



# Annis on MonetDB

Viktor Rosenfeld  
rosenfel@informatik.hu-berlin.de

14. January 2013

Advisors: Prof. Dr. Ulf Leser and Dr. Stefan Manegold



<http://www.flickr.com/photos/karola/3623768629>

1. What is Annis and how is it used?
2. Current implementation on PostgreSQL
3. What are Column-Stores? How can Annis benefit?
4. New implementation on MonetDB and evaluation

1. What is Annis and how is it used?

# What's a corpus?

any principled collection of language

THE WALL STREET JOURNAL.

 die tageszeitung

Falko



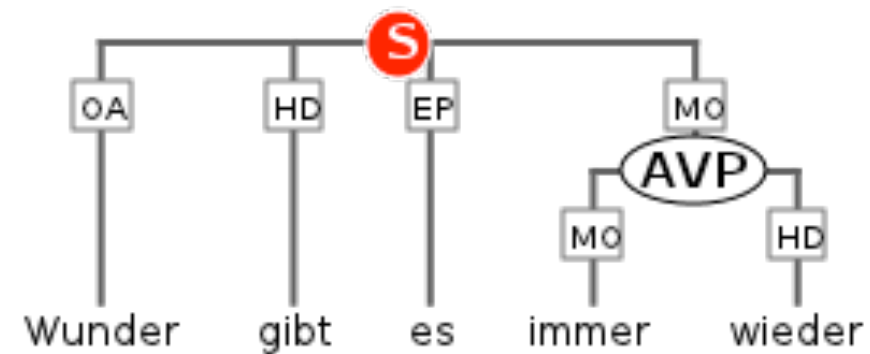


# What's an annotation?

classification and interpretation of the corpus data

additional data to enrich the corpus

Wunder	gibt	es	immer	wieder	!
Wunder	geben	es	immer	wieder	!
Acc.PI.Neut	3.Sg.Pres.Ind	3.Nom.Sg.Neut	--	--	--
NN	VFIN	PPER	ADV	ADV	\$.



Steilpass **Wunder gibt es immer wieder** ! Erst spielen die Dallgower Gemeindevertreter so statisch und verzagt wie die deutsche Abwehrreihe der Fußballkicker . Und dann kommt aus der Tiefe solch ein fulminanter Steilpass , von dem man hofft , dass die Seeburger oder Groß-Glienicker Mitspieler ihn aufnehmen können . Ein Befreiungsschlag ist es allerdings nicht , weil es vorerst keine Gefahr fürs Dallgower Tor gab . Die Seeburger und einige Groß-Glienicker haben den Ball erst zurückgespielt und dann um so drängender wieder gefordert . Nun sollen sie zeigen , wie sie die Chance verwerten . Eine Diskussion , wo künftig die Trainerkabine stehen soll , wäre in der jetzigen Spielsituation verheerend . Und eine Parallele zu den deutschen Grotten-Kickern gibt es immer noch . Auch wenn die Spieler aus den verschiedenen Vereinen zusammengewürfelt sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft " Döberitzer Heide " spielen . Und das heißt gemeinsam und nicht gegeneinander . Ermahnungen von der Seitenlinie , miteinander fair umzugehen und sich nicht beim kleinsten Schubser gegenseitig zu zerfleischen , sind normalerweise überflüssig . Vorerst allerdings hilfreich .

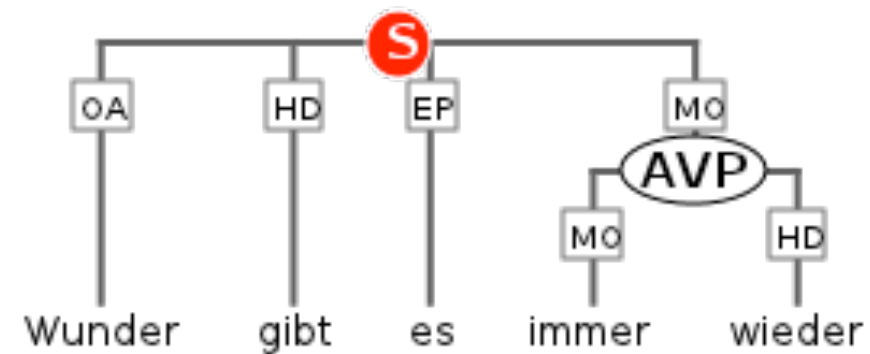
„Steilpass“ – Märkische Allgemeine Zeitung, 12.10.2001  
Potsdam Commentary Corpus (Stede, 2004)

# What's an annotation?

classification and interpretation of the corpus data

additional data to enrich the corpus

Wunder	gibt	es	immer	wieder	!
Wunder	geben	es	immer	wieder	!
Acc.PI.Neut	3.Sg.Pres.Ind	3.Nom.Sg.Neut	--	--	--
NN	VVFIN	PPER	ADV	ADV	\$.



Steilpass Wunder gibt es immer wieder ! Erst spielen die Dallgower Gemeindevertreter so statisch und verzagt wie die deutsche Abwehrreihe der Fußballkicker . Und dann kommt aus der Tiefe solch ein fulminanter Steilpass , von dem man hofft , dass die Seeburger oder Groß-Glienicker Mitspieler ihn aufnehmen können . Ein Befreiungsschlag ist es allerdings nicht , weil es vorerst keine Gefahr fürs Dallgower Tor gab . Die Seeburger und einige Groß-Glienicker haben den Ball erst zurückgespielt und dann um so drängender wieder gefordert . Nun sollen sie zeigen , wie sie die Chance verwerten . Eine Diskussion , wo künftig die Trainerkabine stehen soll , wäre in der jetzigen Spielsituation verheerend . Und eine Parallele zu den deutschen Grotten-Kickern gibt es immer noch . Auch wenn die Spieler aus den verschiedenen Vereinen zusammengewürfelt sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft " Döberitzer Heide " spielen . Und das heißt gemeinsam und nicht gegeneinander . Ermahnungen von der Seitenlinie , miteinander fair umzugehen und sich nicht beim kleinsten Schubser gegenseitig zu zerfleischen , sind normalerweise überflüssig . Vorerst allerdings hilfreich .

„Steilpass“ – Märkische Allgemeine Zeitung, 12.10.2001  
Potsdam Commentary Corpus (Stede, 2004)

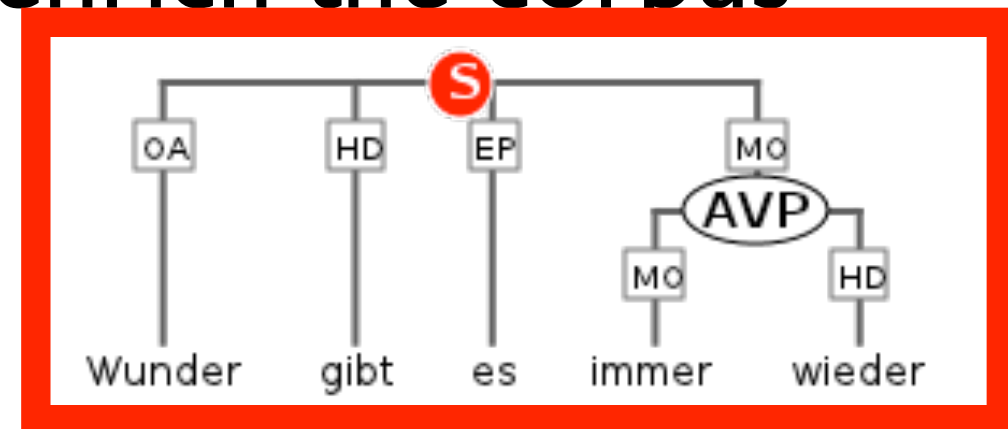


# What's an annotation?

classification and interpretation of the corpus data

additional data to enrich the corpus

Wunder	gibt	es	immer	wieder	!
Wunder	geben	es	immer	wieder	!
Acc.PI.Neut	3.Sg.Pres.Ind	3.Nom.Sg.Neut	--	--	--
NN	VVFIN	PPER	ADV	ADV	\$.



Steilpass **Wunder gibt es immer wieder** ! Erst spielen die Dallgower Gemeindevertreter so statisch und verzagt wie die deutsche Abwehrreihe der Fußballkicker . Und dann kommt aus der Tiefe solch ein fulminanter Steilpass , von dem man hofft , dass die Seeburger oder Groß-Glienicker Mitspieler ihn aufnehmen können . Ein Befreiungsschlag ist es allerdings nicht , weil es vorerst keine Gefahr fürs Dallgower Tor gab . Die Seeburger und einige Groß-Glienicker haben den Ball erst zurückgespielt und dann um so drängender wieder gefordert . Nun sollen sie zeigen , wie sie die Chance verwerten . Eine Diskussion , wo künftig die Trainerkabine stehen soll , wäre in der jetzigen Spielsituation verheerend . Und eine Parallele zu den deutschen Grotten-Kickern gibt es immer noch . Auch wenn die Spieler aus den verschiedenen Vereinen zusammengewürfelt sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft " Döberitzer Heide " spielen . Und das heißt gemeinsam und nicht gegeneinander . Ermahnungen von der Seitenlinie , miteinander fair umzugehen und sich nicht beim kleinsten Schubser gegenseitig zu zerfleischen , sind normalerweise überflüssig . Vorerst allerdings hilfreich .

„Steilpass“ – Märkische Allgemeine Zeitung, 12.10.2001  
Potsdam Commentary Corpus (Stede, 2004)

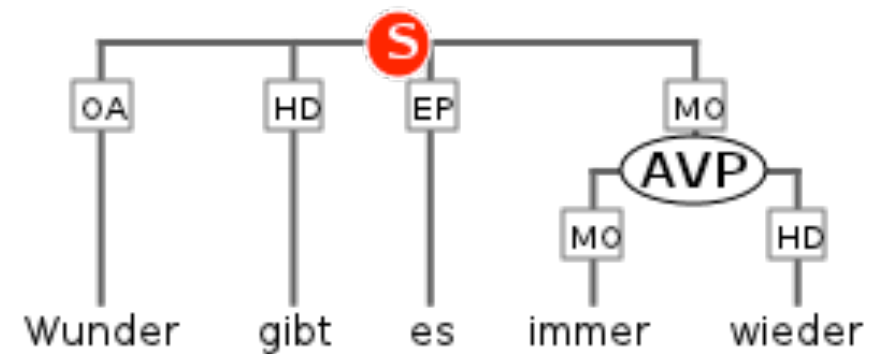


# What's an annotation?

classification and interpretation of the corpus data

additional data to enrich the corpus

Wunder	gibt	es	immer	wieder	!
Wunder	geben	es	immer	wieder	!
Acc.PI.Neut	3.Sg.Pres.Ind	3.Nom.Sg.Neut	--	--	--
NN	VVFIN	PPER	ADV	ADV	\$.



Steilpass **Wunder gibt es immer wieder** ! Erst spielen die Dallgower Gemeindevertreter so statisch und verzagt wie die deutsche Abwehrreihe der Fußballkicker . Und dann kommt aus der Tiefe solch ein fulminanter Steilpass , von dem man hofft , dass die Seeburger oder Groß-Glienicker Mitspieler ihn aufnehmen können . Ein Befreiungsschlag ist es allerdings nicht , weil es vorerst keine Gefahr fürs Dallgower Tor gab . Die Seeburger und einige Groß-Glienicker haben den Ball erst zurückgespielt und dann um so drängender wieder gefordert . Nun sollen sie zeigen , wie sie die Chance verwerten . Eine Diskussion , wo künftig die Trainerkabine stehen soll , wäre in der jetzigen Spielsituation verheerend . Und eine Parallele zu den deutschen Grotten-Kickern gibt es immer noch . Auch wenn die Spieler aus den verschiedenen Vereinen zusammengewürfelt sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft " Döberitzer Heide " spielen . Und das heißt gemeinsam und nicht gegeneinander . Ermahnungen von der Seitenlinie , miteinander fair umzugehen und sich nicht beim kleinsten Schubser gegenseitig zu zerfleischen , sind normalerweise überflüssig . Vorerst allerdings hilfreich .

ANNIS<sup>2</sup> Corpus Search

ANNIS<sup>2</sup> Tutorial logged in as "viktor"

**Search Form**

AnnisQL:

Show Result Query Builder History

Result: 1

**More Corpora**

Name	Texts	Tokens
FalkoEssayL1v2.2	95	70648
FalkoEssayL2v2.2	248	132069
TueBa-D/Z.6.0	2777	975836
b1.aja	117	1452
b2.tangale	49	452
<input checked="" type="checkbox"/> d1.pcc2	2	399
d2.20samplesDEU	22	382
pcc-176	176	33222
ridges.herbology	14	63734
x.ontonotes1.0	597	370789
x.ontonotes1.0_test...	60	38775
x.tiger2	1971	888578

Search Export

Context Left:

Context Right:

Results Per Page:

**Search Result - cat='S' & 'Wunder' & #1 \_1\_ #2 (5, 5)**

Page 1 of 1 | Token Annotations | Show Citation URL | Document Path | Displaying Results 1 - 1 of 1

Path: d1.pcc2 > 4282

Stellpass	Wunder	gibt	es	immer wieder !	Erst	spielen	die	Dallgower
Stellpass	Wunder	geben	es	immer wieder !	erst	spielen	der	Dallgower
Nom.Sg.Masc	Acc.Pl.Neut	3.Sg.Pres.Ind	3.Nom.Sg.Neut	-	-	3.Pl.Pres.Ind	Nom.Pl.Masc	Pos."*.*"
NN	NN	VVFIN	PPER	ADV	ADV	\$ ADV	VVFIN	ART

Constituents (Tree View)

Dependencies (Arches View)

Information Structure (Grid View)

Discourse Referents (Grid View)

Rhetorical Structure (Old Grid View)

Coreference (Discourse View)

Stellpass Wunder gibt es immer wieder ! Erst spielen die Dallgower Gemeindevertreter so statisch und verzagt wie die deutsche Abwehrreihe der Fußballkicker . Und dann kommt aus der Tiefe solch ein fulminanter Stellpass , von dem man hofft , dass die Seeburger oder Groß-Glienicker Mitspieler ihn aufnehmen können . Ein Befreiungsschlag ist es allerdings nicht , weil es vorerst keine Gefahr fürs Dallgower Tor gab . Die Seeburger und einige Groß-Glienicker haben den Ball erst zurückgespielt und dann um so drängender wieder gefordert . Nun sollen sie zeigen , wie sie die Chance verwerten . Eine Diskussion , wo künftig die Trainerkabine stehen soll , wäre in der jetzigen Spielsituation verheerend . Und eine Parallele zu den deutschen Grotten-Kickern gibt es immer noch . Auch wenn die Spieler aus den verschiedenen Vereinen zusammengewürfelt sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft " Döberitzer Heide " spielen . Und das heißt gemeinsam und nicht gegeneinander . Ermahnungen von der Seitenlinie , miteinander fair umzugehen und sich nicht beim kleinsten Schubser gegenseitig zu zerfleischen , sind normalerweise überflüssig . Vorerst allerdings hilfreich .



Query

The screenshot displays the Annis² Corpus Search interface. At the top, the search query is entered as `cat="S" & "Wunder" & #1_#2`. Below the query, a list of corpora is shown, with `d1.pcc2` selected. The search results for `d1.pcc2` are displayed, showing a snippet of text: `Stellpass Wunder gibt es immer wieder ! Erst spielen die Dallgower`. The interface includes various annotation views: Constituents (Tree View), Dependencies (Arches View), and Information Structure (Grid View). The tree view shows the sentence structure with nodes like `S`, `EP`, `AVP`, `NP`, and `NK`. The dependencies view shows arcs between words, such as `adv` and `adv`. The information structure view shows the sentence structure in a grid format.

Annotations

Corpus  
selection



Query

ANNIS<sup>2</sup> Corpus Search

Search Form

AnnisQL: `cat="S" & "Wunder" & #1_#2`

Result: 1

More Corpora

Name	Texts	Tokens
FalkoEssayL1v2.2	95	70648
FalkoEssayL2v2.2	248	132069
TueBa-D/Z.6.0	2777	975836
b1.aja	117	1452
b2.tangale	49	452
<input checked="" type="checkbox"/> d1.pcc2	2	399
d2.20samplesDEU	22	382
pcc-176	176	33222
ridges_herbiology	14	63734
x.ontonotes1.0	597	370789
x.ontonotes1.0_test...	60	38775
x.tiger2	1971	888578

Search Result - cat="S" & "Wunder" & #1\_#2 (5, 5)

Token Annotations

Stellpass Wunder gibt es immer wieder ! Erst spielen die Dallgower

Constituents (Tree View)

```

graph TD
    S((S)) --- OA[OA]
    S --- HD[HD]
    S --- EP[EP]
    S --- MO1[MO]
    S --- AVP((AVP))
    S --- MO2[MO]
    S --- HD2[HD]
    S --- SB[SB]
    S --- NP((NP))
    NP --- NK1[NK]
    NP --- NK2[NK]
  
```

Dependencies (Arches View)

```

graph TD
    W[Wunder] -- obj --> G[gibt]
    G -- adv --> E[es]
    E -- punct --> I[immer]
    I -- adv --> W2[wieder]
    W2 -- punct --> E2[Erst]
    E2 -- adv --> S[spielen]
    S -- det --> D[die]
    D -- attr --> Da[Dallgower]
  
```

Information Structure (Grid View)

Discourse Referents (Grid View)

Rhetorical Structure (Old Grid View)

Coreference (Discourse View)

Context

Context

Results

Stellpass Wunder gibt es immer wieder ! Erst spielen die Dallgower Gemeindevertreter so statisch und verzagt wie aus der Tiefe solch ein fulminanter Stellpass , von dem man hofft , dass die Seeburger oder Groß-Glienicker Mitglieder allerdings nicht , weil es vorerst keine Gefahr fürs Dallgower Tor gab . Die Seeburger und einige Groß-Glienicker wieder gefordert . Nun sollen sie zeigen , wie sie die Chance verwerten . Eine Diskussion , wo künftig die Trainer verheerend . Und eine Parallele zu den deutschen Grotten-Kickern gibt es immer noch . Auch wenn die Spieler an sich daran gewöhnen , dass sie nun in einer Mannschaft " Döberitzer Heide " spielen . Und das heißt gemeinsam miteinander fair umzugehen und sich nicht beim kleinsten Schubser gegenseitig zu zerfleischen , sind normalerweise

Annotations

```

#relation name
#attribute #1_id string
#attribute #1_span string
#attribute #1_tigerioat string
#data
'10870210','Die Jugendlichen wurden somit zum blo\u00DFen Feigenblatt degradiert','S'
'10869648','Nun sollen sie zeigen , wie sie die Chance verwerten','S'
'10869613','Wunder gibt es immer wieder','S'
'10870181','Und aberwitzig dazu','S'
'10870159','Das forderten sie bei der ersten Zossenrunde am Dienstagabend','S'
'10870219','Damit ist eine gro\u00DFe Chance vertan','S'
'10869672','Und das hei\u00DFt gemeinsam und nicht gegeneinander','S'
'10870217','Nicht \u00FCber sondern mit ihnen h\u00E4tten die Politiker reden sollen','S'
'10870156','Die Jugendlichen in Zossen wollen ein Masikoaf\u00E9','S'
'10870222','Vielleicht klappt es bei der n\u00E4chsten Runde Anfang 2002','S'
'10870225','Dann werden auch mehr Jugendliche eingeladen','S'
'10870229','In der Gruppe k\u00FCnnen sie sich hoffentlich mehr Geh\u00F6r verschaffen','S'
'10869681','Vorerst allerdings hilfreich','S'
'10870206','Und d'
'10870237','Die
  
```

Export for statistical analysis

Corpus selection

# Annis query language

cat="S" &  
"Wunder" &  
#1 \_i\_ #2

find a sentence

and find the phrase "Wunder"

the sentence includes the phrase "Wunder"

# Annis query language

cat="S" &  
"Wunder" &  
#1 \_i\_ #2

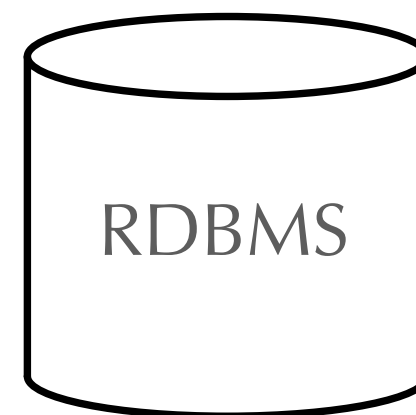
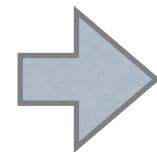
find a sentence

and find the phrase "Wunder"

the sentence includes the phrase "Wunder"



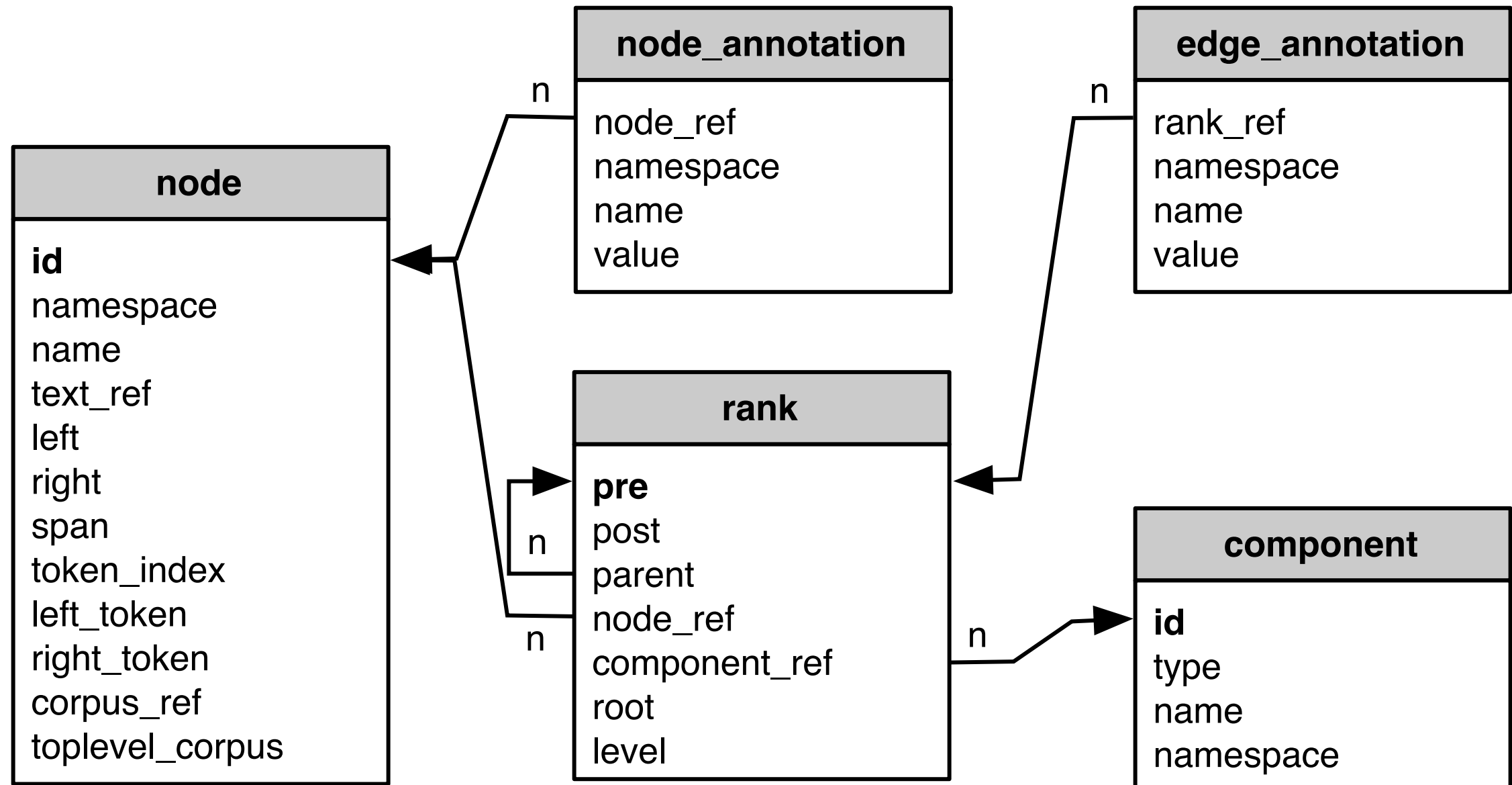
```
SELECT id1, id2  
FROM ...  
WHERE ...
```



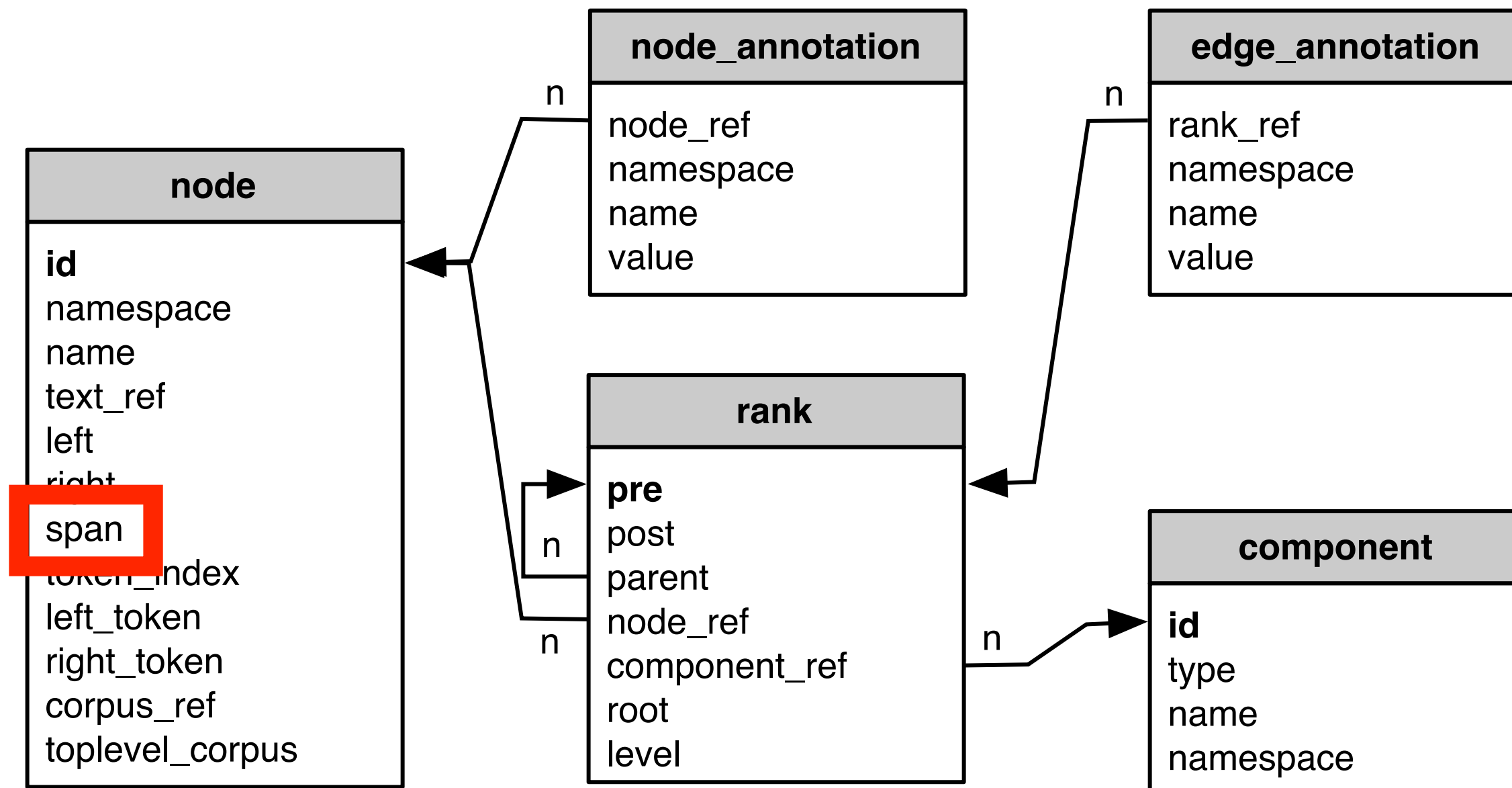


## 2. Current implementation on PostgreSQL

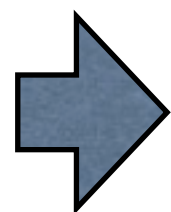
# Database schema



# Example 1: Text search



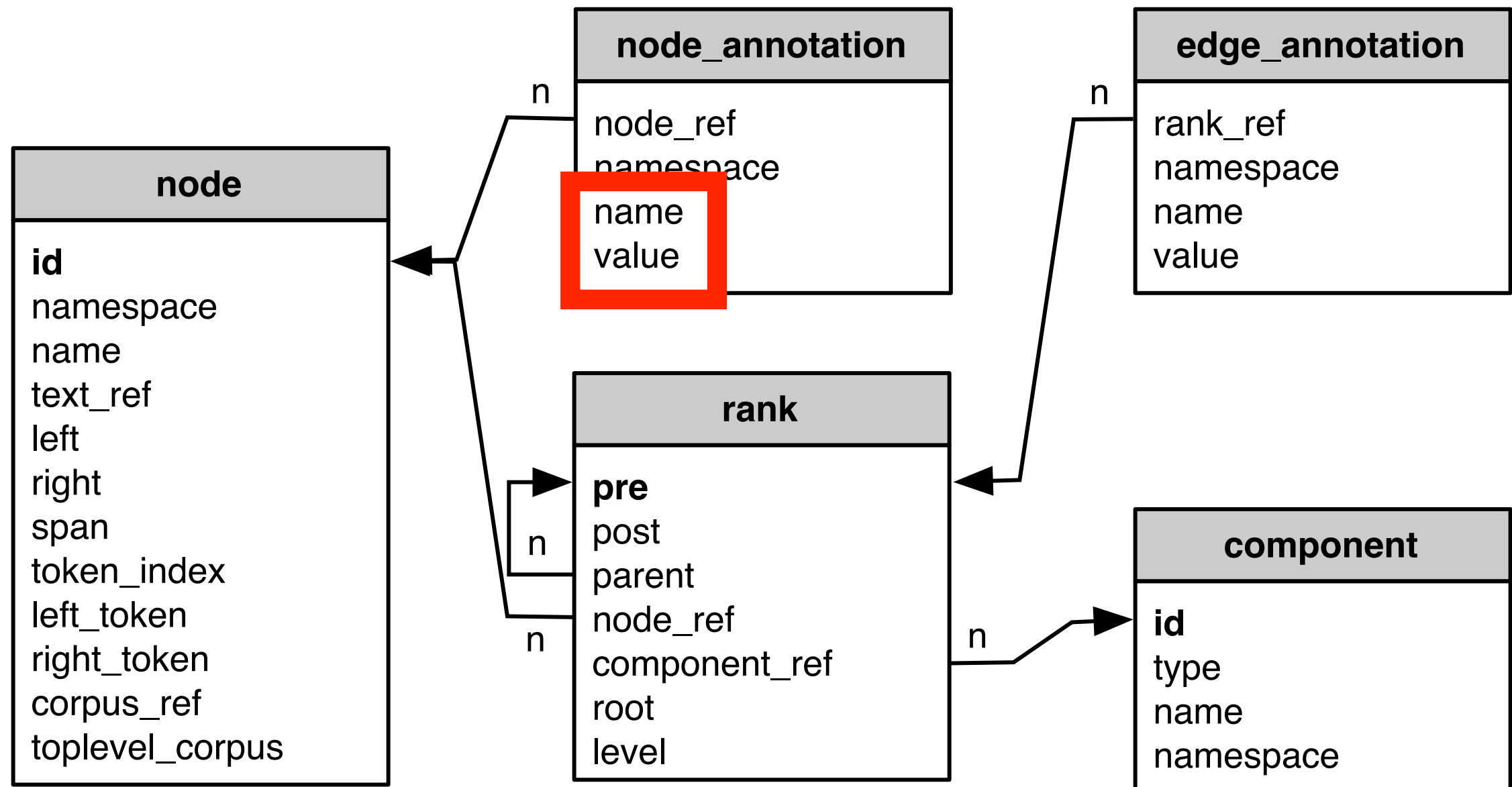
AQL: "Wunder"



SQL: nodeN.span = 'Wunder'



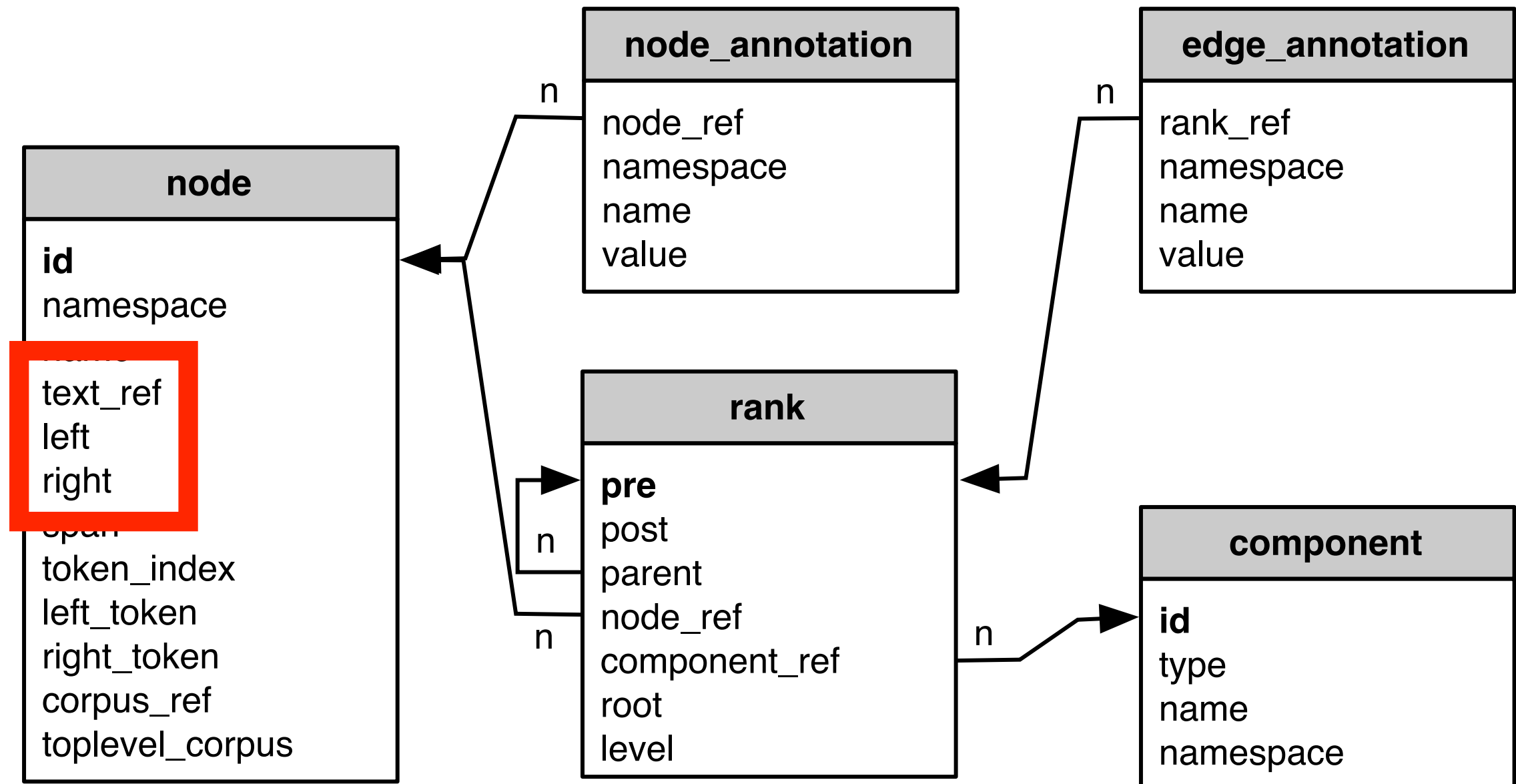
# Example 2: Annotation search



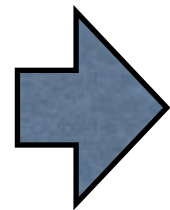
cat="S" →

node\_annotation $N$ .name = 'cat'  
node\_annotation $N$ .value = 'S'

# Example 3: Inclusion operator

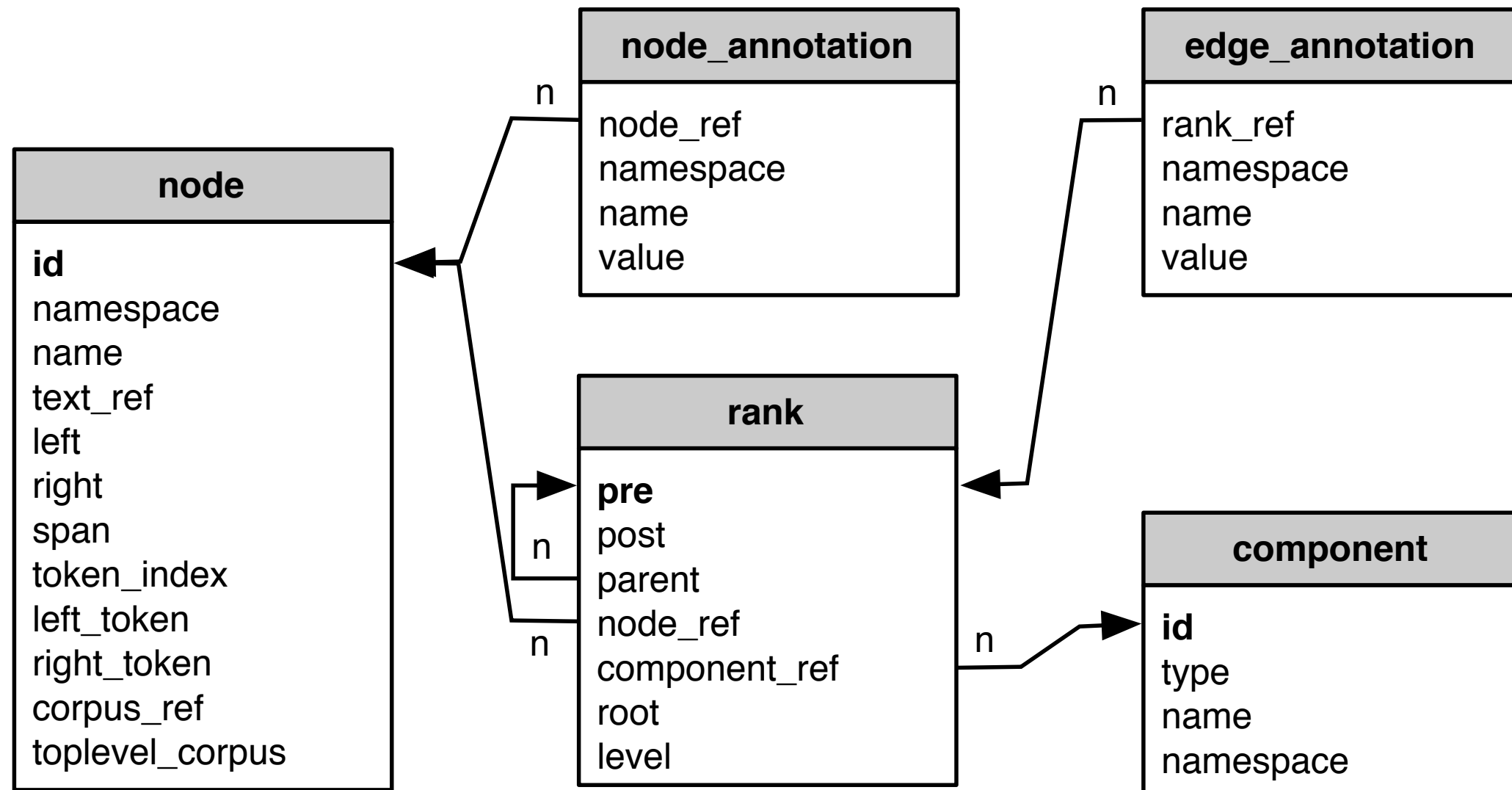


#1 \_i\_ #2



`node1.text_ref = node2.text_ref`  
`node1.right <= node2.right`  
`node1.left >= node2.left`

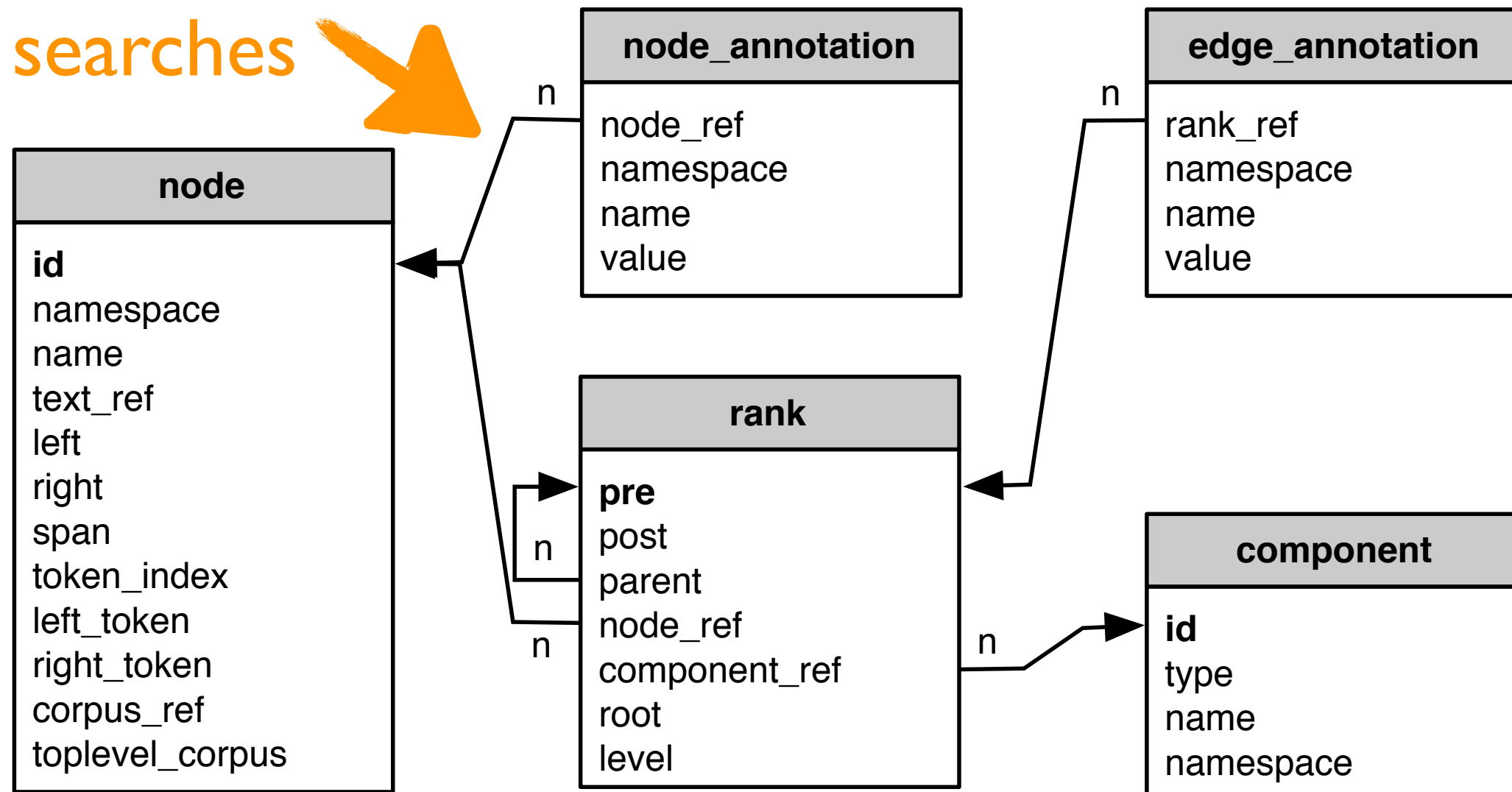
# Many tables – Many joins





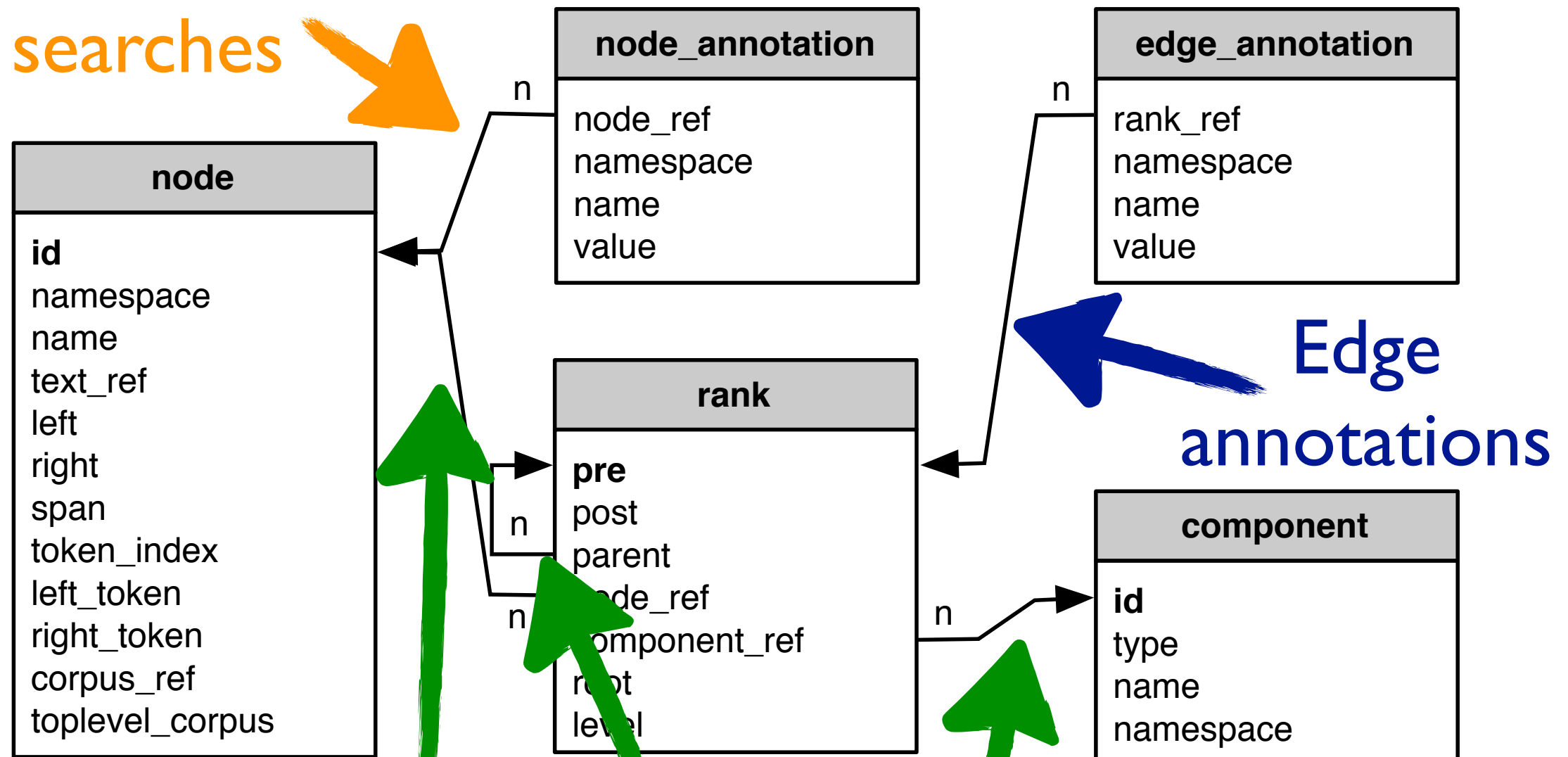
# Many tables – Many joins

Annotation searches



# Many tables – Many joins

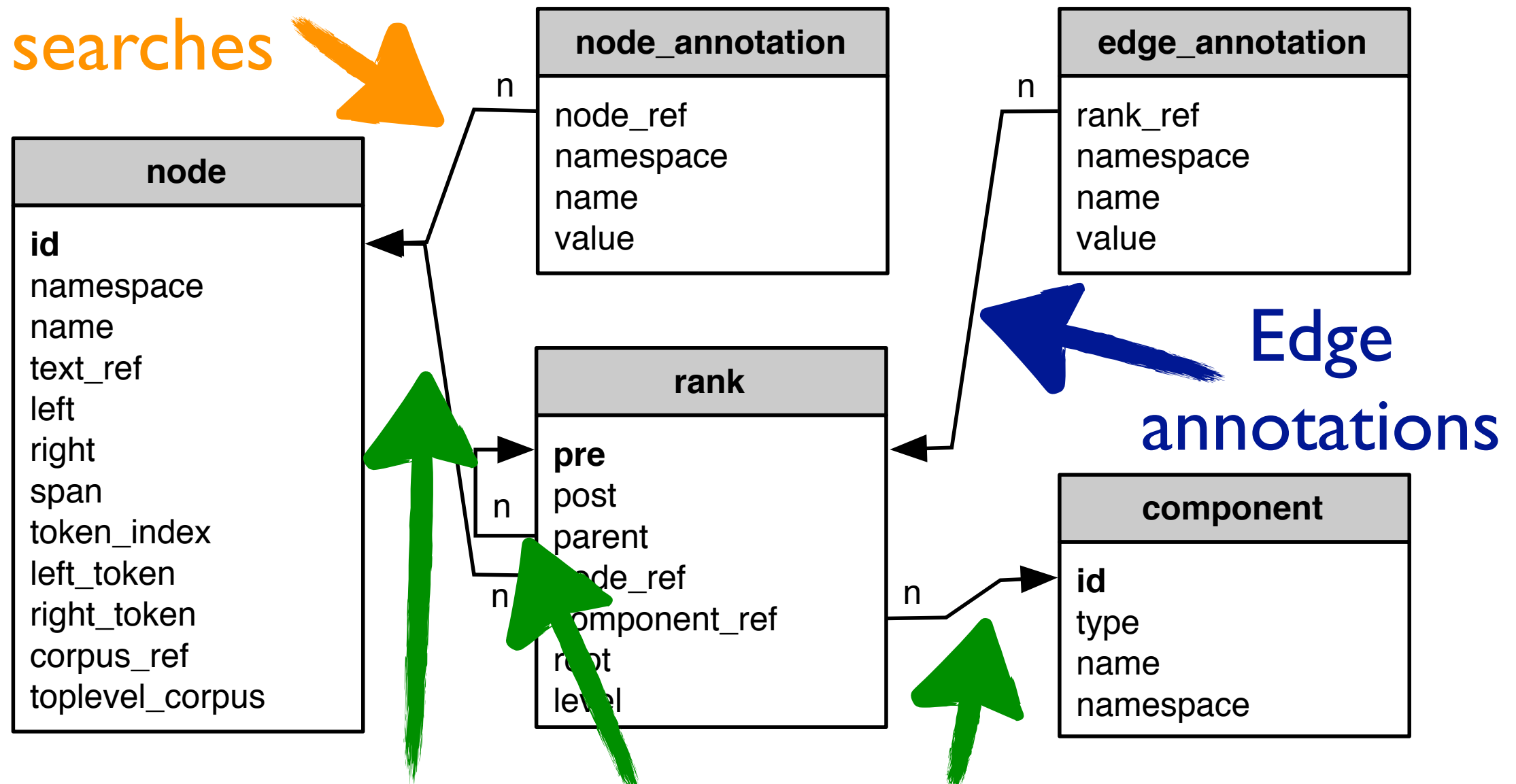
Annotation searches



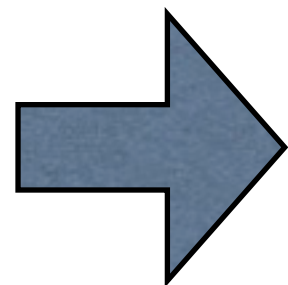
Binary relations on edges

# Many tables – Many joins

Annotation searches

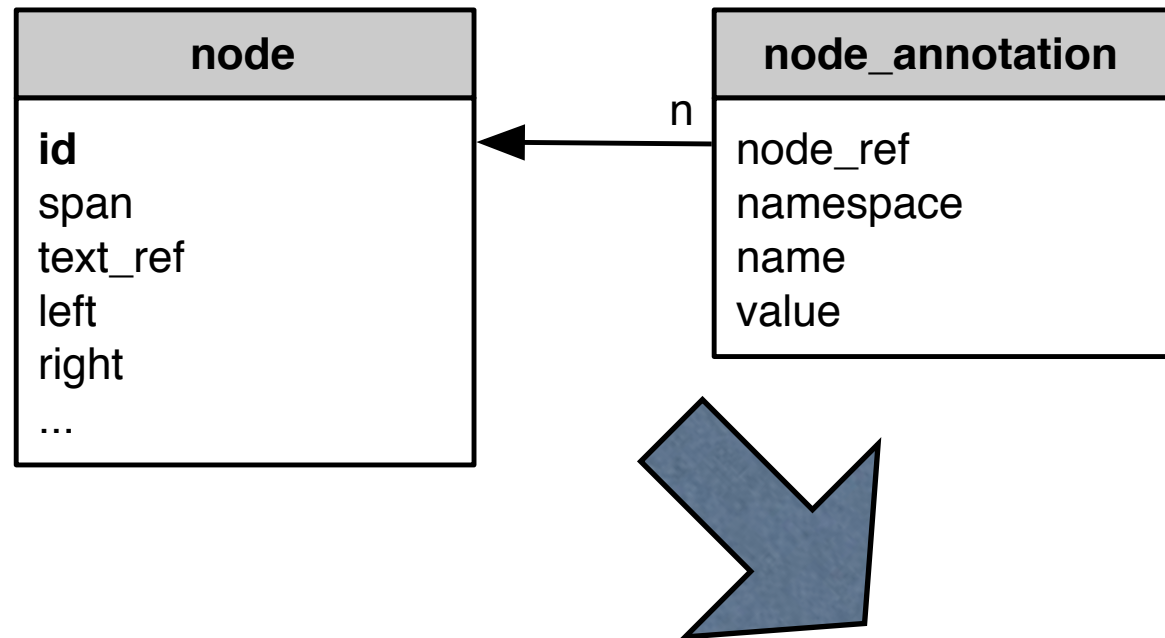


Binary relations on edges



Bad performance on PostgreSQL

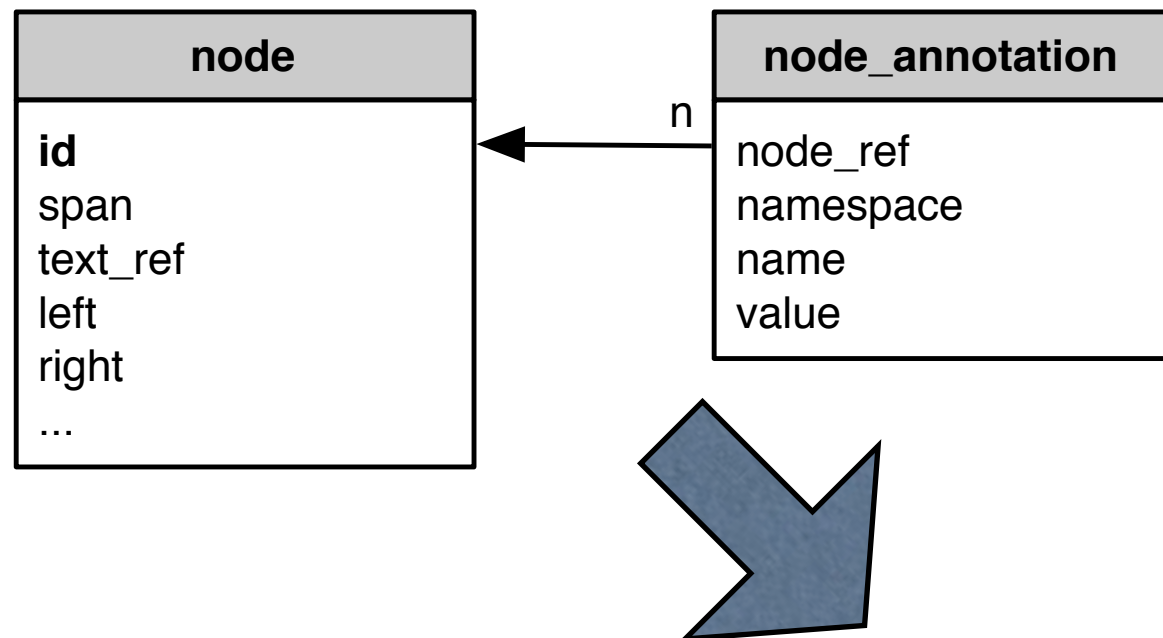
# Solution 1: One big table



id	span	text_ref	left	right	...	na_name	na_value	...
1		1	1	30		cat	S	
2	Wunder	1	1	5		morph	Acc.Pl.Neut	
2	Wunder	1	1	5		pos	NN	
2	Wunder	1	1	5		lemma	Wunder	
...								



# Solution 1: One big table



id	span	text_ref	left	right	...	na_name	na_value	...
1				30		cat	S	
2	Wunder			5		morph	Acc.Pl.Neut	
2	Wunder			5		pos	NN	
2	Wunder			5		lemma	Wunder	
...								

**Pro:** Fewer joins

**Contra:** Increased redundancy, less extensible

# Solution 2: Combined indexes

id	span	text_ref	left	right	...	na_name	na_value	...
1				30		cat	S	
2	Wunder			5		morph	Acc.Pl.Neut	
2	Wunder			5		pos	NN	
2	Wunder			5		lemma	Wunder	
...								

One index over 4 columns

Find nodes

- spanning a certain word,
- in a certain text,
- at a certain position.

cat="S" &

"Wunder" &

#1 \_i\_ #2

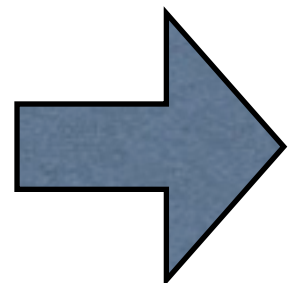
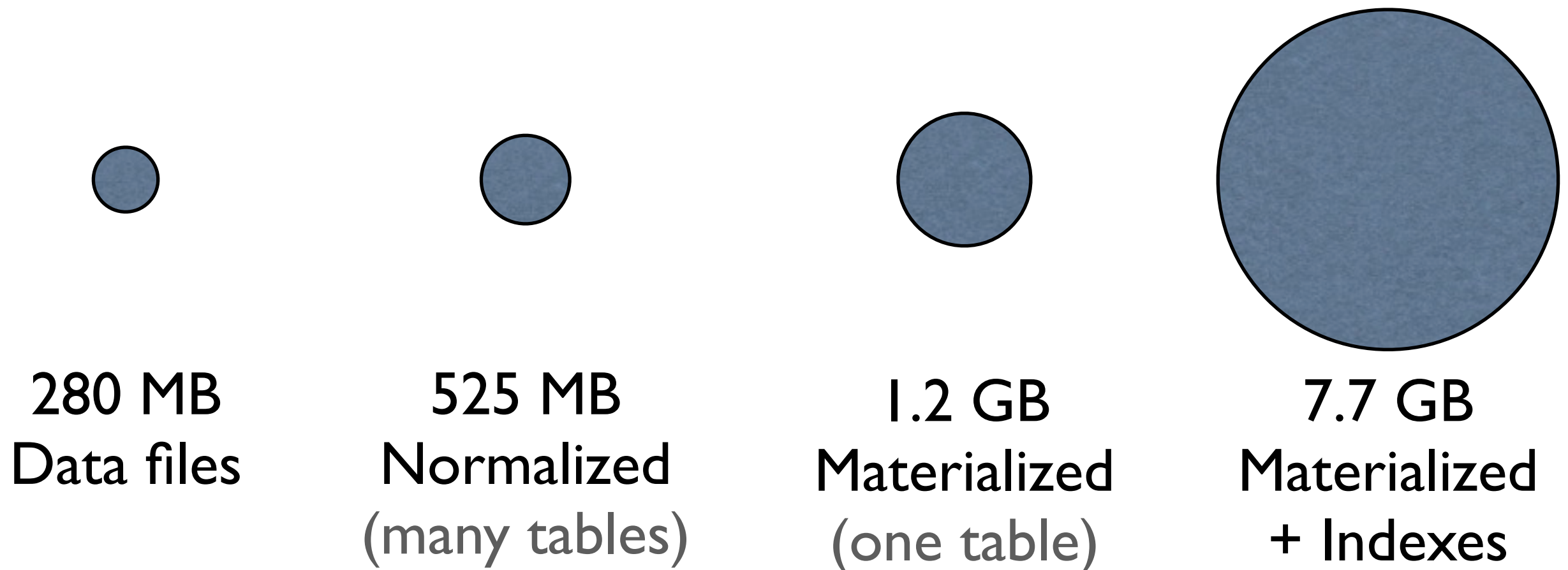
**Pro:** Potentially very fast

**Contra:** Uses lots of disk space

# Disk usage in PostgreSQL

## TIGER Treebank 2.1

ca. 50.000 sentences, 900.000 tokens,  
3 million annotations, 1 million edges



**Increase by factor 15 (or almost 30)**

3. What are Column-Stores? How can Annis benefit?



# What's a Column-Store?

conceptual model

	node_ref	name	value
<b>1</b>	123	pos	VVINF
<b>2</b>	123	lemma	essen
<b>3</b>	456	pos	NN

*table*

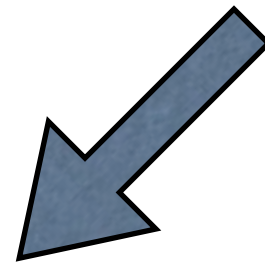
storage model

# What's a Column-Store?

conceptual model

	<b>node_ref</b>	<b>name</b>	<b>value</b>
<b>1</b>	123	pos	VVINF
<b>2</b>	123	lemma	essen
<b>3</b>	456	pos	NN

*table*



storage model

<b>1</b>	123	pos	VVINF
<b>2</b>	123	lemma	essen
<b>3</b>	456	pos	NN

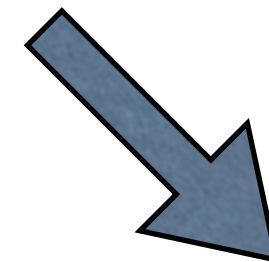
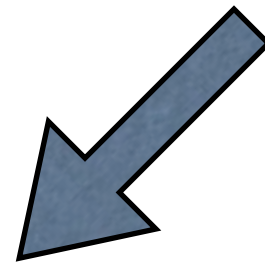
*ROWS*

# What's a Column-Store?

conceptual model

	<b>node_ref</b>	<b>name</b>	<b>value</b>
<b>1</b>	123	pos	VVINF
<b>2</b>	123	lemma	essen
<b>3</b>	456	pos	NN

*table*



storage model

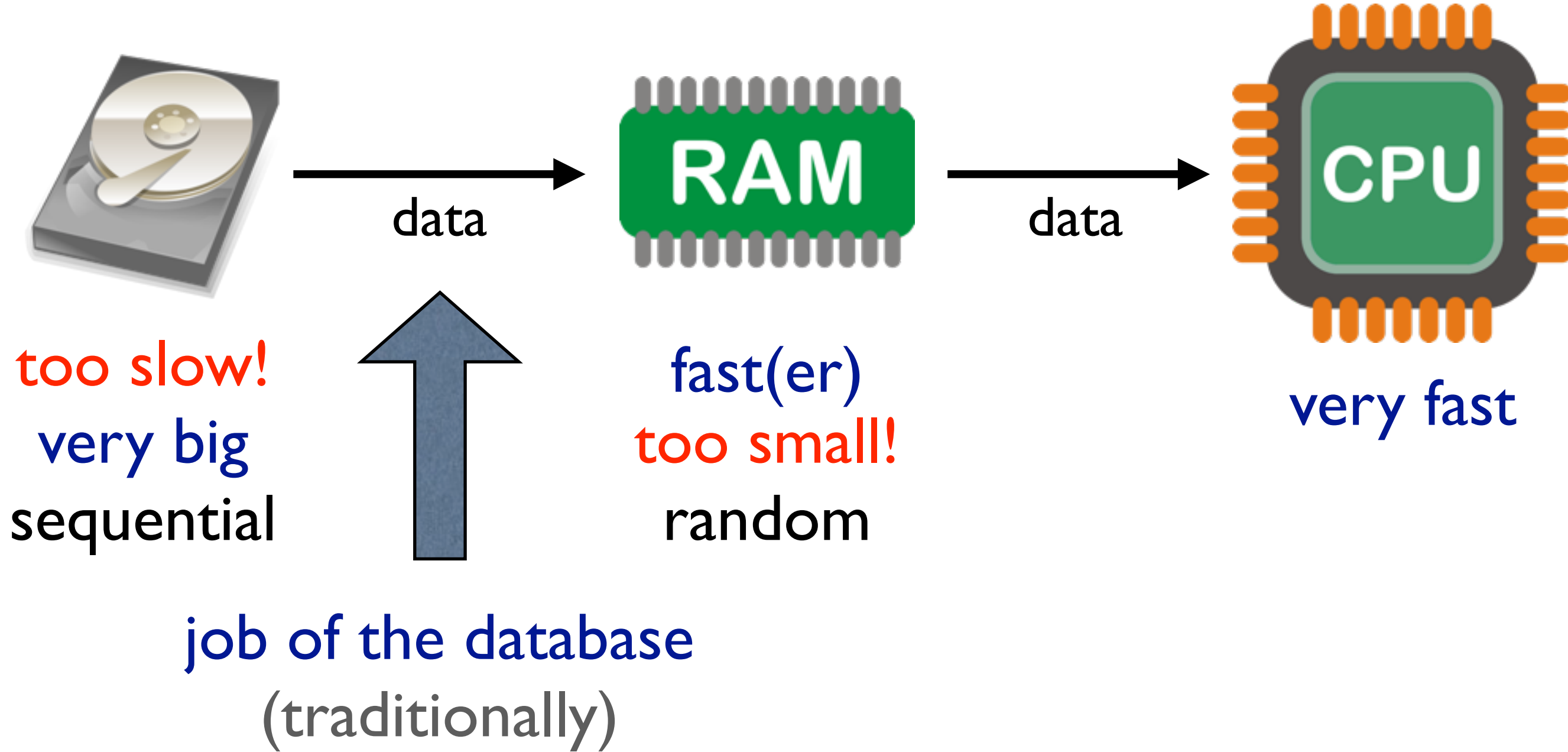
<b>1</b>	123	pos	VVINF
<b>2</b>	123	lemma	essen
<b>3</b>	456	pos	NN

*rows*

<b>node_ref</b>	<b>name</b>	<b>value</b>
123	pos	VVINF
123	lemma	essen
456	pos	NN

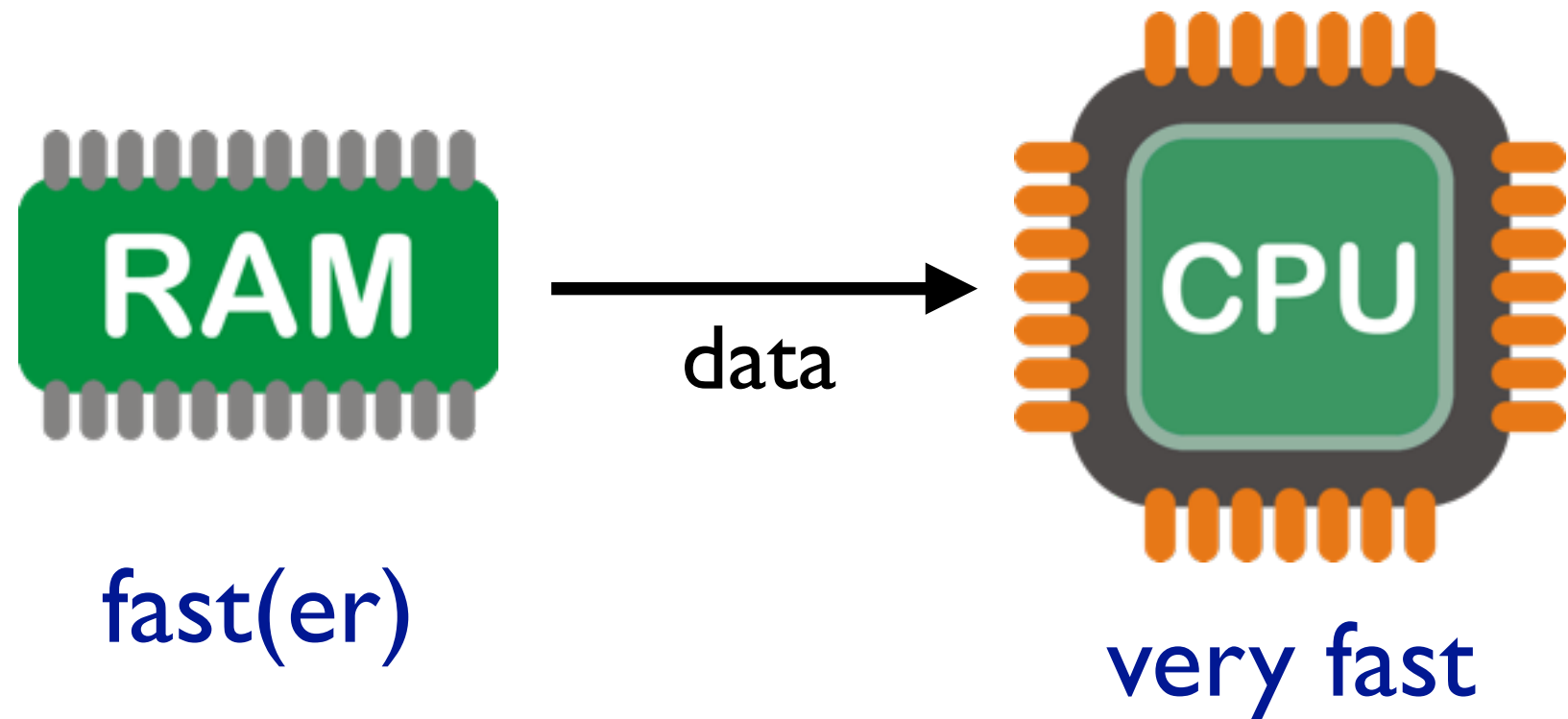
*columns*

# Why Column-Stores? – Why Databases?





# Why Column-Stores? – Why Databases?



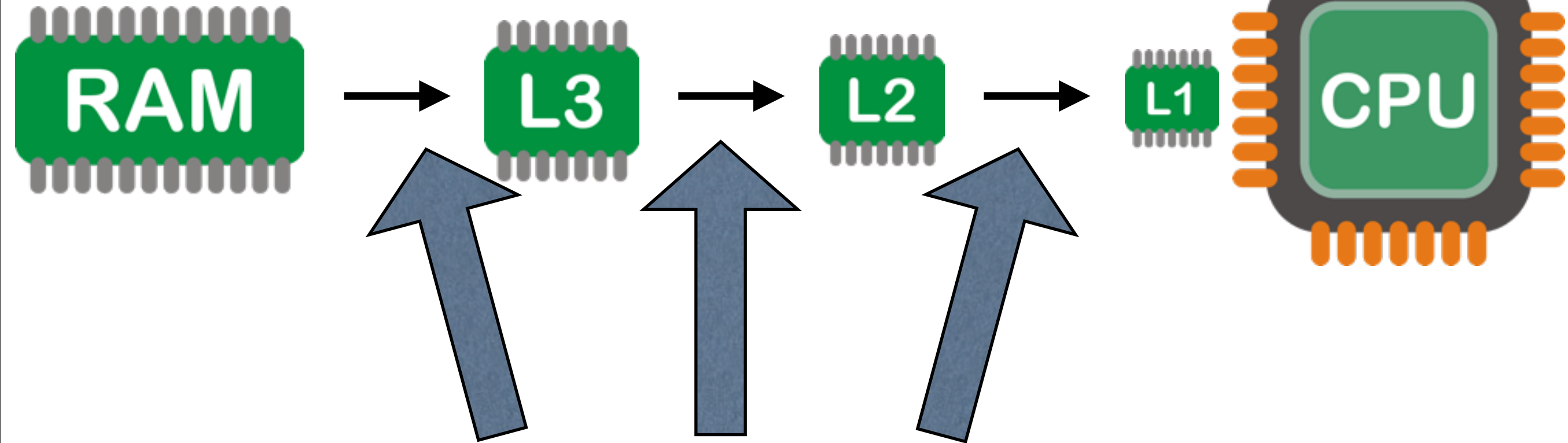
# Caches between RAM and CPU

48 GB  
33 ns

12 MB  
5.4 ns

256 kB  
1.7 ns

32 kB  
1.4 ns



job of the database  
on a modern system  
(among others)

# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

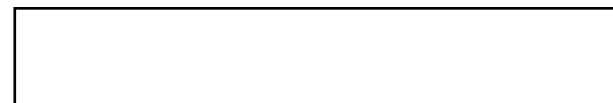
data file:

1	123	pos	VVINF
---	-----	-----	-------

2	123	lemma	essen
---	-----	-------	-------

3	456	pos	NN
---	-----	-----	----

L1 cache:



# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF	2	123	lemma	essen	3	456	pos	NN
---	-----	-----	-------	---	-----	-------	-------	---	-----	-----	----

L1 cache:



1. load first row



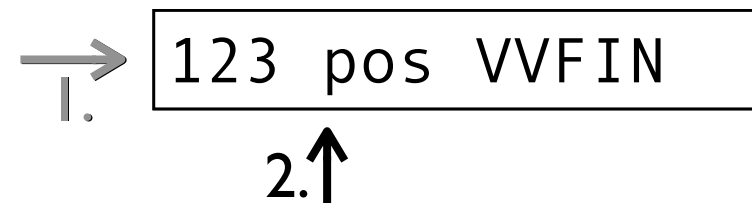
# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

L1 cache:



1. load first row
2. locate *name* attribute

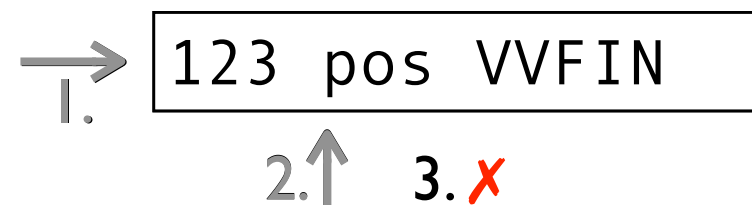
# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

L1 cache:



1. load first row
2. locate *name* attribute
3. test *name* attribute

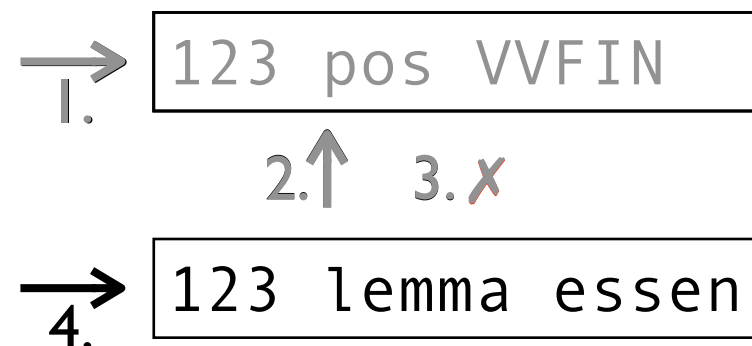
# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

L1 cache:



1. load first row
2. locate *name* attribute
3. test *name* attribute
4. load second row

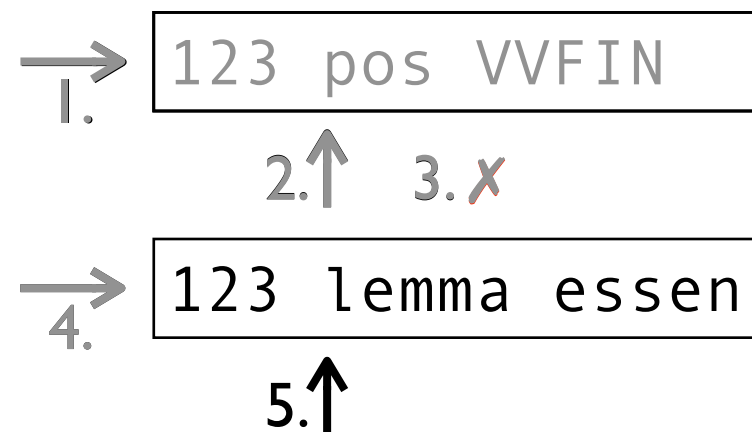
# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

L1 cache:



1. load first row
2. locate *name* attribute
3. test *name* attribute
4. load second row
5. locate *name* attribute

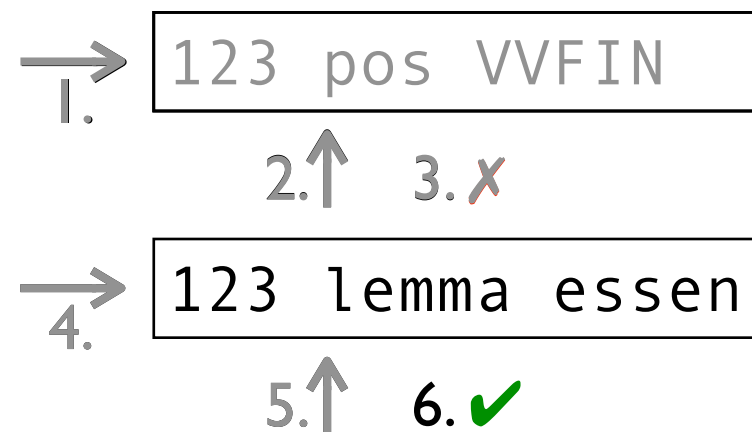
# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

L1 cache:



1. load first row
2. locate *name* attribute
3. test *name* attribute
4. load second row
5. locate *name* attribute
6. test *name* attribute



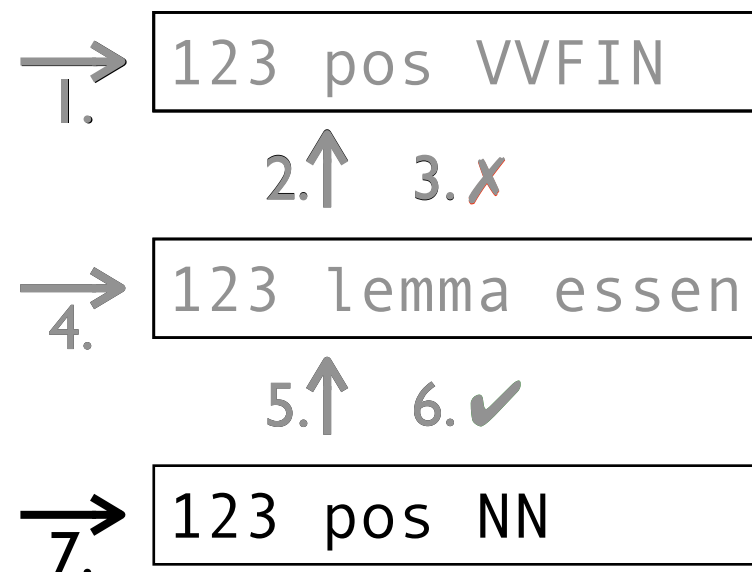
# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

L1 cache:



1. load first row
2. locate *name* attribute
3. test *name* attribute
4. load second row
5. locate *name* attribute
6. test *name* attribute
7. load third row

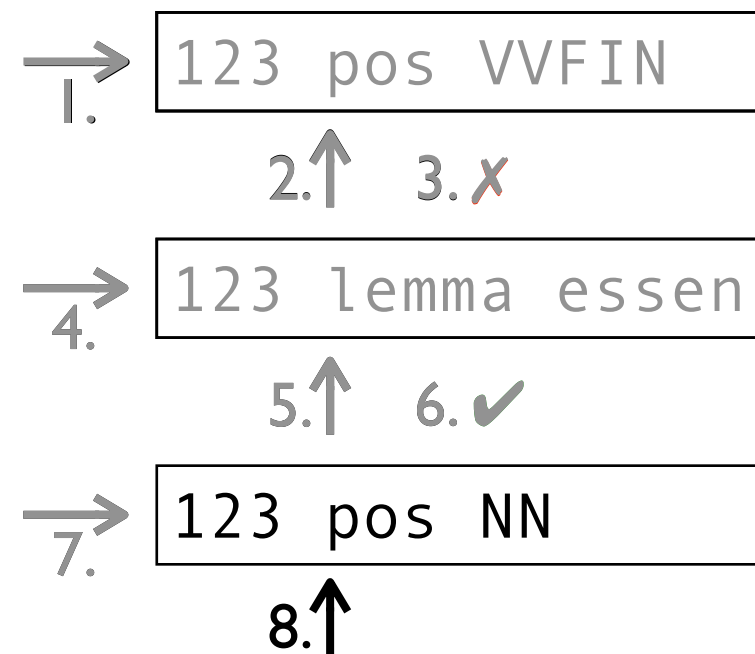
# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

L1 cache:



1. load first row
2. locate *name* attribute
3. test *name* attribute
4. load second row
5. locate *name* attribute
6. test *name* attribute
7. load third row
8. locate *name* attribute

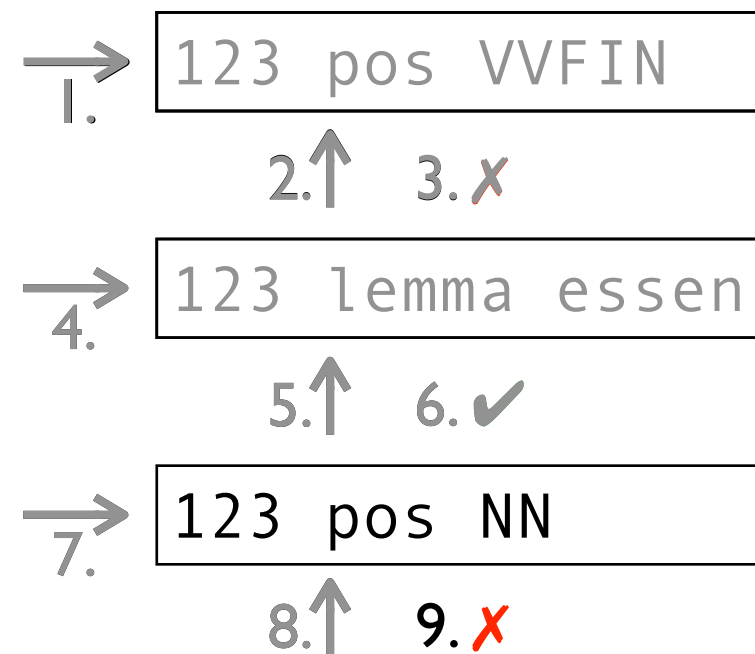
# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

L1 cache:



1. load first row
2. locate *name* attribute
3. test *name* attribute
4. load second row
5. locate *name* attribute
6. test *name* attribute
7. load third row
8. locate *name* attribute
9. test *name* attribute

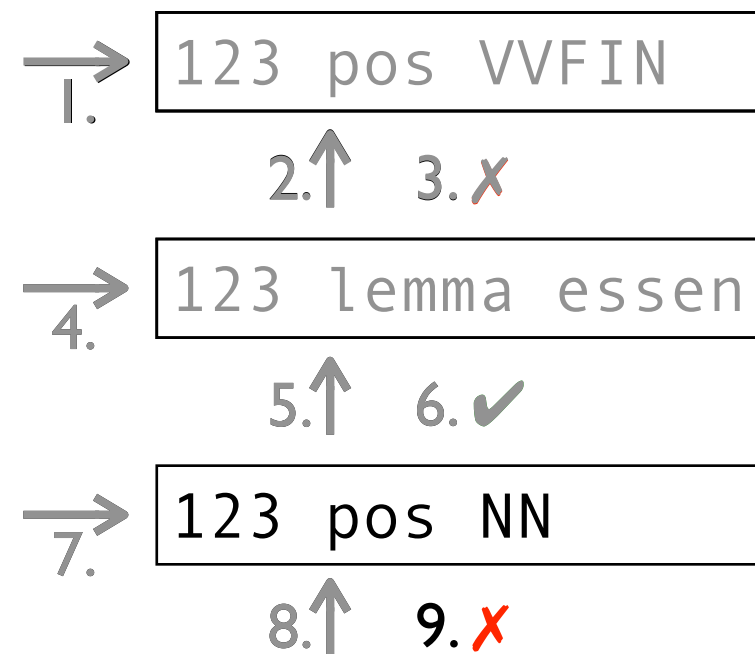
# Cache usage of row layout

query: compare *name* attribute with value 'lemma'

data file:

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

L1 cache:



1. load first row
2. locate *name* attribute
3. test *name* attribute
4. load second row
5. locate *name* attribute
6. test *name* attribute
7. load third row
8. locate *name* attribute
9. test *name* attribute

# Cache usage of column layout

query: compare *name* attribute with value 'lemma'

data file:

node_ref	name	value
123	pos	VVINF
123	lemma	essen
456	pos	NN

L1 cache:



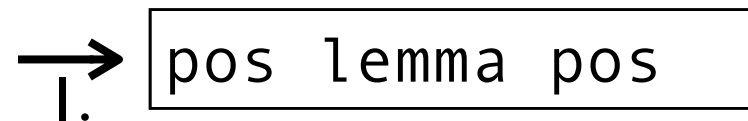
# Cache usage of column layout

query: compare *name* attribute with value 'lemma'

data file:

node_ref	name	value
123	pos	VVINF
123	lemma	essen
456	pos	NN

L1 cache:



1. load *name* column

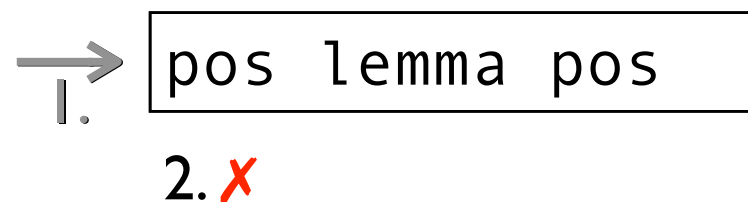
# Cache usage of column layout

query: compare *name* attribute with value 'lemma'

data file:

node_ref	name	value
123	pos	VVINF
123	lemma	essen
456	pos	NN

L1 cache:



1. load *name* column
2. test first *name* attribute

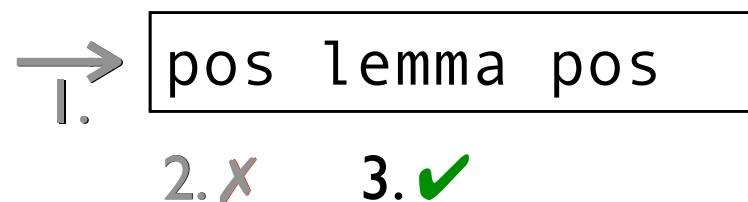
# Cache usage of column layout

query: compare *name* attribute with value 'lemma'

data file:

node_ref	name	value
123	pos	VVINF
123	lemma	essen
456	pos	NN

L1 cache:



1. load *name* column
2. test first *name* attribute
3. test second *name* attribute

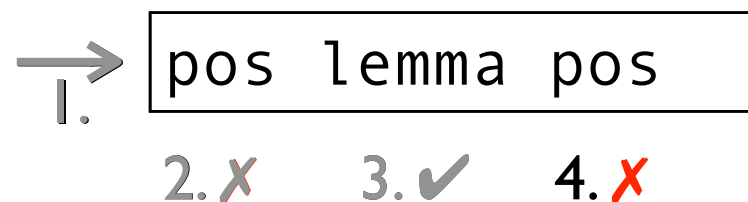
# Cache usage of column layout

query: compare *name* attribute with value 'lemma'

data file:

node_ref	name	value
123	pos	VVINP
123	lemma	essen
456	pos	NN

L1 cache:



1. load *name* column
2. test first *name* attribute
3. test second *name* attribute
4. test third *name* attribute

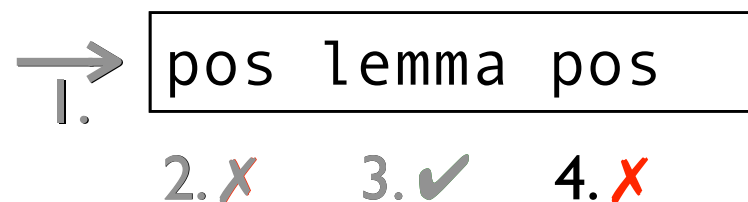
# Cache usage of column layout

query: compare *name* attribute with value 'lemma'

data file:

node_ref	name	value
123	pos	VVINF
123	lemma	essen
456	pos	NN

L1 cache:



1. load *name* column
2. test first *name* attribute
3. test second *name* attribute
4. test third *name* attribute

# Column operations in Annis

## Search terms

- can be indexed

"Wunder"

## Regular expressions

- can often be indexed
- but not always

morph=/.\*\.PL\.Neut/

## Binary operations

- can be indexed
- need many indexes
- slow if there are many index lookups

`_ = _`    `_ i _`    `>`

id	span	text	ref	left	right	...
1					30	
2	Wunder				5	
2	Wunder				5	
2	Wunder				5	
...						



## 4. New implementation on MonetDB and evaluation

# Prototype implementation

## Supported

- Annis 2 Query Language
- COUNT queries

## Not supported

- Annis 3 language features
- ANNOTATE, MATRIX queries
- corpus selection

The screenshot displays the Annis² Corpus Search interface. On the left, the 'Search Form' shows a query: `cat="S" & "Wunder" & #1..#2`. A large green checkmark is overlaid on this query. Below the search form is a table of corpora with columns for Name, Texts, and Tokens. A large red 'X' is overlaid on this table. The main search result area shows the search result for the query: `cat="S" & "Wunder" & #1..#2 (5, 5)`. It displays the path `d1.pcc2 = 4282` and the search results: `Stelpass Wunder gibt es immer wieder ! Erst spielen die Dalglower`. Below the search results are two views: 'Constituents (Tree View)' and 'Dependencies (Arches View)'. Both views show a tree structure for the sentence. A large red 'X' is overlaid on the 'Dependencies (Arches View)' section. At the bottom of the interface, there are options for 'Context Left', 'Context Right', and 'Results Per Page'.

# Realistic test workload

**Corpus:** TIGER Treebank 2.1

## **Queries:**

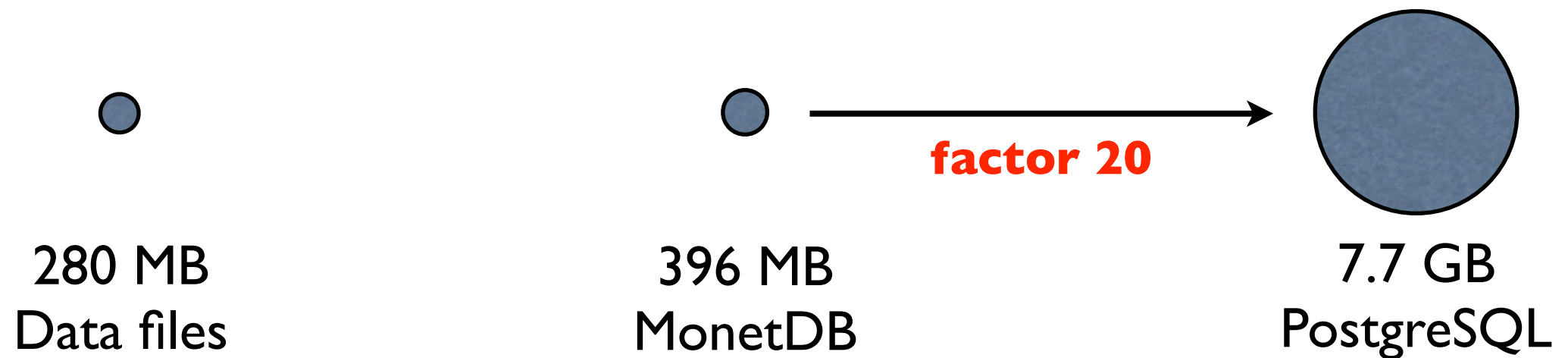
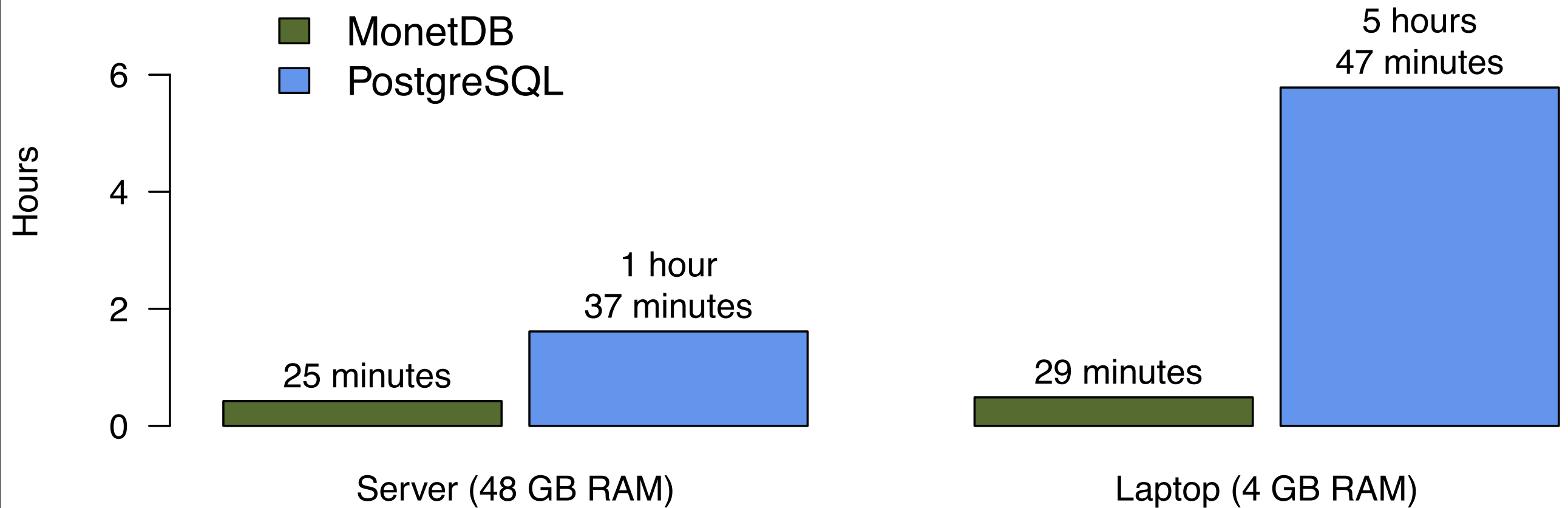
- 3 month query log of Annis instance at the SFB 632
- 337 TIGER queries (224 unique)
- up to 4 search terms
- up to 6 binary operators

## **Random workload:**

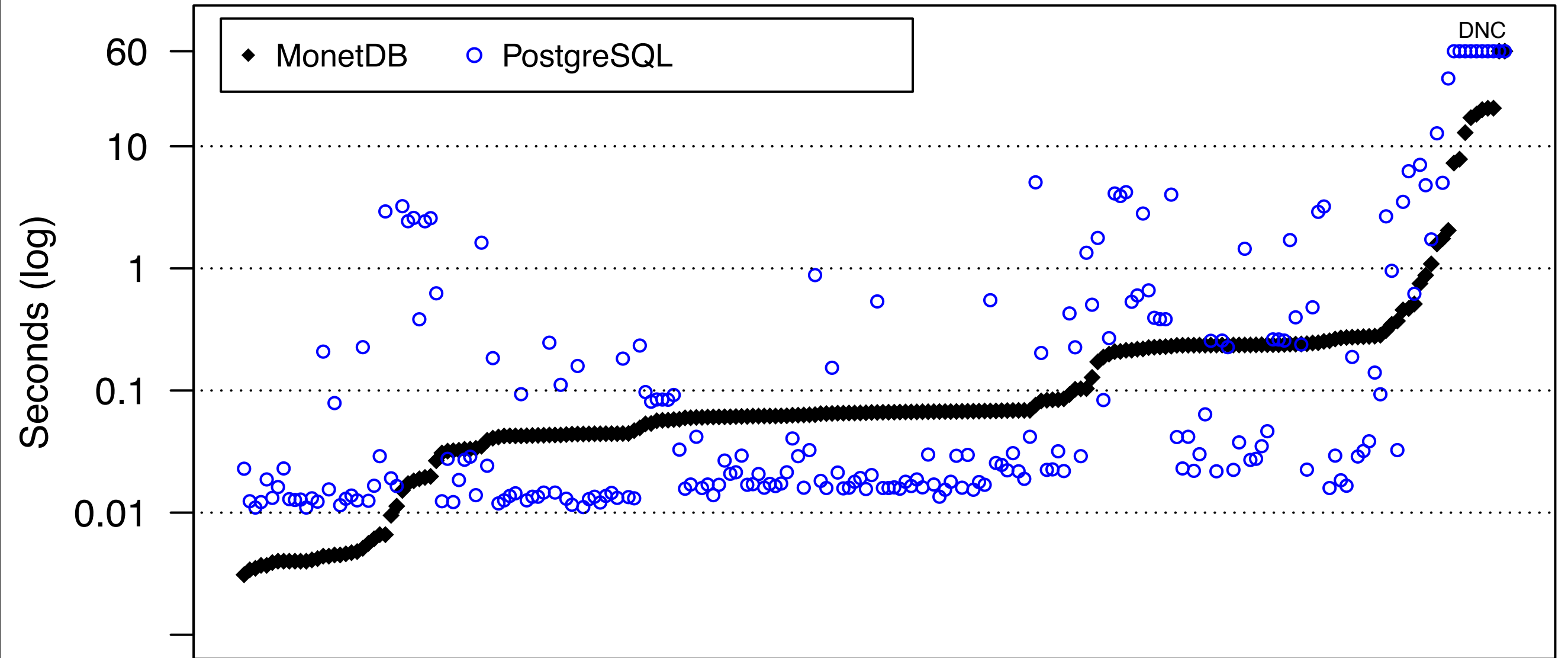
- 10000 queries
- original distribution
- excluded PostgreSQL timeout



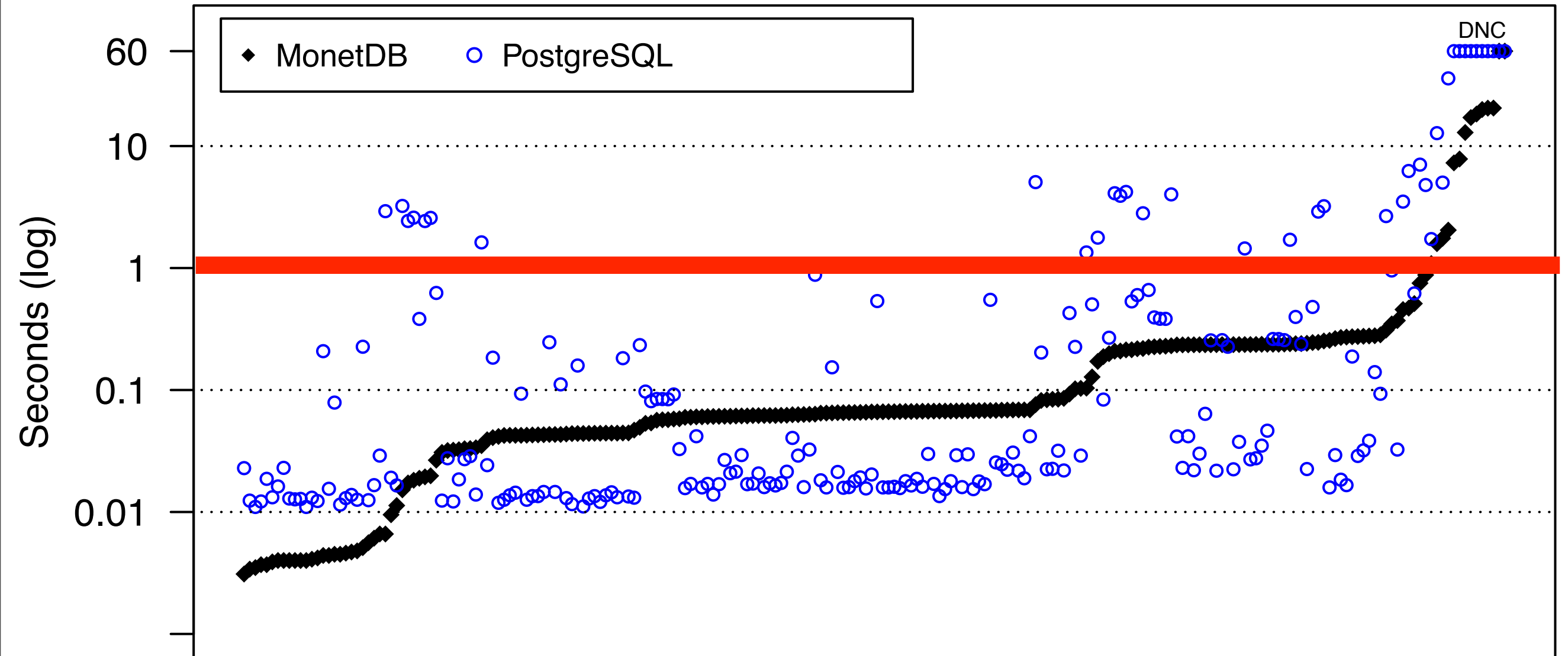
# Workload of 10000 queries



# Individual query performance

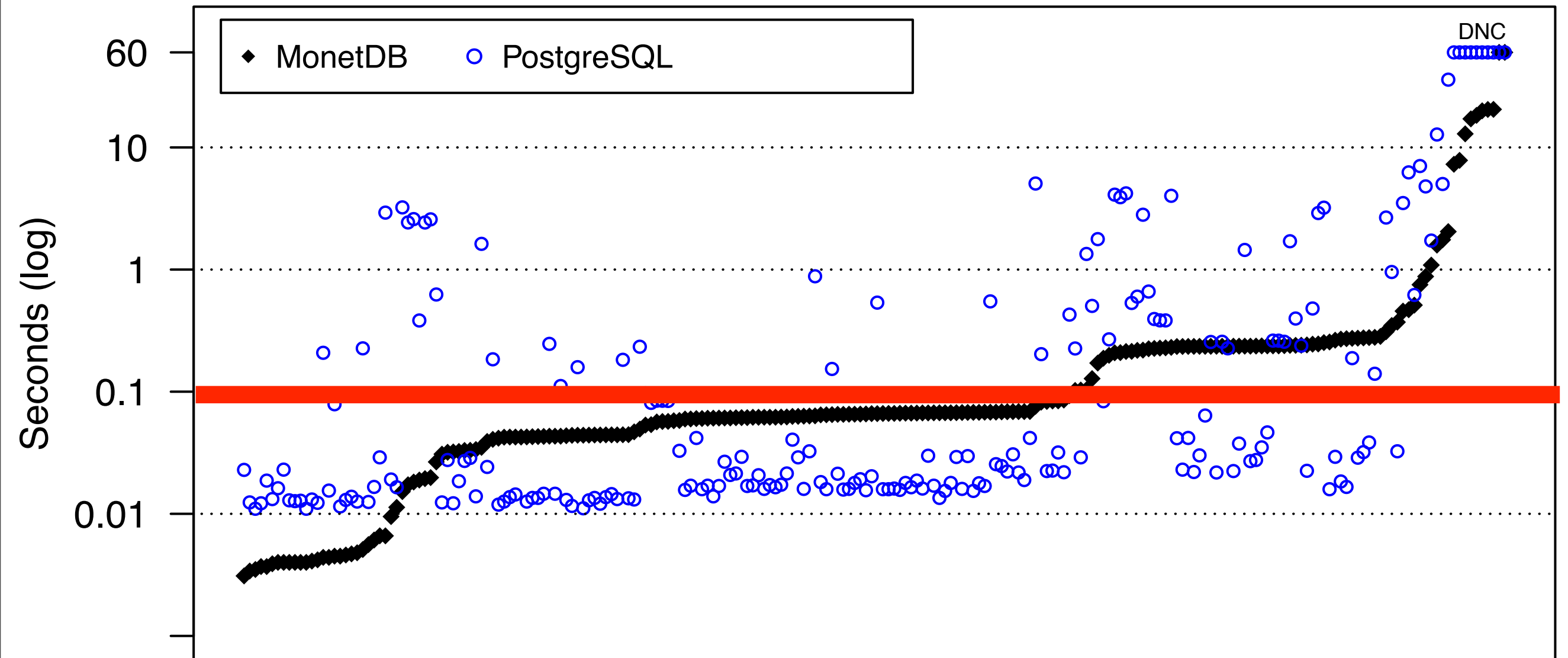


# Individual query performance

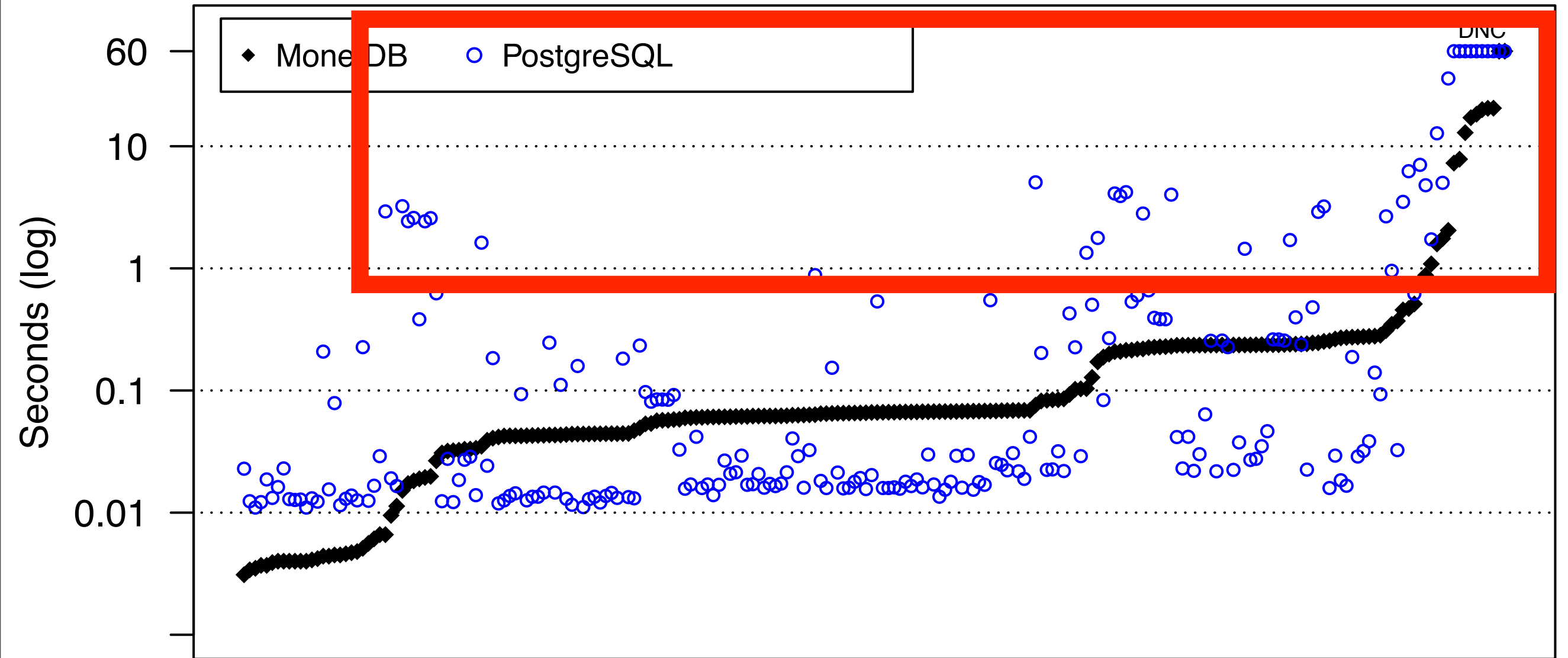




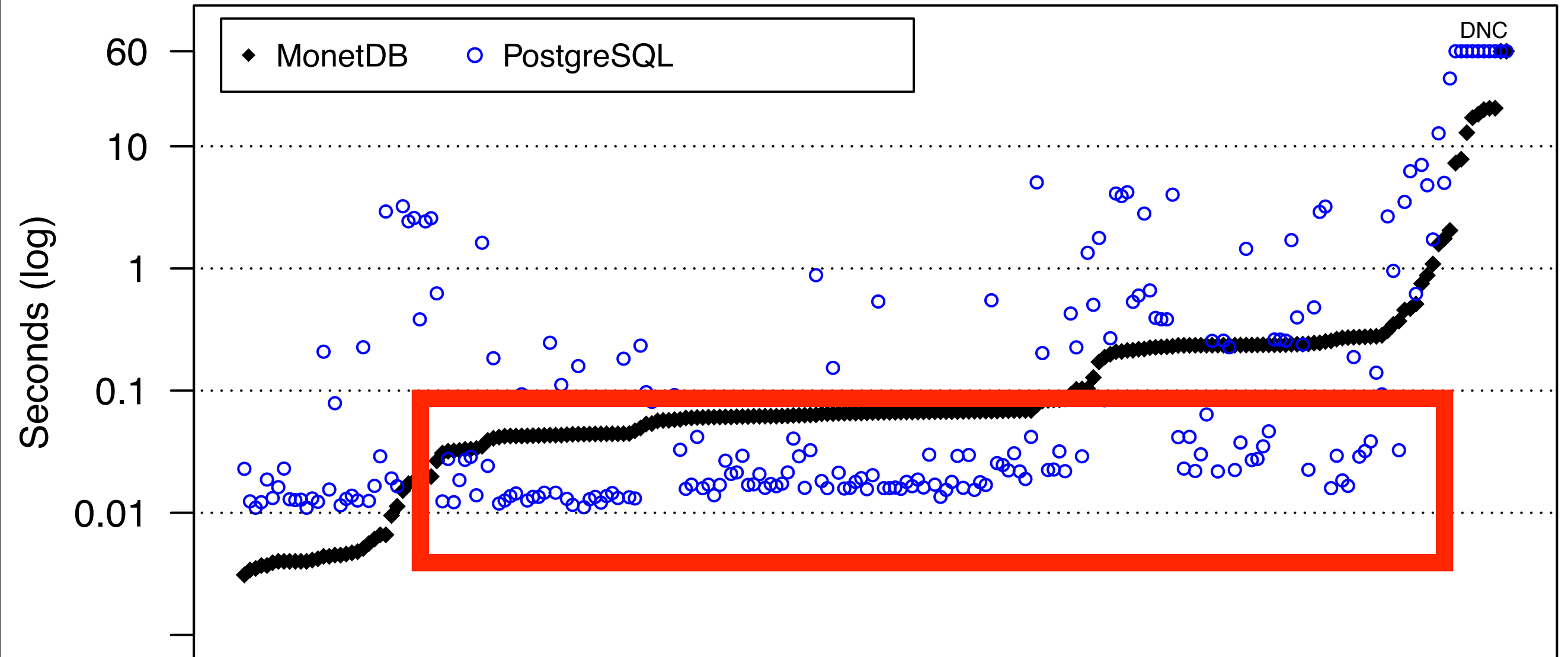
# Individual query performance



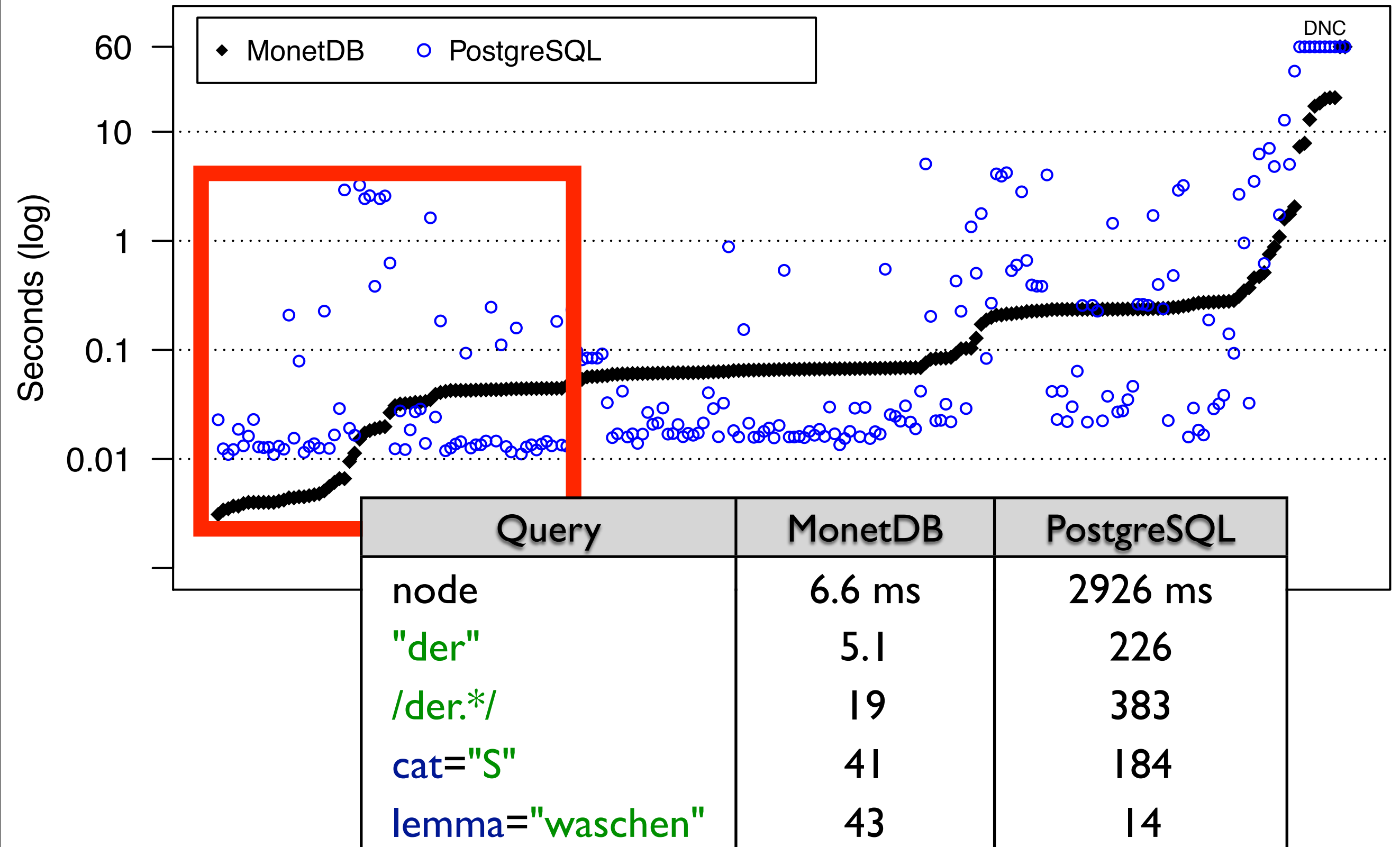
# Individual query performance



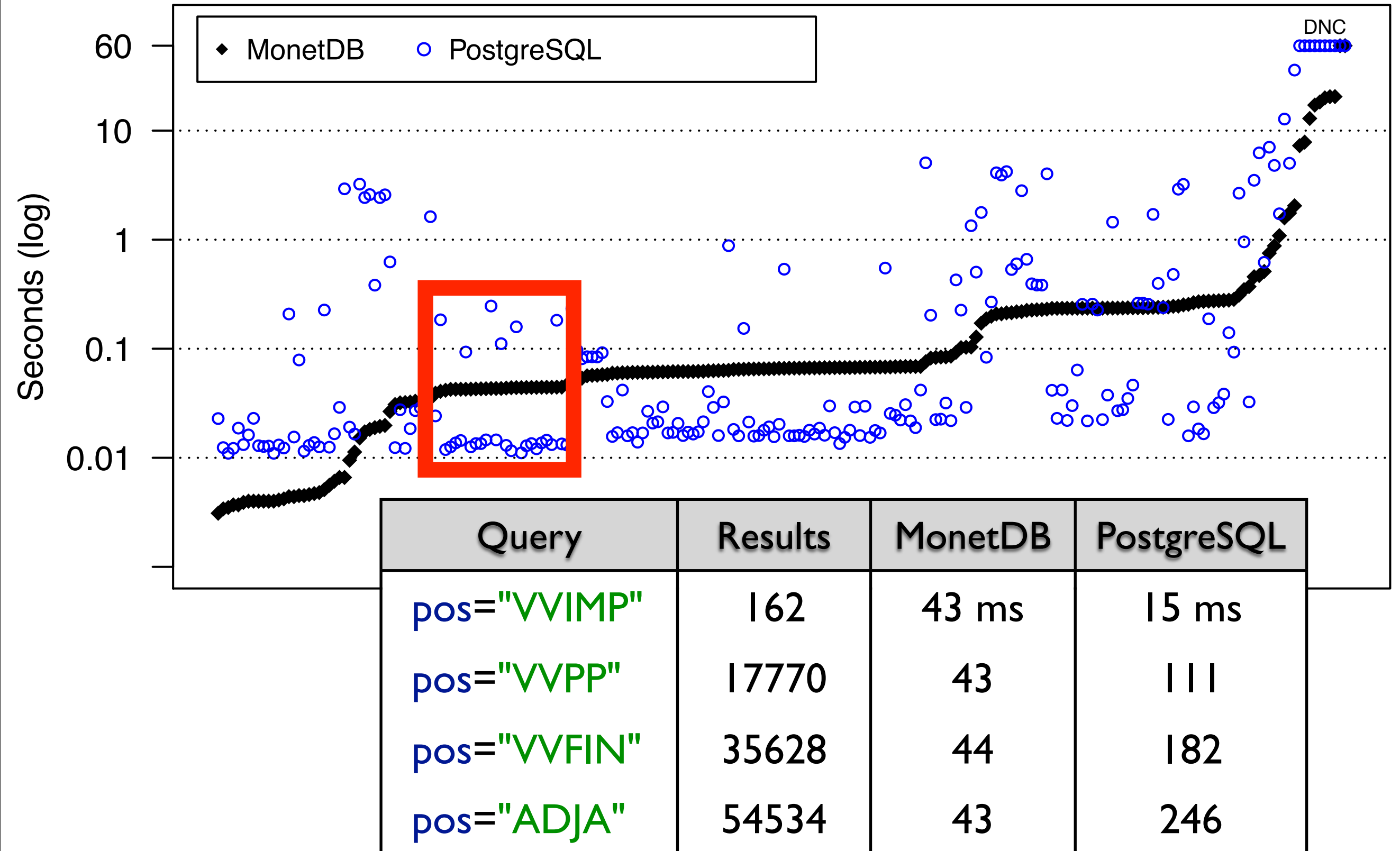
# Individual query performance



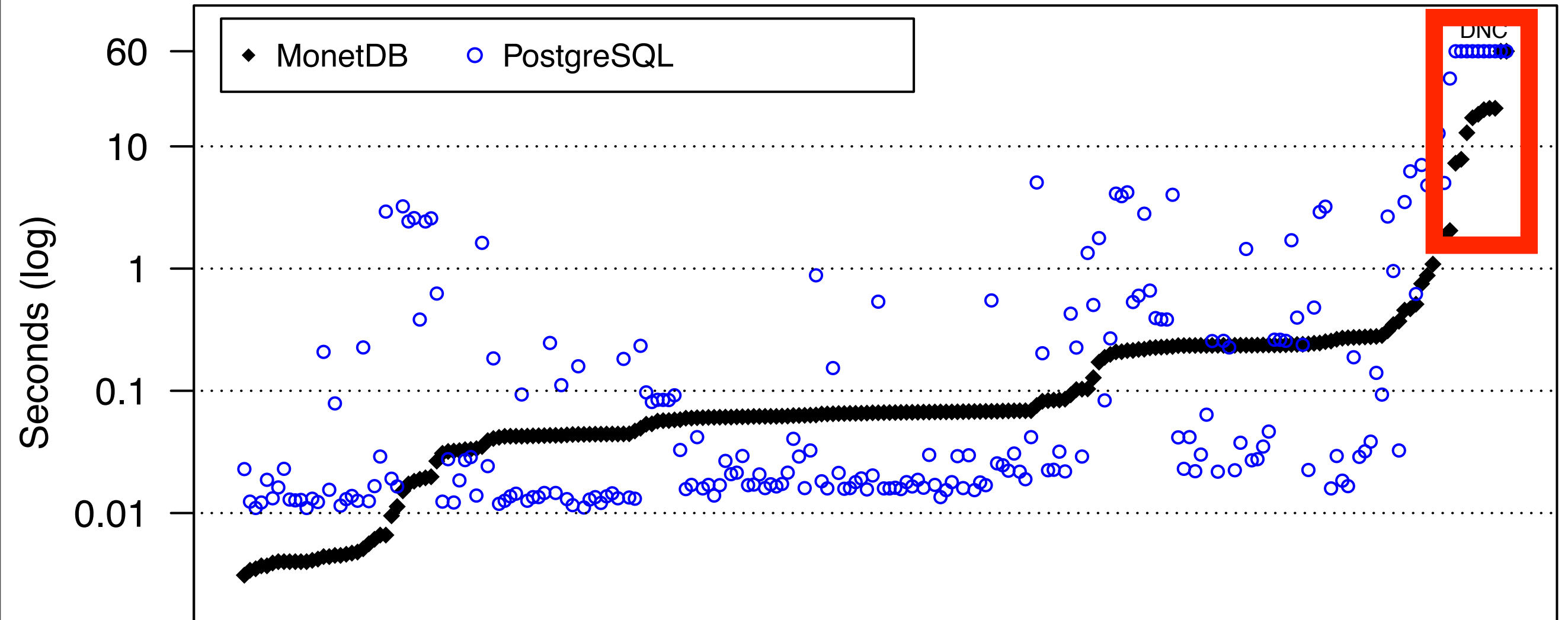
# Simple queries are fast



# Influence of result size



# Queries with millions of results



Query	Results	Monet	PSQL
lemma="müssen" & pos= /VV.* / & pos="\$." & #1 .* #2 & #2 .* #3	4.5 M	2 s	35 s
pos=/VM.* / & pos= /VV.* / & pos=/.*/ & #1 .* #2 & #2 .* #3	384 M	175 s	> 1 h

# Fast regular expressions

regular expression without a fixed prefix  
can't use an index, need to scan the entire column

Query	MonetDB	PostgreSQL
<code>/.*sich.*/</code>	213 ms	4206 ms
<code>/[Kk]ann.*/</code>	219	2812
<code>pos="VVPP" &amp; lemma=/(ge)?kommen/ &amp; #1 _=_ #2</code>	229	383
<code>pos=/N.*/ &amp; /[12][09][0-9][0-9]/ &amp; #1 _=_ #2</code>	246	2902
<code>lemma=/[^äöü]+/ &amp; /.+[äöü].+/ &amp; pos="NN" &amp; #1 _=_ #2 &amp; #2 _=_ #3</code>	469	6246



# Advantages

## **MonetDB**

- better overall performance
- stable query performance
- fast regular expressions
- normalized schema
- greatly reduced disk consumption
- better use of limited resources

## **PostgreSQL**

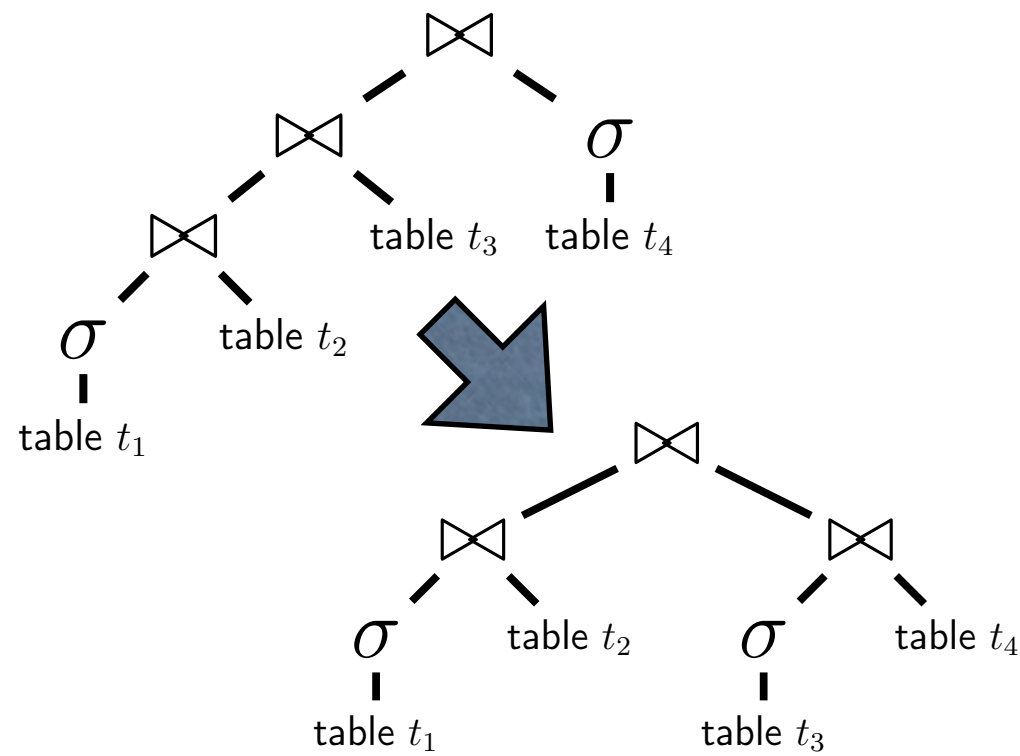
- queries with highly selective search term
- complete implementation
- bug-free SQL processing

# Summary

- prototypical implementation of Annis on MonetDB
- test scenario from an Annis installation in service
- in-depth performance comparison of Annis on MonetDB and PostgreSQL

```
SELECT vielen  
FROM dank;
```

# Port of Annis to MonetDB



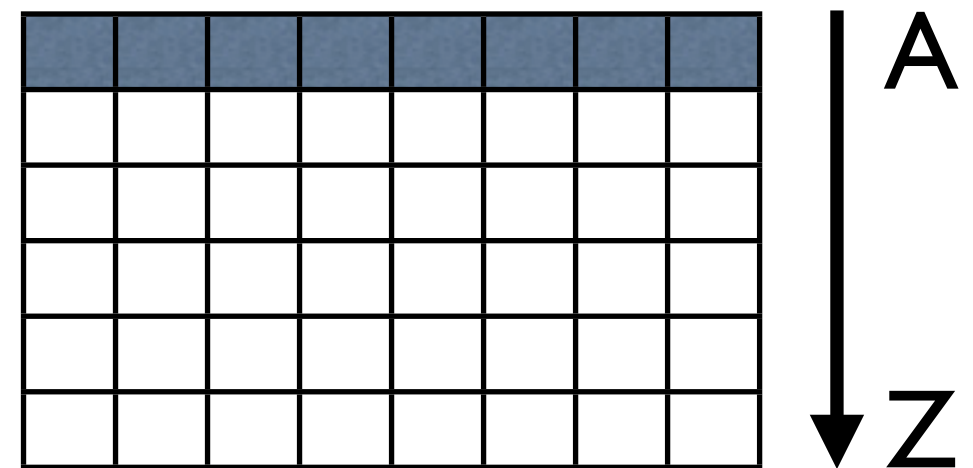
execution plan

**SELECT ~~DISTINCT~~**

SQL optimizations

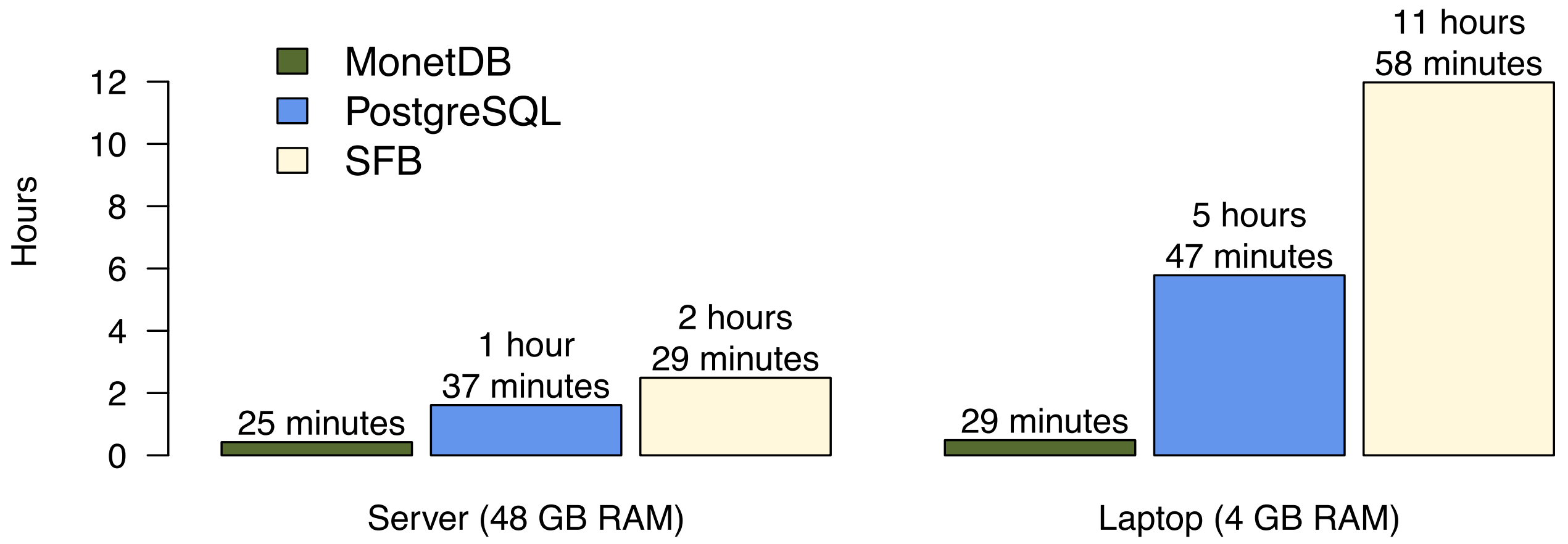
$(a|b)^*$

regular expressions



sorted tables

# PostgreSQL vs. MonetDB



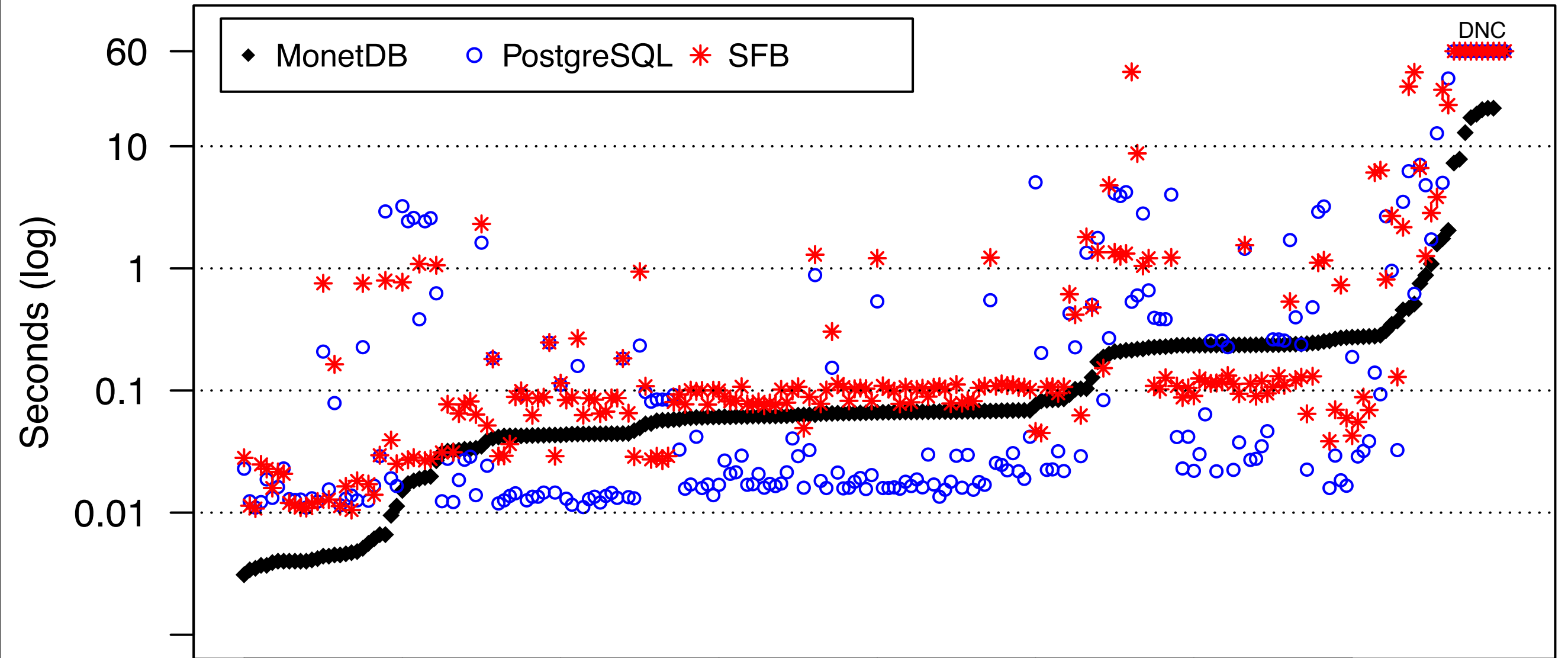
280 MB  
Data files

396 MB  
MonetDB

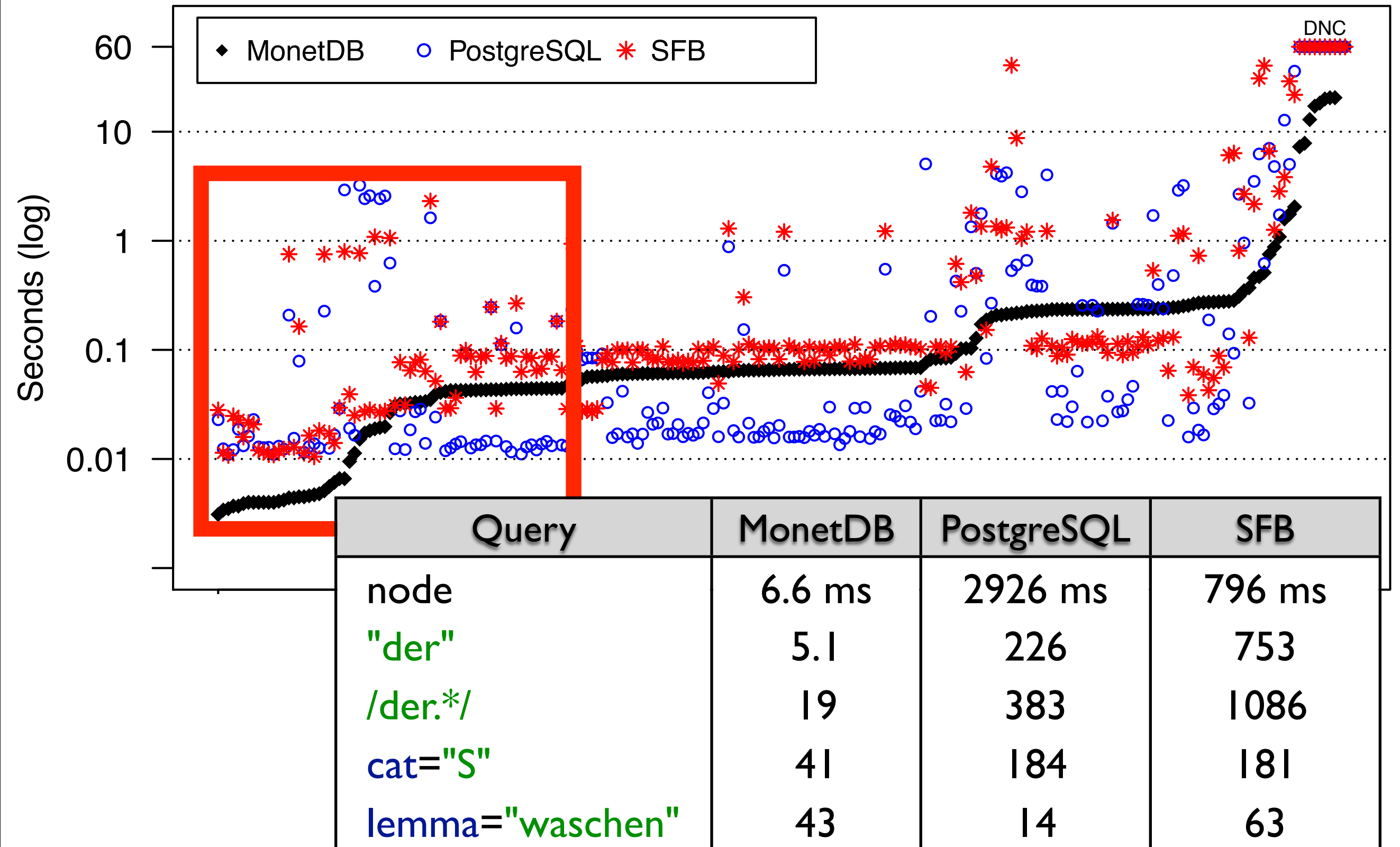
7.7 GB  
Studienarbeit

7.5 GB  
SFB

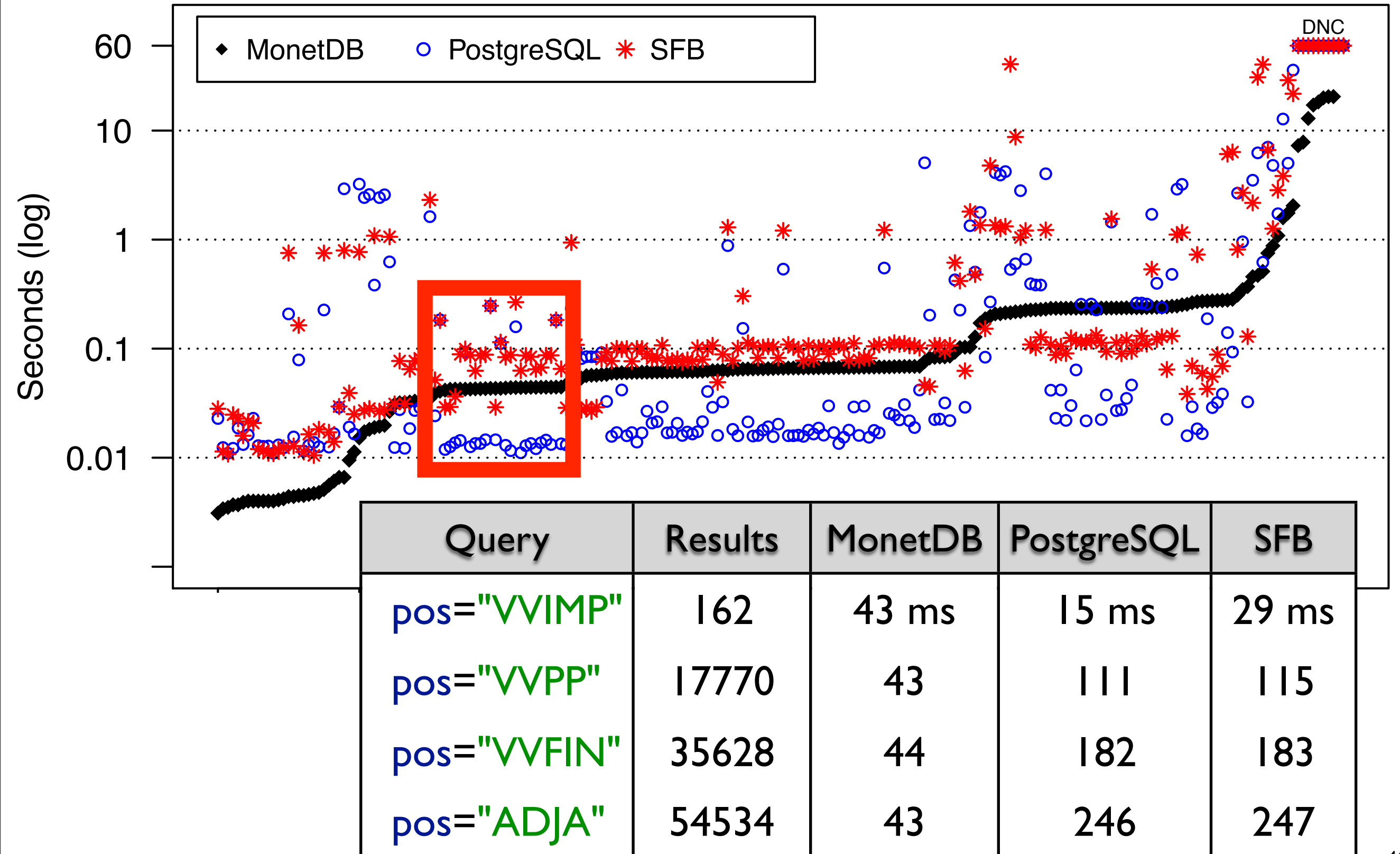
# Individual query performance



# Simple queries

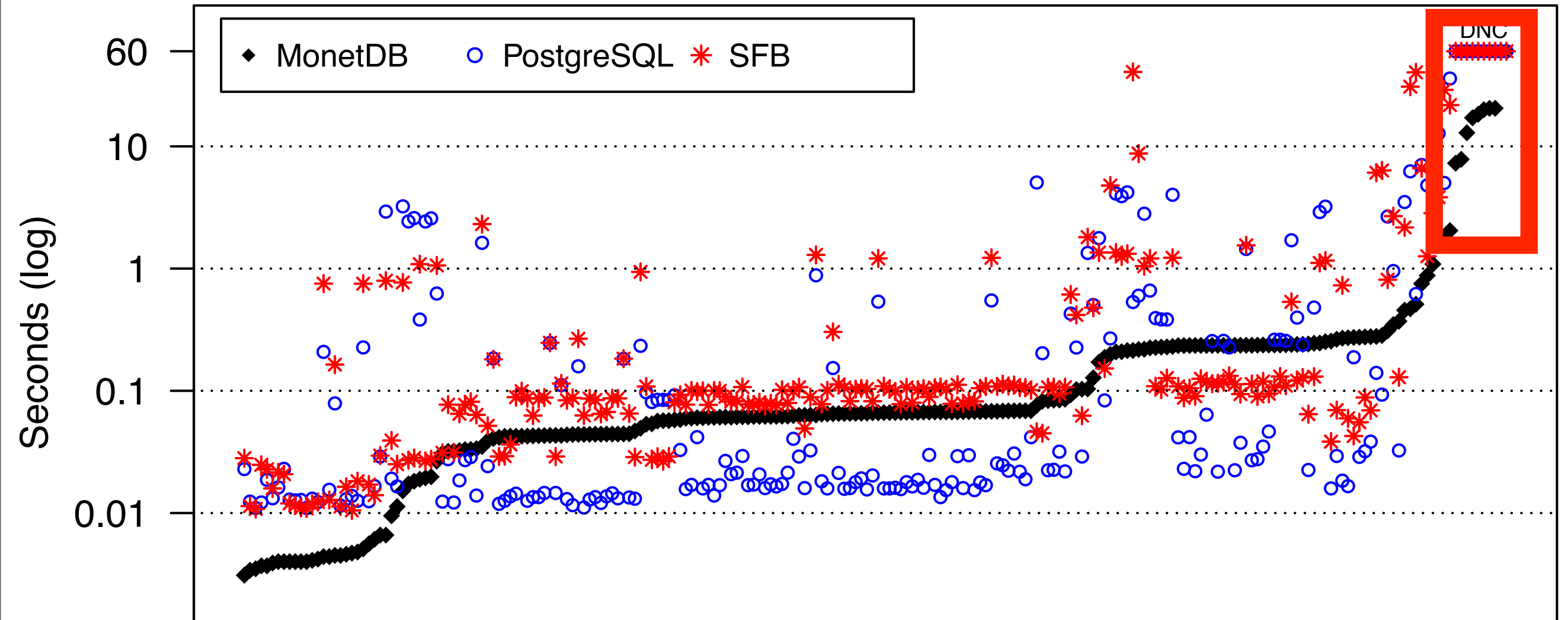


# Influence of result size





# Queries with millions of results



DNC

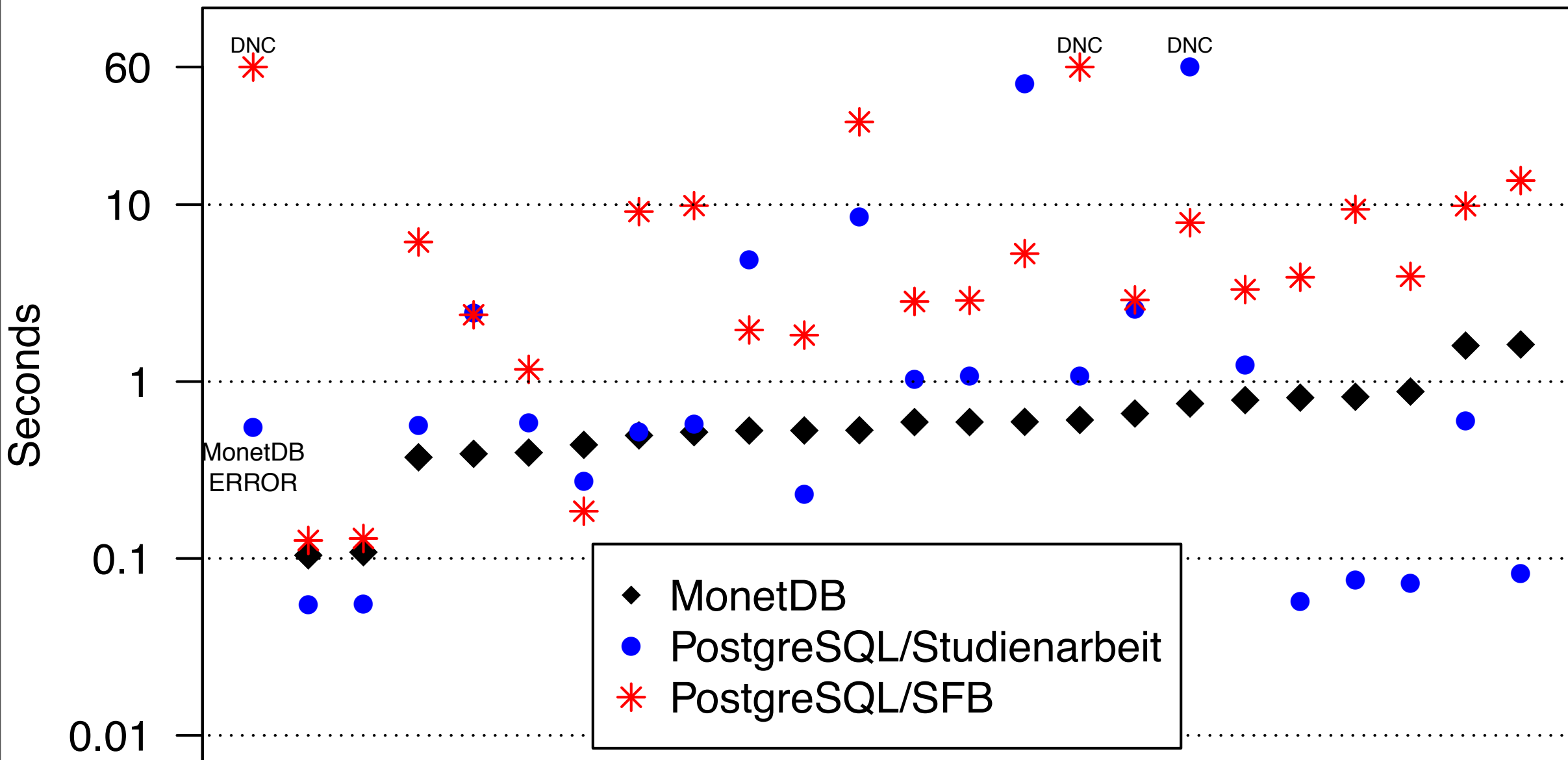
Query	Monet	PSQL	SFB
lemma="müssen" & pos= /VV.* / & pos="\$." & #1 .* #2 & #2 .* #3	2 s	35 s	22 s
pos=/VM.* / & pos= /VV.* / & pos=/.*/ & #1 .* #2 & #2 .* #3	175 s	> 1 h	36 m

# Regular expressions

regular expression without a fixed prefix  
can't use an index, need to scan the entire column

Query	MonetDB	PSQL	SFB
<code>/.*sich.*/</code>	213 ms	4206 ms	1322 ms
<code>/[Kk]ann.*/</code>	219	2812	1047
<code>pos="VVPP" &amp; lemma=/(ge)?kommen/ &amp; #1 _ = _ #2</code>	229	383	127
<code>pos=/N.*/ &amp; /[12][09][0-9][0-9]/ &amp; #1 _ = _ #2</code>	246	2902	1111
<code>lemma=/[^äöü]+/ &amp; /.+[äöü].+/ &amp; pos="NN" &amp; #1 _ = _ #2 &amp; #2 _ = _ #3</code>	469	6246	30803

# Example queries on TüBa-D/Z



1.2 GB  
Data files

1.5 GB  
MonetDB

36 GB  
Studienarbeit

39 GB  
SFB

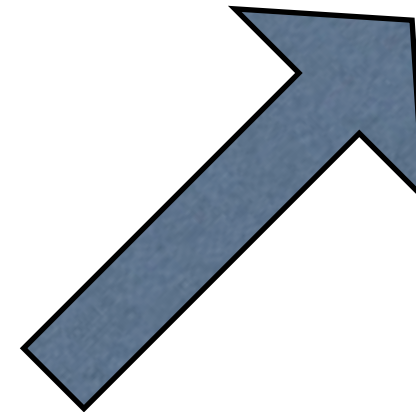
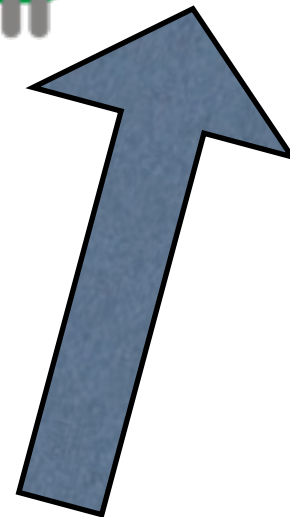
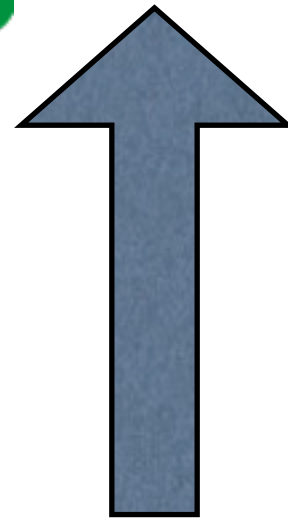
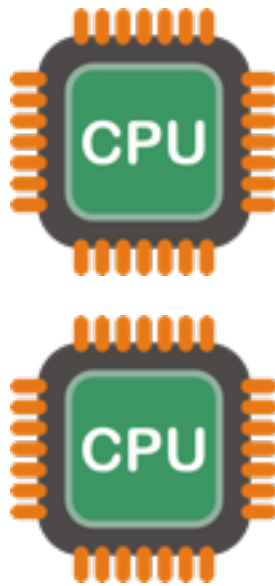
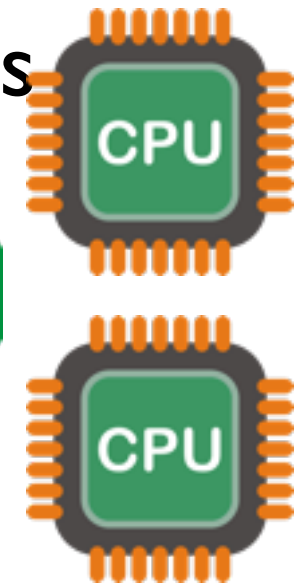
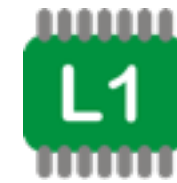
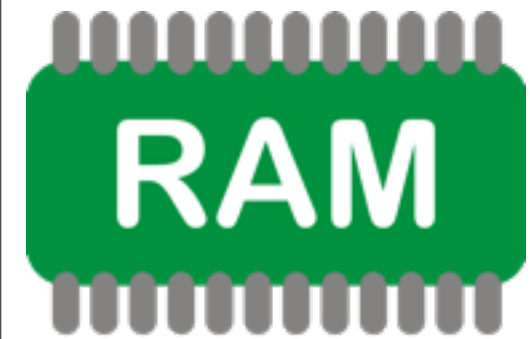
# Multiple CPUs

48 GB  
33 ns

12 MB  
5.4 ns

256 kB  
1.7 ns

32 kB  
1.4 ns



job of the database  
on a modern system  
(among others)

# Algorithms

## Row-Store

1	123	pos	VVINF
2	123	lemma	essen
3	456	pos	NN

## Column-Store

node_ref	name	value
123	pos	VVINF
123	lemma	essen
456	pos	NN

- function calls in a loop  
getNextTuple  
getAttribute('name')  
compareAttribute('lemma')
- pipeline flush
- instruction cache miss

- array lookup in a loop  
column[index] == 'lemma'
- loop unrolling
- no jumps
- instruction locality

**50% of time spent in  
resource stalls**

# History of Column-Stores

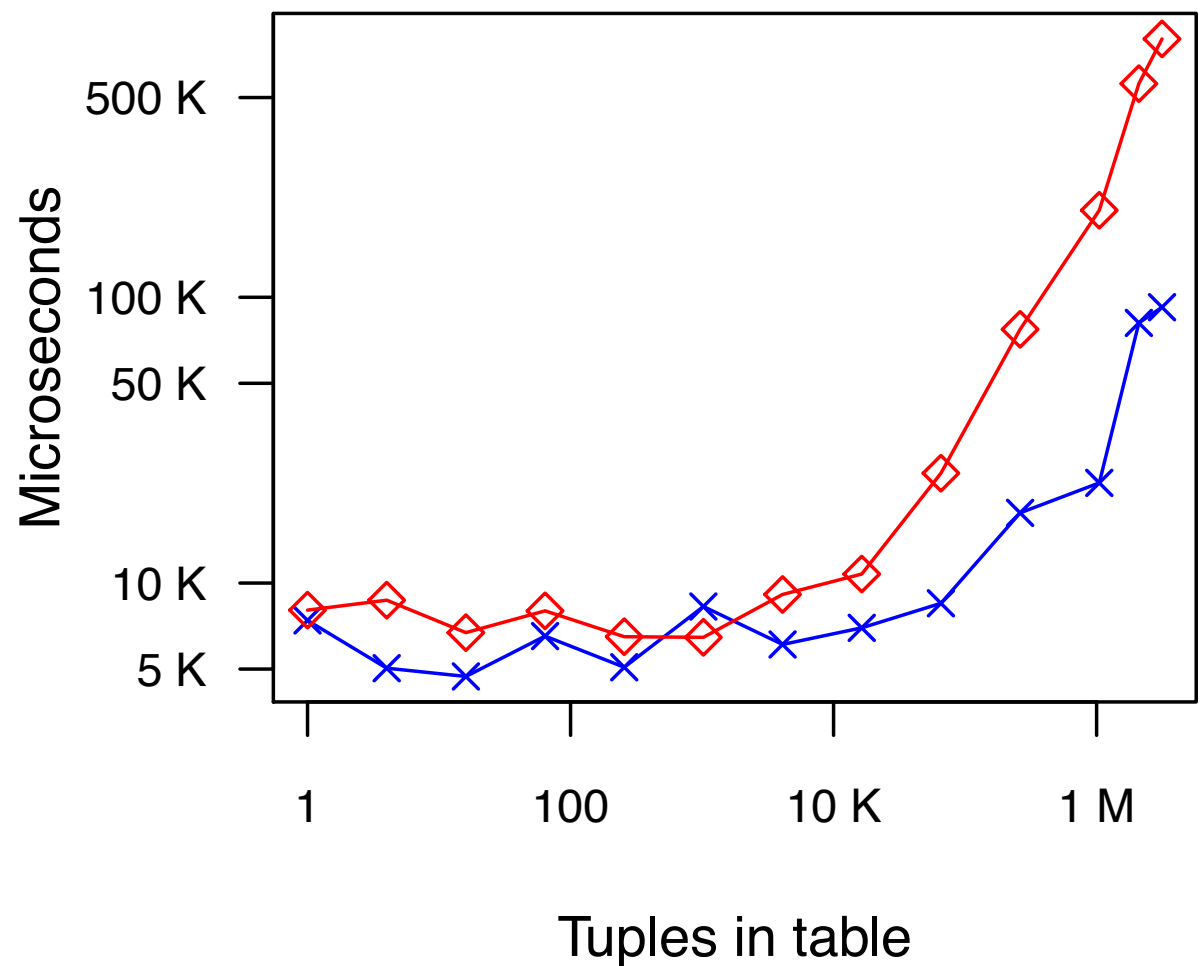
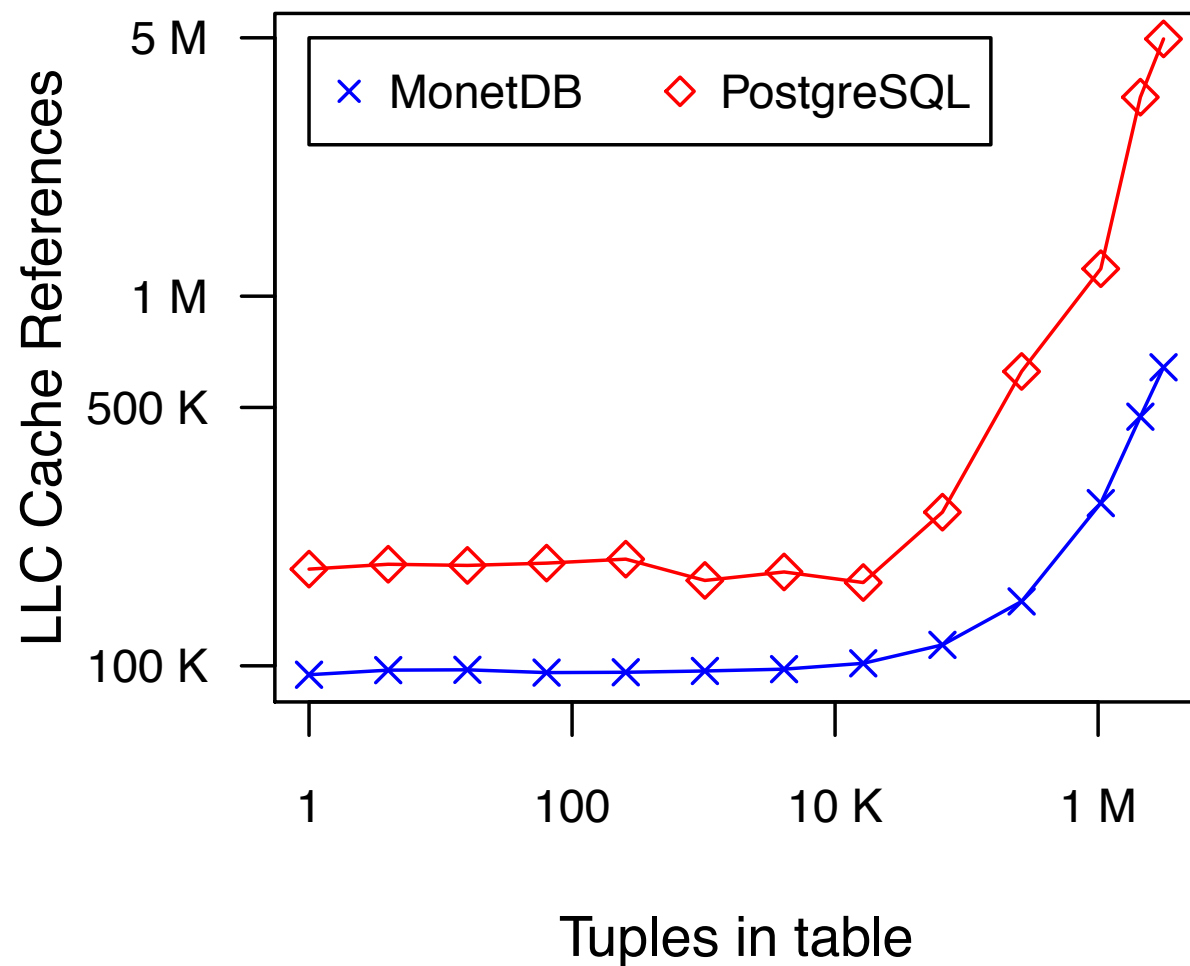
- IBM, 1973: *System R*
- Copeland, 1985: *A Decomposed Storage Model*
- CWI, 1994: *Monet. An Impressionist Sketch of an Advanced Database system*
- 2000s: Column-Stores become mainstream
  - C-Store, Vertica
  - MonetDB/X100, Vectorwise
  - SAP HANA

# MonetDB

- pioneered research into Column-Stores
- active research platform
  - adaptive indexing, *database cracking*
  - partial vectorized bulk-processing, *X100*
  - opportunistic materialization, *recycler*
  - array primitives, *SciQL*
- well-received
  - VLDB 10 Years Best Paper Award, 2009: *Database Architecture Optimized For The New Bottleneck: Memory Access*
  - SIGMOD Jim Gray Doctoral Dissertation Award, 2011: *Database Cracking: Towards Auto-tuning Database Kernels*
  - VLDB Challenges & Vision Best Paper Award, 2011: *The Researcher's Guide to the Data Deluge: Querying a Scientific Database in Just a Few Seconds*



# Cache usage vs. query speed

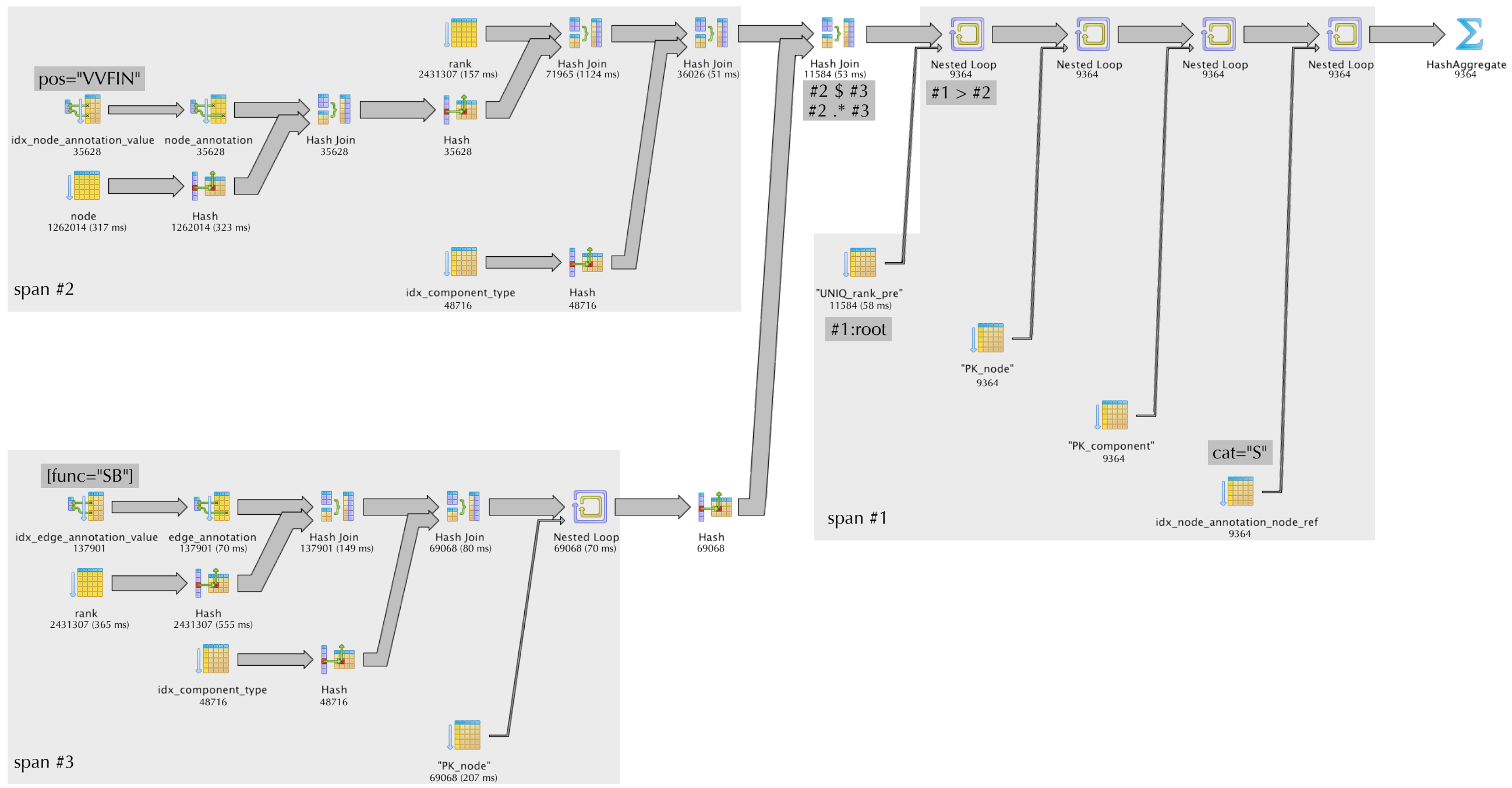


Query: annotation name = 'lemma'

- last level cache references
- indirect measure of processed data
- everything that is processed must pass once through the cache



# Normalized schema on PostgreSQL



# Materialized schema on PostgreSQL

