for whose tokens the authorship was determined and which is denoted in the gold standard. The text plus markup of each of the 40 articles was split into tokens as described in Section 4.3. Out of this tokenized content, for each article, six instances were randomly selected, resulting in a total of 240 tokens. For each of these, the final gold standard contains the revision in which they first appeared (revision of origin) and the starting revision. To assign the correct revision of origin to all of these tokens, we followed three consecutive steps.

STEP 1: Three researchers of the AIFB institute, including the two authors, manually searched the “Revision History”<sup>17</sup> of the respective 40 articles for the origin of each of the 240 tokens in the gold standard independently from each other. No common interpretation of what constitutes a “correct origin” was agreed on beforehand but was entirely up to the individuals. If the researchers initially disagreed on the correct origin of a token, this disagreement could in most cases be resolved. Only in three cases was this not achieved, so that they were excluded from the gold standard and replaced with new randomly selected tokens.

STEP 2: Next, the gold standard was validated by users of the crowdsourcing platform Amazon Mechanical Turk (hereafter called “turkers”).<sup>18</sup> We selected two random tokens for each article in the gold standard. We then created a Human Intelligence Task (HIT) on Mechanical Turk for each of these 80 tokens to be validated by 10 distinct turkers each. We paid 15 US\$ cents and selected turkers

with a past acceptance rate of over 90% and at least 1,000 completed HITs.<sup>19</sup> A HIT was composed of the following elements:

a) A link to a copy of the starting revision of the Wikipedia article with the highlighted token. If the token only appeared in the markup, we represented an excerpt of the markup as a picture next to the front-end text where it appears in the article HTML, explaining to look for it in the markup.

b) A link to the Wikipedia “Difference View” of the revision of origin proposed by the gold standard. It shows which changes the edit introduced that lead to that revision.<sup>20</sup>

c) Detailed instructions explaining how to use the above mentioned pages and a description of what solution was sought. Three different conditions had to be fulfilled by the proposed revision: First, a string equivalent to the token should indeed have been added in that revision (and not only be moved inside the article text). Second, the token added should be the “same” token as highlighted in our gold standard solution. We explicitly left it open to the turkers to interpret what “same” meant to them and gave only one simple, unambiguous example, explaining that not any string matching the gold standard token was looked for but the specific token in the context that it is presented in (e.g., if the token was a specific “and”, we would not be looking for any “and”). The third condition was that the token was actually added in the given revision for the first time and not just reintroduced, e.g., in the course of a vandalism revert. Turkers could choose between one answer option indicating “correct revision”, three choices pointing out various errors and a fifth option with a text box if they had found a revision that was more likely to be the origin of the token. For option 5, we offered a bonus payment of 5 US\$ to propose a better solution than the one presented and gave a detailed manual on how to search the revision history page of a Wikipedia article by hand as well as a list of tools by the Wikimedia community that can be helpful with the task.

RESULTS OF STEP 2: The 800 answers we received as the result of this experiment included 24 answers suggesting a better solution, but none of them fulfilled all three conditions. We therefore reposted these HITs once we assessed them. As these turkers had spent 17 minutes on average for the task and obviously had tried to find a better solution, they were paid bonuses ex-post. Overall, turkers spent from 40 seconds to 13 minutes solving the task, with an average of 4 minutes 49 seconds. Turkers thus spent considerable time assessing the correctness of the presented solutions.<sup>21</sup>

We report on the results aggregating the answers of the “incorrect” options (option 5 was handled as mentioned above). On average, the solutions received 89% agreement. 65 tokens received nine to ten out of ten agreement votes. 12 solutions received 8/10 and three received 7/10 “correct” votes. In the latter cases the disagreeing turkers pointed in 7 of 9 answers at the lack of a matching string being added in the suggested revision, although this was in fact the case.<sup>22</sup> Overall, we consider the result of this experiment to compellingly support the proposed gold standard solutions.

STEP 3: As a further test we ran the WikiWho algorithm as well as the A3 algorithm (in two different variants), as explained in the following Section 5.1.2. For 67 of the 240 tokens in the gold standard at least one of the algorithms produced a result deviating from

<sup>12</sup>Retrieved from: <https://sites.google.com/a/ucsc.edu/luca/the-wikipedia-authorship-project>

<sup>13</sup><http://people.aifb.kit.edu/ffl/wikiwho/>

<sup>14</sup><http://en.wikipedia.org/wiki/Special:Export>

<sup>15</sup>The “random article” feature of Wikipedia was used. Redirect or disambiguation pages were skipped.

<sup>16</sup>Articles under 100 revisions are not challenging for the task. We did sample test-runs with non-crowdsourced test answers and both algorithms scored very close to a precision of 1.0.

<sup>17</sup>Example for the article “Korea”: <https://en.wikipedia.org/w/index.php?title=Korea&action=history>

<sup>18</sup><http://www.mturk.com>

<sup>19</sup>The pay rate was the result of a number of tries with rates at 10 and 13 cents that did not attract enough turkers.

<sup>20</sup>Example diff: <https://en.wikipedia.org/w/index.php?title=Korea&diff=574837201>

<sup>21</sup>This excludes 12 turkers whose HITs were rejected and reposted for obviously incorrect answers, such as choosing option 5 and not reporting a better solution.

<sup>22</sup>We believe this could have been because of particular nature of the respective diffs, where the token was hard to track.

the gold standard. For all of these 67 tokens we set up a Mechanical Turk experiment with the same settings as explained in step 2. In this HIT, however, we presented the turkers with three differing possible revisions of origin and asked them which one was most likely correct or if none of them was (option 4). One of the three solutions was always the gold standard answer and one or two were solutions by one of the algorithms, depending on how many algorithm results disagreed. If only two differing solutions were available, the third one was filled with an incorrect control answer. Answer positions were randomly changed in each HIT.

RESULTS OF STEP 3: 670 single answers were retrieved for the 67 tokens. The general agreement score for the gold standard solution was 81%, with 7/10 or more votes validating the gold standard as correct for 63 tokens. Given the nature of the task and the different possible interpretations, we consider the gold standard to have gained a solid affirmation for these tokens. In four cases, however, the turkers disagreed decisively with the gold standard. In two of these instances, there was complete disagreement over the right solution, while in two other examples four users each endorsed the differing WikiWho and the differing A3 solution, respectively. We therefore removed these tokens from the following evaluation in 5.1.2 since a certain solution for these cases is lacking.<sup>23</sup> The remaining 63 tokens achieved an agreement of 83%.

As a conclusion to these three steps of quality assurance we can assume that the gold standard is sufficiently robust to test algorithm precision against it. We are however publishing the gold standard and encourage the community to assess and expand it further.

### 5.1.2 Measuring the Precision of the Algorithms

After validating the gold standard, WikiWho and A3 algorithms were tested for their ability to correctly detect the revisions of origin for each token. The evaluation metric was precision defined as:  $p = \frac{TP}{TP+FP}$ , where a true positive (TP) means that the authorship label computed by the algorithm is matching the gold standard described in 5.1.1 and otherwise is a false positive (FP).

Three articles in the gold standard from the revisions-size bracket over 10,000 had to be excluded due to technical reasons and are hence exempt from all following experiments to guarantee the same data basis.<sup>24</sup> The remaining 37 articles encompass 218 tokens.

The A3 algorithm we retrieved includes a filter that seems to be intended to remove the Wiki markup that does not appear on the HTML front-end of an article.<sup>25</sup> More important is however that all citations and references get discarded, although they appear in the front-end and can in some cases make up large parts of the article, not to mention their functional importance for the credibility of Wikipedia articles. Hence we ran one variant of the A3 algorithm with this markup filter disabled (henceforth “**A3 MF-OFF**”)<sup>26</sup>, also because our aim was to compare WikiWho to another algorithm that is able to process the entire source text. The unaltered version of the A3 algorithm will be referred to as “**A3 MF-ON**”. A3 MF-ON yielded results for 138 of the 218 tokens as the remaining part was filtered out. We therefore compared its output to the result for the same 138 tokens by WikiWho, as can be seen in the lower part of Table 1. A3 MF-OFF produced output for the whole set and we compared it to the full results of WikiWho, listed in the upper part of Table 1.

<sup>23</sup>We marked these cases in the published gold standard accordingly.

<sup>24</sup>The A3 algorithm did not process these articles despite several intents to resolve the issue. The files were unaltered XML-dumps from the Wikipedia servers. The articles are “Vladimir Putin”, “Apollo11” and “Armenian Genocide”.

<sup>25</sup>The filter is not described in [6].

<sup>26</sup>Apart from this change the settings used in [6] were replicated.

Table 1: Precision comparison of WikiWho and A3

$x \in$	ALL	[10k,∞)	[5k,10k)	[500,5k)	[100,500)
Full sample					
$p$ WikiWho	0.95	0.97	0.93	0.95	0.95
$p$ A3 MF-OFF	0.77	0.77	0.64	0.76	0.87
Gain in $p$ by WikiWho	0.18*	0.20*	0.29*	0.19*	0.08
Available results $n$	218	58	42	58	60
Sample restricted to output of A3 MF-ON ( $n - 80$ )					
$p$ WikiWho (restricted)	0.96	0.97	0.89	1.00	0.95
$p$ A3 MF-ON	0.81	0.69	0.70	0.88	0.95
Gain in $p$ by WikiWho	0.15*	0.28*	0.19	0.12*	0.00
Available results $n$	138	39	27	34	38

Notes:  $n$  = number of tokens,  $k$  = one thousand,  $p$  = precision,  $x$  = number of revisions per article, \* = difference significant at 0.05 (paired t-test)

WikiWho scores at 18% and 15% higher precision overall, respectively for the full and the restricted token sample. As becomes evident from the results, the gain in precision by WikiWho turned out especially high for the two biggest revision-size brackets, while it is lower for the 5000 to 50,000 bracket and much lower and non-existent, respectively, for the smallest-size bracket. On one hand, this seems to indicate that for articles with up to 500 revisions, the difference between the two approaches is negligible and both have a very high precision. Given the long tail of small articles in Wikipedia, this is a very encouraging result. On the other hand, with increasing editing activity and therefore growing number of revisions of an article, it seems to become harder for the A3 algorithm to correctly determine the authorship of certain tokens, while WikiWho can sustain a high level of precision, even for articles with over 10,000 revisions. Given the steady growth of Wikipedia and the size of other revisioned content these approaches might be adaptable to, this is an important aspect of scalability. Moreover, particularly when processing the much “dirtier” Wiki markup, WikiWho seems to have a notable advantage when it comes to precisely determining authorship.

## 5.2 Evaluation of Execution Time

We measured the algorithm time for computing authorship labelling of revisioned content from Wikipedia pages.

### 5.2.1 Experimental Set-up

We used two datasets created by retrieving the full revision history content for each article from the English Wikipedia in XML format.<sup>14</sup> *Dataset 1* was generated by randomly selecting Wiki pages in the article namespace that were no redirects or disambiguation pages. This dataset is comprised of 45,917 revisions in 210 articles, i.e., an average number of 219 revisions per article; the average revision size is 2,968 KB. *Dataset 2* contains the Wiki pages used in the quality evaluation presented in Section 5.1.1. Its articles are larger, with an average number of revisions of 5,952 and an average revision size of 461,522 KB per article. This allowed for some “heavy load” testing. This last dataset is composed of 36 articles with a total of 214,255 revisions.<sup>27</sup>

<sup>27</sup>We excluded again the three articles mentioned in Section 5.1.2 as well as “Jesus”, as it would run over 5 hours for some settings.



We defined execution time as the time elapsed between the point when the algorithm reads the first revision and the point when the authorship labelling of the last revision of a given article is computed. Both algorithms are implemented in Python and the time was measured with the `time.time()` command from the Python library. The experiments were all executed on a dedicated OS X machine with a 2.5 GHz Intel Core i5 processor and 4GB RAM.

### 5.2.2 Algorithm Settings

The A3 algorithm was set according to the configuration presented by de Alfaro and Shavlovsky [6], with sequence length as the rarity function and a threshold equal to 4. The tokenization implemented by A3 uses only whitespaces as delimiters. In addition, A3 employs two types of filters. First, a content ageing filter that limits the number of revisions to be analyzed by excluding the content from old revisions according to the values of the thresholds  $\Delta_N$  and  $\Delta_T$ ;<sup>28</sup> in our experiments, we used the original configuration of the algorithm ( $\Delta_N = 100$ ,  $\Delta_T = 90$ ). Second is the Wiki markup filter, which we discussed in Subsection 4.3. The Wiki markup affects the amount of content to be processed in each iteration and we thus studied the performance of the algorithm with this filter on (**A3 MF-ON**) and disabled (**A3 MF-OFF**).

Regarding the WikiWho vandalism detection mechanisms (cf. Section 4.4), we empirically set up their thresholds by performing tests on a random article sample. In the experiments, the value for the change percentage filter was equal to  $-0.40$ , and the token density was set to 10. This resulted in 0.5% of revisions being filtered. As discussed in Section 4.3, the definition of tokenization units is an important factor that affects the quality as well as the performance of the algorithm. We studied the performance of WikiWho in two variations of tokenization plus one additional setting:

- **WikiWho complex tokenization (CT)**: We implemented the tokenization described in Section 4.3, considering the Wiki markup. This is the original algorithm we propose.
- **WikiWho simple tokenization (ST)**: Tokens are obtained by splitting the raw content using only whitespaces as delimiters. This setting was used to assess which additional load the complex tokenization adds by generating a much higher number of tokens to track.
- **WikiWho ST and content aging filter on (ST/AF-ON)**: We implemented the content aging filter described for A3, with  $\Delta_N = 100$ . This setting and the A3 MF-OFF configuration allow to compare the algorithms under similar conditions.

### 5.2.3 Results

We executed each setting 5 times and report on the average time per article, revision and Kilobyte. The runtime results for all settings are listed in Table 2. Figure 4 plots the time results in relation to increasing total article history size, meaning average revision size times number of revisions.<sup>29</sup> The figures include the function for a fitted linear regression, showing that for the A3 settings the runtime increases with growing article size by a much larger factor than for the WikiWho variants. Although the WikiWho performance seems to be more volatile, we can observe that the behavior of all the settings in both algorithms is consistent in general and increases in a linear or almost linear fashion with an increasing content size. Fluctuations between data points suggest that the execution time is also influenced by other properties of articles, e.g., the amount of content modified from one revision to another.

<sup>28</sup>  $\Delta_N$  limits the processed content to a maximum of  $N$  most recent revisions, while  $\Delta_T$  further filters out revisions older than  $T$  days.

<sup>29</sup> As both values influence the runtime. Text includes Wiki markup.

Table 2: Execution time of algorithm settings

Algorithm setting	Avg. time per article (secs.)	Ratio of runtime (base: ST)	Avg. time per revision (secs.)	Avg. time per KB (secs.)
<i>Dataset 1</i>				
ST	0.84	1:1	0.0038	$2.84 \times 10^{-4}$
CT	1.04	1:1.24	0.0047	$3.49 \times 10^{-4}$
ST/AF-ON	1.32	1:1.57	0.0061	$4.46 \times 10^{-4}$
A3 MF-OFF	14.30	1:17.02	0.0654	$4.82 \times 10^{-3}$
A3 MF-ON	17.69	1:21.05	0.0809	$5.96 \times 10^{-3}$
<i>Dataset 2</i>				
ST	184.97	1:1	0.0322	$4.01 \times 10^{-4}$
CT	284.44	1:1.54	0.0495	$6.16 \times 10^{-4}$
ST/AF-ON	290.97	1:1.57	0.0506	$6.30 \times 10^{-4}$
A3 MF-OFF	2834.37	1:15.32	0.4931	$6.14 \times 10^{-3}$
A3 MF-ON	2559.38	1:13.84	0.4452	$5.55 \times 10^{-3}$

The runtime decrease by WikiWho in contrast to A3 is in the range of one order of magnitude, differing over the settings. The two most comparable settings ST/AF-ON and A3 MF-OFF differ by a factor of 10.83 and 8.80, respectively for the two datasets, while the originally proposed setting CT completes the task in an even shorter time. It appears that the time filter is in fact no accelerator for the WikiWho algorithm, supposedly because it creates more overhead than is saved by not processing older revisions. The A3 algorithm shows the same behavior in *Dataset 1*. Still, for *Dataset 2* the markup filter seems to take effect, presumably because in longer revisions the amount of filtered content is larger.

Overall, WikiWho is able to execute the given task of computing authorship in a very efficient manner and outperforms the A3 algorithm significantly in runtime in all variants. This is possible due to the construction of paragraph and sentence nodes comparable to creating indexes over the text. This allows to efficiently detect large chunks of text that remained unchanged between revisions, vastly reducing the number of necessary comparisons at the token level.

## 5.3 Evaluation of Materialization Size

Since revisioned content is in constant production – particularly in the English Wikipedia, where over 3 Million revisions are created monthly<sup>30</sup> – it might be useful to materialize partial computation in order to allow incremental data processing, i.e., the algorithm can be stopped at a certain point in time and then resume its execution. Therefore, we implemented a JSON serialization mechanism to optionally materialize partial computation. We measured the overhead caused by the serialization in terms of space.

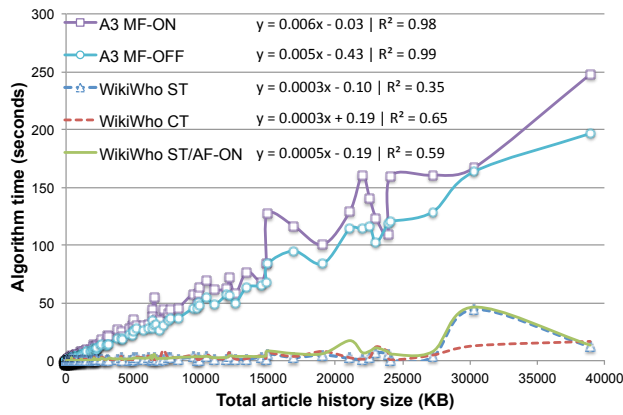
### 5.3.1 Experimental Set-up

We used the articles contained in the two datasets presented in Section 5.2, and serialized the computation of the authorship labels of the whole page history for each article. We compared the serialization mechanism of WikiWho and A3 under similar conditions with the settings ST/AF-ON and MF-OFF, respectively. Both algorithms WikiWho and A3 utilize the `cjson` Python library<sup>31</sup> to implement the (de-)serialization mechanisms. Since we are comparing the algorithms with content aging filter set to  $\Delta_N = 100$ , we report on the size of the JSON serialization in relation to the size of the last  $N = 100$  revisions of each article.

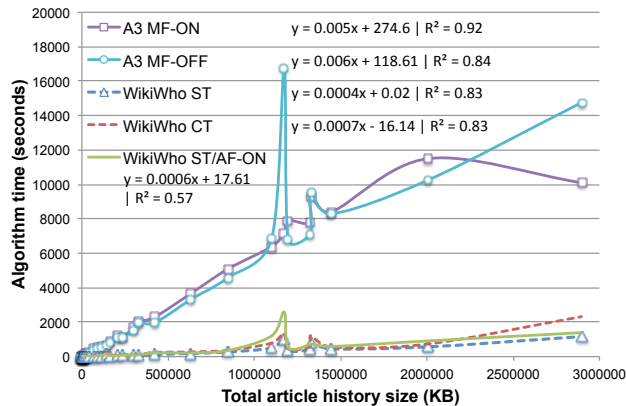
<sup>30</sup> According to the Wikimedia Statistics of June 2013:

<http://stats.wikimedia.org/EN/SummaryEN.htm>

<sup>31</sup> <https://pypi.python.org/pypi/python-cjson>



(a) Performance in *Dataset 1* (Articles randomly selected)



(b) Performance in *Dataset 2* (Articles used in quality evaluation)

Figure 4: Algorithm execution time evaluation for different settings of WikiWho and A3 in *Dataset 1* and *Dataset 2* – the fitted linear functions are denoted by  $y$ , respectively for the data series on the left (fit lines omitted and data points partly omitted for readability)

### 5.3.2 Results

Figure 5 plots the results of the materialization for WikiWho and A3. The behavior of the two algorithms is in general consistent. When the size of the revisioned content increases, the relative size of the serialization decreases exponentially. This suggests that both algorithms efficiently represent redundant content. On average, the size of the serialization is 66% for WikiWho and 56% for the A3 algorithm with respect to the total size of the last 100 revisions. Figure 5 further depicts the percentaged difference of the WikiWho minus the A3 materialization with increasing content size. It shows a volatile behavior with a clear linear trend.

Weighing the cost of storing the results for small articles versus the average time to calculate authorship labels with WikiWho, materializing these results does not bring additional benefits; on the contrary, it incurs on extra space and time. Therefore, the proposed serialization mechanisms should be executed only when the time to compute the authorship labels over the whole article history exceeds a “reasonable” response time, e.g., wait time for end users. Using the “worst case” linear estimation of the originally proposed setting CT for *Dataset 1* (cf. Figure 4), a hypothetical maximum runtime of 5 seconds would allow to process all articles with up to 16,033 KB complete revision history text size without the need for materialization. As far as we can estimate by a random sampling from the Wikipedia database, at least half of all articles in the En-

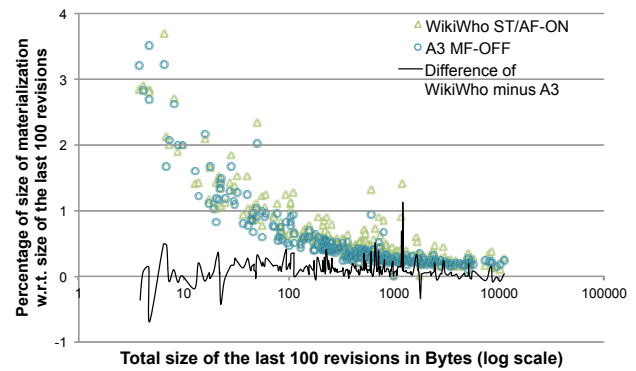


Figure 5: Performance in *Dataset 1* (Wiki pages randomly selected) – article “Rankin County” at  $y$ -axis values 10.69 (WikiWho) and 11.13 (A3) not shown for readability

glish Wikipedia currently stay under this size limit.<sup>32</sup> The needed storage space can of course be further reduced by relaxing the runtime constraint. For articles over this limit, intervals of revisions can be determined when a materialization becomes necessary, although this is beyond the scope of this paper.

## 6. CONCLUSIONS AND FUTURE WORK

In this work we have proposed WikiWho, a solution for the attribution of authorship in revisioned content. We built a graph-based model to represent revisioned content, and provided a formal solution to the authorship problem. In order to measure the quality of WikiWho, we created a gold standard of over 240 tokens from Wikipedia articles, and corroborated it via crowdsourcing. It is, to our expertise, the first gold standard of this kind to measure the precision of authorship attributions on token-level. We compared WikiWho against the state-of-the-art, exceeding it by over 10% on average in precision, and outperforming the baseline execution time by one order of magnitude. Our experimental study confirmed that WikiWho is an effective and efficient solution.

Although in this work we restricted the use of WikiWho to single articles, it is also possible to operate it over several articles in a Wiki, tracking the movement of text between different pages. We used the English Wikipedia as inspiration and testing ground, yet the proposed solution can be understood as a more generally applicable method for revisioned content. We are convinced that many of the assumptions made for our use case also hold true for other Wikis and also for other revisioned content systems.

*Future Work:* We plan to study further techniques to optimize the materialization of intermediate computation. Since each article may show different editing patterns (in terms of size and number of revisions), it is beneficial to adapt the frequency of the serialization routine for each article. We will also look at compression mechanisms that allow to reduce the size of the materialization. Currently, we are establishing an API for WikiWho to allow querying for the authorship of articles from the English and German Wikipedia.<sup>33</sup>

## Acknowledgements

We thank Andriy Rodchenko and Felix Leif Keppmann for their contributions to this research. This work was partially supported by grants from the European Union’s 7th Framework Programme under grant agreement number 257790 (FP7-ICT-2009-5).

<sup>32</sup>[https://wiki.toolserver.org/view/Database\\_access](https://wiki.toolserver.org/view/Database_access)

<sup>33</sup>Will be available at <http://wikiwho.aifb.kit.edu>

## 7. REFERENCES

- [1] T. Adler and L. Alfaro. A content-driven reputation system for the Wikipedia. In *WWW '07: Proceedings of the 16th international conference on World Wide Web*, pages 261–270, 2007.
- [2] T. Adler, K. Chatterjee, L. Alfaro, M. Faella, I. Pye, and V. Raman. Assigning trust to Wikipedia content. In *International Symposium on Wikis*, 2008.
- [3] T. Adler, L. de Alfaro, I. Pye, and V. Raman. Measuring author contributions to the Wikipedia. In *Proceedings of the 4th International Symposium on Wikis, WikiSym '08*, pages 15:1–15:10, New York, NY, USA, 2008. ACM.
- [4] R. Baggen, J. Pedro Correia, K. Schill, and J. Visser. Standardized code quality benchmarking for improving software maintainability. *Software Quality Journal*, 20(2):287–307, 2012.
- [5] R.C. Burns and D.D.E. Long. A linear time, constant space differencing algorithm. In *Performance, Computing, and Communications Conference, 1997. IPCCC 1997., IEEE International*, pages 429–436. IEEE, 1997.
- [6] L. de Alfaro and M. Shavlovsky. Attributing authorship of revisioned content. In *Proceedings of the 22nd international conference on World Wide Web, WWW '13*, pages 343–354, Republic and Canton of Geneva, Switzerland, 2013. International World Wide Web Conferences Steering Committee.
- [7] F. Flöck and A. Rodchenko. Whose article is it anyway?—Detecting authorship distribution in Wikipedia articles over time with WIKIGINI. In *Online proceedings of the Wikipedia Academy 2012*. Wikimedia, 2012.
- [8] A. Forte and A. Bruckman. Why do people write for Wikipedia? Incentives to contribute to open-content publishing. In *Workshop of Sustaining Community: The Role and Design of Incentive Mechanisms in Online Systems. Sanibel Island, FL*, pages 6–9, 2005.
- [9] J. Jones. Patterns of revision in online writing a study of Wikipedia’s featured articles. *Written Communication*, 25(2):262–289, 2008.
- [10] F. L. Keppmann, F. Flöck, A. Adam, E. Simperl, D. Rusu, G. Holz, and A. Metzger. A knowledge diversity dashboard for Wikipedia. In *Proceedings of the ACM WebSci'12*, New York, NY, USA, Juni 2012. ACM.
- [11] A. Kittur, B. Suh, B. A. Pendleton, and E. H. Chi. He says, she says: Conflict and coordination in Wikipedia. In *Proceedings of the SIGCHI conference on Human factors in computing systems, CHI '07*, pages 453–462, New York, NY, USA, 2007. ACM.
- [12] M. Linares-Vasquez, K. Hossen, H. Dang, H. Kagdi, M. Gethers, and D. Poshvanyk. Triaging incoming change requests: Bug or commit history, or code authorship? In *Software Maintenance (ICSM), 2012 28th IEEE International Conference on*, pages 451–460, 2012.
- [13] C. M. Pilato, B. Collins-Sussman, and B. W. Fitzpatrick. *Version control with subversion*. O’Reilly Media, Inc., 2009.
- [14] C. R. Prause. Maintaining fine-grained code metadata regardless of moving, copying and merging. In *Source Code Analysis and Manipulation, 2009. SCAM '09. Ninth IEEE International Working Conference on*, pages 109–118, 2009.
- [15] F. Rahman and P. Devanbu. Ownership, experience and defects: a fine-grained study of authorship. In *33rd International Conference on Software Engineering (ICSE)*, pages 491–500, 2011.
- [16] R. Rivest. The MD5 message-digest algorithm. United States, 1992. RFC Editor, MIT and RSA Data Security, Inc.
- [17] F. B. Viégas, M. Wattenberg, and K. Dave. Studying cooperation and conflict between authors with history flow visualizations. In *Proceedings of the SIGCHI conference on Human factors in computing systems, CHI '04*, pages 575–582, New York, NY, USA, 2004. ACM.