

# Probabilistic Management of OCR Data using an RDBMS

Arun Kumar  
University of Wisconsin-Madison  
arun@cs.wisc.edu

Christopher Ré  
University of Wisconsin-Madison  
chrisre@cs.wisc.edu \*

## ABSTRACT

The digitization of scanned forms and documents is changing the data sources that enterprises manage. To integrate these new data sources with enterprise data, the current state-of-the-art approach is to convert the images to ASCII text using optical character recognition (OCR) software and then to store the resulting ASCII text in a relational database. The OCR problem is challenging, and so the output of OCR often contains errors. In turn, queries on the output of OCR may fail to retrieve relevant answers. State-of-the-art OCR programs, e.g., the OCR powering Google Books, use a probabilistic model that captures many alternatives during the OCR process. Only when the results of OCR are stored in the database, do these approaches discard the uncertainty. In this work, we propose to retain the probabilistic models produced by OCR process in a relational database management system. A key technical challenge is that the probabilistic data produced by OCR software is very large (a single book blows up to 2GB from 400kB as ASCII). As a result, a baseline solution that integrates these models with an RDBMS is over 1000x slower versus standard text processing for single table select-project queries. However, many applications may have quality-performance needs that are in between these two extremes of ASCII and the complete model output by the OCR software. Thus, we propose a novel approximation scheme called STACCATO that allows a user to trade recall for query performance. Additionally, we provide a formal analysis of our scheme's properties, and describe how we integrate our scheme with standard-RDBMS text indexing.

## 1. INTRODUCTION

The mass digitization of books, printed documents, and printed forms is changing the types of data that companies and academics manage. For example, Google Books and

\*This work is supported by the Microsoft Jim Gray Systems Lab, the National Science Foundation (IIS-1054009) and the Office of Naval Research (N000141210041).

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. Articles from this volume were invited to present their results at The 38th International Conference on Very Large Data Bases, August 27th - 31st 2012, Istanbul, Turkey.

*Proceedings of the VLDB Endowment*, Vol. 5, No. 4

Copyright 2011 VLDB Endowment 2150-8097/11/12... \$ 10.00.

their academic partner, the Hathi Trust, have the goal of digitizing all of the world's books to allow scholars to search human knowledge from the pre-Web era. The hope of this effort is that digital access to this data will enable scholars to rapidly mine these vast stores of text for new discoveries.<sup>1</sup> The potential users of this new content are not limited to academics. The market for *enterprise document capture* (scanning of forms) is already in the multibillion dollar range [3]. In many of the applications, the translated data is related to enterprise business data, and so after converting to plain text, the data are stored in an RDBMS [6].

Translating an image of text (e.g., a jpeg) to ASCII is difficult for machines to do automatically. To cope with the huge number of variations in scanned documents, e.g., in spacing of the glyphs and font faces, state-of-the-art approaches for optical character recognition (OCR) use probabilistic techniques. For example, the OCRopus tool from Google Books represents the output of the OCR process as a stochastic automaton called a *finite-state transducer* (FST) that defines a probability distribution over all possible strings that could be represented in the image.<sup>2</sup> An example image and its resulting (simplified) transducer are shown in Figure 1. Each labeled path through the transducer corresponds to a potential string (one multiplies the weights along the path to get the probability of the string). Only to produce the final plain text do current OCR approaches remove the uncertainty. Traditionally, they choose to retain only the single most likely string produced by the FST (called a *maximum a priori estimate* or MAP [1]).

As Google Books demonstrates, the MAP works well for browsing applications. In such applications, one is sensitive to precision (i.e., are the answers I see correct), but one is insensitive to recall (i.e., what fraction of all of the answers in my corpus are returned). But this is not true of all applications: an English professor looking for the earliest dates that a word occurs in a corpus is sensitive to recall [5]. As is an insurance company that wants all insurance claims that were filled in 2010 that mentioned a 'Ford'. This latter query is expressed in SQL in Figure 1(C). In this work, we focus on such single table select-project queries, whose outputs are standard probabilistic RDBMS tables. Using the MAP approach may miss valuable answers. In the example in Figure 1, the most likely string does not contain 'Ford', and so we (erroneously) miss this claim. However, the string 'Ford'

<sup>1</sup>Many repositories of *Digging into Data Challenge* (a large joint effort to bring together social scientists with data analysis) are OCR-based <http://www.diggingintodata.org>.

<sup>2</sup><http://code.google.com/p/ocropus/>.

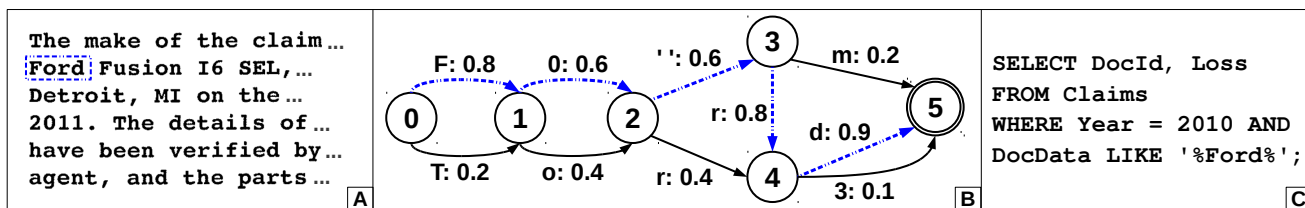


Figure 1: (A) An image of text. (B) A portion of a simple FST resulting from the OCR of the highlighted part of (A). The numbers on the arcs are conditional probabilities of transitioning from one state to another. An emitted string corresponds to a path from states 0 to 5. The string ‘F0 rd’ (highlighted path) has the highest probability,  $0.8 * 0.6 * 0.6 * 0.8 * 0.9 \approx 0.21$ . (C) An SQL query to retrieve loss information that contains ‘Ford’. Using the MAP approach, no claim is found. Using Staccato, a claim is found (with probability 0.12).

does appear (albeit with a lower probability). Empirically, we show that the recall for simple queries on real-world OCR can be as low as 0.3 – and so we may throw away almost 70% of our data if we follow the MAP approach.

To remedy this recall problem, our baseline approach is to store and handle the FSTs as binary large objects inside the RDBMS. As with a probabilistic relational database, the user can then pose questions as if the data are deterministic and it is the job of the system to compute the confidence in its answer. By combining existing open-source tools for transducer composition<sup>3</sup> with an RDBMS, we can then answer queries like that in Figure 1(C). This approach achieves a high quality (empirically, the recall we measured is very close to 1.0, with up to 0.9 precision). Additionally, the enterprise users can ask their existing queries directly on top of the RDBMS data (the query in Figure 1(C) remains unchanged). The downside is that query processing is much slower (up to 1000x slower). While the query processing time for transducers is linear in the data size, the transducers themselves are huge, e.g., a single 200-page book blows up from 400 kB as text to over 2 GB when represented by transducers after OCR. This motivates our central question: “Can we devise an approximation scheme that is somewhere in between these two extremes of recall and performance?”

State-of-the-art OCR tools segment each of the images corresponding to pages in a document into lines using special purpose line-breaking tools. Breaking a single line further into individual words is more difficult (spacing is very difficult to accurately detect). With this in mind, a natural idea to improve the recall of the MAP approach is to retain not only the highest probability string for each line, but instead to retain the  $k$  highest probability strings that appear in each line (called  $k$ -MAP [28, 53]). Indeed, this technique keeps more information around at a linear cost (in  $k$ ) in space and processing time. However, we show that even storing hundreds of paths makes an insignificant jump in the recall of queries.

To combat this problem, we propose a novel approximation scheme called STACCATO, which is our main technical contribution. The main idea is to apply  $k$ -MAP not to the whole line, but to first break the line into smaller chunks which are themselves transducers and apply  $k$ -MAP to each transducer individually. This allows us to store exponentially more alternatives than  $k$ -MAP (exponential in the number of chunks), while using roughly a linear amount more space than the MAP approach. If there is only a single chunk, then STACCATO’s output is equivalent to  $k$ -MAP.

If essentially every possible character is a chunk, then we retain the full FST. Experimentally, we demonstrate that the STACCATO approach *gracefully trades off between performance and recall*. For example, when looking for mentions of laws on a data set that contains scanned acts of the US congress, the MAP approach achieves a recall of 0.28 executing in about 1 second, the full FST approach achieves perfect recall but takes over 2 minutes. An intermediate representation from STACCATO takes around 10 seconds and achieves 0.76 recall. Of course, there is a fundamental trade off between precision and recall. On the same query as above, the MAP has precision 1.0, and the full FST has precision 0.25, while STACCATO achieves 0.73. In general, STACCATO’s precision falls in between the MAP and the full FST.

To understand STACCATO’s approximation more deeply, we conduct a formal analysis, which is our second technical contribution. When constructing STACCATO’s approximation, we ensure two properties (1) each chunk forms a transducer (as opposed to a more general structure), and (2) that the model retains the *unique path property*, i.e., that every string corresponds to a unique path. While both of these properties are satisfied by the transducers produced by OCRopus, neither property is necessary to have a well-defined approximation scheme. Moreover, enforcing these two properties increases the complexity of our algorithm and may preclude some compact approximations. Thus, it is natural to wonder if we can relax these two properties. While we cannot prove that these two conditions are necessary, we show that without these two properties, basic operations become intractable. Without the unique path property, prior work has shown that determining (even approximating) the  $k$ -MAP is intractable for a fixed  $k$  [32]. Even with the unique path property and a fixed set of chunks, we show that essentially the simplest violation of property (1) makes it intractable to construct an approximation even for  $k = 2$  (Theorem 3.1). On the positive side, for any fixed partition, STACCATO retains a set of strings that achieves the highest total probability among approximations that satisfy the above restrictions.

Finally, we describe how to use standard text-indexing techniques to improve query performance. Directly applying an inverted index to transducer data is essentially doomed to failure: the sheer number of terms one would have to index grows exponentially with the length of the document, e.g., an FST for a single line may represent over  $10^{100}$  terms. To combat this, we allow the user to specify a dictionary of terms. We then construct an index of those terms specified in the dictionary. This allows us to process keyword and some regular expressions using standard techniques [14, 52].

<sup>3</sup>OpenFST. <http://www.openfst.org/>

*Outline.* In Section 2, we illustrate our current prototype system to manage OCR data using an RDBMS with an example, and we present a brief background on the use of transducers in OCR. In Section 3, we briefly describe the baseline solutions, and then discuss the main novel technical contributions of this work, viz., the STACCATO approximation scheme and our formal analysis of its properties. In Section 4, we describe our approach for indexing OCR transducer data, which is another technical contribution of this work. In Section 5, we empirically validate that our approach is able to trade off recall for query-runtime performance on several real-world OCR data sets. We validate that our approximation methods can be efficiently implemented, and that our indexing technique provides the expected speedups. In Section 6, we discuss related work.

## 2. PRELIMINARIES

The key functionality that STACCATO provides is to enable users to query OCR data inside an RDBMS as if it were regular text. Specifically, we want to enable the LIKE predicate of SQL on OCR data. We describe STACCATO through an example, followed by a more detailed explanation of its semantics and the formal background.

### 2.1 Using Staccato with OCR

Consider an insurance company that stores loss data with scanned report forms in a table with the following schema:

**Claims**(*DocID*, *Year*, *Loss*, *DocData*)

A document tuple contains an id, the year the form was filed (*Year*), the amount of the loss (*Loss*) and the contents of the report (*DocData*). A simple query that an insurance company may want to ask over the table - “*Get loss amounts of all claims in 2010 where the report mentions ‘Ford’*”. Were *DocData* ASCII text, this could be expressed as an SQL query as follows:

```
SELECT DocID, Loss FROM Claims
WHERE Year = 2010 AND DocData LIKE '%Ford%';
```

If *DocData* is standard text, the semantics of this query is straightforward: we examine each document filed in 2010, and check if it contains the string ‘*Ford*’. The challenge is that instead of a single document, in OCR applications *DocData* represents many different documents (each document is weighted by probability). In STACCATO, we can express this as an SQL query that uses a simple pattern in the LIKE predicate (also in Figure 1(C)). The twist is that the underlying processing must take into account the probabilities from the OCR model.

Formally, STACCATO allows a larger class of queries in the LIKE predicate that can be expressed as deterministic finite automata (DFAs). STACCATO translates the syntax above in to a DFA using standard techniques [27]. As with probabilistic databases [13, 22, 30, 50], STACCATO computes the probability that the document matches the regular expression. STACCATO does this using algorithms from prior work [32, 43]. The result is a probabilistic relation; after this, we can apply probabilistic relational database processing techniques [22, 41, 46]. In this work, we consider only single table select-project queries (joins are handled using the above mentioned techniques).

A critical challenge that STACCATO must address is given a DFA find those documents that are relevant to the query

expressed by the DFA. For a fixed query, the existing algorithms are roughly linear in the size of data that they must process. To improve the runtime of these algorithms, one strategy (that we take) is to reduce the size of the data that must be processed using approximations. The primary contribution of STACCATO is the set of mechanisms that we describe in Section 3 to achieve the trade off of quality and performance by approximating the data. We formally study the properties of our algorithms and describe simple mechanisms to allow the user to set these parameters in Sec. 3.2.

One way to evaluate the query above in the deterministic setting is to scan the string in each report and check for a match. A better strategy may be to use an inverted index to fetch only those documents that contain ‘*Ford*’. In general, this strategy is possible for *anchored* regular expressions [19], which are regular expressions that begin or end with words in the language, e.g. ‘no.(2|3)’ is anchored while ‘(no|num).(2|8)’ is not. STACCATO supports a similar optimization using standard text-indexing techniques. There is, however, one twist: At one extreme, any term may have some small probability of occurring at every location of the document – which renders the index ineffective. Nevertheless, we show that STACCATO is able to provide efficient indexing for anchored regular expressions using a dictionary-based approach.

### 2.2 Background: Stochastic Finite Automata

We formally describe STACCATO’s data model that is based on Stochastic Finite Automata (SFA). This model is essentially identical to the model output by Google’s OCRopus [8, 39].<sup>4</sup> An SFA is a finite state machine that emits strings (e.g., the ASCII conversion of an OCR image). The model is stochastic, which captures the uncertainty in translating the glyphs and spaces to ASCII characters.

At a high level, an SFA over an alphabet  $\Sigma$  represents a discrete probability distribution  $P$  over strings in  $\Sigma^*$ , i.e.,

$$P : \Sigma^* \rightarrow [0, 1] \text{ such that } \sum_{x \in \Sigma^*} P(x) = 1$$

The SFA represents the (finitely many) strings with non-zero probability using an automaton-like structure that we first describe using an example:

**Example 1.** Figure 1 shows an image of text and a simplified SFA created by OCRopus from that data. The SFA is a directed acyclic labeled graph. The graphical structure (i.e., the branching) in the SFA is used by the OCR tool to capture correlations between the emitted letters. Each source-to-sink path (i.e., a path from node 0 to node 5) corresponds to a string with non-zero probability. For example, the string ‘*Ford*’ is one possible path that uses the following sequence of nodes  $0 \rightarrow 1 \rightarrow 2 \rightarrow 4 \rightarrow 5$ . The probability of this string can be found by multiplying the edge weights corresponding to the path:  $0.8 * 0.4 * 0.4 * 0.9 \approx 0.12$ .  $\square$

Formally, we fix an alphabet  $\Sigma$  (in STACCATO, this is the set of ASCII characters). An SFA  $S$  over  $\Sigma$  is a tuple  $S = (V, E, s, f, \delta)$  where  $V$  is a set of nodes,  $E \subseteq V \times V$  is a set of edges such that  $(V, E)$  is a directed acyclic graph, and  $s$

<sup>4</sup>Our prototype uses the same weighted finite state transducer (FST) model that is used by OpenFST and OCRopus. We simplify FST to SFAs here only slightly for presentation. See the full version for more details [34]

(resp.  $f$ ) is a distinguished start (resp. final) node. The function  $\delta$  is a stochastic transition function, i.e.,

$$\delta : E \times \Sigma \rightarrow [0, 1] \text{ s.t. } \sum_{\substack{y:(x,y) \in E \\ \sigma \in \Sigma}} \delta((x, y), \sigma) = 1 \quad \forall x \in V$$

In essence,  $\delta(e, \sigma)$ , where  $e = (x, y)$ , is the conditional probability of transitioning from  $x \rightarrow y$  and emitting  $\sigma$ .

An SFA defines a probability distribution via its labeled paths. A labeled path from  $s$  to  $f$  is denoted by  $p = (e_1, \sigma_1), \dots, (e_N, \sigma_N)$ , where  $e_i \in E$  and  $\sigma_i \in \Sigma$ , corresponding to the string  $\sigma_1 \dots \sigma_n$ , with its probability:<sup>5</sup>

$$\Pr_S[p] = \prod_{i=1}^{|p|} \delta(e_i, \sigma_i)$$

SFAs in OCR satisfy an important property that we call the *unique paths property* that says that any string produced by the SFA with non-zero probability is generated by a unique labeled path through the SFA. We denote by UP the function that takes a string to its unique labeled path. This property guarantees tractability of many important computations over SFAs including finding the highest probability string produced by the SFA [32].

Unlike the example given here, the SFAs produced by Google’s OCRopus are much larger: they contain a weighted arc for every ASCII character. And so, the SFA for a single line can require as much as 600 kB to store.

Queries in STACCATO are formalized in the standard way for probabilistic databases. In this paper, we consider LIKE predicates that contain Boolean queries expressed as DFAs (STACCATO handles non-Boolean queries using algorithms in Kimmelfeld and Ré [32]). Fix an alphabet  $\Sigma$  (the ASCII characters). Let  $q : \Sigma^* \rightarrow \{0, 1\}$  be expressed as DFA and  $x$  be any string. We have  $q(x) = 1$  when  $x$  satisfies the query, i.e., it’s accepted by the DFA. We compute the probability that  $q$  is true; this quantity is denoted  $\Pr[q]$  and is defined by  $\Pr[q] = \sum_{x \in \Sigma^*} q(x) \Pr(x)$  (i.e., simply sum over all possible strings where  $q$  is true). There is a straightforward algorithm based on matrix multiplication to process these queries that is linear in the size of the data and cubic in the number of states of the DFA [43].

### 3. MANAGING SFAS IN AN RDBMS

We start by outlining two baseline approaches that represent the two extremes of query performance and recall. Then, we describe the novel approximation scheme of STACCATO, which enables us to trade performance for recall.

**Baseline Approaches.** We study two baseline approaches:  $k$ -MAP and the FullSFA approach. Fix some  $k \geq 1$ . In the  $k$ -MAP approach we store the  $k$  highest probability strings (simply, top  $k$  strings) generated by each SFA in our databases. We store one tuple per string along with the associated probability. Query processing is straightforward: we process each string using standard text-processing techniques, and then sum the probability of each string (since each string is a disjoint probabilistic event). In the FullSFA approach, we store the entire SFA as a BLOB inside the

<sup>5</sup>Many (including OpenFST) tools use a formalization with log-odds instead of probabilities. It has some intuitive property for graph concepts, e.g., the shortest path corresponds to the most likely string.

RDBMS. To answer a query, we retrieve the BLOB, deserialize it, and then use an open source C++ automata composition library to answer the query [11, 12] and compute all probabilities. Table 1 summarizes the time and space costs for a simple chain SFA (no branching). This table gives an engineer’s intuition about the time and space complexity of the baseline approaches. The factor 16 accounts for the metadata – tuple ID, location in SFA, and probability value (the schema is described in the full version [34]). We also include our proposed approach, STACCATO that depends on a parameter  $m$  (the number of chunks) that we describe below. From the table, we can read that query processing time for STACCATO is essentially linear in  $m$ . Let  $l$  be the length of the document, since  $m \in [1, l]$  query processing time in STACCATO interpolates linearly from the  $k$ -MAP approach to the FullSFA approach.

	$k$ -MAP	FullSFA	STACCATO
QUERY	$lqk$	$lq \Sigma  + q^3(l - 1)$	$lqk + q^3(m - 1)$
SPACE	$lk + 16k$	$l \Sigma  + 16l \Sigma $	$lk + 16mk$

- $l$  : length of the SFA’s strings
- $q$  : # states in the query DFA
- $k$  : # paths parameter in  $k$ -MAP, STACCATO
- $m$  : # chunks in STACCATO ( $1 \leq m \leq l$ )

**Table 1: Space costs and query processing times for a simple chain SFA. The space indicates the number of bytes of storage required.**

### 3.1 Approximating an SFA with Chunks

As mentioned before, the SFAs in OCR are much larger than our example, e.g. one OCR line from a scanned book yielded an SFA of size 600 kB. In turn, the 200-page book blows up to over 2 GB when represented by SFAs. Thus, to answer a query that spans many books in the FullSFA approach, we must read a huge amount of data. This can be a major bottleneck in query processing. To combat this we propose to approximate an SFA with a collection of smaller-sized SFAs (that we call *chunks*). Our goal is to create an approximation that allows us to gracefully tradeoff from the fast-but-low-recall MAP approach to the slow-but-high-recall FullSFA approach.

Recall that the  $k$ -MAP approach is a natural first approximation, wherein we simply store the top- $k$  paths in each of the per-line SFAs. This approach can increase the recall at a linear cost in  $k$ . However, as we demonstrate experimentally, simply increasing  $k$  is insufficient to tradeoff between the two extremes. That is, even for huge values of  $k$  we do not achieve full recall.

Our idea to combat the slow increase of recall starts with the following intuition: the more strings from the SFA we store, the higher our recall will be. We observe that if we store the top  $k$  in each of  $m$  smaller SFAs (that we refer to as ‘chunks’), we effectively store  $k^m$  distinct strings. Thus, increasing the value of  $k$  increases the number of strings polynomially. In contrast, increasing  $m$ , the number of smaller SFAs, increases the number of paths *exponentially*, as illustrated in Figure 2. This observation motivates the idea that to improve quality, we should divide the SFA further. As we demonstrate experimentally, STACCATO achieves the most conceptually important feature of our approximation:

**Algorithm 1: FindMinSFA**

**Inputs:** SFA  $S$  with partial order  $\leq$  on its nodes,  $X \subseteq V$

**while**  $X$  does not form a valid SFA **do**

**if** No unique start node in  $X$  **then**

Compute the least common ancestor of  $X$ , say,  $l$

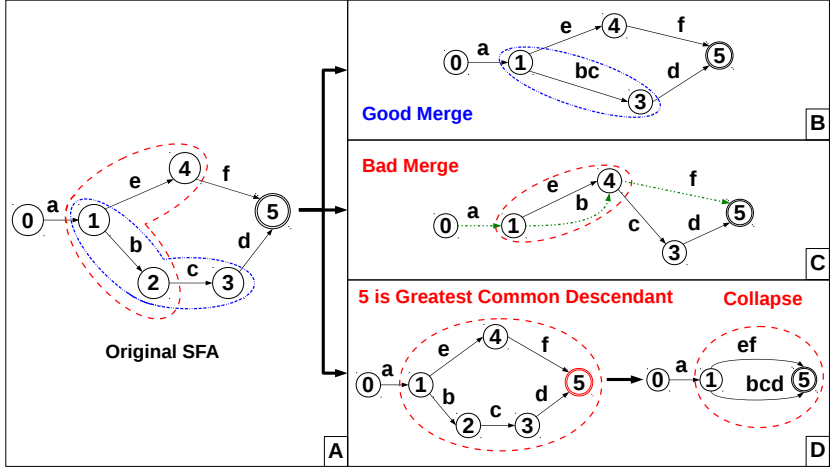
$X \leftarrow X \cup \{y \in V \mid l \leq y \text{ and } \forall x \in X, y \leq x\}$

**if** No unique end node in  $X$  **then**

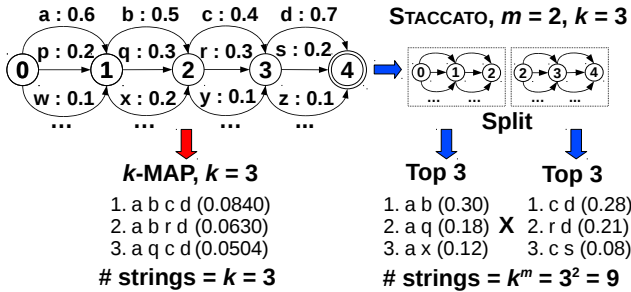
Compute greatest common descendant of  $X$ , say,  $g$

$X \leftarrow X \cup \{y \in V \mid y \leq g \text{ and } \forall x \in X, x \leq y\}$

$\forall e \in E$  s.t. exactly one end-point is in  $X - \{l, g\}$ , add other end-point to  $X$



**Figure 3: Algorithm 1: FindMinSFA.** Illustrating merge and FindMinSFA: (A) Original: The SFA emits two strings:  $ae f$  and  $abcd$ . Two merges considered:  $\{(1,2), (2,3)\}$  (successive edges), and  $\{(1,2), (1,4)\}$  (sibling edges). (B) Good merge: First set gives new edge  $(1,3)$ , emitting  $bc$ . The SFA still emits only  $ae f$  and  $abcd$ . (C) Bad merge: Second set gives new edge  $(1,4)$ , emitting  $e$  and  $b$ . But, the SFA now wrongly emits new strings, e.g.,  $abf$  (dashed lines). (D) Using Algorithm 1 on the second set, the greatest common descendant is obtained (node 5), and the resulting set is collapsed to edge  $(1,5)$ . The SFA now emits only  $ae f$  and  $abcd$ .



**Figure 2: A depiction of conventional Top- $k$  versus Staccato's approximation.**

it allows us to smoothly tradeoff recall for performance. In other words, increasing  $m$  (and  $k$ ) increases the recall at the expense of performance.

**SFA Approximation.** Given an SFA  $S$ , our goal is to find a new SFA  $S'$  that satisfies two competing properties: (1)  $S'$  should be smaller than  $S$ , and (2) the set of strings represented by  $S'$  should contain as many of the high probability strings from  $S$  as possible without containing any strings not in  $S$ .<sup>6</sup> Our technique to approximate the SFA  $S$  is to merge a set of transitions in  $S$  (a ‘chunk’) to produce a new SFA  $S'$ ; then we retain only the top  $k$  transitions on each edge in  $S'$ .

To describe our algorithm, we need some notation. We generalize the definition of SFAs (Section 2) to allow transitions that produce strings (as opposed to single characters). Formally, the transition function  $\delta$  has the type  $\delta : E \times \Sigma^+ \rightarrow [0, 1]$ . Any SFA meets this generalized SFA definition, and so we assume this generalized definition of SFAs for the rest of the section.

<sup>6</sup>This is a type of *sufficient lineage approximation* [44].

Before describing the merging operation formally, we illustrate the challenge in the merging process in Figure 3. Figure 3(A) shows an SFA (without probabilities for readability). We consider two merging operations. First, we have chosen to merge the edges  $(1, 2)$  and  $(2, 3)$  and replaced it with a single edge  $(1, 3)$ . To retain the same strings that are present in the SFA in (A), the transition function must emit the string ‘ $bc$ ’ on the new edge  $(1, 3)$  as illustrated in Figure 3(B). In contrast, if we choose to merge the edges  $(1, 2)$  and  $(1, 4)$ , there is an issue: *no matter what we put on the transition from  $(1, 4)$  we will introduce strings that are not present in the original SFA* (Figure 3(C)). The problem is that the set of nodes  $\{1, 2, 4\}$  do not form an SFA by themselves (there is no unique final node). One could imagine generalizing the definition of SFA to allow richer structures that could capture the correlations between strings, but as we explain in Section 3.2, this approach creates serious technical challenges. Instead, we propose to fix this issue by searching for a minimal SFA  $S'$  that contains this set of nodes (the operation called FINDMINSFA). Then, we replace the nodes in the set with a single edge, retaining only the top  $k$  highest probability strings from  $S'$ . We refer to this operation of replacing  $S'$  with an edge as COLLAPSE. In our example, the result of these operations is illustrated in Figure 3(D).

We describe our algorithm's subroutine FINDMINSFA and then the entire heuristic.

**FindMinSFA.** Given an SFA  $S$  and a set of nodes  $X \subseteq V$ , our goal is to find a SFA  $S'$  whose node set  $Y$  is such that that  $X \subseteq Y$ . We want the set  $Y$  to be minimal in the sense that removing any node  $y \in Y$  causes  $S'$  to violate the SFA property, that is removing  $y$  causes  $S'$  to no longer have a unique start (resp. end) state. Presented in Algorithm 1, our algorithm is based on the observation that the unique start node  $s$  of  $S'$  must come before all nodes in  $X$  in the topological order of the graph (a partial order). Similarly,



the end node  $f$  of the SFA  $S'$  must come after all nodes in  $X$  in topological order. To satisfy these properties, we repeatedly enlarge  $Y$  by computing the start (resp. final node) using the least common ancestor (resp. greatest common descendant) in the DAG. Additionally, we require that any edge in  $S$  that is incident to a node in  $Y$  can be incident to only either  $s$  or  $f$ . (Any node incident to both will be internal to  $S'$ ) If there are no such edges, we are done. Otherwise, for each such edge  $e$ , we include its endpoints in  $Y$  and repeat this algorithm with  $X$  enlarged to  $Y$ . Once we find a suitable set  $Y$ , we replace the set of nodes in the SFA  $S$  with a single edge from  $s$  (the start node of  $S'$ ) to  $f$  (the final node of  $S'$ ). Figure 3(D) illustrates a case when there is no unique end node, and the greatest common descendant has to be computed. More illustrations, covering the other cases, are presented in the full version [34].

**Algorithm Description.** The inputs to our algorithm are the parameters  $k$  (the number of strings retained per edge) and  $m$  (the maximum number of edges that we are allowed to retain in the resulting graph). We describe how a user chooses these parameters in Section 3.2. For now, we focus on the algorithm. At each step, our approximation creates a restricted type of SFA where each edge emits at most  $k$  strings, i.e.,  $\forall e \in E, |\{\sigma \in \Sigma^* \mid \delta(e, \sigma) > 0\}| \leq k$ . When given an SFA not satisfying this property, our algorithm chooses to retain those strings  $\sigma \in \Sigma^*$  with the highest values of  $\delta$  (ties broken arbitrarily). This set can be computed efficiently using the standard Viterbi algorithm [24], which is a dynamic programming algorithm for finding the most likely outputs in probabilistic sequence models, like HMMs. By memoizing the best partial results till a particular state, it can compute the globally optimal results in polynomial time. To compute the top- $k$  results more efficiently, we use an incremental variant by Yen et al [51].

---

**Algorithm 2:** Greedy heuristic over SFA  $S = (V, E)$

---

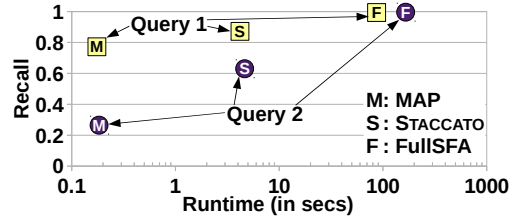
Choose  $\{x, y, z\}$  s.t.  $(x, y), (y, z) \in E$  and maximizing the probability mass of the retained strings.  
 $S \leftarrow \text{COLLAPSE}(\text{FINDMINSFA}(S, \{x, y, z\}))$   
Repeat above steps till  $|E| \leq m$

---

Algorithm 2 summarizes our heuristic: for each triple of nodes  $\{x, y, z\}$  such that  $(x, y), (y, z) \in E$ , we find a minimal containing SFA  $S_{ij}$  by calling  $\text{FINDMINSFA}(\{x, y, z\})$ . We then replace the set of nodes in  $S_{ij}$  by a single edge  $f$  ( $\text{COLLAPSE}$  above). This edge  $f$  keeps only the top- $k$  strings produced by  $S_{ij}$ . Thus, the triple of nodes  $\{x, y, z\}$  generates a candidate SFA. We choose the candidate such that the probability mass of all generated strings is as high as possible (note that since we have thrown away some strings, the total probability mass may be less than 1). Given an SFA we can compute this using the standard sum-product algorithm (a faster incremental variant is actually used in STACCATO). We then continue to recurse until we have reached our goal of finding an SFA that contains fewer than  $m$  edges. A simple optimization (employed by STACCATO) is to cache those candidates we have considered in previous iterations.

While our algorithm is not necessarily optimal, it serves as a proof of concept that our conceptual goal can be achieved. That is, STACCATO offers a knob to tradeoff recall for performance. We describe the experimental setup in more detail

in Section 5, but we illustrate our point with a simple experimental result. Figure 4 plots the recall and runtimes of the two baselines and STACCATO. Here, we have set  $k = 100$  and  $m = 10$ . On these two queries, STACCATO falls in the middle on both recall and performance.



**Figure 4:** Recall - Runtime tradeoff for a keyword query (Query 1) and a regular expression query (Query 2). The parameters are: number of chunks ( $m$ ) = 10, number of paths per chunk ( $k$ ) = 100, and number of answers queried for ( $NumAns$ ) = 100.

### 3.2 Extensions and Analysis

To understand the formal underpinning of our approach, we perform a theoretical analysis. Informally, the first question is: “in what sense is choosing the  $k$ -MAP the best approximation for each chunk in our algorithm?” The second question we ask is to justify our restriction to SFAs as opposed to richer graphical structures in our approximation. We show that  $k$ -MAP in each chunk is no longer the best approximation and that there is likely no simple algorithm (as an underlying problem is NP-complete.)

We formally define the goal of our algorithms. Recall that an SFA  $S$  on  $\Sigma$  represents a probability distribution  $\text{Pr}_S : \Sigma^* \rightarrow [0, 1]$ . Given a set  $X \subseteq \Sigma^*$ , define  $\text{Pr}_S[X] = \sum_{x \in X} \text{Pr}_S[x]$ . All the approximations that we consider emit a subset of strings from the original model. Given an approximation scheme  $\alpha$ , we denote by  $\text{Emit}(\alpha)$  the set of strings that are emitted (retained) by that scheme. All other things being equal, we prefer a scheme  $\alpha$  to  $\alpha'$  whenever

$$\text{Pr}_S[\text{Emit}(\alpha)] \geq \text{Pr}_S[\text{Emit}(\alpha')]$$

That is,  $\alpha$  retains more probability mass than  $\alpha'$ . The formal basis for this choice is a standard statistical measure called the *Kullback-Leibler Divergence* [15], between the original and the approximate probability distributions. In the full version [34], we show that this divergence is lower (which means the approximate distribution is more similar to the original distribution) if the approximation satisfies the above inequality. In other words, a better approximation retains more of the high-probability strings.

We now describe our two main theoretical results. First for SFAs, STACCATO’s approach to choosing the  $k$  highest probability strings in each chunk is optimal. For richer structures than SFAs, finding the optimal approximation is intractable (even if we are given the chunk structure, described below). Showing the first statement is straightforward, while the result about richer structures is more challenging.

**Optimality of  $k$ -MAP for SFAs.** Given a generalized SFA  $S = (V, \delta)$ . Fix  $k \geq 1$ . Let  $\mathbf{S}_{[k]}$  be the set of all SFAs  $(V, \delta')$  that arise from picking  $k$  strings on each edge of  $S$  to store.

That is, for any pair of nodes  $x, y \in V$  the set of strings with non-zero probability has size smaller than  $k$ :

$$|\{\sigma \in \Sigma^* \mid \delta'((x, y), \sigma) > 0\}| \leq k$$

Let  $S_k$  denote an SFA that for each pair  $(x, y) \in V$  chooses the highest probability strings in the model (breaking ties arbitrarily). Then,

PROPOSITION 3.1. *For any  $S' \in \mathbf{S}_{[k]}$ , we have:*

$$\Pr_S[\text{Emit}(S_k)] \geq \Pr_S[\text{Emit}(S')]$$

Since  $S_k$  is selected by STACCATO, we view this as formal justification for STACCATO’s choice.

**Richer Structural Approximation.** We now ask a follow-up question: “If we allow more general partitions (rather than collapsing edges), is  $k$ -MAP still optimal?” To make this precise, we consider a partition of the underlying edges of the SFA into connected components (call that partition  $\Phi$ ). Keeping with our early terminology, an element of the partition is called a *chunk*. In each chunk, we select at most  $k$  strings (corresponding to labeled paths through the chunk). Let  $\alpha : \Phi \times \Sigma^* \rightarrow \{0, 1\}$  be an indicator function such that  $\alpha(\phi, \sigma) = 1$  only if in chunk  $\phi$  we choose string  $\sigma$ . For any  $k \geq 1$ , let  $A_k$  denote the set of all such  $\alpha$ s that picks at most  $k$  strings from each chunk, i.e., for any  $\phi \in \Phi$  we have  $|\{\sigma \in \Sigma^* \mid \alpha(\phi, \sigma) > 0\}| \leq k$ . Let  $\text{Emit}(\alpha)$  be the set of strings emitted by this representation with non-zero probability (all strings that can be created from concatenating paths in the model).

Following the intuition from the SFA case described above, the best  $\alpha$  would select the  $k$ -highest probability strings in each chunk. However, this is not the case. Moreover, we exhibit chunk structures, where finding the optimal choice of  $\alpha$  is NP-hard in the size of the structure. This makes it unlikely that there is any simple description of the optimal approximation.

THEOREM 3.1. *Fix  $k \geq 2$ . The following problem is NP-complete. Given as input  $(S, \Phi, \lambda)$  where  $S$  is an SFA,  $\Phi$  partitions the underlying graph of  $S$ , and  $\lambda \geq 0$ , determine if there exists an  $\alpha \in A_k$  satisfying  $\Pr[\text{Emit}(\alpha)] \geq \lambda$ .*

The above problem remains NP-complete if  $S$  is restricted to satisfy the unique path property and restricted to a binary alphabet. A direct consequence of this theorem is that finding the maximizer is at least NP-hard. We provide the proof of this theorem in the full version [34]. The proof includes a detailed outline of a reduction from a matrix multiplication-related problem that is known to be hard. The reduction is by a gadget construction that encodes matrix multiplication as SFAs. Each chunk has at most 2 nodes in either border (as opposed to an SFA which has a single start and final node). This is about the weakest violation of the SFA property that we can imagine, and suggests to us that the SFA property is critical for tractable approximations.

**Automated Construction of STACCATO.** Part of our goal is to allow knobs to trade recall for performance on a per application basis, but setting the correct values for  $m$  and  $k$  may be unintuitive for users. To reduce the burden on the user, we devise a simple parameter tuning heuristic that maximizes query performance, while achieving acceptable

recall. To measure recall, the user provides a set of labeled examples and representative queries. The user specifies a quality constraint (average recall for the set of queries) and a size constraint (storage space as percentage of the original dataset size). The goal is to find a pair of parameters  $(m, k)$  that satisfies both these constraints. We note that the size of the data is a function of  $(m, k)$  (see Table 1), which along with the size constraint helps us express  $k$  in terms of  $m$  (or vice versa). We empirically observed that for a fixed size, a smaller  $m$  usually yields faster query performance than a smaller  $k$ , which suggests that we need to minimize the value of  $m$  to maximize query performance. Our method works as follows: we pick a given value of  $m$ , then calculate the corresponding  $k$  that lies on the size constraint boundary. Given the resulting  $(m, k)$  pair, we compute the STACCATO approximation of the dataset and estimate the average recall. This problem is now a one-dimensional search problem: our goal is to find the smallest  $m$  that satisfies the recall constraint. We solve this using essentially a binary search. If infeasible, the user relaxes one of the constraints and repeats the above method. We experimentally validated this tuning method and compared it with an exhaustive search on the parameter space. The results are discussed in Section 5.5.

## 4. INVERTED INDEXING

To speedup keywords and anchored regex queries on standard ASCII text, a popular technique is to use standard inverted-indexing [14]. While indexing  $k$ -MAP data is pretty straightforward, the FullSFA is difficult. The reason is that the FullSFA encodes exponentially many strings in its length, and so indexing all strings for even a moderate-sized SFA is hopeless. Figure 5 shows the size of the index obtained (in number of *postings* [14], in our case line number item pairs) when we try to directly index the STACCATO text of a single SFA (one OCR line).

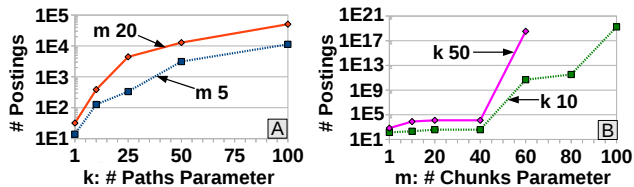


Figure 5: Number of postings (in logscale) from directly indexing one SFA. (A) Fix  $m$ , vary  $k$ . (B) Fix  $k$ , vary  $m$ . In (B), for  $k = 50$ , the number of postings overflows the 64-bit representation beyond  $m = 60$ .

Figure 5 shows an exponential blowup with  $m$  – which is not surprising as we store exponentially more paths with increasing  $m$ . Our observation is that many of these exponentially many terms are useless to applications. Thus, to extend the reach of indexing, we apply a standard technique. We use a dictionary of terms input by the user, and construct the index only for these terms [20]. These terms may be extracted from a known clean text corpus or from other sources like an English dictionary. Our construction algorithm builds a DFA from the dictionary of terms, and runs a slight modification of the SFA composition algorithm [27] with the data to find the start locations of all terms (details of the modification are in the full version [34]). The running time of the algorithm is linear in the size of the dictionary.

**Projection.** In traditional text processing, given the length of the keyword and the offset of a match, we can read only that small portion of the document to process the query. We extend this idea to STACCATO by finding a small portion of the SFA that is needed to answer the query – an operation that we call *projection*. Given a term  $t$  of length  $u$ , we obtain start locations of  $t$  from the postings. For each start location, we compute an (over)estimate of the nodes that we must process to obtain the term  $t$ . More precisely, we want the descendant nodes in the DAG that can be reached by a directed path from the start location that contains  $u$  or fewer edges (we find such nodes using a breadth-first search). This gives us a set of nodes that we must retrieve, which is often much smaller than the entire SFA.

We empirically show that even a simple indexing scheme as above can be used by STACCATO to speedup keyword and anchored regular expression queries by over an order of magnitude versus a filescan-based approach. This validates our claim that indexing is possible for OCR transducers, and opens the possibility of adapting more advanced indexing techniques to improve the runtime speedups.

## 5. EXPERIMENTAL EVALUATION

We experimentally verify that the STACCATO approach can gracefully tradeoff between performance and quality. We also validate that our modifications to standard inverted indexing allow us to speedup query answering.

Dataset	No. of Pages	No. of SFAs	Size as:	
			SFAs	Text
Cong. Acts (CA)	38	1590	533MB	90kB
English Lit. (LT)	32	1211	524MB	78kB
DB Papers (DB)	16	627	359MB	54kB

**Table 2: Dataset Statistics.** Each SFA represents one line of a scanned page.

**Datasets Used.** We use three real-world datasets from domains where document digitization is growing. Congress Acts (CA) is a set of scans of acts of the U.S. Congress, obtained from The Hathi Trust [9]. English Literature (LT) is a set of scans of an English literature book, obtained from the JSTOR Archive [10]. Database Papers (DB) is a set of papers that we scanned ourselves to simulate a setting where an organization would scan documents for in-house usage. All the scan images were converted to SFAs using the OCRopus tool [8]. Each line of each document is represented by one SFA. We created a manual ground truth for these documents. The relevant statistics of these datasets are shown in Table 2. In order to study the scalability of the approaches on much larger datasets, we used a 100 GB dataset obtained from Google Books [7].

**Experimental Setup.** The three approaches were implemented in C++ using PostgreSQL 9.0.3. The current implementation is single threaded so as to assess the impact of the approximation. All experiments are run on Intel Core-2 E6600 machines with 2.4 GHz CPU, 4 GB RAM, running Linux 2.6.18-194. The runtimes are averaged over 7 runs. The notation for the parameters is summarized in Table 3. We set  $NumAns = 100$ , which is greater than the number of answers in the ground truth for all reported queries. If

Symbol	Description
$k$	# Paths Parameter ( $k$ -MAP, STACCATO)
$m$	# Chunks Parameter (STACCATO)
$NumAns$	# Answers queried for

**Table 3: Notations for Parameters**

STACCATO finds fewer matches than  $NumAns$ , it may return fewer answers.  $NumAns$  affects precision, and we do sensitivity analysis for  $NumAns$  in the full version [34].

### 5.1 Quality - Performance Tradeoff (Filescan)

We now present the detailed quality and performance results for queries run with a full filescan. The central technical claim of this paper is that STACCATO bridges the gap from the low-recall-but-fast MAP to the high-recall-but-slow FullSFA. To verify this claim, we measured the recall and performance of 21 queries on the three datasets. We formulated these queries based on our discussions with practitioners in companies and researchers in the social sciences who work with real-world OCR data. Table 4 presents a subset of these results (the rest are presented in the full version of this paper [34]).

Query	MAP	$k$ -MAP	FullSFA	STACCATO
<b>Precision/Recall</b>				
CA1	1.00/0.79	1.00/0.79	0.14/1.00	1.00/0.79
CA2	1.00/0.28	1.00/0.52	0.25/1.00	0.73/0.76
LT1	0.96/0.87	0.96/0.90	0.92/1.00	0.97/0.91
LT2	0.78/0.66	0.76/0.66	0.31/0.97	0.44/0.81
DB1	0.93/0.75	0.90/0.92	0.67/0.99	0.90/0.96
DB2	0.96/0.76	0.96/0.76	0.33/1.00	0.91/0.97
<b>Runtime (in seconds)</b>				
CA1	0.17	0.75	86.72	2.87
CA2	0.18	0.84	150.35	3.36
LT1	0.13	0.19	83.78	1.98
LT2	0.14	0.24	155.45	2.88
DB1	0.07	0.29	40.73	0.75
DB2	0.07	0.33	619.31	0.86

**Table 4: Recall and runtime results across datasets.** The keyword queries are – CA1: ‘*President*’, LT1: ‘*Brinkmann*’ and DB1: ‘*Trio*’. The regex queries are – CA2: ‘*U.S.C. 2\d\d\d*’, LT2: ‘*19\d\d, \d\d*’ and DB2: ‘*Sec(\x)\*\d*’. Here,  $\backslash x$  is any character and  $\backslash d$  is any digit. The number of ground truth matches are – CA1: 28, LT1: 92, DB1: 68, CA2: 55, LT2: 32 and DB2: 33. The parameter setting here is:  $k = 25$ ,  $m = 40$ ,  $NumAns = 100$ .

We classify the kinds of queries to *keywords* and *regular expressions*. The intuition is that keyword queries are likely to achieve higher recall on  $k$ -MAP compared to more complex queries that contain spaces, special characters, and wildcards. Table 4 presents the recall and runtime results for six queries – one keyword and one regular expression (regex) query per dataset. Table 4 confirms that indeed there are intermediate points in our approximation that have faster runtimes than FullSFA (even up to two orders of magnitude), while providing higher quality than  $k$ -MAP.

We would like the tradeoff of quality for performance to be smooth as we vary  $m$  and  $k$ . To validate that our approximation can support this, we present two queries, a keyword



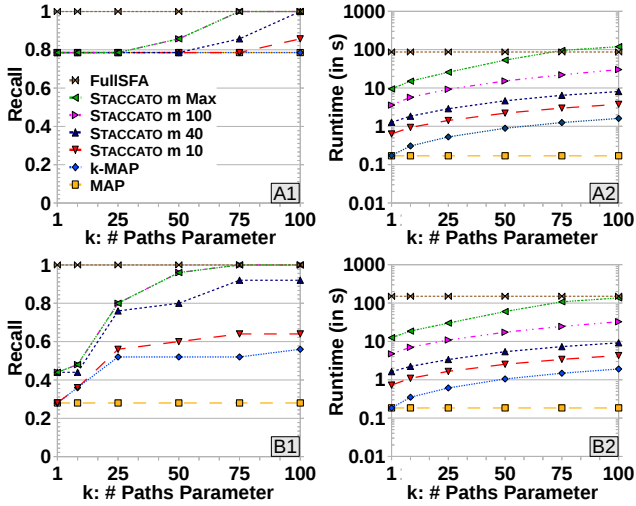


Figure 6: Recall and Runtime variations with  $k$ , for different values of  $m$ , on two queries: (A) ‘*President*’ (keyword), and (B) ‘*U.S.C. 2\d\d\d*’ (regex). The  $\d$  is short for  $(0|1|\dots|9)$ . The runtimes are in logscale.  $NumAns$  is set to 100. Recall that  $m$  is the number of chunks parameter and  $NumAns$  is the number of answers queried for.

and a regex, on the Congress Acts dataset (described below). To demonstrate this point, we vary  $k$  (the number of paths) for several values of  $m$  (the number of chunks) and plot the results in Figure 6. Given an SFA,  $m$  takes values from 1 to the number of the edges in the SFA (the latter being the nominal parameter setting ‘Max’). When  $m = 1$ , STACCATO is equivalent to  $k$ -MAP. Note that the state-of-the-art in our comparison is essentially the MAP approach ( $k$ -MAP with  $k = 1$ , or STACCATO with  $m = 1, k = 1$ ), which is what is employed by Google Books.

**Keyword Queries.** In Figures 6 (A1) and (A2), we see the recall and performance behavior of running a keyword query (here ‘*President*’) in STACCATO for various combinations of  $k$  and  $m$ . We observe that the recall of  $k$ -MAP is high (0.8) but not perfect and in (A2)  $k$ -MAP is efficient (0.1s) to answer the query. Further, as we increase  $k$  there is essentially no change in recall (the running time does increase by an order of magnitude). We verified that the reason is that the top- $k$  paths change in only a small set of locations – and so no new occurrences of the string ‘*President*’ are found. In contrast, the FullSFA approach achieves perfect recall, but it takes over 3 orders of magnitude longer to process the query. As we can see from the plots, for the STACCATO approach, the recall improves as we increase  $m$  – with corresponding slowdowns in query time. We believe that our approach is promising because of the gradual tradeoff of running time for quality. The fact that the  $k$ -MAP recall does not increase substantially with  $k$ , and does not manage to achieve the recall of FullSFA even for large  $k$  underscores the need for finer-grained partition, which is what STACCATO does.

**Regular Expressions.** Figures 6 (B1) and (B2) present the results for a more sophisticated regex query that looks for a congressional code (‘*U.S.C. 2\d\d\d*’) referenced in the text.

As the figure shows, this more sophisticated query has much lower recall for the MAP approach, and increases slowly with increasing  $k$ . Again, we see the same tradeoff that the FullSFA approach is orders of magnitude slower than  $k$ -MAP, but achieves perfect recall. Here, we see that the STACCATO approach does well: there are substantial (but smooth) jumps in quality as we increase  $k$  and  $m$ , going all the way from MAP to FullSFA. This suggests that more sophisticated queries benefit from our scheme more, which is an encouraging first step to enable applications to do rich analytics over such data.

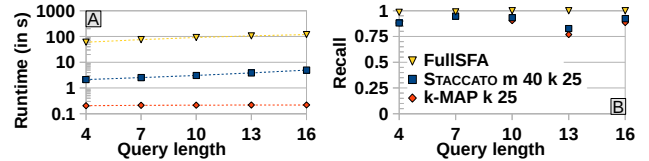


Figure 7: Impact of Query Length on (A) Runtime and (B) Recall.  $NumAns$ , the number of answers queried for, is set to 100.

**Query Efficiency.** To assess the impact of query length on recall and runtime, we plot the two for a set of keyword queries of increasing length in Figure 7. We observe that the runtimes increase polynomially but slowly for all the approaches, while no clear trends exist for the recall. We saw similar results with regular expression queries, and discuss the details in the full version [34].

We also studied the impact of  $m$  and  $k$  on precision (and F-1 score), and observed that the precision of STACCATO usually falls in between  $k$ -MAP and FullSFA (but F-1 of STACCATO can be better than both in some cases). Similar to the recall-runtime tradeoff, STACCATO also manages to gracefully tradeoff on precision and recall. Due to space constraints, these results are discussed in the full version [34].

## 5.2 Staccato Construction Time

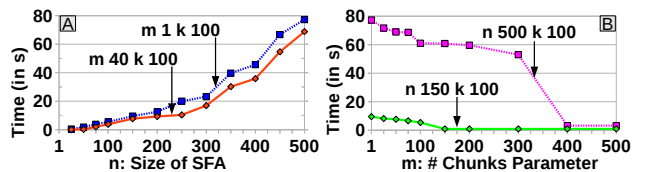


Figure 8: (A) Variation of Staccato approximation runtimes with the size of the SFA ( $n =$  number of nodes + edges) fixing  $m$  and  $k$ . (B) Sensitivity of the runtimes to  $m$ , fixing  $n$  and  $k$ . Recall that  $m$  is the number of chunks parameter and  $k$  is the number of paths parameter.

We now investigate the runtime of the STACCATO’s approximation algorithm. The runtime depends on the size of the input SFA data as well as  $m$  and  $k$ . We first fix  $m$  and  $k$ , then we plot the construction time for SFAs of varying size (number of nodes) from the CA dataset (Figure 8(A)). Overall, we can see that the algorithm runs efficiently – even in our unoptimized implementation. As this is an offline process, speed may not be critical for some applications. Also,

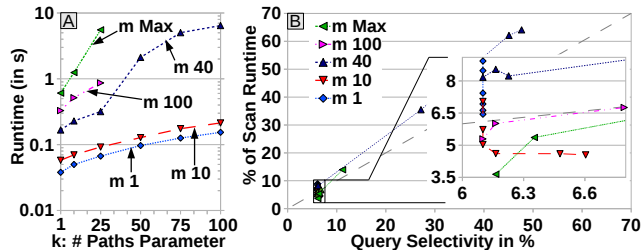
this computation is embarrassingly parallel (across SFAs). We used Condor [2] to run the STACCATO construction on all the SFAs in the three datasets, for all of the above parameters. This process completed in approximately 11 hours.

To study the sensitivity of the construction time to  $m$ , we select a fixed SFA from the CA dataset (Figure 8(B)). When  $m \geq |E|$ , the algorithm picks each transition as a block, and terminates. But when  $m = 300 < |E|$ , the algorithm computes several candidate merges, leading to a sudden spike in the runtime. There onwards, the runtime varies almost linearly with decreasing  $m$ . However, there are some spikes in the middle. We verified that the spikes arise since the ‘FindMinSFA’ operation has to fix merged chunks not satisfying the SFA property, thus causing the variation to be less smooth. We also verified that the runtime was linear in  $k$ , fixing the SFA and  $m$  (see full version [34]). In general, a linear runtime in  $k$  is not guaranteed since the chunk structure obtained during merging may not be similar across  $k$ , for a given SFA and  $m$ .

### 5.3 Inverted Indexing

We now verify that standard inverted indexing can be made to work on SFAs. We implement the index as a relational table with a B+-tree on top of it. More efficient inverted indexing implementations are possible, and so our results are an upperbound on indexing performance. However, this prototype serves to illustrate our main technical point that indexing is possible for such data.

A dictionary of about 60,000 terms from a freely available dictionary [4] was converted to a prefix-trie automaton, and used for index construction. While parsing the query, we ascertain if the given regex contains a left-anchor term. If so, we look up the anchor in the index to obtain the postings, and retrieve the data to employ query processing on them.



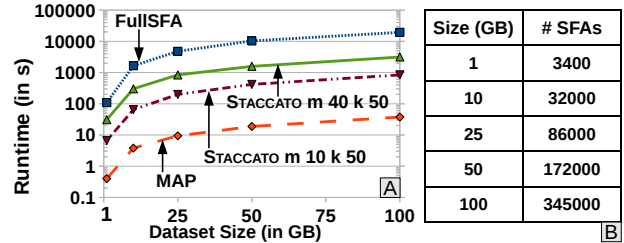
**Figure 9: (A) Total Runtimes, and (B) Fractional Runtimes Vs Selectivity for the query ‘Public Law (8|9)\d’, using the inverted index with the left anchor term ‘public’. Runtimes are in logscale. Recall that  $m$  is the number of chunks parameter.**

Figure 9 shows the results for a fixed length left anchored regex on the CA data set that is anchored by a word in the dictionary (here, ‘Public’). We omit some combinations ( $m = 100, \text{Max}$  and  $k = 50, 75, 100$ ) since their indexes had nearly 100% selectivity for all queries that we consider, rendering them useless. The first plot shows the sensitivity of the total runtimes to  $m$  and  $k$ . Mostly, there is a linear trend with  $k$ , except for a spike at  $m = 40, k = 50$ . To understand this behavior, we plot the runtime, as a percentage of the filescan runtime, against the selectivity of the term in the index. Ideally, the points should lie on the  $Y = X$  line, or slightly above it. For the lowest values of  $m$  and  $k$ , the

relative speedup is slightly lowered by the index lookup overhead. But as  $k$  increases, the query processing dominates, and hence the speedup improves, though selectivity changes only slightly. For higher  $m$ , the projection overhead lowers the speedup, and as  $k$  goes up, the selectivity shoots up, increasing the runtime. Overall, we see that dictionary-based indexing provides substantial speedups in many cases.

### 5.4 Scalability

To understand the feasibility of our approaches on larger amounts of data, we now study how the runtimes scale with increasing dataset sizes. We use a set of 8 scanned books from Google Books [7] and use OCRopus to obtain the SFAs. The total size of the SFA dataset is 100 GB.



**Figure 10: (A) Filescan runtimes (logscale) against the dataset size for MAP, FullSFA and Staccato with two parameter settings. (B) Number of SFAs in the respective datasets.**

Figure 10 shows the scalability results for a regex query. The filescans for FullSFA, MAP and STACCATO all scale linearly in the dataset size. Overall, the filescan runtimes are in the order a few hours for FullSFA. The runtimes are one to two orders of magnitude lower for STACCATO, depending on the parameters, and about three orders of magnitude lower for MAP. We also verified that indexing over this data provides further speedup (subject to query selectivity) as shown before. One can speedup query answering in all of the approaches by partitioning the dataset across multiple machines (or even using multiple disks). Thus, to scale to much larger corpora (say, millions of books), we plan to investigate the use of parallel data processing frameworks to attack this problem.

### 5.5 Automated Parameter Tuning

We now empirically demonstrate the parameter tuning method on a labeled set of 1590 SFAs (from the CA dataset), and a set of 5 queries (both keywords and regular expressions). The size constraint is chosen as 10% and the recall constraint is chosen as 0.9. We use increments of 5 for both  $m$  and  $k$ . Based on the tuning method described in Section 3.2, we obtain the following size equation:  $20mk + 58k = 45540$ , and the resultant parameter estimates of  $m = 45, k = 45$ , with a recall of 0.91. We then performed an exhaustive search on the parameter space to obtain the optimal values subject to the same constraints. Figure 11 shows the surface plots of the size and the recall obtained by varying  $m$  and  $k$ . The optimal values obtained are:  $m = 35, k = 80$ , again with a recall of 0.91. The difference in the parameter values arises primarily because the tuning method overestimated the size at this location. Nevertheless, we see that the tuning method provides parameter estimates satisfying the user requirements.

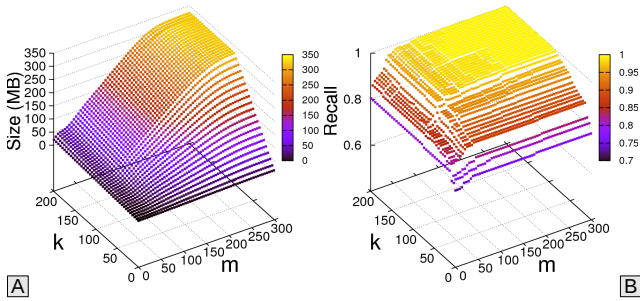


Figure 11: 3-D plots showing the variation of (A) the size of the approximated dataset (in MB), and (B) the average recall obtained. Recall that  $m$  is the number of chunks parameter and  $k$  is the number of paths parameter.

## 6. RELATED WORK

Transducers are widely used in the OCR and speech communities [11, 39] and mature open-source tools exist to process in-memory transducers [12]. For example we use a popular open-source tool, OCRopus [8], from Google Books that provides well-trained language models and outputs transducers. See Mohri et al. [39] for a discussion of why transducers are well-suited to represent the uncertainty for OCR. In the same work, Mohri et al. also describe speech data. We experimented with speech data, but we were hampered by the lack of high quality open-source speech recognizer toolkits. Using the available toolkits, we found that the language quality from open source speech recognizers is substantially below commercial quality.

The Lahar system [37, 43] manages Hidden Markov Models (HMMs) as *Markovian streams* inside an RDBMS and allows querying them with SQL-like semantics. In contrast to an HMM [42] that requires that all strings be of the same length, transducers are able to encode strings of different lengths. This is useful in OCR, since identifying spaces between words is difficult, and this uncertainty is captured by the branching in the SFA [39]. Our work drew inspiration from the empirical study of work of approximation trade-offs from Letchner et al. [37]. Directly relevant to this work is the recent theoretical results of Kimelfeld and Ré [32], who studied the problem of evaluating transducers as queries over uncertain sequence data modeled using Hidden Markov Models [42, 43]. STACCATO represents both the data and query by transducers which simplifies the engineering of our system.

Transducers are a graphical representation of probability models which makes them related to graphical models. Graphical models have been a hot topic in the database research community. Kanagal et al. [30] handle general graphical models. Wang et al. [49] also process Conditional Random Fields (CRFs) [35]. Though transducers can be viewed as a specialized directed graphical model, the primary focus of our work here is on the application of transducers to OCR in the domain of content management and the approximations that are critical to achieve good performance. However, our work is similar in spirit to these in that we too want to enable SQL-like querying of probabilistic OCR data inside an RDBMS.

Probabilistic graphical models have been successfully applied to various kinds of sequential data including OCR [17], RFID [43], speech [38], etc. Various models have been studied in both the machine learning and data management communities [21, 29, 30, 43, 49].

Many approximation schemes for probabilistic models have been studied [28, 37]. We built on the technique  $k$ -MAP [1], which is particularly relevant to us. Essentially, the idea is to infer the top  $k$  most likely results from the model and keep only those around. Another popular type of approximation is based on *mean-field theory*, where the intuition is that we replace complex dependencies (say in a graphical model) with their average (in some sense) [48]. Both mean-field theory and our approach share a common formal framework: minimizing KL-divergence. For a good overview of various probabilistic graphical models, approximation and inference techniques, we refer the reader to the excellent book by Wainwright and Jordan [48].

Gupta and Sarawagi [25] devise efficient approximation schemes to represent the outputs of a CRF, viz., labeled segmentations of text, in a probabilistic database. They partition the space of segmentations (i.e., the outputs) using boolean constraints on the output segment labels, and then structurally merge the partitions to a pre-defined count using Expectation Maximization, without any enumeration. Thus, their final partitions are disjoint sets of full-row outputs (‘horizontally’ partitioned). Both their approach and STACCATO use KL-divergence to measure the goodness of approximation. However, STACCATO is different in that we partition the underlying structure of the model (‘vertically’ partitioned). They also consider soft-partitioning approaches to overcome the limitations of disjoint partitioning. It is interesting future work to adapt such ideas for our problem, and compare with STACCATO’s approach.

Probabilistic databases have been studied in several recent projects (e.g., ORION [18], Trio [45], MystiQ [22], Sprout [41], and MayBMS [13]). Our work is complementary to these efforts: the queries we consider can produce probabilistic data that can be ingested by many of the above systems, while the above systems focus on querying restricted models (e.g., U-Relations or BIDs). We also use model-based views [23] to expose the results of query-time inference over the OCR transducers to applications.

The OCR, speech and IR communities have explored error correction techniques as well as approximate retrieval schemes [16, 26, 40]. However, prior work primarily focus on keyword search over plain-text transcriptions. STACCATO can benefit from these approaches and is orthogonal to our goal of integrating OCR data into an RDBMS. In contrast, we advocate retaining the uncertainty in the transcription.

Many authors have explored indexing techniques for probabilistic data [31, 33, 36, 47]. Letchner et al. [36] design new indexes for RFID data stored in an RDBMS as Markovian streams. Kanagal et al. [31] consider indexing correlated probabilistic streams using tree partitioning algorithms and describe a new technique called *shortcut potentials* to speedup query answering. Kimura et al. [33] propose a new *uncertain primary index* that clusters heap files according to uncertain attributes. Singh et al. [47] consider indexing categorical data and propose an R-tree based index as well as a probabilistic inverted index. Our work focuses on the challenges that content models like OCR raise for integrating indexing with an RDBMS.

## 7. CONCLUSION AND FUTURE WORK

We present our prototype system, STACCATO, that integrates a probabilistic model for OCR into an RDBMS. We demonstrated that it is possible to devise an approximation scheme that trades query runtime performance for result quality (in particular, increased recall). The technical contributions are a novel approximation scheme and a formal analysis of this scheme. Additionally, we showed how to adapt standard text-indexing schemes to OCR data, while retaining more answers.

Our future work is in two main directions. Firstly, we aim to extend STACCATO to handle larger data sets and more sophisticated querying (e.g., using aggregation with a probabilistic RDBMS, sophisticated indexing, parallel processing etc.). Secondly, we aim to extend our techniques to more types of content-management data such as speech transcription data. Interestingly, transducers provide a unifying formal framework for both transcription processes. Our initial experiments with speech data suggest that similar approximations techniques may be useful. This direction is particularly exciting to us: it is a first step towards unifying RDBMS and content-management systems, two multibillion dollar industries.

## 8. REFERENCES

- [1] A Brief Introduction to Graphical Models and Bayesian Networks. <http://www.cs.ubc.ca/~murphyk/Bayes/bayes.html>.
- [2] Condor high-throughput computing system. <http://www.cs.wisc.edu/condor/>.
- [3] Content Management Systems. <http://www.cmswire.com/>.
- [4] Corncob List. <http://www.mieliestronk.com/wordlist.html>.
- [5] Digital humanities by UW's Prof. Witmore. <http://winedarksea.org>.
- [6] ExperVision Inc. <http://www.expervision.com/>.
- [7] Google Books. <http://books.google.com/>.
- [8] OCRopus open source OCR system. <http://code.google.com/p/ocropus>.
- [9] The Hathi Trust. <http://www.hathitrust.org/>.
- [10] The JSTOR Archive. <http://www.jstor.org/>.
- [11] C. Allauzen, M. Mohri, and M. Saraclar. General indexation of weighted automata - application to spoken utterance retrieval. In *Workshop on Interdisciplinary Approaches to Speech Indexing and Retrieval (HLT/NAACL)*, pages 33–40, 2004.
- [12] C. Allauzen, M. Riley, J. Schalkwyk, W. Skut, and M. Mohri. Openfst: A general and efficient weighted finite-state transducer library. In *CJAA*, pages 11–23, 2007.
- [13] L. Antova, C. Koch, and D. Olteanu. Maybms: Managing incomplete information with probabilistic world-set decompositions. In *ICDE*, pages 1479–1480, 2007.
- [14] R. A. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval*. Addison-Wesley Longman Publishing Co., Inc., 1999.
- [15] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag New York, Inc., 2006.
- [16] J. Callan, W. B. Croft, and S. M. Harding. The inquiry retrieval system. In *DEXA*, pages 78–83, 1992.
- [17] M. Y. Chen, A. Kundu, and J. Zhou. Off-line handwritten word recognition using a hidden markov model type stochastic network. *Pattern Anal. Mach. Intell.*, 16:481–496, May 1994.
- [18] R. Cheng, D. V. Kalashnikov, and S. Prabhakar. Evaluating probabilistic queries over imprecise data. In *SIGMOD*, pages 551–562, 2003.
- [19] J. Cho and S. Rajagopalan. A fast regular expression indexing engine. In *ICDE*, pages 419–430, 2001.
- [20] R. Cole, L.-A. Gottlieb, and M. Lewenstein. Dictionary matching and indexing with errors and don't cares. In *STOC*, pages 91–100, 2004.
- [21] R. G. Cowell, A. P. Dawid, S. L. Lauritzen, and D. J. Spiegelhalter. *Probabilistic Networks and Expert Systems: Exact Computational Methods for Bayesian Networks*. Springer, 2007.
- [22] N. N. Dalvi and D. Suciu. Efficient query evaluation on probabilistic databases. In *VLDB*, pages 864–875, 2004.
- [23] A. Deshpande and S. Madden. Mauvedb: supporting model-based user views in database systems. In *SIGMOD*, pages 73–84, 2006.
- [24] J. Forney, G. D. The viterbi algorithm. *Proc. IEEE*, 61:268–278, 1973.
- [25] R. Gupta and S. Sarawagi. Creating probabilistic databases from information extraction models. In *VLDB*, pages 965–976, 2006.
- [26] S. Harding, W. B. Croft, and C. Weir. Probabilistic retrieval of ocr degraded text using n-grams. In *ECDDL*, pages 345–359, 1997.
- [27] J. E. Hopcroft, R. Motwani, and J. D. Ullman. *Introduction to Automata Theory, Languages, and Computation (3rd Edition)*. Addison-Wesley Longman Publishing Co., Inc., 2006.
- [28] F. Jensen and S. Andersen. Approx. in bayesian belief universes for knowledge-based systems. In *UAI*, pages 162–169, 1990.
- [29] M. I. Jordan. *Learning in graphical models*. MIT Press, 1999.
- [30] B. Kanagal and A. Deshpande. Efficient query evaluation over temporally correlated probabilistic streams. In *ICDE*, pages 1315–1318, 2009.
- [31] B. Kanagal and A. Deshpande. Indexing correlated probabilistic databases. In *SIGMOD*, pages 455–468, 2009.
- [32] B. Kimelfeld and C. Ré. Transducing markov sequences. In *PODS*, pages 15–26, 2010.
- [33] H. Kimura, S. Madden, and S. B. Zdonik. Upi: A primary index for uncertain databases. *PVLDB*, 3(1):630–637, 2010.
- [34] A. Kumar and C. Ré. Probabilistic management of ocr data using an rdbms. *UW-CS-Technical Report*, 2011. available from [http://www.cs.wisc.edu/hazy/staccato/papers/HazyOCR\\_TR.pdf](http://www.cs.wisc.edu/hazy/staccato/papers/HazyOCR_TR.pdf).
- [35] J. Lafferty. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. Morgan Kaufmann, 2001.
- [36] J. Letchner, C. Ré, M. Balazinska, and M. Philipose. Access methods for markovian streams. In *ICDE*, pages 246–257, 2009.
- [37] J. Letchner, C. Ré, M. Balazinska, and M. Philipose. Approximation trade-offs in markovian stream processing: An empirical study. In *ICDE*, pages 936–939, 2010.
- [38] S. E. Levinson, L. R. Rabiner, and M. M. Sondhi. An introduction to the application of the theory of probabilistic functions of a markov process to automatic speech recognition. *Bell Systems Technical Journal*, 62:1035–1074, 1983.
- [39] M. Mohri. Finite-state transducers in language and speech processing. *Computational Linguistics*, 23(2):269–311, 1997.
- [40] S. Mori, H. Nishida, and H. Yamada. *Optical character recognition*. John Wiley & Sons, Inc., 1999.
- [41] D. Olteanu, J. Huang, and C. Koch. Sprout: Lazy vs. eager query plans for tuple-independent probabilistic databases. In *ICDE*, pages 640–651, 2009.
- [42] L. R. Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. In *Proc. of IEEE*, pages 257–286, 1989.
- [43] C. Ré, J. Letchner, M. Balazinska, and D. Suciu. Event queries on correlated probabilistic streams. In *SIGMOD*, pages 715–728, 2008.
- [44] C. Ré and D. Suciu. Approximate lineage for probabilistic databases. *PVLDB*, 1(1):797–808, 2008.
- [45] A. D. Sarma, O. Benjelloun, A. Halevy, and J. Widom. Working models for uncertain data. *ICDE*, pages 7–18, 2006.
- [46] A. D. Sarma, M. Theobald, and J. Widom. Exploiting lineage for confidence computation in uncertain and probabilistic databases. In *ICDE*, pages 1023–1032, 2008.
- [47] S. Singh, C. Mayfield, S. Prabhakar, R. Shah, and S. Hambrusch. Indexing uncertain categorical data. In *ICDE*, pages 616–625, 2007.
- [48] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends of Machine Learning*, 1, 2008.
- [49] D. Z. Wang, E. Michelakis, M. N. Garofalakis, and J. M. Hellerstein. Bayesstore: managing large, uncertain data repositories with probabilistic graphical models. *PVLDB*, 1(1):340–351, 2008.
- [50] J. Widom. Trio: A system for integrated management of data, accuracy, and lineage. In *CIDR*, pages 262–276, 2005.
- [51] J. Y. Yen. Finding the k shortest loopless paths in a network. In *Management Science*, 1971.
- [52] J. Zobel, A. Moffat, and R. Sacks-davis. An efficient indexing technique for full-text database systems. In *VLDB*, 1992.
- [53] A. Zymnis, S. Boyd, and D. Gorinevsky. Relaxed maximum a posteriori fault identification. *Signal Process.*, 89, June 2009.