

# Document Spanners — A Brief Overview of Concepts, Results, and Recent Developments

Markus L. Schmid and Nicole Schweikardt

Humboldt-Universität zu Berlin

Gems of PODS

PODS 2022: 41st Symposium on Principles of Database Systems

# Introduction

**Goal:** extract structured information from texts

**Idea:** query texts like you would query a relational database

# Introduction

**Goal:** extract structured information from texts

**Idea:** query texts like you would query a relational database

Timeline:

**Since 2005:** development of IBM's SystemT:  
a rule-based information extraction system with an SQL-like declarative language called AQL (Annotation Query Language)

# Introduction

**Goal:** extract structured information from texts

**Idea:** query texts like you would query a relational database

Timeline:

- Since 2005:** development of IBM's SystemT:  
a rule-based information extraction system with an SQL-like declarative language called AQL (Annotation Query Language)
- 2010:** release of SystemT as a commercial IBM product

# Introduction

**Goal:** extract structured information from texts

**Idea:** query texts like you would query a relational database

Timeline:

**Since 2005:** development of IBM's SystemT:  
a rule-based information extraction system with an SQL-like declarative language called AQL (Annotation Query Language)

**2010:** release of SystemT as a commercial IBM product

**PODS 2013:** introduction of the framework of **document spanners** as a formalisation of AQL:

## Spanners: A Formal Framework for Information Extraction

Ronald Fagin  
IBM Research – Almaden  
San Jose, CA, USA

Benny Kimelfeld  
IBM Research – Almaden  
San Jose, CA, USA

Frederick Reiss  
IBM Research – Almaden  
San Jose, CA, USA

Stijn Vansummeren  
Université Libre de  
Bruxelles (ULB)  
Bruxelles, Belgium

# Introduction

**Goal:** extract structured information from texts

**Idea:** query texts like you would query a relational database

Timeline:

**Since 2005:** development of IBM's SystemT:  
a rule-based information extraction system with an SQL-like declarative language called AQL (Annotation Query Language)

**2010:** release of SystemT as a commercial IBM product

**PODS 2013:** introduction of the framework of **document spanners** as a formalisation of AQL:

## Spanners: A Formal Framework for Information Extraction

Ronald Fagin  
IBM Research – Almaden  
San Jose, CA, USA

Benny Kimelfeld  
IBM Research – Almaden  
San Jose, CA, USA

Frederick Reiss  
IBM Research – Almaden  
San Jose, CA, USA

Stijn Vansummeren  
Université Libre de  
Bruxelles (ULB)  
Bruxelles, Belgium

**since then:** in-depth investigation of this framework

# Overview

Introduction

Concepts

Results

Recent Developments

# Overview

Introduction

Concepts

Results

Recent Developments



# Spanners: map documents to relations of spans

# Spanners: map documents to relations of spans

$D = \text{aggagccagc}$

x	y	z
$[2, 5\rangle$	$[4, 7\rangle$	$[1, 10\rangle$
$[3, 5\rangle$	$[5, 8\rangle$	$[4, 7\rangle$
$[1, 3\rangle$	$[4, 6\rangle$	$[2, 4\rangle$

# Spanners: map documents to relations of spans

a span

D = a g g a g c c a g c

x	y	z
[2, 5]	[4, 7]	[1, 10]
[3, 5]	[5, 8]	[4, 7]
[1, 3]	[4, 6]	[2, 4]

# Spanners: map documents to relations of spans

a span

D = aggagccagc

x	y	z
[2, 5]	[4, 7]	[1, 10]
[3, 5]	[5, 8]	[4, 7]
[1, 3]	[4, 6]	[2, 4]

# Spanners: map documents to relations of spans

a span

D = a g g a g c c a g c

x	y	z
[2, 5]	[4, 7]	[1, 10]
[3, 5]	[5, 8]	[4, 7]
[1, 3]	[4, 6]	[2, 4]

# Spanners: map documents to relations of spans

$D = \text{aggagccagc}$

x	y	z
$[2, 5\rangle$	$[4, 7\rangle$	$[1, 10\rangle$
$[3, 5\rangle$	$[5, 8\rangle$	$[4, 7\rangle$
$[1, 3\rangle$	$[4, 6\rangle$	$[2, 4\rangle$

# Spanners: map documents to relations of spans

$D = \text{aggagccagc}$

In SystemT / AQL: 2-stage approach

1. extract a relation from a document

regular expressions  
with capture variables

x	y	z
[2, 5]	[4, 7]	[1, 10]
[3, 5]	[5, 8]	[4, 7]
[1, 3]	[4, 6]	[2, 4]

# Spanners: map documents to relations of spans

$D = \text{aggagccagc}$

In SystemT / AQL: 2-stage approach

1. extract a relation from a document

x	y	z
[2, 5)	[4, 7)	[1, 10)
[3, 5)	[5, 8)	[4, 7)
[1, 3)	[4, 6)	[2, 4)

regular expressions  
with capture variables

finite automata  
(with marker symbols)



# Spanners: map documents to relations of spans

$D = \text{aggagccagc}$

In SystemT / AQL: 2-stage approach

1. extract a relation from a document

x	y	z
$[2, 5\rangle$	$[4, 7\rangle$	$[1, 10\rangle$
$[3, 5\rangle$	$[5, 8\rangle$	$[4, 7\rangle$
$[1, 3\rangle$	$[4, 6\rangle$	$[2, 4\rangle$

regular expressions  
with capture variables

finite automata  
(with marker symbols)

2. manipulate this relation by  
algebraic operations

standard relational operators:

union  $\cup$

projection  $\pi$

join  $\bowtie$

# Spanners: map documents to relations of spans

$D = \text{aggagccagc}$

In SystemT / AQL: 2-stage approach

1. extract a relation from a document

x	y	z
[2, 5]	[4, 7]	[1, 10]
[3, 5]	[5, 8]	[4, 7]
[1, 3]	[4, 6]	[2, 4]

regular expressions  
with capture variables

finite automata  
(with marker symbols)

2. manipulate this relation by  
algebraic operations

standard relational operators:

union  $\cup$

projection  $\pi$

join  $\bowtie$

text-centric operators:

string-equality selection  $\zeta^=$

containment

# Spanners: map documents to relations of spans

D = **a g g a g c c a g c**

In SystemT / AQL: 2-stage approach

1. extract a relation from a document

regular expressions  
with capture variables

$R =$

x	y	z
[2, 5]	[4, 7]	[1, 10]
[3, 5]	[5, 8]	[4, 7]
[1, 3]	[4, 6]	[2, 4]

finite automata  
(with marker symbols)

2. manipulate this relation by  
algebraic operations

$\zeta_{x,y}^{\equiv}(R)$

selects tuples of  $R$  where the sub-strings described by  $x$  and  $y$  are equal

standard relational operators:

union  $\cup$

projection  $\pi$

join  $\bowtie$

text-centric operators:

string-equality selection  $\zeta^{\equiv}$

containment

# Spanners: map documents to relations of spans

$D = \text{aggagccagc}$

In SystemT / AQL: 2-stage approach

1. extract a relation from a document

x	y	z
$[2, 5\rangle$	$[4, 7\rangle$	$[1, 10\rangle$
$[3, 5\rangle$	$[5, 8\rangle$	$[4, 7\rangle$
$[1, 3\rangle$	$[4, 6\rangle$	$[2, 4\rangle$

regular expressions  
with capture variables

finite automata  
(with marker symbols)

2. manipulate this relation by  
algebraic operations

standard relational operators:

union  $\cup$

projection  $\pi$

join  $\bowtie$

text-centric operators:

string-equality selection  $\zeta^=$

containment

conflict resolution:

resolve undesired overlapping spans

# Core Spanners: map documents to relations of spans

$D = \text{aggagccagc}$

In SystemT / AQL: 2-stage approach

1. extract a relation from a document

x	y	z
$[2, 5\rangle$	$[4, 7\rangle$	$[1, 10\rangle$
$[3, 5\rangle$	$[5, 8\rangle$	$[4, 7\rangle$
$[1, 3\rangle$	$[4, 6\rangle$	$[2, 4\rangle$

regular expressions  
with capture variables

finite automata  
(with marker symbols)

2. manipulate this relation by  
algebraic operations

standard relational operators:

union  $\cup$

projection  $\pi$

join  $\bowtie$

text-centric operators:

string-equality selection  $\zeta^=$

containment

conflict resolution:

resolve undesired overlapping spans

# Core-Simplification Lemma

[Fagin, Kimelfeld, Reiss, Vansummeren, PODS'13]

Every core spanner can be represented by a single variable-set automaton, followed by string-equality selections and then by a projection.

## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $\bowtie^x$  and  $\triangleleft^x$ .

## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $\langle^x$  and  $\rangle^x$ .

document	span tuple		subword-marked word
	$x$	$y$	
agcca	$[1, 2)$	$[3, 4)$	$\langle^x a \rangle^x g \rangle^y c \langle^y ca$



## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $x \triangleright$  and  $\triangleleft^x$ .

document	span tuple		subword-marked word
	$x$	$y$	
agcca	$[1, 2)$	$[3, 4)$	$x \triangleright a \triangleleft^x g \ y \triangleright c \triangleleft^y ca$
agcca	$[1, 2)$	$[4, 5)$	$x \triangleright a \triangleleft^x gc \ y \triangleright c \triangleleft^y a$
agcca	$[1, 2)$	$[3, 5)$	$x \triangleright a \triangleleft^x g \ y \triangleright cc \triangleleft^y a$

## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $x \triangleright$  and  $\triangleleft^x$ .

document	span tuple		subword-marked word
	$x$	$y$	
agcca	$[1, 2)$	$[3, 4)$	$x \triangleright a \triangleleft^x g \triangleright^y c \triangleleft^y ca$
agcca	$[1, 2)$	$[4, 5)$	$x \triangleright a \triangleleft^x gc \triangleright^y c \triangleleft^y a$
agcca	$[1, 2)$	$[3, 5)$	$x \triangleright a \triangleleft^x g \triangleright^y cc \triangleleft^y a$
$e(w)$	$st(w)$		$w$

## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $x \triangleright$  and  $\triangleleft^x$ .

document	span tuple		subword-marked word
	$x$	$y$	
agcca	$[1, 2)$	$[3, 4)$	$x \triangleright a \triangleleft^x g y \triangleright c \triangleleft^y ca$
agcca	$[1, 2)$	$[4, 5)$	$x \triangleright a \triangleleft^x gc y \triangleright c \triangleleft^y a$
agcca	$[1, 2)$	$[3, 5)$	$x \triangleright a \triangleleft^x g y \triangleright cc \triangleleft^y a$
$\epsilon(w)$	$st(w)$		$w$

**Observation:**

- Every subword-marked language  $L$  represents a spanner  $\llbracket L \rrbracket$ :  
 $\llbracket L \rrbracket(D) = \{st(w) : w \in L, \epsilon(w) = D\}$ .

## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $x \triangleright$  and  $\triangleleft^x$ .

document	span tuple		subword-marked word
	$x$	$y$	
agcca	$[1, 2)$	$[3, 4)$	$x \triangleright a \triangleleft^x g y \triangleright c \triangleleft^y ca$
agcca	$[1, 2)$	$[4, 5)$	$x \triangleright a \triangleleft^x gc y \triangleright c \triangleleft^y a$
agcca	$[1, 2)$	$[3, 5)$	$x \triangleright a \triangleleft^x g y \triangleright cc \triangleleft^y a$
$\epsilon(w)$	$st(w)$		$w$

**Observation:**

- Every subword-marked language  $L$  represents a spanner  $\llbracket L \rrbracket$ :  
 $\llbracket L \rrbracket(D) = \{st(w) : w \in L, \epsilon(w) = D\}$ .
- For every spanner  $S$  there is a subword-marked language  $L$  with  $S = \llbracket L \rrbracket$ .

## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $x \triangleright$  and  $\triangleleft^x$ .

document	span tuple		subword-marked word
	$x$	$y$	
agcca	$[1, 2]$	$[3, 4]$	$x \triangleright a \triangleleft^x g y \triangleright c \triangleleft^y ca$
agcca	$[1, 2]$	$[4, 5]$	$x \triangleright a \triangleleft^x gc y \triangleright c \triangleleft^y a$
agcca	$[1, 2]$	$[3, 5]$	$x \triangleright a \triangleleft^x g y \triangleright cc \triangleleft^y a$
$\epsilon(w)$	$st(w)$		$w$

**Observation:**

- Every subword-marked language  $L$  represents a spanner  $\llbracket L \rrbracket$ :  
 $\llbracket L \rrbracket(D) = \{st(w) : w \in L, \epsilon(w) = D\}$ .
- For every spanner  $S$  there is a subword-marked language  $L$  with  $S = \llbracket L \rrbracket$ .

**Regular spanners:** defined by regular languages of subword-marked words

## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $x \triangleright$  and  $\triangleleft^x$ .

document	span tuple		subword-marked word
	$x$	$y$	
agcca	$[1, 2]$	$[3, 4]$	$x \triangleright a \triangleleft^x g y \triangleright c \triangleleft^y ca$
agcca	$[1, 2]$	$[4, 5]$	$x \triangleright a \triangleleft^x gc y \triangleright c \triangleleft^y a$
agcca	$[1, 2]$	$[3, 5]$	$x \triangleright a \triangleleft^x g y \triangleright cc \triangleleft^y a$
$\epsilon(w)$	$st(w)$		$w$

**Observation:**

- Every subword-marked language  $L$  represents a spanner  $\llbracket L \rrbracket$ :  
 $\llbracket L \rrbracket(D) = \{st(w) : w \in L, \epsilon(w) = D\}$ .
- For every spanner  $S$  there is a subword-marked language  $L$  with  $S = \llbracket L \rrbracket$ .

**Regular spanners:** defined by regular languages of subword-marked words (equivalent to the *variable-set automata* of [Fagin et al, PODS'13])

## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $x \triangleright$  and  $\triangleleft^x$ .

document	span tuple		subword-marked word
	$x$	$y$	
agcca	$[1, 2]$	$[3, 4]$	$x \triangleright a \triangleleft^x g y \triangleright c \triangleleft^y ca$
agcca	$[1, 2]$	$[4, 5]$	$x \triangleright a \triangleleft^x gc y \triangleright c \triangleleft^y a$
agcca	$[1, 2]$	$[3, 5]$	$x \triangleright a \triangleleft^x g y \triangleright cc \triangleleft^y a$
$\epsilon(w)$	$st(w)$		$w$

**Observation:**

- Every subword-marked language  $L$  represents a spanner  $\llbracket L \rrbracket$ :  
 $\llbracket L \rrbracket(D) = \{st(w) : w \in L, \epsilon(w) = D\}$ .
- For every spanner  $S$  there is a subword-marked language  $L$  with  $S = \llbracket L \rrbracket$ .

**Regular spanners:** defined by regular languages of subword-marked words (equivalent to the *variable-set automata* of [Fagin et al, PODS'13])

**Context-free spanners:** def. by context-free languages of subword-marked words

## Representing spanners by languages of subword-marked words

**Markers:** For every variable  $x$  we use meta-symbols (markers)  $x^{\triangleright}$  and  $\triangleleft^x$ .

document	span tuple		subword-marked word
	$x$	$y$	
agcca	$[1, 2]$	$[3, 4]$	$x^{\triangleright} a \triangleleft^x g y^{\triangleright} c \triangleleft^y ca$
agcca	$[1, 2]$	$[4, 5]$	$x^{\triangleright} a \triangleleft^x gc y^{\triangleright} c \triangleleft^y a$
agcca	$[1, 2]$	$[3, 5]$	$x^{\triangleright} a \triangleleft^x g y^{\triangleright} cc \triangleleft^y a$
$\epsilon(w)$	$st(w)$		$w$

**Observation:**

- Every subword-marked language  $L$  represents a spanner  $\llbracket L \rrbracket$ :  
 $\llbracket L \rrbracket(D) = \{ st(w) : w \in L, \epsilon(w) = D \}$ .
- For every spanner  $S$  there is a subword-marked language  $L$  with  $S = \llbracket L \rrbracket$ .

**Regular spanners:** defined by regular languages of subword-marked words (equivalent to the *variable-set automata* of [Fagin et al, PODS'13])

**Context-free spanners:** def. by context-free languages of subword-marked words (equivalent to the *extraction grammars* of [Peterfreund, ICDT'21])



# Overview

Introduction

Concepts

Results

Recent Developments

# Enumerating Results of Spanners

**Goal:** given spanner  $S$  and document  $D$ , compute  $S(D)$

$D = \text{aggagccagc}$

$S(D) =$

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

# Enumerating Results of Spanners

**Goal:** given spanner  $S$  and document  $D$ , compute  $S(D)$

$D = \text{aggagccagc}$

$S(D) =$

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

**Preprocessing phase:** compute a suitable data structure

# Enumerating Results of Spanners

**Goal:** given spanner  $S$  and document  $D$ , compute  $S(D)$

$D = \text{aggagccagc}$

$S(D) =$

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

**Preprocessing phase:** compute a suitable data structure

**Enumeration phase:** use this data structure to enumerate, without repetition, all tuples in  $S(D)$

# Enumerating Results of Spanners

**Goal:** given spanner  $S$  and document  $D$ , compute  $S(D)$

$D = \text{aggagccagc}$

$S(D) =$

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

**Preprocessing phase:** compute a suitable data structure

**linear preprocessing time:**  $O(|D|)$

**Enumeration phase:** use this data structure to enumerate, without repetition, all tuples in  $S(D)$

# Enumerating Results of Spanners

**Goal:** given spanner  $S$  and document  $D$ , compute  $S(D)$

$D = \text{aggagccagc}$

$S(D) =$

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

**Preprocessing phase:** compute a suitable data structure

**linear preprocessing time:**  $O(|D|)$ , i.e.,  $f(S) \cdot O(|D|)$

**Enumeration phase:** use this data structure to enumerate, without repetition, all tuples in  $S(D)$

# Enumerating Results of Spanners

**Goal:** given spanner  $S$  and document  $D$ , compute  $S(D)$

$D = \text{aggagccagc}$

$S(D) =$

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

**Preprocessing phase:** compute a suitable data structure

**linear preprocessing time:**  $O(|D|)$ , i.e.,  $f(S) \cdot O(|D|)$

**Enumeration phase:** use this data structure to enumerate, without repetition, all tuples in  $S(D)$

**constant delay:**  $O(1)$

# Enumerating Results of Spanners

**Goal:** given spanner  $S$  and document  $D$ , compute  $S(D)$

$D = \text{aggagccagc}$

$S(D) =$

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

**Preprocessing phase:** compute a suitable data structure

**linear preprocessing time:**  $O(|D|)$ , i.e.,  $f(S) \cdot O(|D|)$

**Enumeration phase:** use this data structure to enumerate, without repetition, all tuples in  $S(D)$

**constant delay:**  $O(1)$ , i.e.,  $f(S)$



# Enumerating Results of Spanners

**Goal:** given spanner  $S$  and document  $D$ , compute  $S(D)$

$D = \text{aggagccagc}$

$S(D) =$

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

**Preprocessing phase:** compute a suitable data structure

**linear preprocessing time:**  $O(|D|)$ , i.e.,  $f(S) \cdot O(|D|)$

**Enumeration phase:** use this data structure to enumerate, without repetition, all tuples in  $S(D)$

**constant delay:**  $O(1)$ , i.e.,  $f(S)$

**Results for regular spanners:**

- [Florenzano, Riveros, Ugarte, Vansummeren, Vrgoc, PODS'18]:  
linear preprocessing  $f(S) \cdot O(|D|)$  and constant delay  $f(S)$

# Enumerating Results of Spanners

**Goal:** given spanner  $S$  and document  $D$ , compute  $S(D)$

$D = \text{aggagccagc}$

$S(D) =$

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

**Preprocessing phase:** compute a suitable data structure

**linear preprocessing time:**  $O(|D|)$ , i.e.,  $f(S) \cdot O(|D|)$

**Enumeration phase:** use this data structure to enumerate, without repetition, all tuples in  $S(D)$

**constant delay:**  $O(1)$ , i.e.,  $f(S)$

**Results for regular spanners:**

- [Florenzano, Riveros, Ugarte, Vansummeren, Vrgoc, PODS'18]: linear preprocessing  $f(S) \cdot O(|D|)$  and constant delay  $f(S)$
- [Amarilli, Bourhis, Mengel, Niewerth, ICDT'19]: improvement to  $f(S) = \text{poly}(|M|)$  where  $S$  is represented by a variable-set automaton  $M$

# Decision Problems for Spanners

Query evaluation problems:

# Decision Problems for Spanners

Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

# Decision Problems for Spanners

Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

# Decision Problems for Spanners

Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

Static analysis problems:

# Decision Problems for Spanners

Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .



# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

# Decision Problems for Spanners

Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.<sup>1</sup>

---

<sup>1</sup>i.e.: **balanced** brackets  $\{x \triangleright, \triangleleft x : x \in \mathcal{X}\}$  in subword-marked words.

not balanced:  $x \triangleright y \triangleright \triangleleft x \triangleleft y$       balanced:  $x \triangleright y \triangleright \triangleleft y \triangleleft x \ z \triangleright \triangleleft z$

# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.

Known:

# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.

## Known:

computational complexity of these problems for regular spanners & core spanners

# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.

## Known:

computational complexity of these problems for regular spanners & core spanners

## In short:

- regular spanners: mild complexity

# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.

## Known:

computational complexity of these problems for regular spanners & core spanners

## In short:

- regular spanners: mild complexity
- core spanners: intractable

# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.

**Known:** **Satisfiability** is

# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.

**Known:** **Satisfiability** is

- in  $O(|M|)$  if  $S$  is a regular spanner represented by a var.-set automaton  $M$



# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.

**Known:** **Satisfiability** is

- in  $O(|M|)$  if  $S$  is a regular spanner represented by a var.-set automaton  $M$
- PSpace-complete for core spanners  $S$  [Freydenberger, Holldack, ICDT'16]

# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.

**Known:** **Satisfiability** is

- in  $O(|M|)$  if  $S$  is a regular spanner represented by a var.-set automaton  $M$
  - PSpace-complete for core spanners  $S$  [Freydenberger, Holldack, ICDT'16]
- Idea:** intersection non-emptiness of regular languages is PSpace-hard

# Decision Problems for Spanners

## Query evaluation problems:

**Testing:** Given spanner  $S$ , document  $D$ , span tuple  $t$ , decide if  $t \in S(D)$ .

**NonEmptiness:** Given spanner  $S$ , document  $D$ , decide if  $S(D) \neq \emptyset$ .

## Static analysis problems:

**Satisfiability:** Given spanner  $S$ , decide if there is a document  $D$  with  $S(D) \neq \emptyset$ .

**Equivalence:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) = S_2(D)$  for every document  $D$ .

**Containment:** Given spanners  $S_1, S_2$ , decide if  
 $S_1(D) \subseteq S_2(D)$  for every document  $D$ .

**Hierarchicality:** Given spanner  $S$ , decide if  $S$  is hierarchical.

**Known:** **Satisfiability** is

- in  $O(|M|)$  if  $S$  is a regular spanner represented by a var.-set automaton  $M$
- PSpace-complete for core spanners  $S$  [Freydenberger, Holldack, ICDT'16]

**Idea:** intersection non-emptiness of regular languages is PSpace-hard

$$r_1 \cap \dots \cap r_n \neq \emptyset \iff \zeta_{x_1, \dots, x_n}^{\equiv} ({}^{x_1} \triangleright r_1 \triangleleft^{x_1} \dots {}^{x_n} \triangleright r_n \triangleleft^{x_n}) \text{ is satisfiable}$$

# Overview

Introduction

Concepts

Results

Recent Developments

# Recent Developments

- Evaluating spanners on compressed representations of documents

# Recent Developments

- Evaluating spanners on compressed representations of documents

- Beyond regular spanners:

Search for classes of spanners that

- are more expressive than the regular spanners
- but avoid the high complexity of the class of core spanners

Beyond regular spanners:

Refl-Spanners

# Ref-Words and Ref-Languages

A ref-word over  $\Sigma$  and  $\mathcal{X}$  is ...

... a subword-marked word  $w$  over  $\Sigma' = \mathcal{X} \cup \Sigma \cup \{x_{\triangleright}, x_{\triangleleft} : x \in \mathcal{X}\}$  such that

$$\forall x \in \mathcal{X} : w = w_1 x w_2 \implies w_1 = v_1 x_{\triangleright} v_2 x_{\triangleleft} v_3 \text{ and } x \notin v_2.$$



# Ref-Words and Ref-Languages

A ref-word over  $\Sigma$  and  $\mathcal{X}$  is ...

... a subword-marked word  $w$  over  $\Sigma' = \Sigma \cup \{\triangleright^x, \triangleleft^x : x \in \mathcal{X}\}$  such that

$$\forall x \in \mathcal{X} : w = w_1 x w_2 \implies w_1 = v_1 \triangleright^x v_2 \triangleleft^x v_3 \text{ and } x \notin v_2.$$

## Example

ref-words	subword-marked words, but no ref-words
$ag \triangleright^x ag \triangleleft^x c \triangleright^y xaa \triangleleft^y y$	

# Ref-Words and Ref-Languages

A ref-word over  $\Sigma$  and  $\mathcal{X}$  is ...

... a subword-marked word  $w$  over  $\Sigma' = \Sigma \cup \{\triangleright^x, \triangleleft^x : x \in \mathcal{X}\}$  such that

$$\forall x \in \mathcal{X} : w = w_1 x w_2 \implies w_1 = v_1 \triangleright^x v_2 \triangleleft^x v_3 \text{ and } x \notin v_2.$$

## Example

ref-words	subword-marked words, but no ref-words
$ag \triangleright^x ag \triangleleft^x c \triangleright^y xaa \triangleleft^y y$	
$a \triangleright^x ag \triangleright^y ag \triangleleft^x a \triangleleft^y xy$	

# Ref-Words and Ref-Languages

A ref-word over  $\Sigma$  and  $\mathcal{X}$  is ...

... a subword-marked word  $w$  over  $\Sigma' = \Sigma \cup \{\triangleright^x, \triangleleft^x : x \in \mathcal{X}\}$  such that

$$\forall x \in \mathcal{X} : w = w_1 x w_2 \implies w_1 = v_1 \triangleright^x v_2 \triangleleft^x v_3 \text{ and } x \notin v_2.$$

## Example

ref-words	subword-marked words, but no ref-words
$ag \triangleright^x ag \triangleleft^x c \triangleright^y xaa \triangleleft^y y$	$axg \triangleright^x ag \triangleleft^x c \triangleright^y xaa \triangleleft^y y$
$a \triangleright^x ag \triangleright^y ag \triangleleft^x a \triangleleft^y xy$	

# Ref-Words and Ref-Languages

A ref-word over  $\Sigma$  and  $\mathcal{X}$  is ...

... a subword-marked word  $w$  over  $\Sigma' = \mathcal{X} \cup \Sigma \cup \{^x\triangleright, \triangleleft^x : x \in \mathcal{X}\}$  such that

$$\forall x \in \mathcal{X} : w = w_1 x w_2 \implies w_1 = v_1 ^x\triangleright v_2 \triangleleft^x v_3 \text{ and } x \notin v_2.$$

## Example

ref-words	subword-marked words, but no ref-words
$ag ^x\triangleright ag \triangleleft^x c ^y\triangleright xaa \triangleleft^y y$	$axg ^x\triangleright ag \triangleleft^x c ^y\triangleright xaa \triangleleft^y y$
$a ^x\triangleright ag ^y\triangleright ag \triangleleft^x a \triangleleft^y xy$	$aa ^x\triangleright ag \triangleleft^x c ^y\triangleright ya \triangleleft^y$

# Ref-Words and Ref-Languages

A ref-word over  $\Sigma$  and  $\mathcal{X}$  is ...

... a subword-marked word  $w$  over  $\Sigma' = \Sigma \cup \{\triangleright^x, \triangleleft^x : x \in \mathcal{X}\}$  such that

$$\forall x \in \mathcal{X} : w = w_1 x w_2 \implies w_1 = v_1 \triangleright^x v_2 \triangleleft^x v_3 \text{ and } x \notin v_2.$$

## Example

ref-words	subword-marked words, but no ref-words
$ag \triangleright^x ag \triangleleft^x c \triangleright^y xaa \triangleleft^y y$	$axg \triangleright^x ag \triangleleft^x c \triangleright^y xaa \triangleleft^y y$
$a \triangleright^x ag \triangleright^y ag \triangleleft^x a \triangleleft^y xy$	$aa \triangleright^x ag \triangleleft^x c \triangleright^y ya \triangleleft^y$

A [ref-language](#) is a set of ref-words.

# Dereferencing Ref-Words

## Dereference-Function

Let  $w$  be a ref-word (over  $\Sigma$  and  $\mathcal{X}$ ).

$\partial(w)$  = replace all  $x$  by whatever is in  $\langle x \rangle \dots \langle x \rangle$ .

# Dereferencing Ref-Words

## Dereference-Function

Let  $w$  be a ref-word (over  $\Sigma$  and  $\mathcal{X}$ ).

$\partial(w)$  = replace all  $x$  by whatever is in  $\langle x \rangle \dots \langle x \rangle$ .

## Example

$\langle x_1 \rangle a \langle x_2 \rangle aa \langle x_1 \rangle cx_1 \langle x_3 \rangle ac \langle x_2 \rangle x_2 ax_1 \langle x_3 \rangle$

# Dereferencing Ref-Words

## Dereference-Function

Let  $w$  be a ref-word (over  $\Sigma$  and  $\mathcal{X}$ ).

$\partial(w)$  = replace all  $x$  by whatever is in  $\langle x \rangle \dots \langle x \rangle$ .

## Example

$\langle x_1 \rangle a \langle x_2 \rangle aa \langle x_1 \rangle c \langle x_1 \rangle \langle x_3 \rangle ac \langle x_2 \rangle x_2 a x_1 \langle x_3 \rangle$



# Dereferencing Ref-Words

## Dereference-Function

Let  $w$  be a ref-word (over  $\Sigma$  and  $\mathcal{X}$ ).

$\partial(w)$  = replace all  $x$  by whatever is in  $\langle x \rangle \dots \langle x \rangle$ .

## Example

$\langle x_1 \rangle a \langle x_2 \rangle aa \langle x_1 \rangle caaa \langle x_3 \rangle ac \langle x_2 \rangle x_2aaaa \langle x_3 \rangle$

# Dereferencing Ref-Words

## Dereference-Function

Let  $w$  be a ref-word (over  $\Sigma$  and  $\mathcal{X}$ ).

$\partial(w)$  = replace all  $x$  by whatever is in  $\langle x \rangle \dots \langle x \rangle$ .

## Example

$\langle x_1 \rangle a \langle x_2 \rangle aa \langle x_1 \rangle caaa \langle x_3 \rangle ac \langle x_2 \rangle x_2 aaaa \langle x_3 \rangle$

# Dereferencing Ref-Words

## Dereference-Function

Let  $w$  be a ref-word (over  $\Sigma$  and  $\mathcal{X}$ ).

$\partial(w)$  = replace all  $x$  by whatever is in  $\langle^{x_1} \rangle \dots \langle^{x_n} \rangle$ .

## Example

$\langle^{x_1} \rangle a \langle^{x_2} \rangle aa \langle^{x_1} \rangle caaa \langle^{x_3} \rangle ac \langle^{x_2} \rangle \mathbf{aacaaaa} caaaa \langle^{x_3} \rangle$

# Dereferencing Ref-Words

## Dereference-Function

Let  $w$  be a ref-word (over  $\Sigma$  and  $\mathcal{X}$ ).

$\partial(w)$  = replace all  $x$  by whatever is in  $x \triangleright \dots \triangleleft x$ .

## Example

$x_1 \triangleright a \ x_2 \triangleright aa \ \triangleleft^{x_1} caaa \ x_3 \triangleright ac \ \triangleleft^{x_2} aacaaaacaaaa \triangleleft^{x_3}$

ref-word	subword-marked word	document	span-tuple
$w$	$\partial(w)$	$\epsilon(\partial(w))$	$st(\partial(w))$

# Dereferencing Ref-Words

## Dereference-Function

Let  $w$  be a ref-word (over  $\Sigma$  and  $\mathcal{X}$ ).

$\partial(w)$  = replace all  $x$  by whatever is in  $x_{\triangleright} \dots \triangleleft_x$ .

## Example

$x_1_{\triangleright} a x_2_{\triangleright} aa \triangleleft_{x_1} caaa x_3_{\triangleright} ac \triangleleft_{x_2} aacaaaacaaaa \triangleleft_{x_3}$

ref-word	subword-marked word	document	span-tuple
$w$	$\partial(w)$	$\epsilon(\partial(w))$	$st(\partial(w))$

## Example

$w$  =  $x_1_{\triangleright} a x_2_{\triangleright} aa \triangleleft_{x_1} cx_1 x_3_{\triangleright} ac \triangleleft_{x_2} x_2ax_1 \triangleleft_{x_3}$   
 $\partial(w)$  =  $x_1_{\triangleright} a x_2_{\triangleright} aa \triangleleft_{x_1} caaa x_3_{\triangleright} ac \triangleleft_{x_2} aacaaaacaaaa \triangleleft_{x_3}$   
 $\epsilon(\partial(w))$  =  $aaacaaaacacaaaacaaaa$   
 $st(\partial(w))$  =  $([1, 4], [2, 10], [8, 22])$

# Refl-Spanners

## Refl-Spanners

A spanner  $S$  is a refl-spanner if  $S = \llbracket \partial(L) \rrbracket$  for a regular ref-language  $L$ .

# Refl-Spanners

## Refl-Spanners

A spanner  $S$  is a refl-spanner if  $S = \llbracket \partial(L) \rrbracket$  for a regular ref-language  $L$ .

## Example

Regular spanner:  $r = \Sigma^* x \triangleright (a \vee g)^+ \triangleleft^x \Sigma^* c \Sigma^* y \triangleright (a \vee g)^+ \triangleleft^y \Sigma^*$

String-equ. sel.:  $\overline{\zeta_{x,y}}$

# Refl-Spanners

## Refl-Spanners

A spanner  $S$  is a refl-spanner if  $S = \llbracket \partial(L) \rrbracket$  for a regular ref-language  $L$ .

## Example

Regular spanner:  $r = \Sigma^* x \triangleright (a \vee g)^+ \triangleleft^x \Sigma^* c \Sigma^* y \triangleright (a \vee g)^+ \triangleleft^y \Sigma^*$

String-equ. sel.:  $\overline{\zeta_{x,y}}$



Refl-spanner:  $r = \Sigma^* x \triangleright (a \vee g)^+ \triangleleft^x \Sigma^* c \Sigma^* y \triangleright \times \triangleleft^y \Sigma^*$



# Expressive Power of Refl-Spanners

## Theorem (Schmid, Schweikardt, ICDT'21)

- Every regular spanner is a refl-spanner.

# Expressive Power of Refl-Spanners

## Theorem (Schmid, Schweikardt, ICDT'21)

- Every regular spanner is a refl-spanner.
- Every reference-bounded<sup>1</sup> refl-spanner is a core spanner.

---

<sup>1</sup>no unbounded number of references, like in  $x \triangleright (a \vee g)^* \triangleleft^x x^*$

# Expressive Power of Refl-Spanners

## Theorem (Schmid, Schweikardt, ICDT'21)

- Every regular spanner is a refl-spanner.
- Every reference-bounded<sup>1</sup> refl-spanner is a core spanner.
- Core spanners with non-overlapping string-equality selections can be expressed by refl-spanners with span fusions.<sup>2</sup>

<sup>1</sup>no unbounded number of references, like in  $x \triangleright (a \vee g)^* \triangleleft^x x^*$

<sup>2</sup> $\bigcup_{\{x_1, x_2\} \rightarrow z} ($

$x_1$	$x_2$	$x_3$	$z$	$x_3$
$[1, 4]$	$[4, 7]$	$[1, 10]$	$[1, 7]$	$[1, 10]$
$[3, 5]$	$[5, 8]$	$[4, 7]$	$[3, 8]$	$[4, 7]$
$[1, 3]$	$[3, 10]$	$[2, 4]$	$[1, 10]$	$[2, 4]$

$) =$

# Complexity of Evaluation of Refl-Spanners

Problem	Regular sp.	Refl-sp. <sup>4</sup>	Core sp. <sup>6</sup>
Testing	$O( w  M  \log  \mathcal{X} )$	$O( w  M  \log  \mathcal{X} )$	NP-c
NonEmptiness	$O( w  M )$	NP-c	NP-h
Satisfiability	$O( M )$	$O( M )$	PSpace-c
Containment	PSpace-c <sup>3</sup>	PSpace-c <sup>5</sup>	undec.
Equivalence	PSpace-c <sup>3</sup>	PSpace-c <sup>5</sup>	undec.
Hierarchicality	$O( M  \mathcal{X} ^3)$	$O( M  \mathcal{X} ^3)$	PSpace-c

---

<sup>3</sup>Maturana et al., PODS'18.

<sup>4</sup>Schmid, Schweikardt, ICDT'21.

<sup>5</sup>For strongly reference extracting refl-spanners.

<sup>6</sup>Freydenberger, Holldack, ICDT'16.

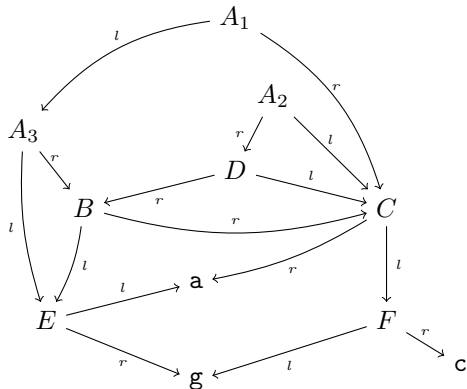
Evaluating spanners on  
compressed representations of documents

# SLP-Represented Document Databases

$$\text{DDB} = \left\{ \underbrace{\text{agaggcagca}}_{D_1}, \underbrace{\text{gcagcaaggca}}_{D_2}, \underbrace{\text{agaggca}}_{D_3} \right\}$$

# SLP-Represented Document Databases

$$\text{DDB} = \left\{ \underbrace{\text{agaggcagca}}_{D_1}, \underbrace{\text{gcagcaaggca}}_{D_2}, \underbrace{\text{agaggca}}_{D_3} \right\}$$



$$A_1 \rightarrow A_3 C$$

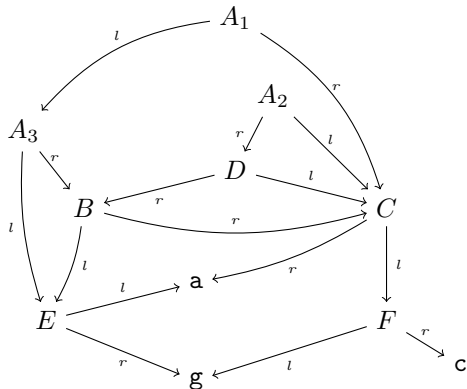
$$A_2 \rightarrow CD$$

$$C \rightarrow Fa$$

$$\vdots$$

# SLP-Represented Document Databases

$$\text{DDB} = \left\{ \underbrace{\text{agaggcagca}}_{D_1}, \underbrace{\text{gcagcaaggca}}_{D_2}, \underbrace{\text{agaggca}}_{D_3} \right\}$$



$$\begin{aligned} A_1 &\rightarrow A_3 C \\ A_2 &\rightarrow CD \\ C &\rightarrow Fa \\ &\vdots \end{aligned}$$

$$A_1 \rightarrow A_3 C \rightarrow EBC \rightarrow \dots \rightarrow D_1$$

$$A_2 \rightarrow \dots \rightarrow D_2$$



# Complex Document Editing

## CDE-Expression

An expression over documents, using typical text-editor copy-and-paste operations.

# Complex Document Editing

## CDE-Expression

An expression over documents, using typical text-editor copy-and-paste operations.

## Example

$$e = \text{concat}(D_1, \text{insert}(D_3, \text{extract}(D_7, 5, 21), 12))$$

# Complex Document Editing

## CDE-Expression

An expression over documents, using typical text-editor copy-and-paste operations.

## Example

$$e = \text{concat}(D_1, \text{insert}(D_3, \text{extract}(D_7, 5, 21), 12))$$

## Main Result (Schmid, Schweikardt, PODS'21, PODS'22):

- Represent a document database DDB by a suitably balanced SLP  $\mathcal{S}$ .
- For a given regular spanner  $R$  use time  $O(|\mathcal{S}|)$  to build a data structure that enables to efficiently enumerate the result of  $R$  on each document in DDB
- add the document  $D_e$  described by a CDE-expression  $e$  into this data structure in time  $\approx O(\log |D_e|)$

## More Precisely: Main Result

### Theorem<sup>7</sup> (Schmid, Schweikardt, PODS'21, PODS'22)

Let DDB be represented by a strongly balanced<sup>8</sup> SLP  $S$ , and let  $Q$  be a finite set of regular document spanners.

**Enumeration:** In time  $O(|Q| \cdot |S|)$ , we can compute a data structure such that, for any  $R \in Q$  and  $D \in \text{DDB}$ , we can enumerate  $\llbracket R \rrbracket(D)$  with delay  $O(\log |D|)$ .

---

<sup>7</sup>All bounds in data complexity.

<sup>8</sup>A property that implies that paths are bounded logarithmically in the size of the data.

## More Precisely: Main Result

### Theorem<sup>7</sup> (Schmid, Schweikardt, PODS'21, PODS'22)

Let DDB be represented by a strongly balanced<sup>8</sup> SLP  $S$ , and let  $Q$  be a finite set of regular document spanners.

**Enumeration:** In time  $O(|Q| \cdot |S|)$ , we can compute a data structure such that, for any  $R \in Q$  and  $D \in \text{DDB}$ , we can enumerate  $\llbracket R \rrbracket(D)$  with delay  $O(\log |D|)$ .

**Updates:** Given a CDE-expression  $e$  over DDB, we can extend  $S$  into a strongly balanced SLP  $S'$  for

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

in time  $O(|e| \cdot \log d)$ ,

where  $d$  is the maximum length of any intermediate document represented by a subexpression of  $e$ .

We can update the data structure in time  $O(|Q| \cdot |e| \cdot \log d)$  such that afterwards, upon input of any  $D \in \text{DDB}'$  and any  $R \in Q$ , we can enumerate  $\llbracket R \rrbracket(D)$  with delay  $O(\log |D|)$ .

---

<sup>7</sup>All bounds in data complexity.

<sup>8</sup>A property that implies that paths are bounded logarithmically in the size of the data.

Thank you very much  
for your attention!