

Spanner Evaluation over SLP-Compressed Documents

Markus L. Schmid, Nicole Schweikardt

HU Berlin, Germany

PODS 2021

Document Spanners

Document Spanners

- ▶ $\Sigma = \{a, b, c, \dots\}$ is an alphabet.
- ▶ $\mathcal{X} = \{x, y, z, \dots\}$ is a set of variables.
- ▶ D is a document over Σ .

Document Spanners

- ▶ $\Sigma = \{a, b, c, \dots\}$ is an alphabet.
- ▶ $\mathcal{X} = \{x, y, z, \dots\}$ is a set of variables.
- ▶ D is a document over Σ .

$D = a b b a b c c a b c$

Document Spanners

- ▶ $\Sigma = \{a, b, c, \dots\}$ is an alphabet.
- ▶ $\mathcal{X} = \{x, y, z, \dots\}$ is a set of variables.
- ▶ D is a document over Σ .

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

\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

Document Spanners

- ▶ $\Sigma = \{a, b, c, \dots\}$ is an alphabet.
- ▶ $\mathcal{X} = \{x, y, z, \dots\}$ is a set of variables.
- ▶ D is a document over Σ .

$D = a b b a b c c a b c$

\Rightarrow

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

Document Spanners

- ▶ $\Sigma = \{a, b, c, \dots\}$ is an alphabet.
- ▶ $\mathcal{X} = \{x, y, z, \dots\}$ is a set of variables.
- ▶ D is a document over Σ .

$D = abbabcabc$

\Rightarrow

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

Document Spanners

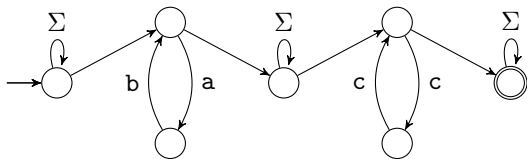
- ▶ $\Sigma = \{a, b, c, \dots\}$ is an alphabet.
- ▶ $\mathcal{X} = \{x, y, z, \dots\}$ is a set of variables.
- ▶ D is a document over Σ .

$D =$ `abbabccabc`

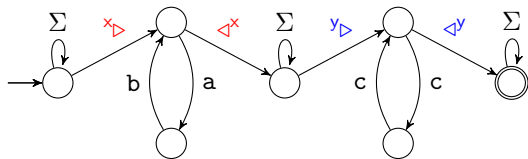
\implies

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

Regular (Document) Spanners



Regular (Document) Spanners

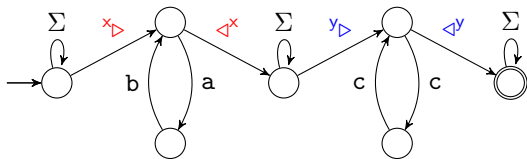


Meta-symbols for variables $x, y, \dots \in \mathcal{X}$:

$\triangleright^x \dots \triangleleft^x$ (start and end position of span extracted by x),

$\triangleright^y \dots \triangleleft^y$ (start and end position of span extracted by y).

Regular (Document) Spanners

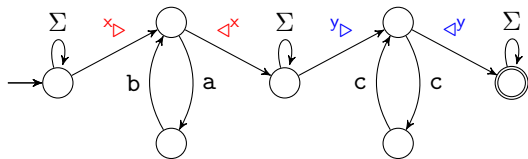


Meta-symbols for variables $x, y, \dots \in \mathcal{X}$:

$\triangleright^x \dots \triangleleft^x$ (start and end position of span extracted by x),

$\triangleright^y \dots \triangleleft^y$ (start and end position of span extracted by y).

Regular (Document) Spanners



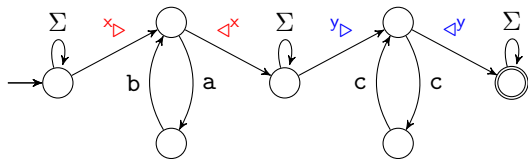
Meta-symbols for variables $x, y, \dots \in \mathcal{X}$:

$x \triangleright \dots \triangleleft x$ (start and end position of span extracted by x),

$y \triangleright \dots \triangleleft y$ (start and end position of span extracted by y).

$D = \text{abbababccca}$

Regular (Document) Spanners



Meta-symbols for variables $x, y, \dots \in \mathcal{X}$:

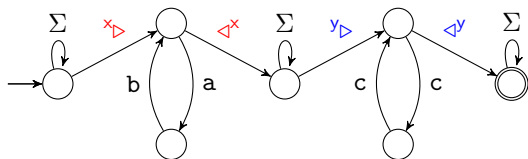
$\triangleright^x \dots \triangleleft^x$ (start and end position of span extracted by x),

$\triangleright^y \dots \triangleleft^y$ (start and end position of span extracted by y).

$D = \text{abbababccca}$

$\text{abb}^{\triangleright^x} \text{ab}^{\triangleleft^x} \text{ab}^{\triangleright^y} \text{cc}^{\triangleleft^y} \text{ca}$

Regular (Document) Spanners



Meta-symbols for variables $x, y, \dots \in \mathcal{X}$:

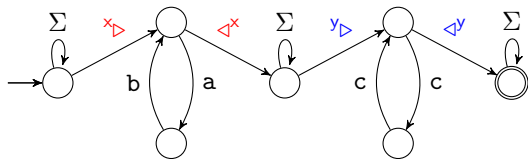
$x \triangleright \dots \triangleleft x$ (start and end position of span extracted by x),

$y \triangleright \dots \triangleleft y$ (start and end position of span extracted by y).

$D = \text{abbababccca}$

$\text{abb}x \triangleright \text{ab} \triangleleft x \text{ab}y \triangleright \text{cc} \triangleleft y \text{ca} \implies ([4, 6], [8, 10])$

Regular (Document) Spanners



Meta-symbols for variables $x, y, \dots \in \mathcal{X}$:

$x \triangleright \dots \triangleleft x$ (start and end position of span extracted by x),

$y \triangleright \dots \triangleleft y$ (start and end position of span extracted by y).

$D = \text{abbababccca}$

$\text{abb}x \triangleright \text{ab} \triangleleft x \text{ab}y \triangleright \text{cc} \triangleleft y \text{ca} \implies ([4, 6], [8, 10])$

$\text{abb}x \triangleright \text{abab} \triangleleft x c^y \triangleright \text{cc} \triangleleft y a \implies ([4, 8], [9, 11])$

\vdots

Regular Spanners – Notations

$\llbracket M \rrbracket$ denotes the spanner represented by an NFA M .

Regular Spanners – Notations

$\llbracket M \rrbracket$ denotes the spanner represented by an NFA M .

$\llbracket M \rrbracket(D)$ denotes the span-relation extracted from a document D .

Regular Spanners – Notations

$\llbracket M \rrbracket$ denotes the spanner represented by an NFA M .

$\llbracket M \rrbracket(D)$ denotes the span-relation extracted from a document D .

A spanner S is a *regular spanner* if $S = \llbracket M \rrbracket$ for some NFA M .

Results About Regular Spanners

Introduced by Fagin et al. PODS 2013, JACM 2015.

Results About Regular Spanners

Introduced by Fagin et al. PODS 2013, JACM 2015.

Since then intensely studied; many positive results.

Results About Regular Spanners

Introduced by Fagin et al. PODS 2013, JACM 2015.

Since then intensely studied; many positive results.

A major result: linear preprocessing and constant delay enumeration (Florenzano et al. PODS 2018, Amarilli et al. ICDT 2019).

Approach of this Paper

Spanner Evaluation over Compressed Documents

- Input: A spanner represented by an NFA M ,
a document D given in a **compressed form*** \mathcal{S} .
- Task: Evaluate M on D (e. g., model checking, computing or enumerating $\llbracket M \rrbracket(D)$)... but **without** decompressing \mathcal{S} .

*Compression Scheme: Straight-Line Programs (SLPs).

Straight-Line Programs

Straight-Line Programs

Straight-Line Program

A straight-line program for document D is a context-free grammar \mathcal{S} that describes the language $\{D\}$.

Straight-Line Programs

Example

Let S have rules

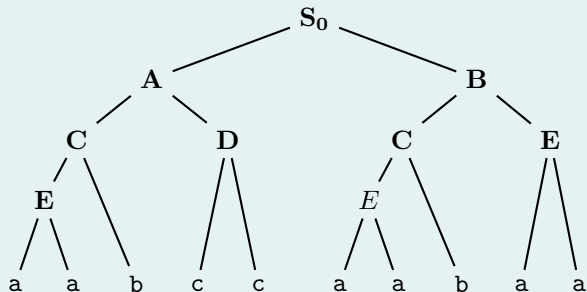
$$\begin{array}{lll} S_0 \rightarrow AB, & A \rightarrow CD, & B \rightarrow CE, \\ C \rightarrow Eb, & D \rightarrow cc, & E \rightarrow aa \end{array}$$

Straight-Line Programs

Example

Let S have rules

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



Straight-Line Programs

Example

Let S have rules

$$S_0 \rightarrow AB,$$

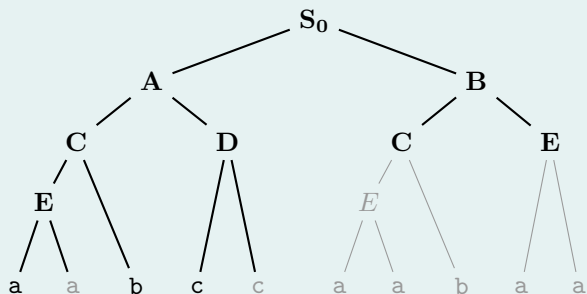
$$A \rightarrow CD,$$

$$B \rightarrow CE,$$

$$C \rightarrow Eb,$$

$$D \rightarrow cc,$$

$$E \rightarrow aa$$



Good Properties of SLPs

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

- ▶ Exponential compression rates.

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).

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).

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.

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.

Research Task

Spanner Evaluation over SLP-Compressed Documents

Input:	A spanner represented by an NFA M , an SLP \mathcal{S} for a document D .
Non-emptiness:	Check whether $\llbracket M \rrbracket(D) \neq \emptyset$.
Model Checking:	Check whether $t \in \llbracket M \rrbracket(D)$ for a span-tuple t .
Computation:	Compute $\llbracket M \rrbracket(D)$.
Enumeration:	Enumerate $\llbracket M \rrbracket(D)$.

Results

Results

Theorem (Data Complexity)

Non-emptiness: $O(\text{size}(\mathcal{S}))$

Model Checking: $O(\text{size}(\mathcal{S}))$

Computation: $O(\text{size}(\mathcal{S}) \cdot \text{size}(\llbracket M \rrbracket(D)))$

Enumeration: preprocessing time $O(\text{size}(\mathcal{S}))$ and
delay $O(\log(|D|))$.

Results

Theorem (Data Complexity)

Non-emptiness:	$O(\text{size}(\mathcal{S}))$
Model Checking:	$O(\text{size}(\mathcal{S}))$
Computation:	$O(\text{size}(\mathcal{S}) \cdot \text{size}(\llbracket M \rrbracket(D)))$
Enumeration:	preprocessing time $O(\text{size}(\mathcal{S}))$ and delay $O(\log(D))$.

Two remarks about combined-complexity:

- ▶ Sets of markers (“ $\{x_{\triangleright}, \triangleleft^y, z_{\triangleright}\}$ ”) as arc labels of the NFA (a.k.a. *extended variable-set automata*), which makes the NFA larger.

Results

Theorem (Data Complexity)

Non-emptiness: $O(\text{size}(\mathcal{S}))$

Model Checking: $O(\text{size}(\mathcal{S}))$

Computation: $O(\text{size}(\mathcal{S}) \cdot \text{size}(\llbracket M \rrbracket(D)))$

Enumeration: preprocessing time $O(\text{size}(\mathcal{S}))$ and delay $O(\log(|D|))$.

Two remarks about combined-complexity:

- ▶ Sets of markers (“ $\{x_{\triangleright}, \triangleleft^y, z_{\triangleright}\}$ ”) as arc labels of the NFA (a.k.a. *extended variable-set automata*), which makes the NFA larger.
- ▶ For the enumeration result, we require the NFA also to be deterministic.

Proof Sketches

Non-Emptiness, Model-Checking and Computation

Follows (non-trivial) from known results about the regular membership problem for SLP-compressed words:

Non-emptiness: $O(\text{size}(\mathcal{S}))$

Model Checking: $O(\text{size}(\mathcal{S}))$

Computation: $O(\text{size}(\mathcal{S}) \cdot \text{size}(\llbracket M \rrbracket(D)))$

Non-Emptiness, Model-Checking and Computation

Follows (non-trivial) from known results about the regular membership problem for SLP-compressed words:

Non-emptiness: $O(\text{size}(\mathcal{S}))$

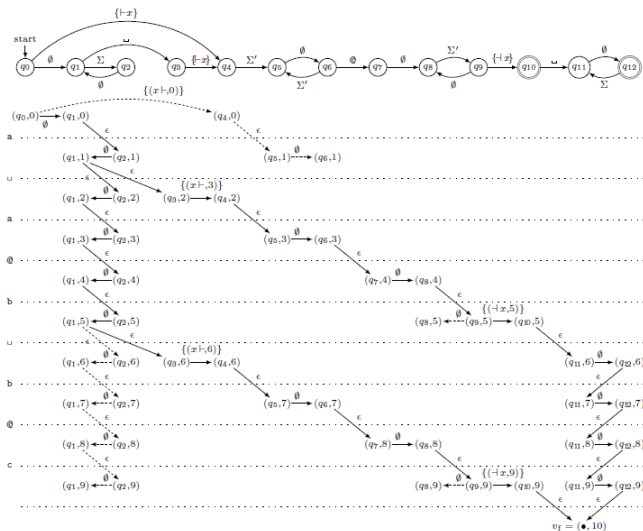
Model Checking: $O(\text{size}(\mathcal{S}))$

Computation: $O(\text{size}(\mathcal{S}) \cdot \text{size}(\llbracket M \rrbracket(D)))$

In the following: Sketch for Enumeration!

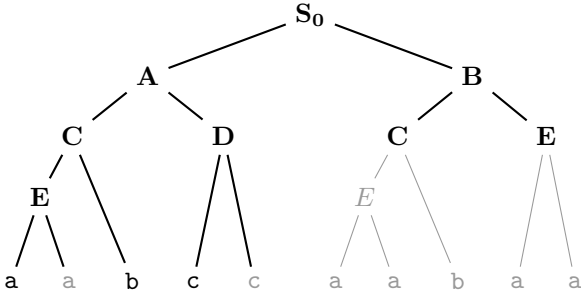
Non-Compressed Enumeration

(Florenzano et al. PODS 2018, Amarilli et al. ICDT 2019)

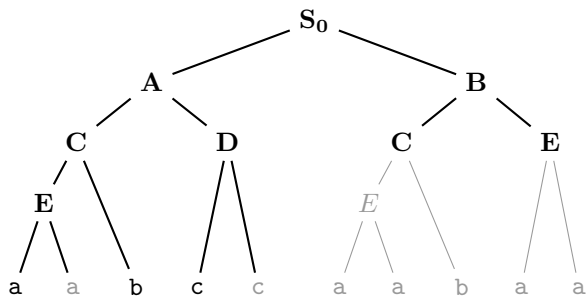


[Amarilli et al. ICDT 2019, SIGMOD Rec., 2020]

Marking SLPs

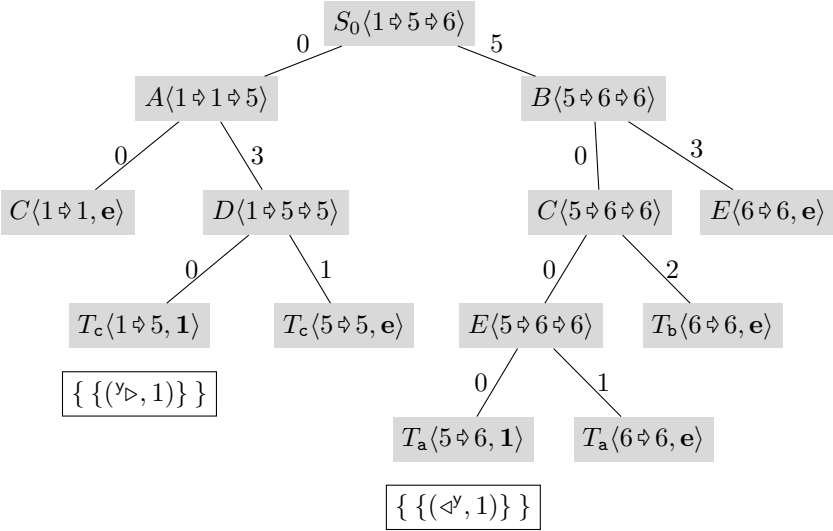


Marking SLPs



\Rightarrow enumerate *partially decompressed* SLPs.

Enumerating Partially Decompressed SLPs



Balancing SLPs

SLP Balancing Theorem, Ganardi, Jez and Lohrey, FOCS 2019:

Theorem

Any given SLP \mathcal{S} can be *balanced** in linear time.

* $\text{depth}(\mathcal{S}) = O(\log(|D|))$.

Future Work

Dynamic setting with updates!

Thank you very much for your attention.