

Query Evaluation Over
SLP-Represented Document Databases
With Complex Document Editing

Markus L. Schmid, Nicole Schweikardt

HU Berlin, Germany

PODS 2022

The Data Model

Document Databases

- ▶ Finite alphabet $\Sigma = \{a, b, c, \dots\}$.
- ▶ *Documents*: strings over Σ , e. g., $D = \text{abaacbca}$.
- ▶ *Document Databases*: Sets of documents, e. g.,

$$\begin{aligned} \text{DDB} &= \{D_1, D_2, D_3\} \\ &= \{\text{ababbcabca}, \text{bcabcaabbca}, \text{ababbca}\}. \end{aligned}$$

Information Extraction

Evaluation Task

Given a document database DDB,
evaluate a given query q over some $D \in \text{DDB}$.

Information Extraction

Evaluation Task

Given a document database DDB,
evaluate a given query q over some $D \in \text{DDB}$.

Prominent choice as query class:

Document spanners

Information Extraction

Evaluation Task

Given a document database DDB,
evaluate a given query q over some $D \in \text{DDB}$.

Prominent choice as query class:

Document spanners

$D = \text{a b b a b c c a b c}$

Information Extraction

Evaluation Task

Given a document database DDB,
evaluate a given query q over some $D \in \text{DDB}$.

Prominent choice as query class:

Document spanners

$D = \text{abbabccabc}$

\implies

x	y	z
$[2, 5\rangle$	$[4, 7\rangle$	$[1, 10\rangle$
$[3, 5\rangle$	$[5, 8\rangle$	$[4, 7\rangle$
$[1, 3\rangle$	$[3, 10\rangle$	$[2, 4\rangle$
\vdots	\vdots	\vdots

Information Extraction

Evaluation Task

Given a document database DDB,
evaluate a given query q over some $D \in \text{DDB}$.

Prominent choice as query class:

Document spanners

$D = a \mathbf{b b a} b c c a b c$

\implies

x	y	z
$[2, 5\rangle$	$[4, 7\rangle$	$[1, 10\rangle$
$[3, 5\rangle$	$[5, 8\rangle$	$[4, 7\rangle$
$[1, 3\rangle$	$[3, 10\rangle$	$[2, 4\rangle$
\vdots	\vdots	\vdots

Information Extraction

Evaluation Task

Given a document database DDB,
evaluate a given query q over some $D \in \text{DDB}$.

Prominent choice as query class:

Document spanners

$D = \text{abb} \mathbf{abc} \text{cabc}$

\implies

x	y	z
$[2, 5\rangle$	$[4, 7\rangle$	$[1, 10\rangle$
$[3, 5\rangle$	$[5, 8\rangle$	$[4, 7\rangle$
$[1, 3\rangle$	$[3, 10\rangle$	$[2, 4\rangle$
\vdots	\vdots	\vdots

Information Extraction

Evaluation Task

Given a document database DDB,
evaluate a given query q over some $D \in \text{DDB}$.

Prominent choice as query class:

Document spanners

$D =$ `abbabccabc`

\implies

x	y	z
[2, 5)	[4, 7)	[1, 10)
[3, 5)	[5, 8)	[4, 7)
[1, 3)	[3, 10)	[2, 4)
⋮	⋮	⋮

Compressed Strings

- ▶ In practice, strings have many redundancies (\leadsto are highly compressible).
- ▶ Many practical (dictionary based) compression schemes for strings exist.
- ▶ Good compression rates, low (near linear-time) running times.

Algorithmics on Compressed Strings

String-Problem P

Input: A string w .

Task: Solve P for w .

Running time: $f(|w|)$.

Algorithmics on Compressed Strings

String-Problem P (compressed)

Input: A **compressed form** S of string w .

Task: Solve P for w (**without explicitly constructing** w).

Running time: $f(|S|)$.

Algorithmics on Compressed Strings

String-Problem P (compressed)

Input: A **compressed form** S of string w .

Task: Solve P for w (**without explicitly constructing w**).

Running time: $f(|S|)$.

uncompressed	vs.	compressed
$f(w)$	vs.	$g(S)$

Algorithmics on Compressed Strings

String-Problem P (compressed)

Input: A **compressed form** S of string w .

Task: Solve P for w (**without explicitly constructing w**).

Running time: $f(|S|)$.

uncompressed	vs.	compressed
--------------	-----	------------

$f(w)$	vs.	$g(S)$
----------	-----	----------

$O(w)$	vs.	$\text{poly}(S)$
----------	-----	--------------------

Algorithmics on Compressed Strings

String-Problem P (compressed)

Input: A **compressed form** S of string w .

Task: Solve P for w (**without explicitly constructing w**).

Running time: $f(|S|)$.

uncompressed	vs.	compressed
$f(w)$	vs.	$g(S)$
$O(w)$	vs.	$\text{poly}(S)$
$O(w)$	vs.	$\text{polylog}(w)$

Straight-Line Programs (SLPs)

Main idea: Represent a string w by a context-free grammar \mathcal{S} (in chomsky normal form) for language $\{w\}$.

Straight-Line Programs (SLPs)

Main idea: Represent a string w by a context-free grammar \mathcal{S} (in chomsky normal form) for language $\{w\}$.

Example

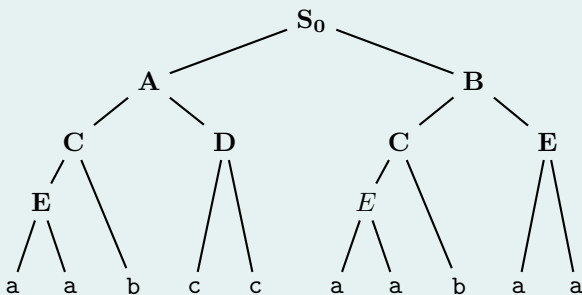
$$\begin{array}{lll} S_0 \rightarrow AB, & A \rightarrow CD, & B \rightarrow CE, \\ C \rightarrow Eb, & D \rightarrow cc, & E \rightarrow aa \end{array}$$

Straight-Line Programs (SLPs)

Main idea: Represent a string w by a context-free grammar \mathcal{S} (in chomsky normal form) for language $\{w\}$.

Example

$$\begin{array}{lll} S_0 \rightarrow AB, & A \rightarrow CD, & B \rightarrow CE, \\ C \rightarrow Eb, & D \rightarrow cc, & E \rightarrow aa \end{array}$$

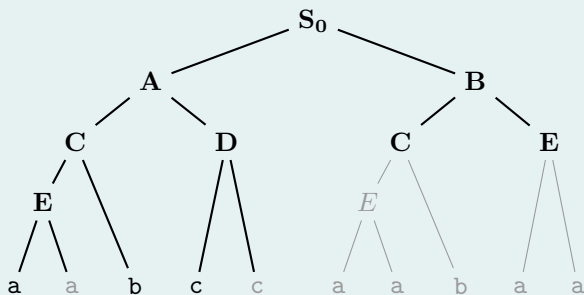


Straight-Line Programs (SLPs)

Main idea: Represent a string w by a context-free grammar \mathcal{S} (in chomsky normal form) for language $\{w\}$.

Example

$$\begin{array}{lll} S_0 \rightarrow AB, & A \rightarrow CD, & B \rightarrow CE, \\ C \rightarrow Eb, & D \rightarrow cc, & E \rightarrow aa \end{array}$$



Good Properties of SLPs

SLPs are intensely researched in TCS and many things are known:

- ▶ Exponential compression rates.
- ▶ SLPs are mathematically easy to handle (\Rightarrow good for theoretical considerations).
- ▶ High practical relevance (SLPs cover many practically applied dictionary-based compression schemes).
- ▶ Many approximations and heuristics exist that efficiently compute small SLPs.
- ▶ SLPs are suitable for algorithmics on compressed strings: comparison, pattern matching, membership in a regular language, retrieving subwords, etc.

Straight-Line Programs for Document Databases

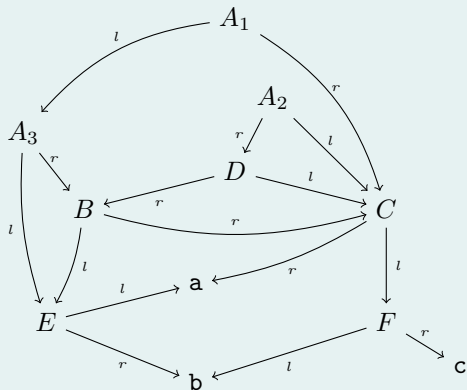
Example

$$\text{DDB} = \{\underbrace{\text{ababbcabca}}_{D_1}, \underbrace{\text{bcabcaabbca}}_{D_2}, \underbrace{\text{ababbca}}_{D_3}\}.$$

Straight-Line Programs for Document Databases

Example

DDB = { $\underbrace{\text{ababbcbca}}_{D_1}$, $\underbrace{\text{bcabcaabbca}}_{D_2}$, $\underbrace{\text{ababbca}}_{D_3}$ }.



$A_1 \rightarrow A_3 C$

$A_2 \rightarrow CD$

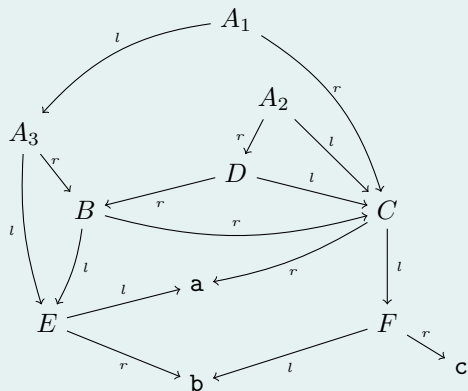
$C \rightarrow Fa$

\vdots

Straight-Line Programs for Document Databases

Example

$DDB = \{\underbrace{ababbcabca}_{D_1=\mathcal{D}(A_1)}, \underbrace{bcabcaabbca}_{D_2=\mathcal{D}(A_2)}, \underbrace{ababbca}_{D_3=\mathcal{D}(A_3)}\}.$



$A_1 \rightarrow A_3 C$
 $A_2 \rightarrow CD$
 $C \rightarrow Fa$
 \vdots

Spanner Evaluation Over SLP-Compressed Documents

What about database theory (i. e., information extraction)?

Theorem (S., Schweikardt, PODS'21)

Let D be a document represented by an SLP \mathcal{S} and let M be a document spanner. The set $\llbracket M \rrbracket(D)$ can be enumerated with preprocessing time $O(|\mathcal{S}|)$ and delay $O(\log |D|)$.¹

¹In data-complexity.

Spanner Evaluation Over SLP-Compressed Documents

What about database theory (i. e., information extraction)?

Theorem (S., Schweikardt, PODS'21)

Let D be a document represented by an SLP \mathcal{S} and let M be a document spanner. The set $\llbracket M \rrbracket(D)$ can be enumerated with preprocessing time $O(|\mathcal{S}|)$ and delay $O(\log |D|)$.¹

Remark

Delay $O(\log |D|)$ independent of the size $|\mathcal{S}|$ is possible since \mathcal{S} is **balanced** in the preprocessing.

(**Balanced** SLPs will play a central role in the following).

¹In data-complexity.

This Paper: Dynamic Setting

Basic Setting

Given:

- ▶ A document database DDB represented by an SLP \mathcal{S} .
- ▶ Data structures for enumerating spanners M_1, M_2, \dots on documents of DDB.

This Paper: Dynamic Setting

Basic Setting

Given:

- ▶ A document database DDB represented by an SLP \mathcal{S} .
- ▶ Data structures for enumerating spanners M_1, M_2, \dots on documents of DDB.

Task:

- ▶ “Update” DDB by directly updating \mathcal{S} .

This Paper: Dynamic Setting

Basic Setting

Given:

- ▶ A document database DDB represented by an SLP \mathcal{S} .
- ▶ Data structures for enumerating spanners M_1, M_2, \dots on documents of DDB.

Task:

- ▶ “Update” DDB by directly updating \mathcal{S} .

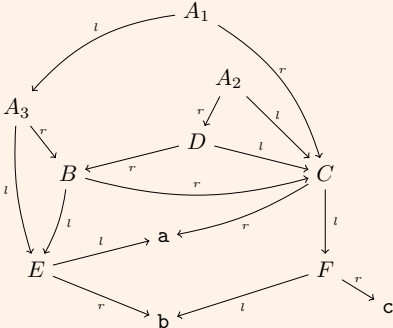
Additional conditions:

- ▶ Do not decompress \mathcal{S} .
- ▶ Also update the data structures for enumeration.
- ▶ Maintain the “**balancedness**” property of \mathcal{S} .

Balanced SLPs

Notations for Balancedness

Let A be a node of an SLP \mathcal{S} that represents a document database.



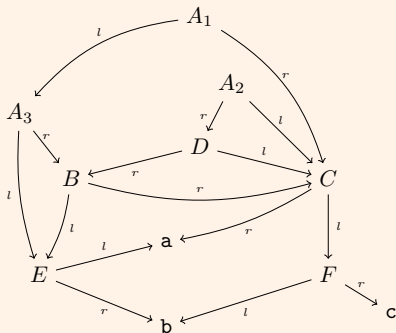
Balanced SLPs

Notations for Balancedness

Let A be a node of an SLP \mathcal{S} that represents a document database.

$\text{ord}(A)$ is the length of a longest path from A to a sink node.

(E. g., $\text{ord}(A_3) = 4$)



Balanced SLPs

Notations for Balancedness

Let A be a node of an SLP \mathcal{S} that represents a document database.

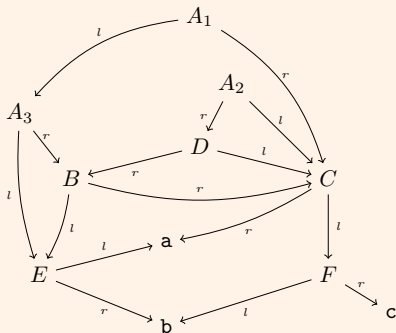
$\text{ord}(A)$ is the length of a longest path from A to a sink node.

(E. g., $\text{ord}(A_3) = 4$)

Let $A \rightarrow BC$.

$\text{bal}(A) = \text{ord}(B) - \text{ord}(C)$.

(E. g., $\text{bal}(A_3) = -2$)



Balanced SLPs

Notations for Balancedness

Let A be a node of an SLP \mathcal{S} that represents a document database.

$\text{ord}(A)$ is the length of a longest path from A to a sink node.

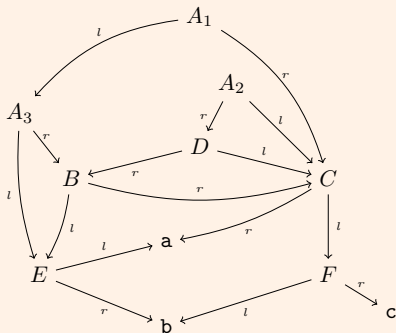
(E. g., $\text{ord}(A_3) = 4$)

Let $A \rightarrow BC$.

$\text{bal}(A) = \text{ord}(B) - \text{ord}(C)$.

(E. g., $\text{bal}(A_3) = -2$)

A is c -**shallow** (for a constant $c \in \mathbb{N}$) if $\text{ord}(A) \leq c \cdot \log |\mathcal{D}(A)|$.



Balanced SLPs

Notations for Balancedness

Let A be a node of an SLP \mathcal{S} that represents a document database.

$\text{ord}(A)$ is the length of a longest path from A to a sink node.

(E. g., $\text{ord}(A_3) = 4$)

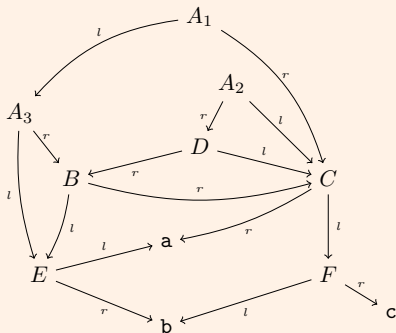
Let $A \rightarrow BC$.

$\text{bal}(A) = \text{ord}(B) - \text{ord}(C)$.

(E. g., $\text{bal}(A_3) = -2$)

A is **c -shallow** (for a constant $c \in \mathbb{N}$) if $\text{ord}(A) \leq c \cdot \log |\mathcal{D}(A)|$.

A is **balanced** if $\text{bal}(A) \in \{-1, 0, 1\}$.



Balanced SLPs

Notations for Balancedness

Let A be a node of an SLP \mathcal{S} that represents a document database.

$\text{ord}(A)$ is the length of a longest path from A to a sink node.

(E. g., $\text{ord}(A_3) = 4$)

Let $A \rightarrow BC$.

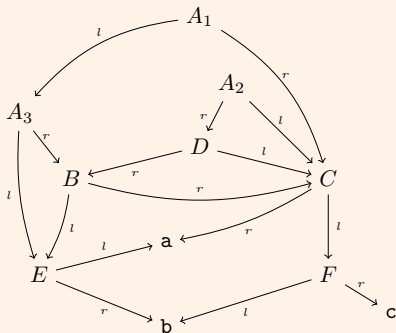
$\text{bal}(A) = \text{ord}(B) - \text{ord}(C)$.

(E. g., $\text{bal}(A_3) = -2$)

A is **c -shallow** (for a constant $c \in \mathbb{N}$) if $\text{ord}(A) \leq c \cdot \log |\mathcal{D}(A)|$.

A is **balanced** if $\text{bal}(A) \in \{-1, 0, 1\}$.

\mathcal{S} is **c -shallow** or **strongly balanced** if all nodes are **c -shallow** or **balanced**, respectively.



Balanced SLPs

Important Results About SLP-Balancing

A given SLP \mathcal{S} for a single document D ...

... can be made c -shallow in time $O(|\mathcal{S}|)$.

(Ganardi, Jez, Lohrey, JACM 2021)

Balanced SLPs

Important Results About SLP-Balancing

A given SLP \mathcal{S} for a single document D ...

... can be made c -shallow in time $O(|\mathcal{S}|)$.

(Ganardi, Jez, Lohrey, JACM 2021)

... can be made strongly balanced in time $O(|\mathcal{S}| \cdot \log |D|)$.

(Rytter, TCS 2003)

Balanced SLPs

Important Results About SLP-Balancing

A given SLP \mathcal{S} for a single document D ...

... can be made c -shallow in time $O(|\mathcal{S}|)$.

(Ganardi, Jez, Lohrey, JACM 2021)

... can be made strongly balanced in time $O(|\mathcal{S}| \cdot \log |D|)$.

(Rytter, TCS 2003)

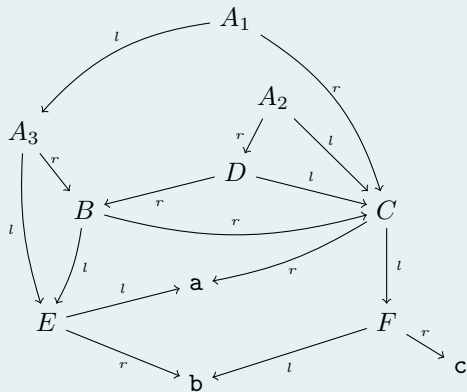
... can in general not be made strongly balanced without a size increase by a factor of $\Omega(\log |D|)$.

(Ganardi, ESA 2021)

Updating SLP-Represented Document Databases

Example

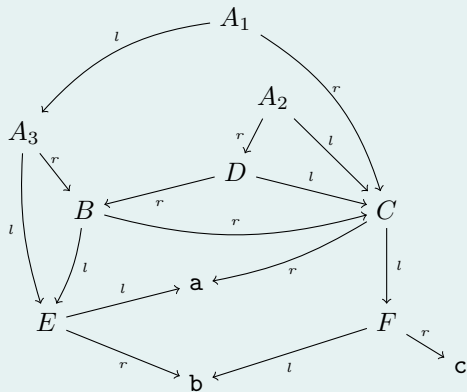
DDB = {ababbcabca, bcabcaabbca, ababbca}.



Updating SLP-Represented Document Databases

Example

DDB = {ababbcabca, bcabcaabbca, ababbca}.

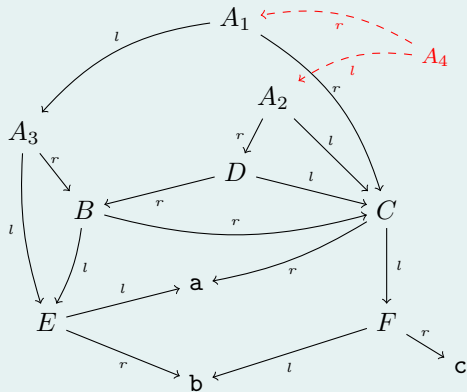


Add document
 $D_4 = D_2 D_1$

Updating SLP-Represented Document Databases

Example

DDB = {ababbcabca, bcabcaabbca, ababbca}.

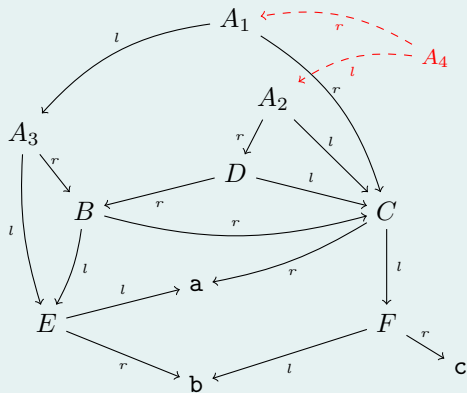


Add document
 $D_4 = D_2 D_1$

Updating SLP-Represented Document Databases

Example

DDB = {ababbcabca, bcabcaabbca, ababbca}.



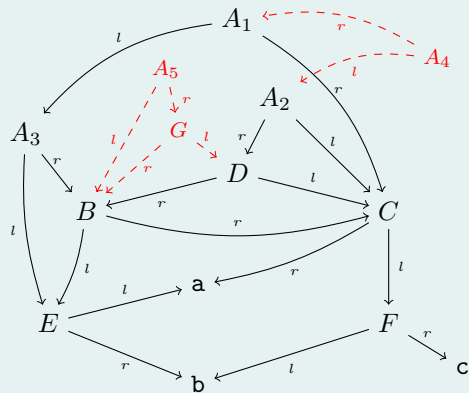
Add document
 $D_4 = D_2 D_1$

Add document
 $D_5 = \mathcal{D}(B) \mathcal{D}(D) \mathcal{D}(B)$

Updating SLP-Represented Document Databases

Example

DDB = {ababbcabca, bcabcaabbca, ababbca}.



Add document
 $D_4 = D_2 D_1$

Add document
 $D_5 = \mathcal{D}(B) \mathcal{D}(D) \mathcal{D}(B)$

Complex Document Editing

Notation

Let D be a document and $i, j \in \{1, 2, \dots, |D|\}$.
 $D[i..j]$ is D 's factor from position i to position j .

Complex Document Editing

Notation

Let D be a document and $i, j \in \{1, 2, \dots, |D|\}$.
 $D[i..j]$ is D 's factor from position i to position j .

Basic CDE-Algebra Operations

$$\begin{aligned}\text{concat}(D, D') &= D \cdot D', \\ \text{extract}(D, i, j) &= D[i..j], \\ \text{delete}(D, i, j) &= D[1..i-1] \cdot D[j+1..|D|], \\ \text{insert}(D, D', k) &= D[1..k-1] \cdot D' \cdot D[k..|D|], \\ \text{copy}(D, i, j, k) &= D[1..k-1] \cdot D[i..j] \cdot D[k..|D|].\end{aligned}$$

Complex Document Editing

Notation

Let D be a document and $i, j \in \{1, 2, \dots, |D|\}$.
 $D[i..j]$ is D 's factor from position i to position j .

Basic CDE-Algebra Operations

$$\begin{aligned}\text{concat}(D, D') &= D \cdot D', \\ \text{extract}(D, i, j) &= D[i..j], \\ \text{delete}(D, i, j) &= D[1..i-1] \cdot D[j+1..|D|], \\ \text{insert}(D, D', k) &= D[1..k-1] \cdot D' \cdot D[k..|D|], \\ \text{copy}(D, i, j, k) &= D[1..k-1] \cdot D[i..j] \cdot D[k..|D|].\end{aligned}$$

CDE-Expressions

Expressions over DDB's documents using the basic CDE-algebra operations.

E. g., $\varphi = \text{concat}(D_1, \text{insert}(D_3, \text{extract}(D_7, 5, 21), 12))$

Our Main Result

CDE Extension Theorem

Let DDB be represented by a strongly balanced SLP \mathcal{S} , let φ be a CDE-expression over DDB.

We can construct a strongly balanced SLP \mathcal{S}' for

$$\text{DDB} \cup \{\text{eval}(\varphi)\}$$

in time

$$O(|\varphi|^2 + |\varphi| \cdot \log(d_{\max})),$$

where $d_{\max} = \{|D| \mid D \in \text{DDB}\}$.

Our Main Result

Corollary

Let DDB be represented by a strongly balanced SLP \mathcal{S} .

Let \mathcal{M} be a class of regular spanners such that every $M \in \mathcal{M}$ can be enumerated on every $D \in \text{DDB}$ with delay $O(\log(|D|))$.

Given a CDE-expression φ over DDB, we can construct in time

$$O(|\mathcal{M}| \cdot (|\varphi|^2 + |\varphi| \cdot \log(d_{\max})))$$

a strongly balanced SLP \mathcal{S}' for

$$\text{DDB}' = \text{DDB} \cup \{\text{eval}(\varphi)\}$$

and data structures such that we can now enumerate every $M \in \mathcal{M}$ on every $D \in \text{DDB}'$ with delay $O(\log(|D|))$.

Proof Ideas

Proof Roadmap

- ▶ For each of the basic CDE-operation, prove a lemma that “handles” this operation.
(I. e., “For nodes B and C , create a node that derives $\text{concat}(\mathcal{D}(B), \mathcal{D}(C))$ ”, “For node B and i, j , create a node that derives $\text{extract}(\mathcal{D}(B), i, j)$ ”, etc.)

Proof Ideas

Proof Roadmap

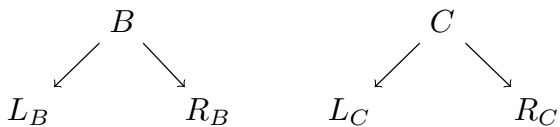
- ▶ For each of the basic CDE-operation, prove a lemma that “handles” this operation.
(I. e., “For nodes B and C , create a node that derives $\text{concat}(\mathcal{D}(B), \mathcal{D}(C))$ ”, “For node B and i, j , create a node that derives $\text{extract}(\mathcal{D}(B), i, j)$ ”, etc.)
- ▶ For a given CDE-expression φ , we can then apply these lemmas “bottom-up” (along φ 's syntax tree).

Proof Ideas

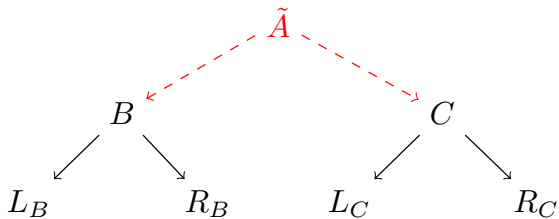
Proof Roadmap

- ▶ For each of the basic CDE-operation, prove a lemma that “handles” this operation.
(I. e., “For nodes B and C , create a node that derives $\text{concat}(\mathcal{D}(B), \mathcal{D}(C))$ ”, “For node B and i, j , create a node that derives $\text{extract}(\mathcal{D}(B), i, j)$ ”, etc.)
- ▶ For a given CDE-expression φ , we can then apply these lemmas “bottom-up” (along φ 's syntax tree).
- ▶ Lemmas for $\text{concat}(\cdot, \cdot)$ and $\text{extract}(\cdot, \cdot, \cdot)$ are sufficient (since all CDE-operations can be defined by applications of these two operations).

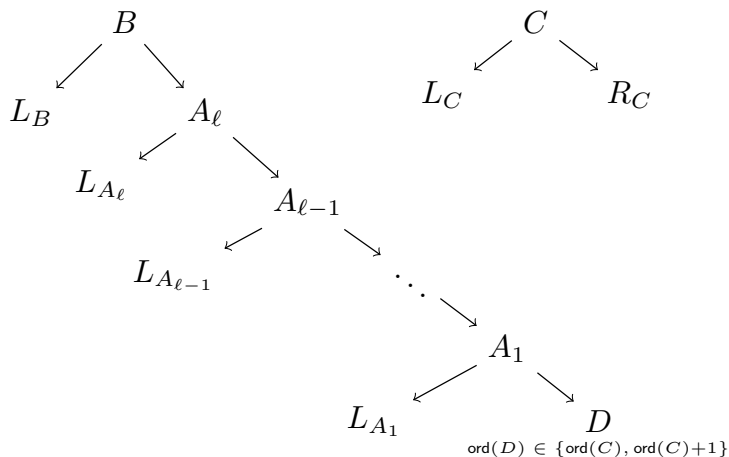
Proof Sketch for Concatenation Lemma



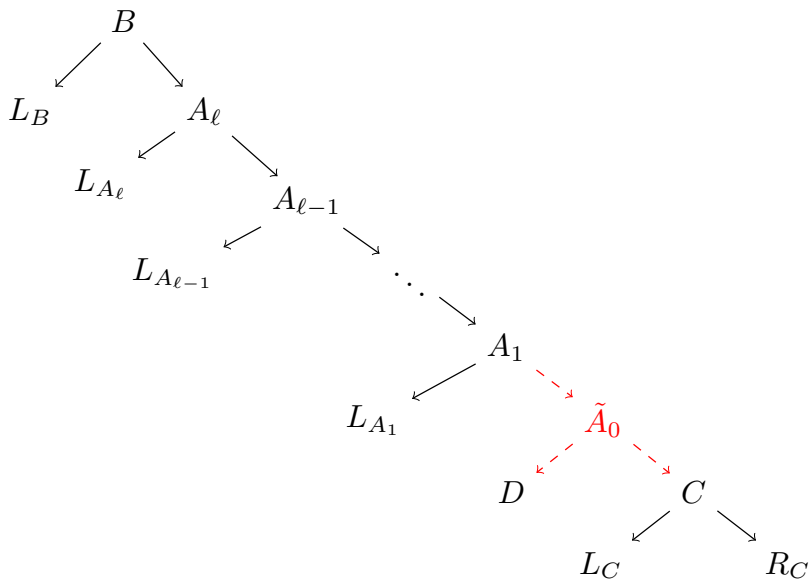
Proof Sketch for Concatenation Lemma



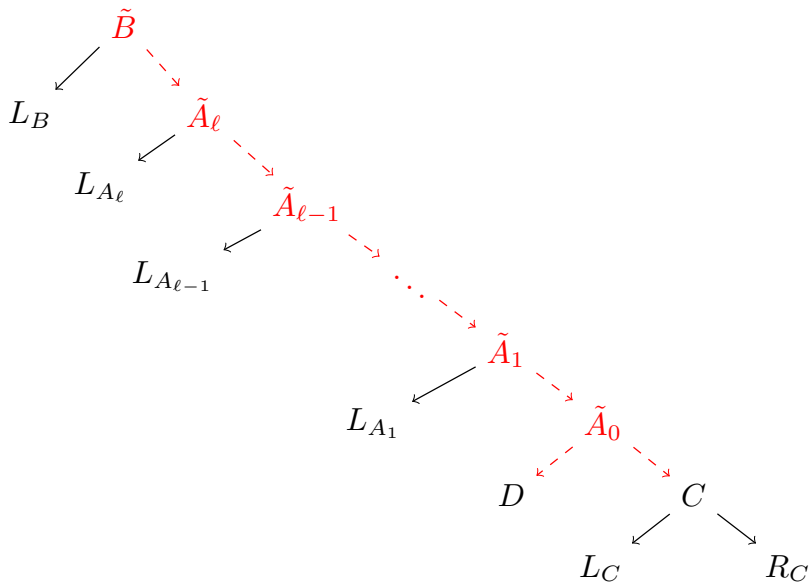
Proof Sketch for Concatenation



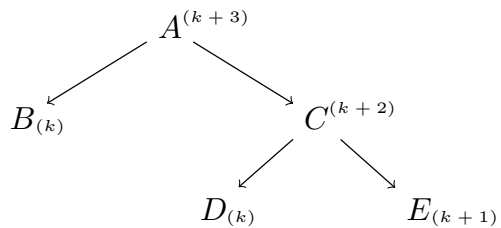
Proof Sketch for Concatenation



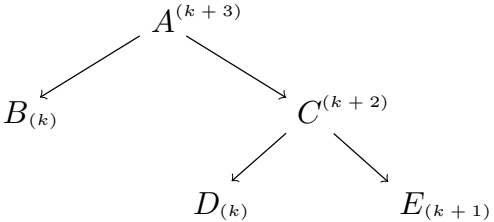
Proof Sketch for Concatenation



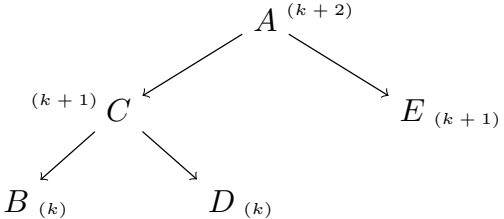
Rotations



Rotations



\Downarrow



Related Question Discussed in the Paper

- ▶ Adding a completely new document (in compressed form or plain text).
- ▶ Building an SLP-represented document database from scratch.
- ▶ Retrieving SLP-represented documents for the spans extracted by a spanner.
- ▶ Overall setting works for other evaluation problems like testing, non-emptiness, etc.
- ▶ Extends to general finite transducers as extractors.

Thank you very much for your attention.