

Computing n -Gram Statistics in MapReduce

Klaus Berberich
Max Planck Institute for Informatics
Saarbrücken, Germany
kberberi@mpi-inf.mpg.de

Srikanta Bedathur
Indraprastha Institute of Information Technology
New Delhi, India
bedathur@iiitd.ac.in

ABSTRACT

Statistics about n -grams (i.e., sequences of contiguous words or other tokens in text documents or other string data) are an important building block in information retrieval and natural language processing. In this work, we study how n -gram statistics, optionally restricted by a maximum n -gram length and minimum collection frequency, can be computed efficiently harnessing MapReduce for distributed data processing. We describe different algorithms, ranging from an extension of word counting, via methods based on the APRIORI principle, to a novel method SUFFIX- σ that relies on sorting and aggregating suffixes. We examine possible extensions of our method to support the notions of maximality/closedness and to perform aggregations beyond occurrence counting. Assuming Hadoop as a concrete MapReduce implementation, we provide insights on an efficient implementation of the methods. Extensive experiments on The New York Times Annotated Corpus and ClueWeb09 expose the relative benefits and trade-offs of the methods.

Categories and Subject Descriptors

H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing

General Terms

Algorithms, Experimentation

Keywords

n -Grams, MapReduce

1. INTRODUCTION

Applications in various fields including information retrieval [12, 46] and natural language processing [13, 18, 39] rely on statistics about n -grams (i.e., sequences of contiguous words in text documents or other string data) as an important building block. Google and Microsoft have made available n -gram statistics computed on parts of the Web.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

EDBT/ICDT'13 March 18 - 22 2013, Genoa, Italy
Copyright 2013 ACM 978-1-4503-1597-5/13/03 ...\$15.00.

While certainly a valuable resource, one limitation of these datasets is that they only consider n -grams consisting of up to five words. With this limitation, there is no way to capture idioms, quotations, poetry, lyrics, and other types of named entities (e.g., products, books, songs, or movies) that typically consist of more than five words and are crucial to applications including plagiarism detection, opinion mining, and social media analytics.

MapReduce has gained popularity in recent years both as a programming model and in its open-source implementation Hadoop. It provides a platform for distributed data processing, for instance, on web-scale document collections. MapReduce imposes a rigid programming model, but treats its users with features such as handling of node failures and an automatic distribution of the computation. To make most effective use of it, problems need to be cast into its programming model, taking into account its particularities.

In this work, we address the problem of efficiently computing n -gram statistics on MapReduce platforms. We allow for a restriction of the n -gram statistics to be computed by a maximum length σ and a minimum collection frequency τ . Only n -grams consisting of up to σ words and occurring at least τ times in the document collection are thus considered.

While this can be seen as a special case of frequent sequence mining, our experiments on two real-world datasets show that MapReduce adaptations of APRIORI-based methods [38, 45] do not perform well – in particular when long and/or less frequent n -grams are of interest. In this light, we develop our novel method SUFFIX- σ that is based on ideas from string processing. Our method makes thoughtful use of MapReduce's grouping and sorting functionality. It keeps the number of records that have to be sorted by MapReduce low and exploits their order to achieve a compact main-memory footprint, when determining collection frequencies of all n -grams considered.

We also describe possible extensions of our method. This includes the notions of maximality/closedness, known from frequent sequence mining, that can drastically reduce the amount of n -gram statistics computed. In addition, we investigate to what extent our method can support aggregations beyond occurrence counting, using n -gram time series, recently made popular by Michel et al. [32], as an example.

Contributions made in this work include:

- a *novel method* SUFFIX- σ to compute n -gram statistics that has been specifically designed for MapReduce;
- a detailed account on *efficient implementation* and *possible extensions* of SUFFIX- σ (e.g., to consider maximal/closed n -grams or support other aggregations);

- a *comprehensive experimental evaluation* on The New York Times Annotated Corpus (1.8 million news articles from 1987–2007) and ClueWeb09-B (50 million web pages crawled in 2009), as two large-scale real-world document collections, comparing our method against state-of-the-art competitors and investigating their trade-offs.

SUFFIX- σ outperforms its best competitor in our experiments by up to a *factor 12x* when long and/or less frequent n -grams are of interest. Otherwise, it is at least on par with its competitors.

Organization. Section 2 introduces our model. Section 3 details on methods to compute n -gram statistics based on prior ideas. Section 4 introduces our method SUFFIX- σ . Aspects of efficient implementation are addressed in Section 5. Possible extensions of SUFFIX- σ are sketched in Section 6. Our experiments are the subject of Section 7. In Section 8, we put our work into context, before concluding in Section 9.

2. PRELIMINARIES

We now introduce our model, establish our notation, and provide some technical background on MapReduce.

2.1 Data Model

Our methods operate on sequences of terms (i.e., words or other textual tokens) drawn from a vocabulary \mathcal{V} . We let \mathcal{S} denote the universe of all sequences over \mathcal{V} . Given a sequence $\mathbf{s} = \langle s_0, \dots, s_{n-1} \rangle$ with $s_i \in \mathcal{V}$, we refer to its length as $|\mathbf{s}|$, write $\mathbf{s}[i..j]$ for the subsequence $\langle s_i, \dots, s_j \rangle$, and let $\mathbf{s}[i]$ refer to the element s_i . For two sequences \mathbf{r} and \mathbf{s} , we let $\mathbf{r}||\mathbf{s}$ denote their concatenation. We say that

- \mathbf{r} is a *prefix* of \mathbf{s} ($\mathbf{r} \triangleright \mathbf{s}$) iff

$$\forall 0 \leq i < |\mathbf{r}| : \mathbf{r}[i] = \mathbf{s}[i]$$

- \mathbf{r} is a *suffix* of \mathbf{s} ($\mathbf{r} \triangleleft \mathbf{s}$) iff

$$\forall 0 \leq i < |\mathbf{r}| : \mathbf{r}[i] = \mathbf{s}[|\mathbf{s}| - |\mathbf{r}| + i]$$

- \mathbf{r} is a *subsequence* of \mathbf{s} ($\mathbf{r} \diamond \mathbf{s}$) iff

$$\exists 0 \leq j < |\mathbf{s}| : \forall 0 \leq i < |\mathbf{r}| : \mathbf{r}[i] = \mathbf{s}[i + j]$$

and capture how often \mathbf{r} occurs in \mathbf{s} as

$$f(\mathbf{r}, \mathbf{s}) = |\{0 \leq j < |\mathbf{s}| \mid \forall 0 \leq i < |\mathbf{r}| : \mathbf{r}[i] = \mathbf{s}[i + j]\}|.$$

To avoid confusion, we use the following convention: When referring to sequences of terms having a specific length k , we will use the notion k -gram or indicate the considered length by alluding to, for instance, 5-grams. The notion n -gram, as found in the title, will be used when referring to variable-length sequences of terms.

As an input, all methods considered in this work receive a document collection \mathcal{D} consisting of sequences of terms as documents. Our focus is on determining how often n -grams occur in the document collection. Formally, the *collection frequency* of an n -gram \mathbf{s} is defined as

$$cf(\mathbf{s}) = \sum_{\mathbf{d} \in \mathcal{D}} f(\mathbf{s}, \mathbf{d}).$$

Alternatively, one could consider the document frequency of n -grams as the total number of documents that contain a specific n -gram. While this corresponds to the notion of *support* typically used in frequent sequence mining, it is

less common for natural language applications. However, all methods presented below can easily be modified to produce document frequencies instead.

2.2 MapReduce

MapReduce, as described by Dean and Ghemawat [17], is a programming model and an associated runtime system at Google. While originally proprietary, the MapReduce programming model has been widely adopted in practice and several implementations exist. In this work, we rely on Hadoop [1] as a popular open-source MapReduce platform. The objective of MapReduce is to facilitate distributed data processing on large-scale clusters of commodity computers. MapReduce enforces a functional style of programming and lets users express their tasks as two functions

```
map()      : (k1,v1) -> list<(k2,v2)>
reduce()   : (k2, list<v2>) -> list<(k3,v3)>
```

that consume and emit key-value pairs. Between the `map`- and `reduce`-phase, the system sorts and groups the key-value pairs emitted by the `map`-function. The partitioning of key-value pairs (i.e., how they are assigned to cluster nodes) and their sort order (i.e., in which order they are seen by the `reduce`-function on each cluster node) can be customized, if needed for the task at hand. For detailed introductions to working with MapReduce and Hadoop, we refer to Lin and Dyer [29] as well as White [41].

3. METHODS BASED ON PRIOR IDEAS

With our notation established, we next describe three methods based on prior ideas to compute n -gram statistics in MapReduce. Before delving into their details, let us state the problem that we address in more formal terms:

Given a document collection \mathcal{D} , a minimum collection frequency τ , a maximum length σ , our objective is to identify all n -grams \mathbf{s} with their collection frequency $cf(\mathbf{s})$, for which $cf(\mathbf{s}) \geq \tau$ and $|\mathbf{s}| \leq \sigma$ hold.

We thus assume that n -grams are only of interest to the task at hand, if they occur at least τ times in the document collection, coined *frequent* in the following, and consist of at most σ terms. Consider, as an example task, the construction of n -gram language models [46], for which one would only look at n -grams up to a specific length and/or resort to back-off models [24] to obtain more robust estimates for n -grams that occur less than specific number of times.

The problem statement above can be seen as a special case of frequent sequence mining that considers only contiguous sequences of single-element itemsets. We believe this to be an important special case that warrants individual attention and allows for an efficient solution in MapReduce, as we show in this work. A more elaborate comparison to existing research on frequent sequence mining is part of Section 8.

To ease our explanations below, we use the following running example, considering a collection of three documents:

```

d1 = <a x b x x>
d2 = <b a x b x>
d3 = <x b a x b>
```

With parameters $\tau = 3$ and $\sigma = 3$, we expect as output

```

<a> : 3   <b> : 5   <x> : 7
<a x> : 3  <x b> : 4
<a x b> : 3
```

from any method, when applied to this document collection.

Algorithm 1: NAÏVE

```

// Mapper
1 map(long did, seq d) begin
2   for b = 0 to |d| - 1 do
3     for e = b to min(b + σ - 1, |d| - 1) do
4       emit(seq d[b..e], long did)

// Reducer
1 reduce(seq s, list<long> l) begin
2   if |l| ≥ τ then
3     emit(seq s, int |l|)

```

3.1 Naïve Counting

One of the example applications of MapReduce, given by Dean and Ghemawat [17] and also used in many tutorials, is word counting, i.e., determining the collection frequency of every word in the document collection. It is straightforward to adapt word counting to consider variable-length n -grams instead of only unigrams and discard those that occur less than τ times. Pseudo code of this method, which we coin NAÏVE, is given in Algorithm 1.

In the `map`-function, the method emits all n -grams of length up to σ for a document together with the document identifier. If an n -gram occurs more than once, it is emitted multiple times. In the `reduce`-phase, the collection frequency of every n -gram is determined and, if it exceeds τ , emitted together with the n -gram itself.

Interestingly, apart from minor optimizations, this is the method that Brants et al. [13] used for training large-scale language models at Google, considering n -grams up to length five. In practice, several tweaks can be applied to improve this simple method including local pre-aggregation in the `map`-phase (e.g., using a combiner in Hadoop). Implementation details of this kind are covered in more detail in Section 5. The potentially vast number of emitted key-value pairs that needs to be transferred and sorted, though, remains a shortcoming.

In the worst case, when $\sigma \geq |\mathbf{d}|$, NAÏVE emits $\mathcal{O}(|\mathbf{d}|^2)$ key-value pairs for a document \mathbf{d} , each consuming $\mathcal{O}(|\mathbf{d}|)$ bytes, so that the method transfers $\mathcal{O}(|\mathbf{d}|^3)$ bytes between the `map`- and `reduce`-phase. Complementary to that, we can determine the number of key-value pairs emitted based on the n -gram statistics. NAÏVE emits a total of $\sum_{\mathbf{s} \in \mathcal{S}: |\mathbf{s}| \leq \sigma} cf(\mathbf{s})$ key-value pairs, each of which consumes $\mathcal{O}(|\mathbf{s}|)$ bytes.

3.2 Apriori-Based Methods

How can one do better than the naïve method just outlined? One idea is to exploit the APRIORI principle, as described by Agrawal et al. [9] in their seminal paper on identifying frequent itemsets and follow-up work on frequent pattern mining [10, 37, 38, 45]. Cast into our setting, the APRIORI principle states that

$$\mathbf{r} \diamond \mathbf{s} \Rightarrow cf(\mathbf{r}) \geq cf(\mathbf{s})$$

holds for any two sequences \mathbf{r} and \mathbf{s} , i.e., the collection frequency of a sequence \mathbf{r} is an upper bound for the collection frequency of any supersequence \mathbf{s} . In what follows, we describe two methods that make use of the APRIORI principle to compute n -gram statistics in MapReduce.

APRIORI-SCAN

The first APRIORI-based method APRIORI-SCAN, like the original APRIORI algorithm [9] and GSP [38], performs mul-

Algorithm 2: APRIORI-SCAN

```

int k = 1
repeat
  hashset<int[]> dict = load(output-(k - 1))

  // Mapper
  1 map(long did, seq d) begin
  2   for b = 0 to |d| - k do
  3     if k = 1 ∨
  4       (contains(dict, d[b..(b + k - 2)]) ∧
  5        contains(dict, d[(b + 1)..(b + k - 1)])) then
  6       emit(seq d[b..(b + k - 1)], long did)

  // Reducer
  1 reduce(seq s, list<long> l) begin
  2   if |l| ≥ τ then
  3     emit(seq s, int |l|)

  k += 1
until isEmpty(output-(k - 1)) ∨ k = σ + 1;

```

iple scans over the input data. During the k -th scan the method determines k -grams that occur at least τ times in the document collection. To this end, it exploits the output from the previous scan via the APRIORI principle to prune the considered k -grams. In the k -th scan, only those k -grams are considered whose two constituent $(k - 1)$ -grams are known to be frequent. Unlike GSP, which first generates all potentially frequent sequences as candidates, APRIORI-SCAN considers only sequences that actually occur in the document collection. The method terminates after σ scans or when a scan does not produce any output.

Algorithm 2 shows how the method can be implemented in MapReduce. The outer `repeat`-loop controls the execution of multiple MapReduce jobs, each of which performs one distributed parallel scan over the input data. In the k -th iteration, and thus the k -th scan of the input data, the method considers all k -grams from an input document in the `map`-function, but discards those that have a constituent $(k - 1)$ -gram that is known to be infrequent. This pruning is done leveraging the output from the previous iteration that is kept in a dictionary. In the `reduce`-function, analogous to NAÏVE, collection frequencies of k -grams are determined and output if above the minimum collection frequency τ . After σ iterations or once an iteration does not produce any output, the method terminates, which is safe since the APRIORI principle guarantees that no longer n -gram can occur τ or more times in the document collection.

When applied to our running example, in its third scan of the input data, APRIORI-SCAN emits in the `map`-phase for every document \mathbf{d}_i only the key-value pair $(\langle \mathbf{a} \mathbf{x} \mathbf{b} \rangle, \mathbf{d}_i)$, but discards other trigrams (e.g., $\langle \mathbf{b} \mathbf{x} \mathbf{x} \rangle$) that contain an infrequent bigram (e.g., $\langle \mathbf{x} \mathbf{x} \rangle$).

When implemented in MapReduce, every iteration corresponds to a separate job that needs to be run and comes with its administrative fix cost (e.g., for launching and finalizing the job). Another challenge in APRIORI-SCAN is the implementation of the dictionary that makes the output from the previous iteration available and accessible to cluster nodes. This dictionary can either be implemented locally, so that every cluster node receives a replica of the previous iteration's output (e.g., implemented using the distributed cache in Hadoop), or by loading the output from the previous iteration into a shared dictionary (e.g., implemented using a distributed key-value store), which can then be accessed remotely by cluster nodes. Either way, to make

lookups in the dictionary efficient, significant main memory at cluster nodes is required.

An apparent shortcoming of APRIORI-SCAN is that it has to scan the entire input data in every iteration. Thus, although typically only few frequent n -grams are found in later iterations, the cost of an iteration depends on the size of the input data. The number of iterations needed, on the other hand, is determined by the parameter σ or the length of the longest frequent n -gram.

In the worst case, when $\sigma \geq |\mathbf{d}|$ and $cf(\mathbf{d}) \geq \tau$, APRIORI-SCAN emits $\mathcal{O}(|\mathbf{d}|^2)$ key-value pairs per document \mathbf{d} , each consuming $\mathcal{O}(|\mathbf{d}|)$ bytes, so that the method transfers $\mathcal{O}(|\mathbf{d}|^3)$ bytes between the `map`- and `reduce`-phase. Again, we provide a complementary analysis based on the actual n -gram statistics. To this end, let

$$\mathcal{S}_{NP} = \{s \in \mathcal{S} \mid \forall r \in \mathcal{S} : (r \neq s \wedge r \diamond s) \Rightarrow cf(r) \geq \tau\}$$

denote the set of sequences that cannot be pruned based on the APRIORI principle, i.e., whose true subsequences all occur at least τ times in the document collection. APRIORI-SCAN emits a total of $\sum_{s \in \mathcal{S}_{NP}: |s| \leq \sigma} cf(s)$ key-value pairs, each of which amounts to $\mathcal{O}(|s|)$ bytes. Obviously, $\mathcal{S}_{NP} \subseteq \mathcal{S}$ holds, so that APRIORI-SCAN emits at most as many key-value pairs as NAIVE. Its concrete gains, though, depend on the value of τ and characteristics of the document collection.

APRIORI-INDEX

The second APRIORI-based method APRIORI-INDEX does not repeatedly scan the input data but incrementally builds an inverted index of frequent n -grams from the input data as a more compact representation. Operating on an index structure as opposed to the original data and considering n -grams of increasing length, it resembles SPADE [45] when breadth-first traversing the sequence lattice.

Pseudo code of APRIORI-INDEX is given in Algorithm 3. In its first phase, the method constructs an inverted index with positional information for all frequent n -grams up to length K (cf. `Mapper #1` and `Reducer #1` in the pseudo code). In its second phase, to identify frequent n -grams beyond that length, APRIORI-INDEX harnesses the output from the previous iteration. Thus, to determine a frequent k -gram (e.g., $\langle \mathbf{b a x} \rangle$), the method joins the posting lists of its constituent $(k-1)$ -grams (i.e., $\langle \mathbf{b a} \rangle$ and $\langle \mathbf{a x} \rangle$). In MapReduce, this can be accomplished as follows (cf. `Mapper #2` and `Reducer #2` in the pseudo code): The `map`-function emits for every frequent $(k-1)$ -gram two key-value pairs. The frequent $(k-1)$ -gram itself along with its posting list serves in both as a value. As keys the prefix and suffix of length $(k-2)$ are used. In the pseudo code, the method keeps track of whether the key is a prefix or suffix of the sequence in the value by using the `r-seq` and `l-seq` subtypes. The `reduce`-function identifies for a specific key all compatible sequences from the values, joins their posting lists, and emits the resulting k -gram along with its posting list if its collection frequency is at least τ . Two sequences are compatible and must be joined, if one has the current key as a prefix, and the other has it as a suffix. In its nested `for`-loops, the method considers all compatible combinations of sequences. This second phase of APRIORI-INDEX can be seen as a distributed candidate generation and pruning step.

Applied to our running example and assuming $K = 2$, the method only sees one pair of compatible sequences with their posting lists for the key $\langle \mathbf{x} \rangle$ in its third iteration, namely:

Algorithm 3: APRIORI-Index

```

int k = 1
repeat
  if k ≤ K then
    // Mapper #1
    map(long did, seq d) begin
      1  hashmap<seq, int[]> pos = {}
      2  for b = 0 to |d| - 1 do
      3    | add(get(pos, d[b..(b+k-1)]), b)
      4    | for seq s : keys(pos) do
      5    |   emit(seq s, posting(did, get(pos,s)))
      6
    // Reducer #1
    1  reduce(seq s, list<posting> l) begin
    2  | if cf(l) ≥ τ then
    3  |   emit(seq s, list<posting> l)
  else
    // Mapper #2
    1  map(seq s, list<posting> l) begin
    2  | emit(seq s[0..|s| - 2],
    3  |   (r-seq, list<posting>)(s, l))
    4  | emit(seq s[1..|s| - 1],
    5  |   (l-seq, list<posting>)(s, l))
    // Reducer #2
    1  reduce(seq s, list<(seq, list<posting>)> l) begin
    2  | for (l-seq, list<posting>)(m, lm) : l do
    3  |   for (r-seq, list<posting>)(n, ln) : l do
    4  |     list<posting> lj = join(lm, ln)
    5  |     if cf(lj) ≥ τ then
    6  |       seq j = m || ⟨n[|n| - 1]⟩
    7  |       emit(seq j, list<posting> lj)
    k += 1
until isEmpty(output-(k-1)) ∨ k = min(σ, K);

```

$$\begin{aligned} \langle \mathbf{a x} \rangle & : \langle \mathbf{d}_1 : [0], \mathbf{d}_2 : [1], \mathbf{d}_3 : [2] \rangle \\ \langle \mathbf{x b} \rangle & : \langle \mathbf{d}_1 : [1], \mathbf{d}_2 : [2], \mathbf{d}_3 : [0, 3] \rangle . \end{aligned}$$

By joining those, APRIORI-INDEX obtains the only frequent 3-gram with its posting list

$$\langle \mathbf{a x b} \rangle : \langle \mathbf{d}_1 : [0], \mathbf{d}_2 : [1], \mathbf{d}_3 : [2] \rangle .$$

For all $k < K$, it would be enough to determine only collection frequencies, as opposed to, positional information of n -grams. While a straightforward optimization in practice, we opted for simpler pseudo code. When implemented as described in Algorithm 3, the method produces an inverted index with positional information, which can be used to quickly determine the locations of a specific frequent n -gram.

One challenge when implementing APRIORI-INDEX is that the number and size of posting-list values seen for a specific key can become large in practice. Moreover, to join compatible sequences, these posting lists have to be buffered, and a scalable implementation must deal with the case when this is not possible in the available main memory. This can, for instance, be accomplished by storing posting lists temporarily in a disk-resident key-value store.

The number of iterations needed by APRIORI-INDEX is determined by the parameter σ or the length of the longest frequent n -gram. Since every iteration, as for APRIORI-SCAN, corresponds to a separate MapReduce job, a non-negligible administrative fix cost is incurred.

In the worst case, when $\sigma \geq |\mathbf{d}|$ and $cf(\mathbf{d}) \geq \tau$, APRIORI-INDEX emits $\mathcal{O}(|\mathbf{d}|^2)$ key-value pairs per document \mathbf{d} , each consuming $\mathcal{O}(|\mathbf{d}|)$ bytes, so that $\mathcal{O}(|\mathbf{d}|^3)$ bytes are transferred between the `map`- and `reduce`-phase. We assume $K <$

σ for the complementary analysis. In its first K iterations, APRIORI-INDEX emits $\sum_{\mathbf{s} \in \mathcal{S}: |\mathbf{s}| \leq K} df(\mathbf{s})$ key-value pairs, where $df(\mathbf{s}) \leq cf(\mathbf{s})$ refers to the document frequency of the n -gram \mathbf{s} , as mentioned in Section 2. Each key-value pair consumes $\mathcal{O}(cf(\mathbf{s}))$ bytes. To analyze the following iterations, let

$$\mathcal{S}_F = \{\mathbf{s} \in \mathcal{S} \mid cf(\mathbf{s}) \geq \tau\}$$

denote the set of frequent n -grams that occur at least τ times. APRIORI-INDEX emits a total of

$$2 \cdot |\{\mathbf{s} \in \mathcal{S}_F \mid K \leq |\mathbf{s}| < \sigma\}|$$

key-value pairs, each of which consumes $\mathcal{O}(cf(\mathbf{s}))$ bytes. Like for APRIORI-SCAN, the concrete gains depend on the value of τ and characteristics of the document collection.

4. SUFFIX SORTING & AGGREGATION

As already argued, the methods presented so far suffer from either excessive amounts of data that need to be transferred and sorted, requiring possibly many MapReduce jobs, or a high demand for main memory at cluster nodes. Our novel method SUFFIX- σ avoids these deficiencies: It requires a single MapReduce job, transfers only a modest amount of data, and requires little main memory at cluster nodes.

Consider again what the `map`-function in the NAIVE approach emits for document \mathbf{d}_3 from our running example. Emitting key-value pairs for all of the n -grams $\langle \mathbf{b a x} \rangle$, $\langle \mathbf{b a} \rangle$, and $\langle \mathbf{b} \rangle$ is clearly wasteful. The key observation here is that the latter two are subsumed by the first one and can be obtained as its prefixes. Suffix arrays [31] and other string processing techniques exploit this very idea.

Based on this observation, it is safe to emit key-value pairs only for a subset of the n -grams contained in a document. More precisely, it is enough to emit at every position in the document a single key-value pair with the suffix starting at that position as a key. These suffixes can further be truncated to length σ – hence the name of our method.

To determine the collection frequency of a specific n -gram \mathbf{r} , we have to determine how many of the suffixes emitted in the `map`-phase are prefixed by \mathbf{r} . To do so correctly using only a single MapReduce job, we must ensure that all relevant suffixes are seen by the same reducer. This can be accomplished by partitioning suffixes based on their first term only, as opposed to, all terms therein. It is thus guaranteed that a single reducer receives all suffixes that begin with the same term. This reducer is then responsible for determining the collection frequencies of all n -grams starting with that term. One way to accomplish this would be to enumerate all prefixes of a received suffix and aggregate their collection frequencies in main memory (e.g., using a hashmap or a prefix tree). Since it is unknown whether an n -gram is represented by other yet unseen suffixes from the input, it cannot be emitted early along with its collection frequency. Bookkeeping is thus needed for many n -grams and requires significant main memory.

How can we reduce the main-memory footprint and emit n -grams with their collection frequency early on? The key idea is to exploit that the order in which key-value pairs are sorted and received by reducers can be influenced. SUFFIX- σ sorts key-value pairs in reverse lexicographic order of their

Algorithm 4: SUFFIX- σ

```

// Mapper
1 map(long did, seq d) begin
2   for b = 0 to |d| - 1 do
3     emit(seq d[b..min(b + sigma - 1, |d| - 1)], long did)

// Reducer
stack<int> terms = empty
stack<int> counts = empty
1 reduce(seq s, list<long> l) begin
2   while lcp(s, seq(terms)) < len(terms) do
3     if peek(counts) >= tau then
4       emit(seq seq(terms), int peek(counts))
5     pop(terms)
6     push(counts, pop(counts) + pop(counts))
7   if len(terms) = |s| then
8     push(counts, pop(counts) + |l|)
   else
10    for i = lcp(s, seq(terms)) to |s| - 1 do
11      push(terms, s[i])
12      push(counts, (i == |s| - 1 ? |l| : 0))

1 cleanup() begin
2   reduce(seq empty, list<long> empty)

// Partitioner
1 partition(seq s) begin
2   return hashCode(s[0]) mod R

// Comparator
1 compare(seq r, seq s) begin
2   for b = 0 to min(|r|, |s|) - 1 do
3     if r[b] < s[b] then
4       return +1 // r > s
5     else if r[b] > s[b] then
6       return -1 // r < s
7   return |s| - |r|

```

suffix key, formally defined as follows for sequences \mathbf{r} and \mathbf{s} :

$$\mathbf{r} < \mathbf{s} \Leftrightarrow (|\mathbf{r}| > |\mathbf{s}| \wedge \mathbf{s} \triangleright \mathbf{r}) \vee \\ \exists 0 \leq i < \min(|\mathbf{r}|, |\mathbf{s}|) : \mathbf{r}[i] > \mathbf{s}[i] \wedge \forall 0 \leq j < i : \mathbf{r}[j] = \mathbf{s}[j].$$

To see why this is useful, recall that each suffix from the input represents all n -grams that can be obtained as its prefixes. Let \mathbf{s} denote the current suffix from the input. The reverse lexicographic order guarantees that we can safely emit any n -gram \mathbf{r} such that $\mathbf{r} < \mathbf{s}$, since no yet unseen suffix from the input can represent \mathbf{r} . Conversely, at this point, the only n -grams for which we have to do bookkeeping, since they are represented both by the current suffix \mathbf{s} and potentially by yet unseen suffixes, are the prefixes of \mathbf{s} . We illustrate this observation with our running example. The reducer responsible for suffixes starting with \mathbf{b} receives:

$$\begin{aligned} \langle \mathbf{b x x} \rangle & : \langle \mathbf{d}_1 \rangle \\ \langle \mathbf{b x} \rangle & : \langle \mathbf{d}_2 \rangle \\ \langle \mathbf{b a x} \rangle & : \langle \mathbf{d}_2, \mathbf{d}_3 \rangle \\ \langle \mathbf{b} \rangle & : \langle \mathbf{d}_3 \rangle. \end{aligned}$$

When seeing the third suffix $\langle \mathbf{b a x} \rangle$, we can immediately finalize the collection frequency of the n -gram $\langle \mathbf{b x} \rangle$ and emit it, since no yet unseen suffix can have it as a prefix. On the contrary, the n -grams $\langle \mathbf{b} \rangle$ and $\langle \mathbf{b a} \rangle$ cannot be emitted, since yet unseen suffixes from the input may have them as a prefix.

Building on this observation, we can do efficient bookkeeping for prefixes of the current suffix \mathbf{s} only and lazily aggregate their collection frequencies using two stacks. On the first stack *terms*, we keep the terms constituting \mathbf{s} . The sec-

ond stack *counts* keeps one counter per prefix of *s*. Between invocations of the **reduce**-function, we maintain two invariants. First, the two stacks have the same size *m*. Second,

$$\sum_{j=i}^{m-1} counts[j]$$

reflects how often the *n*-gram

$$\langle terms[0], \dots, terms[i] \rangle$$

has been seen so far in the input. To maintain these invariants, when processing a suffix *s* from the input, we first synchronously pop elements from both stacks until the contents of *terms* form a prefix of *s*. Before each pop operation, we emit the contents of *terms* and the top element of *counts*, if the latter is above our minimum collection frequency τ . When popping an element from *counts*, its value is added to the new top element. Following that, we update *terms*, so that its contents equal the suffix *s*. For all but the last term added, a zero is put on *counts*. For the last term, we put the frequency of *s*, reflected by the length of its associated document-identifier list value, on *counts*. Figure 1 illustrates how the states of the two stacks evolve, as the above example input is processed.

x 1		x 2		
x 0	x 2	a 0		
b 0	b 0	b 2	b 4	- -

Figure 1: Suffix- σ 's bookkeeping illustrated

Pseudo code of SUFFIX- σ is given in Algorithm 4. The **map**-function emits for every document all its suffixes truncated to length σ if possible. The **reduce**-function reads suffixes in reverse lexicographic order and performs the bookkeeping using two separate stacks for *n*-grams (*terms*) and their collection frequencies (*counts*), as described above. The function **seq()** returns the *n*-gram corresponding to the entire *terms* stack. The function **lcp()** returns the length of the longest common prefix that two *n*-grams share. In addition, Algorithm 4 contains a **partition**-function ensuring that suffixes are assigned to one of *R* reducers solely based on their first term, as well as, a **compare**-function that ensures the reverse lexicographic order of input suffixes in the **map**-phase. When implemented in Hadoop, these two functions would materialize as a custom partitioner class and a custom comparator class. Finally, **cleanup()** is a method invoked once, when all input has been seen.

SUFFIX- σ emits $\mathcal{O}(|d|)$ key-value pairs per document *d*. Each of these key-value pairs consumes $\mathcal{O}(|d|)$ bytes in the worst case when $\sigma \geq |d|$. The method thus transfers $\mathcal{O}(|d|^2)$ bytes between the **map**- and **reduce**-phase. For every term occurrence in the document collection, SUFFIX- σ emits exactly one key-value pair, so that in total $\sum_{s \in S: |s|=1} cf(s)$ key-value pairs are emitted, each consuming $\mathcal{O}(\sigma)$ bytes.

5. EFFICIENT IMPLEMENTATION

Having described the different methods at a conceptual level, we now provide details on aspects of their implementation, which we found to have a significant impact on performance in practice:

Document Splits. Collection frequencies of individual terms (i.e., unigrams) can be exploited to drastically reduce

required work by splitting up every document at infrequent terms that it contains. Thus, assuming that *z* is an infrequent term given the current value of τ , we can split up a document like $\langle c b a z b a c \rangle$ into the two shorter sequences $\langle c b a \rangle$ and $\langle b a c \rangle$. Again, this is safe due to the APRIORI principle, since no frequent *n*-gram can contain *z*. All methods profit from this – for large values of σ in particular.

Sequence Encoding. It is inefficient to operate on documents in a textual representation. As a one-time preprocessing, we therefore convert our document collections, so that they are represented as a dictionary, mapping terms to term identifiers, and one integer term-identifier sequence for every document. We assign identifiers to terms in descending order of their collection frequency to optimize compression. From there on, our implementation internally only deals with arrays of integers. Whenever serialized for transmission or storage, these are compactly represented using variable-byte encoding [42]. This also speeds up sorting, since *n*-grams can now be compared using integer operations as opposed to operations on strings, thus requiring generally fewer machine instructions. Compact sequence encoding benefits all methods – in particular APRIORI-SCAN with its repeated scans of the document collection.

Key-Value Store. For APRIORI-SCAN and APRIORI-INDEX, reducers potentially buffer a lot of data, namely, the dictionary of frequent $(k-1)$ -grams or the set of posting lists to be joined. Our implementation keeps this data in main memory as long as possible. Otherwise, it migrates the data into a disk-resident key-value store (Berkeley DB Java Edition [3]). Most main memory is then used for caching, which helps APRIORI-SCAN in particular, since lookups of frequent $(k-1)$ -grams typically hit the cache.

Hadoop-Specific Optimizations that we use in our implementation include local aggregation (cf. Mapper #1 in Algorithm 3), Hadoop's distributed cache facility, raw comparators to avoid deserialization and object instantiation, as well as other best practices (e.g., described in [41]).

How easy to implement are the methods presented in previous sections? While hard to evaluate systematically, we still want to address this question based on our own experience. NAÏVE is the clear winner here. Implementations of the APRIORI-based methods, as explained in Section 3, require various tweaks (e.g., the use of a key-value store) to make them work. SUFFIX- σ does not require any of those and, when Hadoop is used as a MapReduce implementation, can be implemented using only on-board functionality. Our code is available – details provided at the end of this paper.

6. EXTENSIONS

In this section, we describe how SUFFIX- σ can be extended to consider only maximal/closed *n*-grams and thus produce a more compact result. Moreover, we explain how it can support aggregations beyond occurrence counting, using *n*-gram time series, recently made popular by [32], as an example.

6.1 Maximality & Closedness

The number of *n*-grams that occur at least τ times in the document collection can be huge in practice. To reduce it, we can adopt the notions of maximality and closedness common in frequent pattern mining. Formally, an *n*-gram *r* is *maximal*, if there is no *n*-gram *s* such that *r* \diamond *s* and $cf(s) \geq \tau$. Similarly, an *n*-gram *r* is *closed*, if no *n*-gram *s* exists such that *r* \diamond *s* and $cf(r) = cf(s) \geq \tau$. The sets of maximal

or closed n -grams are subsets of all n -grams that occur at least τ times. Omitted n -grams can be reconstructed – for closedness even with their accurate collection frequency.

SUFFIX- σ can be extended to produce maximal or closed n -grams. Recall that our method processes suffixes in reverse lexicographic order in its `reduce`-function. Let \mathbf{r} denote the last n -gram emitted. For maximality, we only emit the next n -gram \mathbf{s} , if it is no prefix of \mathbf{r} (i.e., $\neg(\mathbf{s} \triangleright \mathbf{r})$). For closedness, we only emit \mathbf{s} , if it is no prefix of \mathbf{r} or if it has a different collection frequency (i.e., $\neg(\mathbf{s} \triangleright \mathbf{r} \wedge cf(\mathbf{s}) = cf(\mathbf{r}))$). In our example, the reducer responsible for term \mathbf{a} receives

$$\langle \mathbf{a} \times \mathbf{b} \rangle : \langle \mathbf{d}_1, \mathbf{d}_2, \mathbf{d}_3 \rangle$$

and, both for maximality and closedness, emits only the n -gram $\langle \mathbf{a} \times \mathbf{b} \rangle$ but none of its prefixes. With this extension, we thus emit only prefix-maximal or prefix-closed n -grams, whose formal definitions are analogous to those of maximality and closedness above but replace \diamond by \triangleright . In our example, we still emit $\langle \mathbf{x} \mathbf{b} \rangle$ and $\langle \mathbf{b} \rangle$ on the reducers responsible for terms \mathbf{x} and \mathbf{b} , respectively. For maximality, as subsequences of $\langle \mathbf{a} \times \mathbf{b} \rangle$, these n -grams must be omitted. We achieve this by means of an additional post-filtering MapReduce job. As input, the job consumes the output produced by SUFFIX- σ with the above extensions. In its `map`-function, n -grams are reversed (e.g., $\langle \mathbf{a} \times \mathbf{b} \rangle$ becomes $\langle \mathbf{b} \times \mathbf{a} \rangle$). These reversed n -grams are partitioned based on their first term and sorted in reverse lexicographic order, reusing ideas from SUFFIX- σ . In the `reduce`-function, we apply the same filtering as described above to keep only prefix-maximal or prefix-closed reversed n -grams. Before emitting a reversed n -gram, we restore its original order by reversing it. In our example, the reducer responsible for \mathbf{b} receives

$$\begin{aligned} \langle \mathbf{b} \times \mathbf{a} \rangle &: 3 \\ \langle \mathbf{b} \times \mathbf{x} \rangle &: 4 \\ \langle \mathbf{b} \rangle &: 5 \end{aligned}$$

and, for maximality, only emits $\langle \mathbf{a} \times \mathbf{b} \rangle$. In summary, we obtain maximal or closed n -grams by first determining prefix-maximal or prefix-closed n -grams and, after that, identifying the suffix-maximal or suffix-closed among them.

6.2 Beyond Occurrence Counting

Our focus so far has been on determining collection frequencies of n -grams, i.e., counting their occurrences in the document collection. One can move beyond occurrence counting and aggregate other information about n -grams, e.g.:

- build an *inverted index* that records for every n -gram how often or where it occurs in individual documents;
- compute *statistics based on meta-data of documents* (e.g., timestamp or location) that contain a n -gram.

In the following, we concentrate on the second type of aggregation and, as a concrete instance, consider the computation of n -gram time series. Here, the objective is to determine for every n -gram a time series whose observations reveal how often the n -gram occurs in documents published, e.g., in a specific year. SUFFIX- σ can be extended to produce such n -gram time series as follows: In the `map`-function we emit every suffix along with the document identifier and its associated timestamp. In the `reduce`-function, the *counts* stack is replaced by a stack of time series, which we aggregate lazily. When popping an element from the stack, instead

of adding counts, we add time series observations. In the same manner, we can compute other statistics based on the occurrences of an n -gram in documents and their associated meta-data. While these could also be computed by an extension of NAÏVE, the benefit of using SUFFIX- σ is that the required document meta-data is transferred only per suffix of a document, as opposed to, per contained n -gram.

7. EXPERIMENTAL EVALUATION

We conducted comprehensive experiments to compare the different methods and understand their relative benefits and trade-offs. Our findings from these experiments are the subject of this section.

7.1 Setup & Implementation

Cluster Setup. All experiments were run on a local cluster consisting of ten Dell R410 server-class computers, each equipped with 64 GB of main memory, two Intel Xeon X5650 6-core CPUs, and four internal 2 TB SAS 7,200 rpm hard disks configured as a bunch-of-disks. Debian GNU/Linux 5.0.9 (Lenny) was used as an operating system. Machines in the cluster are connected via 1 GBit Ethernet. We use Cloudera CDH3u0 as a distribution of Hadoop 0.20.2 running on Oracle Java 1.6.0_26. One of the machines acts a master and runs Hadoop’s namenode and jobtracker; the other nine machines are configured to run up to ten map tasks and ten reduce tasks in parallel. To restrict the number of map/reduce slots, we employ a capacity-constrained scheduler pool in Hadoop. When we state that n map/reduce slots are used, our cluster executes up to n map tasks and n reduce tasks in parallel. Java virtual machines to process tasks are always launched with 4 GB heap space.

Implementation. All methods are implemented in Java (JDK 1.6) applying the optimizations described in Section 5 to the extent possible and sensible for each of them.

Methods. We compare the methods NAÏVE, APRIORI-SCAN, APRIORI-INDEX, and SUFFIX- σ in our experiments. For APRIORI-INDEX, we set $K = 4$, so that the method directly computes collection frequencies of n -grams having length four or less. We found this to be the best-performing parameter setting in a series of calibration experiments.

Measures. For our experiments, in the following, we report as performance measures:

- wallclock time* as the total time elapsed between launching a method and receiving the final result (possibly involving multiple Hadoop jobs),
- bytes transferred* as the total amount of data transferred between `map`- and `reduce`-phase(s) (obtained from Hadoop’s `MAP_OUTPUT_BYTES` counter),
- # records* as the total number of key-value pairs transferred and sorted between `map`- and `reduce`-phase(s) (obtained from Hadoop’s `MAP_OUTPUT_RECORDS` counter).

For APRIORI-SCAN and APRIORI-INDEX, measures (b) and (c) are aggregates over all Hadoop jobs launched. All measurements reported are based on single runs and were performed with exclusive access to the Hadoop cluster, i.e., without concurrent activity by other jobs, services, or users.

7.2 Datasets

We use two publicly-available real-world datasets for our experiments, namely:

	NYT	C09
# documents	1, 830, 592	50, 221, 915
# term occurrence	1, 049, 440, 645	21, 404, 321, 682
# distinct terms	345, 827	979, 935
# sentences	55, 362, 552	1, 257, 357, 167
sentence length (mean)	18.96	17.02
sentence length (stddev)	14.05	17.56
dataset size raw (GB)	3.23	246.84
dataset size encoded (GB)	2.07	41.50

Table 1: Dataset characteristics

- **The New York Times Annotated Corpus** [7] consisting of more than 1.8 million newspaper articles from the period 1987–2007 (NYT);
- **ClueWeb09-B** [6], as a well-defined subset of the ClueWeb09 corpus of web documents, consisting of more than 50 million web documents in English language that were crawled in 2009 (CW).

These two are extremes: NYT is a well-curated, relatively clean, longitudinal corpus, i.e., documents therein have a clear structure, use proper language with few typos, and cover a long time period. CW is a “World Wild Web” corpus, i.e., documents therein are highly heterogeneous in structure, content, and language.

For NYT a document consists of the newspaper article’s title and body. To make CW more handleable, we use boilerplate detection as described by Kohlschütter et al. [25] and implemented in boilerpipe’s [4] default extractor, to identify the core content of documents. On both datasets, we use OpenNLP [2] to detect sentence boundaries in documents. Sentence boundaries act as barriers, i.e., we do not consider n -grams that span across sentences in our experiments. As described in Section 5, in a pre-processing step, we convert both datasets into sequences of integer term-identifiers. The term dictionary is kept as a single text file; documents are spread as key-value pairs of 64-bit document identifier and content integer array over a total of 256 binary files. Table 1 reports sizes of the two datasets (before and after boilerplate removal and encoding) and summarizes their characteristics.

7.3 Output Characteristics

Let us first look at the n -gram statistics that (or, parts of which) we expect as output from all methods. To this end, for both document collections, we determine all n -grams that occur at least five times (i.e., $\tau = 5$ and $\sigma = \infty$). We bin n -grams into 2-dimensional buckets of exponential width, i.e., the n -gram \mathbf{s} with collection frequency $cf(\mathbf{s})$ goes into bucket (i, j) where $i = \lfloor \log_{10} |\mathbf{s}| \rfloor$ and $j = \lfloor \log_{10} cf(\mathbf{s}) \rfloor$. Figure 2 reports the number of n -grams per bucket.

The figure reveals that the distribution is biased toward short and less frequent n -grams. Consequently, as we lower the value of τ , all methods have to deal with a drastically increasing number of n -grams. What can also be seen from Figure 2 is that, in both datasets, n -grams exist that are very long, containing hundred or more terms, and occur more than ten times in the document collection. Examples of long n -grams that we see in the output include ingredient lists of recipes (e.g., `. . . 1 tablespoon cooking oil. . .`) and chess openings (e.g., `e4 e5 2 nf3. . .`) in NYT; in CW they include web spam (e.g., `travel tips san miguel tourism san`

`miguel transport san miguel. . .`) as well as error messages and stack traces from web servers and other software (e.g., `. . . php on line 91 warning. . .`) that also occur within user discussions in forums. For the APRIORI-based methods, such long n -grams are unfavorable, since they require many iterations to identify them.

7.4 Use Cases

As a first experiment, we investigate how the methods perform for parameter settings chosen to reflect two typical use cases, namely, *training a language model* and *text analytics*. For the first use case, we set $\tau = 10$ on NYT and $\tau = 100$ on CW, as relatively low minimum collection frequencies, in combination with $\sigma = 5$. The n -gram statistics made public by Google [5], as a comparison, were computed with parameter settings $\tau = 40$ and $\sigma = 5$ on parts of the Web. For the second use case, we choose $\sigma = 100$, as a relatively high maximum sequence length, combined with $\tau = 100$ on NYT and $\tau = 1,000$ on CW. The idea in the analytics use case is to identify recurring fragments of text (e.g., quotations or idioms) to be analyzed further (e.g., their spread over time).

Figure 3 reports wallclock-time measurements obtained for these two use cases with 64 map/reduce slots. For our language-model use case, SUFFIX- σ outperforms APRIORI-SCAN as the best competitor by a *factor 3x* on both datasets. For our analytics use case, we see a *factor 12x* improvement over APRIORI-INDEX as the best competitor on NYT; on CW SUFFIX- σ still outperforms the next best APRIORI-SCAN by a *factor 1.5x*. Measurements for NAÏVE on CW in are missing, since the method did not complete in reasonable time.

7.5 Varying Minimum Collection Frequency

Our second experiment studies how the methods behave as we vary the minimum collection frequency τ . We use a maximum length $\sigma = 5$ and apply all methods to the entire datasets. Measurements are performed using 64 map/reduce slots and reported in Figure 4.

We observe that for high minimum collection frequencies, SUFFIX- σ performs as well as the best competitor APRIORI-SCAN. For low minimum collection frequencies, it significantly outperforms the other methods. Both APRIORI-based method show steep increases in wallclock time as we lower the minimum collection frequency – especially when we reach the lowest value of τ on each document collection. This is natural, because for both methods the work that has to be done in the k -th iteration depends on the number of $(k - 1)$ -grams output in the previous iteration, which have to be joined or kept in a dictionary, as described in Section 3. As observed in Figure 2 above, the number of k -grams grows drastically as we decrease the value of τ . When looking at the number of bytes and the number of records transferred, we see analogous behavior. For low values of τ , SUFFIX- σ transfers significantly less data than its competitors.

7.6 Varying Maximum Length

In this third experiment, we study the methods’ behavior as we vary the maximum length σ . The minimum collection frequency is set as $\tau = 100$ for NYT and $\tau = 1,000$ for CW to reflect their different scale. Measurements are performed on the entire datasets with 64 map/reduce slots and reported in Figure 5. Measurements for $\sigma > 5$ are missing for NAÏVE on CW, since the method did not finish within reasonable time for those parameter settings.

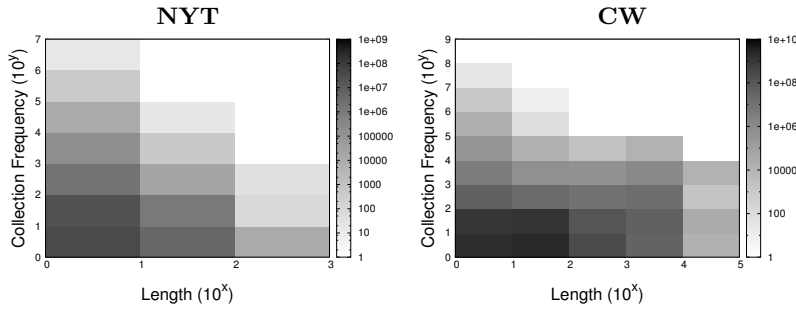


Figure 2: Output characteristics as # of n -grams s with $cf(s) \geq 5$ per n -gram length and collection frequency

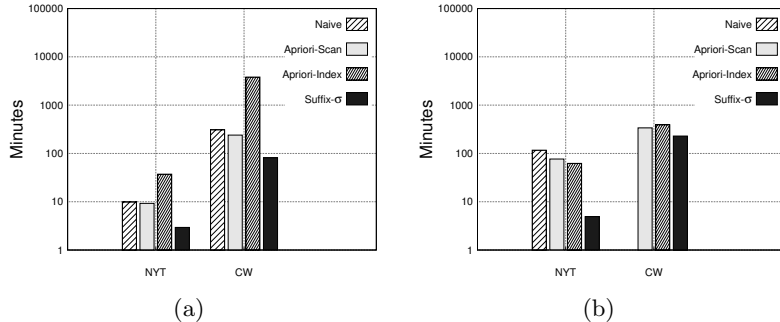


Figure 3: Wallclock times in minutes for (a) *training a language model* ($\sigma = 5$, NYT: $\tau = 10$ / CW: $\tau = 100$) and (b) *text analytics* ($\sigma = 100$, NYT: $\tau = 100$ / CW: $\tau = 1,000$) as two typical use cases

SUFFIX- σ is on par with the best-performing competitor on CW, when considering n -grams of length up to 50. For $\sigma = 100$, it outperforms the next best APRIORI-SCAN by a factor 1.5x. On NYT, Suffix- σ consistently outperforms all competitors by a wide margin. When we increase the value of σ , the APRIORI-based methods need to run more Hadoop jobs, so that their wallclock times keep increasing. For NAÏVE and SUFFIX- σ , on the other hand, we observe a saturation of wallclock times. This is expected, since these methods have to do additional work only for input sequences longer than σ consisting of terms that occur at least τ times in the document collection. When looking at the number of bytes and the number of records transferred, we observe a saturation for NAÏVE for the reason mentioned above. For SUFFIX- σ only the number of bytes saturates, the number of records transferred is constant, since it depends only on the minimum collection frequency τ . Further, we see that SUFFIX- σ consistently transfers fewest records.

7.7 Scaling the Datasets

Next, we investigate how the methods react to changes in the scale of the datasets. To this end, both from NYT and CW, we extract smaller datasets that contain a random 25%, 50%, or 75% subset of the documents. Again, the minimum collection frequency is set as $\tau = 100$ for NYT and $\tau = 1,000$ for CW. The maximum length is set as $\sigma = 5$. Wallclock times are measured using 64 map/reduce slots.

From Figure 6, we observe that NAÏVE handles additional data equally well on both datasets. The other methods' scalability is comparable to that of NAÏVE on CW, as can be seen from their almost-identical slopes. On NYT, in contrast, APRIORI-SCAN, APRIORI-INDEX, and SUFFIX- σ cope slightly better with additional data than NAÏVE. This is due to the different characteristics of the two datasets.

7.8 Scaling Computational Resources

Our final experiment explores how the methods behave as we scale computational resources. Again, we set $\tau = 100$ for NYT and $\tau = 1,000$ for CW. All methods are applied to the 50% samples of documents from the collections. We vary the number of map/reduce slots as 16, 32, 48, and 64. The number of cluster nodes remains constant in this experiment, since we cannot add/remove machines to/from the cluster due to organizational restrictions. We thus only vary the amount of parallel work every machine can do; their total number remains constant throughout this experiment.

We observe from Figure 7 that all methods show comparable behavior as we make additional computational resources available. Or, put differently, all methods make equally effective use of them. What can also be observed across all methods is that the gains of adding more computational resources are diminishing – because of mappers and reducers competing for shared devices such as hard disks and network interfaces. This phenomenon is more pronounced on NYT than CW, since methods take generally less time on the smaller dataset, so that competition for shared devices is fiercer and has no chance to level out over time.

Summary

What we see in our experiments is that SUFFIX- σ outperforms its competitors when long and/or less frequent n -grams are considered. Even otherwise, when the focus is on short and/or very frequent n -grams, SUFFIX- σ performs never significantly worse than the other methods. It is hence robust and can handle a wide variety of parameter choices. To substantiate this, consider that SUFFIX- σ could compute statistics about arbitrary-length n -grams that occur at least five times (i.e., $\tau = 5$ and $\sigma = \infty$), as reported in Figure 2, in less than six minutes on NYT and six hours on CW.

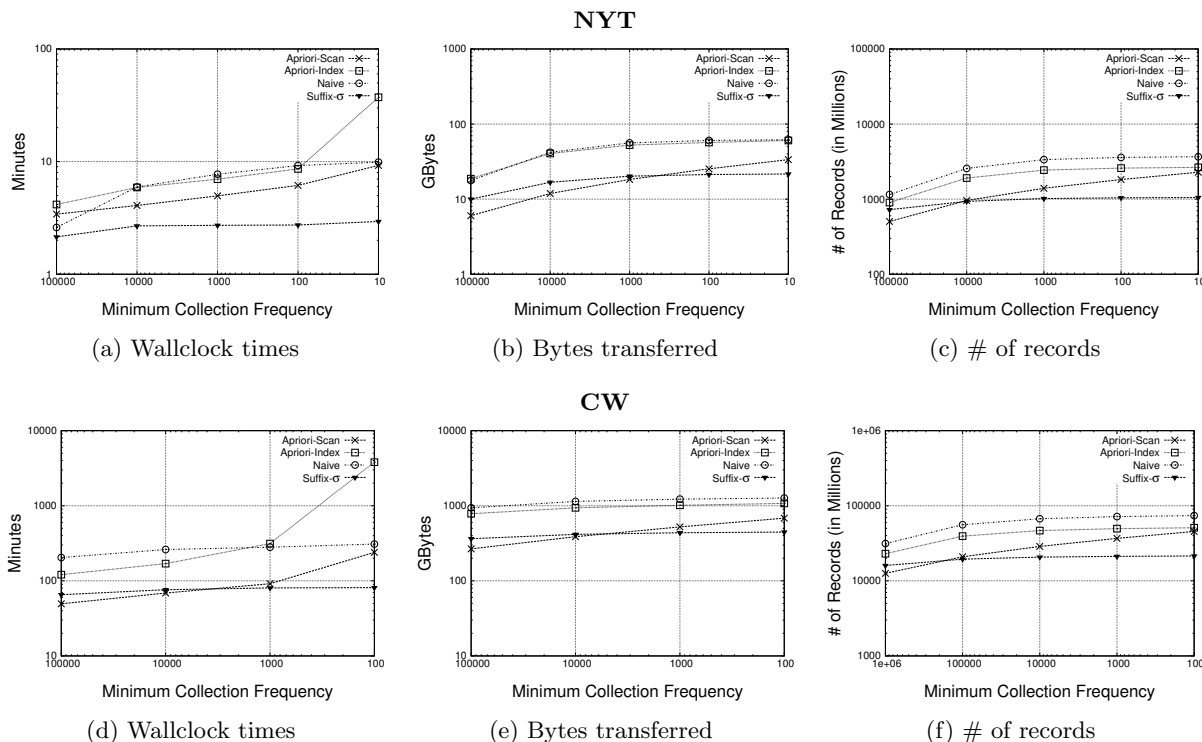


Figure 4: Varying the minimum collection frequency τ

8. RELATED WORK

We now discuss the connection between this work and existing literature, which can broadly be categorized into:

Frequent Pattern Mining goes back to the seminal work by Agrawal et al. [8] on identifying frequent itemsets in customer transactions. While the APRIORI algorithm described therein follows a candidate generation & pruning approach, Han et al. [21] have advocated pattern growth as an alternative approach. To identify frequent sequences, which is a problem closer to our work, the same kinds of approaches can be used. Agrawal and Srikant [10, 38] describe candidate generation & pruning approaches; Pei et al. [37] propose a pattern-growth approach. SPADE by Zaki [45] also generates and prunes candidates but operates on an index structure as opposed to the original data. Parallel methods for frequent pattern mining have been devised both for distributed-memory [19] and shared-memory machines [36, 44]. Little work exists that assumes MapReduce as a model of computation. Li et al. [26] describe a pattern-growth approach to mine frequent itemsets in MapReduce. Huang et al. [22] sketch an approach to maintain frequent sequences while sequences in the database evolve. Their approach is not applicable in our setting, since it expects input sequences to be aligned (e.g. based on time) and only supports document frequency. For more detailed discussions, we refer to Ceglar and Roddick [14] for frequent itemset mining, Mabroukeh and Ezeife [30] for frequent sequence mining, and Han et al. [20] for frequent pattern mining in general.

Natural Language Processing & Information Retrieval. Given their role in NLP, multiple efforts [11, 15, 18, 23, 39] have looked into n -gram statistics computation. While these approaches typically consider document collections of modest size, recently Lin et al. [27] and Nguyen et al. [34] targeted web-scale data. Among the aforemen-

tioned work, Huston et al. [23] is closest to ours, also focusing on less frequent n -grams and using a cluster of machines. However, they only consider n -grams consisting of up to eleven words and do not provide details on how their methods can be adapted to MapReduce. Yamamoto and Church [43] augment suffix arrays, so that the collection frequency of substrings in a document collection can be determined efficiently. Bernstein and Zobel [12] identify long n -grams as a means to spot co-derivative documents. Brants et al. [13] and Wang et al. [40] describe the n -gram statistics made available by Google and Microsoft, respectively. Zhai [46] gives details on the use of n -gram statistics in language models. Michel et al. [32] demonstrated recently that n -gram time series are powerful tools to understand the evolution of culture and language.

MapReduce Algorithms. Several efforts have looked into how specific problems can be solved using MapReduce, including all-pairs document similarity [28], processing relational joins [35], coverage problems [16], content matching [33]. However, no existing work has specifically addressed computing n -gram statistics in MapReduce.

9. CONCLUSIONS

In this work, we have presented SUFFIX- σ , a novel method to compute n -gram statistics using MapReduce as a platform for distributed data processing. Our evaluation on two real-world datasets demonstrated that SUFFIX- σ outperforms MapReduce adaptations of APRIORI-based methods significantly, in particular when long and/or less frequent n -grams are considered. Otherwise, SUFFIX- σ is robust, performing at least on par with the best competitor. We also argued that our method is easier to implement than its competitors, having been designed with MapReduce in mind.

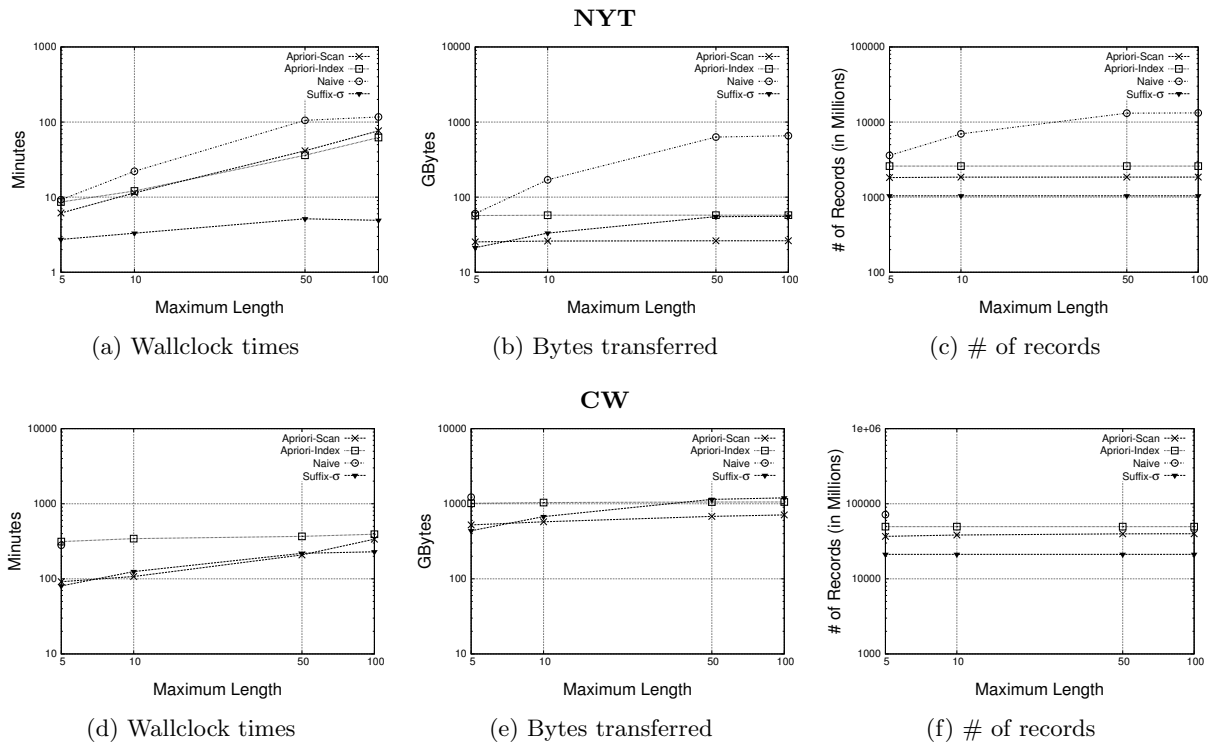


Figure 5: Varying the maximum length σ

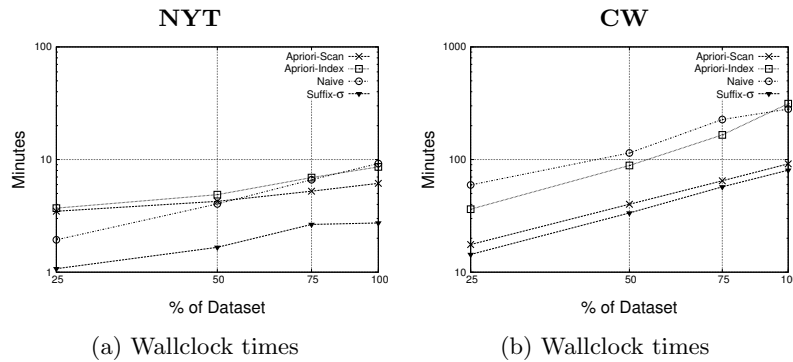


Figure 6: Scaling the datasets

Finally, we established our method’s versatility by showing that it can be extended to produce maximal/closed n -grams and perform aggregations beyond occurrence counting.

Acknowledgments

(i) This work is supported by the 7th Framework IST programme of the European Union through the focused research project (STREP) on Longitudinal Analytics of Web Archive data (LAWA) under contract no. 258105. (ii) Calibration experiments for this work were run on the Amazon Elastic Compute Cloud (EC2) and possible thanks to an Amazon Web Service Research Grant.

Code

<http://www.mpi-inf.mpg.de/~kberberi/edbt2013/>

10. REFERENCES

- [1] Apache Hadoop <http://hadoop.apache.org/>.
- [2] Apache OpenNLP <http://opennlp.apache.org/>.
- [3] Berkeley DB Java Edition <http://www.oracle.com/products/berkeleydb/>.
- [4] Boilerpipe <http://code.google.com/p/boilerpipe/>.
- [5] Google n -Gram Corpus <http://googleresearch.blogspot.de/2006/08/all-our-n-gram-are-belong-to-you.html>.
- [6] The ClueWeb09 Dataset <http://lemurproject.org/clueweb09>.
- [7] The New York Times Annotated Corpus <http://corpus.nytimes.com>.
- [8] R. Agrawal et al. Mining association rules between sets of items in large databases. *SIGMOD* 1993

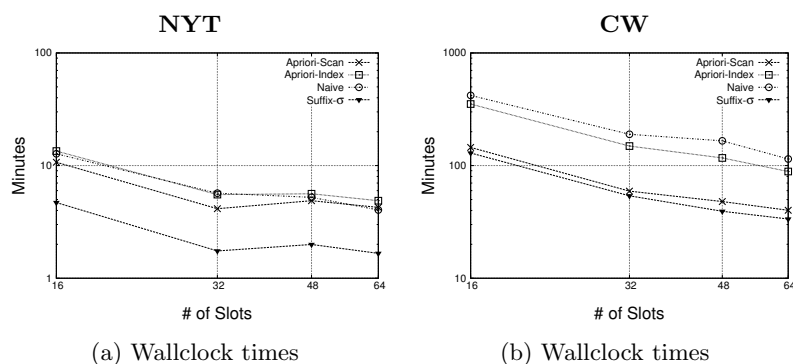


Figure 7: Scaling computational resources

- [9] R. Agrawal and R. Srikant. Fast algorithms for mining association rules in large databases. *VLDB* 1994.
- [10] R. Agrawal and R. Srikant. Mining sequential patterns. *ICDE* 1995.
- [11] S. Banerjee and T. Pedersen. The design, implementation, and use of the ngram statistics package. *CICLing* 2003.
- [12] Y. Bernstein and J. Zobel. Accurate discovery of co-derivative documents via duplicate text detection. *Inf. Syst.*, 31(7):595–609, 2006.
- [13] T. Brants et al. Large language models in machine translation. *EMNLP-CoNLL* 2007.
- [14] A. Ceglar and J. F. Roddick. Association mining. *ACM Comput. Surv.*, 38(2), 2006.
- [15] H. Ceylan and R. Mihalcea. An efficient indexer for large n-gram corpora. *ACL* 2011.
- [16] F. Chierichetti et al. Max-cover in map-reduce. *WWW* 2010.
- [17] J. Dean and S. Ghemawat. Mapreduce: Simplified data processing on large clusters. *OSDI* 2004.
- [18] M. Federico et al. Irstlm: an open source toolkit for handling large scale language models. *INTERSPEECH* 2008.
- [19] V. Guralnik and G. Karypis. Parallel tree-projection-based sequence mining algorithms. *Parallel Computing*, 30(4):443–472, 2004.
- [20] J. Han et al. Frequent pattern mining: current status and future directions. *DMKD*, 15(1):55–86, 2007.
- [21] J. Han et al. Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *DMKD*, 8(1):53–87, 2004.
- [22] J.-W. Huang et al. Dpsp: Distributed progressive sequential pattern mining on the cloud. *PAKDD* 2010.
- [23] S. Huston et al. Efficient indexing of repeated n-grams. *WSDM* 2011.
- [24] S. Katz. Estimation of probabilities from sparse data for the language model component of a speech recognizer. *ASSP*, 35(3):400–401, 1987.
- [25] C. Kohlschütter et al. Boilerplate detection using shallow text features. *WSDM* 2010.
- [26] H. Li et al. Pfp: parallel fp-growth for query recommendation. *RecSys* 2008.
- [27] D. Lin et al. New tools for web-scale n-grams. *LREC* 2010.
- [28] J. Lin. Brute force and indexed approaches to pairwise document similarity comparisons with mapreduce. *SIGIR* 2009.
- [29] J. Lin and C. Dyer. *Data-Intensive Text Processing with MapReduce*. Morgan & Claypool Publishers, 2010.
- [30] N. R. Mabroukeh and C. I. Ezeife. A taxonomy of sequential pattern mining algorithms. *ACM Comput. Surv.*, 43(1):3, 2010.
- [31] U. Manber and E. W. Myers. Suffix arrays: A new method for on-line string searches. *SIAM J. Comput.*, 22(5):935–948, 1993.
- [32] J.-B. Michel et al. Quantitative Analysis of Culture Using Millions of Digitized Books. *Science*, 2010.
- [33] G. D. F. Morales et al. Social content matching in mapreduce. *PVLDB*, 4(7):460–469, 2011.
- [34] P. Nguyen et al. Msrlm: a scalable language modeling toolkit. MSR-TR-2007-144
- [35] A. Okcan and M. Riedewald. Processing theta-joins using mapreduce. *SIGMOD* 2011.
- [36] S. Parthasarathy et al. Parallel data mining for association rules on shared-memory systems. *Knowl. Inf. Syst.*, 3(1):1–29, 2001.
- [37] J. Pei et al. Mining sequential patterns by pattern-growth: The prefixspan approach. *TKDE* 2004.
- [38] R. Srikant and R. Agrawal. Mining sequential patterns: Generalizations and performance improvements. *EDBT* 1996.
- [39] A. Stolcke. Srilm - an extensible language modeling toolkit. *INTERSPEECH* 2002.
- [40] K. Wang et al. An Overview of Microsoft Web N-gram Corpus and Applications. *NAACL-HLT* 2010.
- [41] T. White. *Hadoop: The Definitive Guide*. O’Reilly Media, Inc., 2nd edition, 2010.
- [42] I. H. Witten et al. *Managing Gigabytes: Compressing and Indexing Documents and Images, Second Edition*. Morgan Kaufmann, 1999.
- [43] M. Yamamoto and K. W. Church. Using suffix arrays to compute term frequency and document frequency for all substrings in a corpus. *Comput. Linguist.*, 27:1–30, March 2001.
- [44] M. J. Zaki. Parallel sequence mining on shared-memory machines. *J. Parallel Distrib. Comput.*, 61(3):401–426, 2001.
- [45] M. J. Zaki. Spade: An efficient algorithm for mining frequent sequences. *ML*, 42(1/2):31–60, 2001.
- [46] C. Zhai. Statistical language models for information retrieval: a critical review. *FTIR*, 2:137–213, 2008.