

Index support for regular expression search

Alexander Korotkov
PGCon 2012, Ottawa

The background features a complex, abstract pattern of overlapping, flowing lines in various shades of teal, light blue, and white. These lines create a sense of movement and depth, resembling liquid or smoke. A solid black horizontal band is positioned across the middle of the image, serving as a backdrop for the text.

Introduction

What is regular expressions?

Regular expressions are:

- powerful tool for text processing
- based on formal language theory
- expressing same class of “languages” as finite automata

So... automata

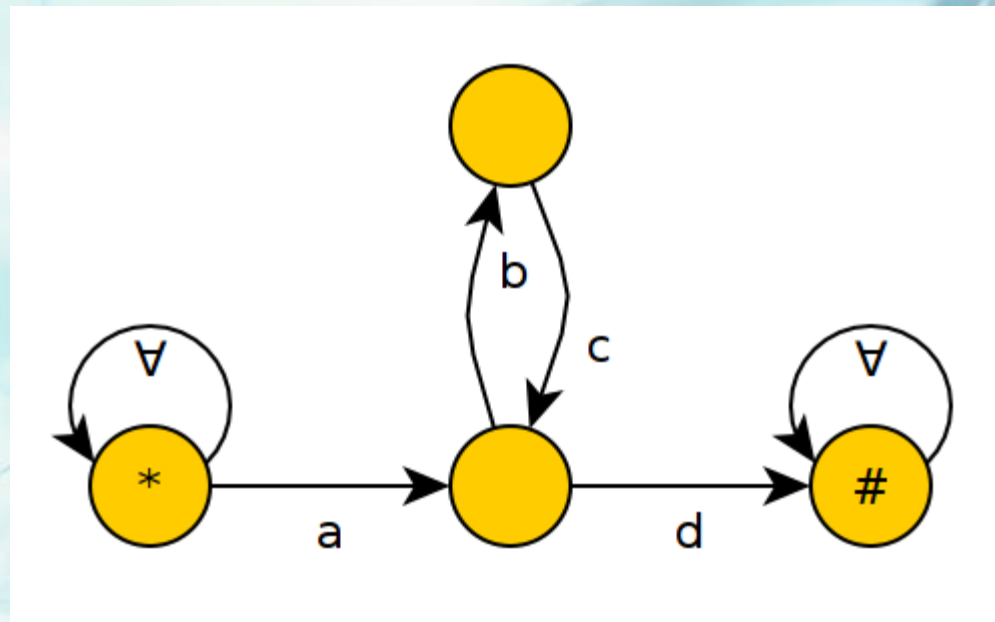
- Regular expression can be transformed into automaton.
- Moreover, such transformation is really used by regex engines

So... automata

- Automaton is a graph which vertices are "states" and which arcs are labeled by characters.
- Automaton "reads" string if you can type that string by a traversal from "initial" state to "final" state

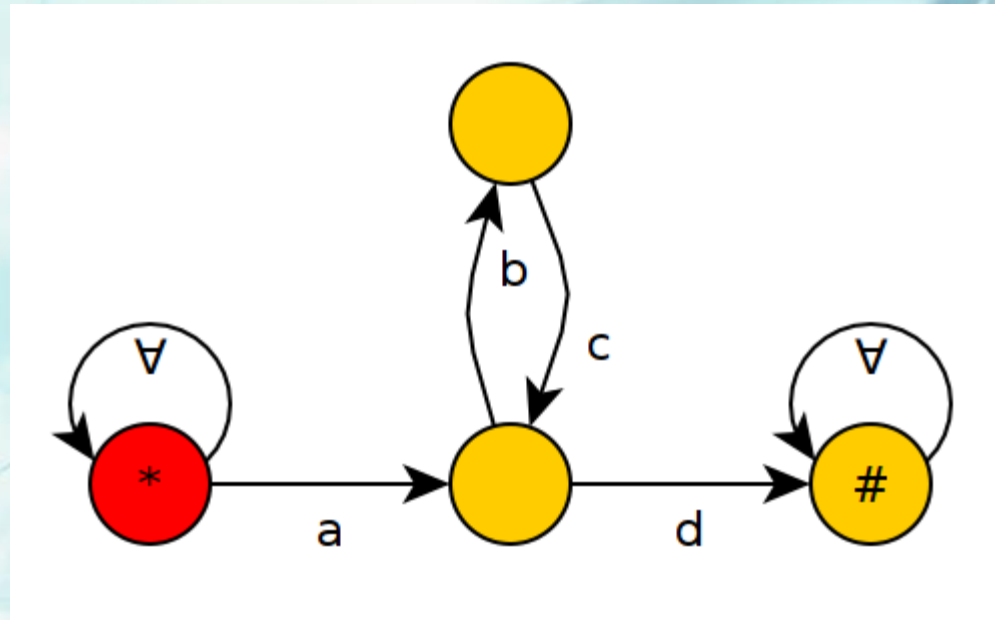
Example

/a(bc)*d/
becomes



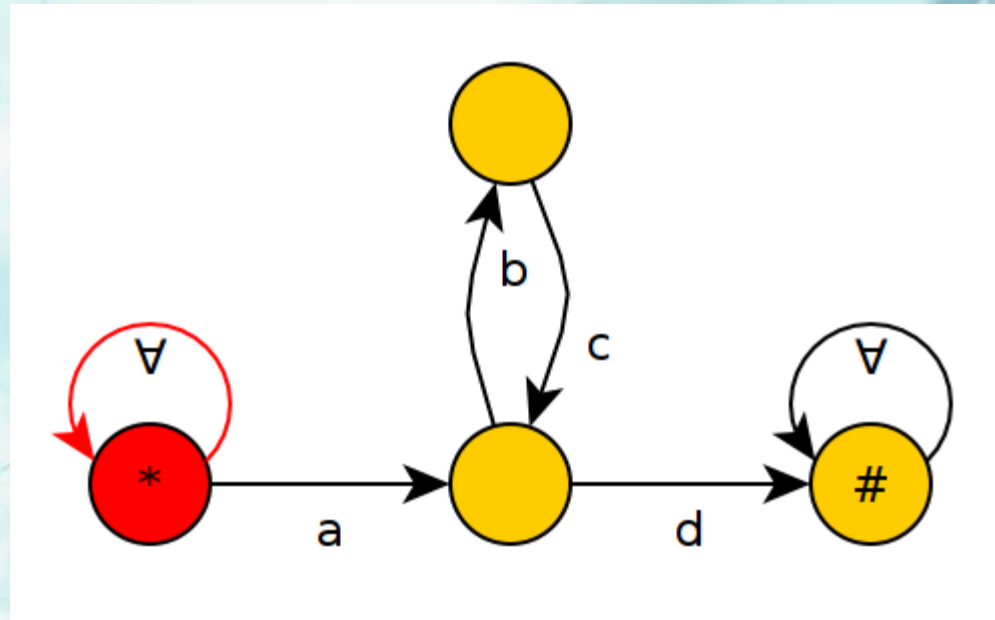
Example

xyzabcbcdxyz



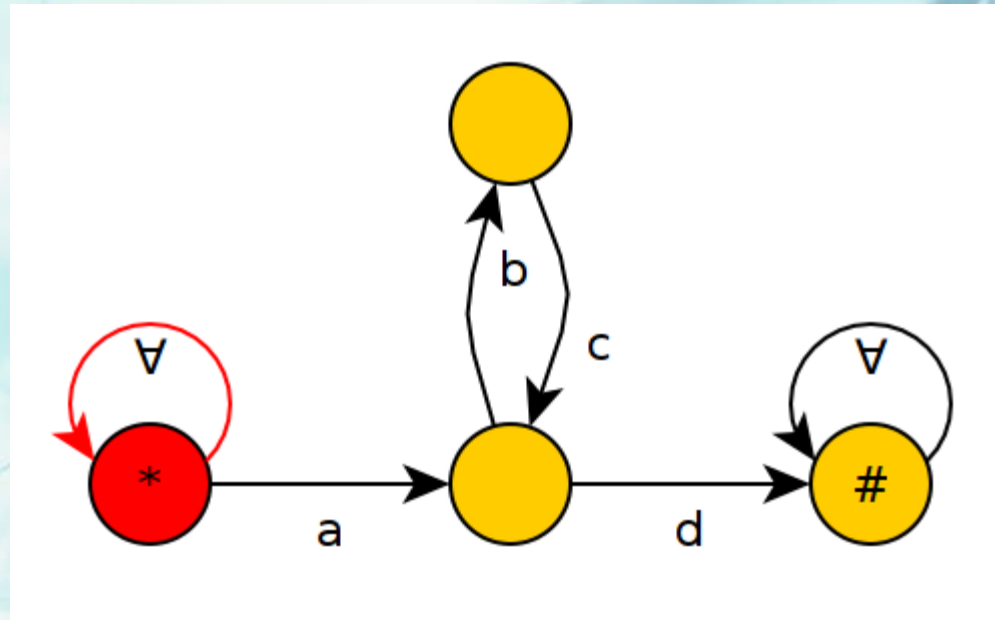
Example

x y z a b c b c d x y z



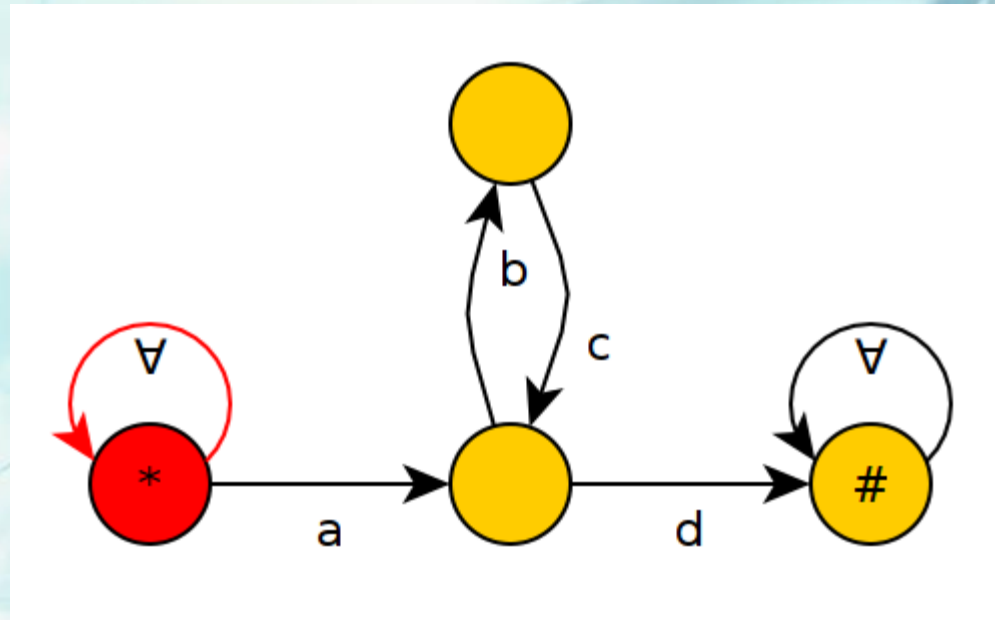
Example

xyzabc**bc**cdxyz



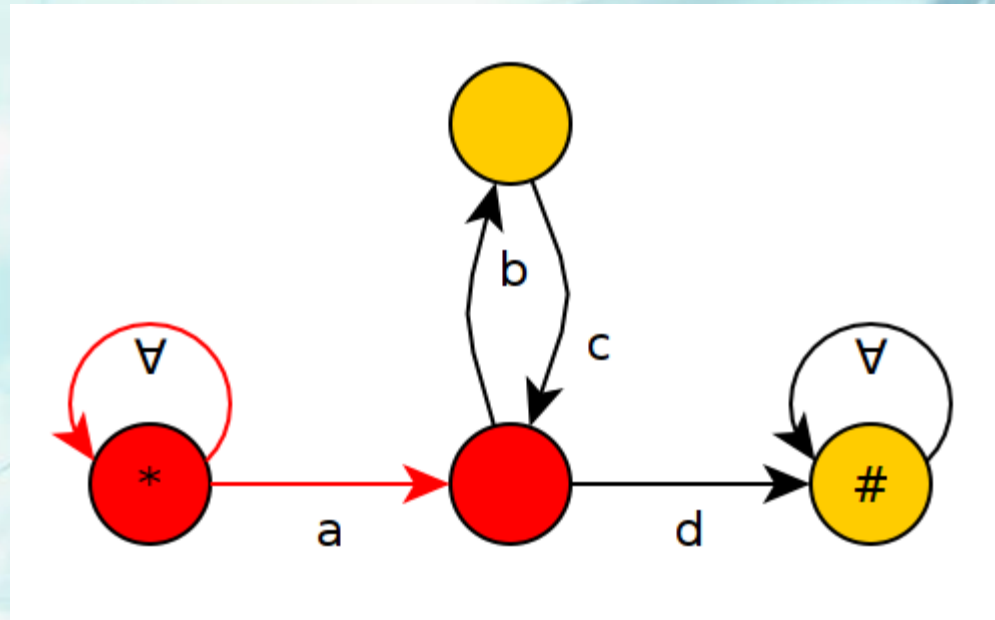
Example

xy**z**abc**b**cdxyz



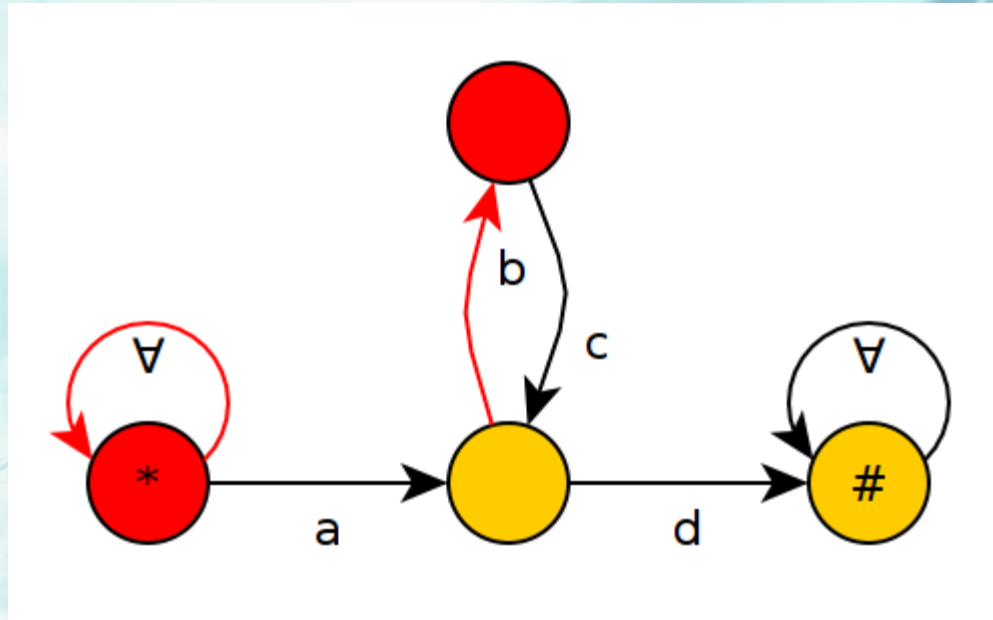
Example

xyz**a**bcbcdxyz



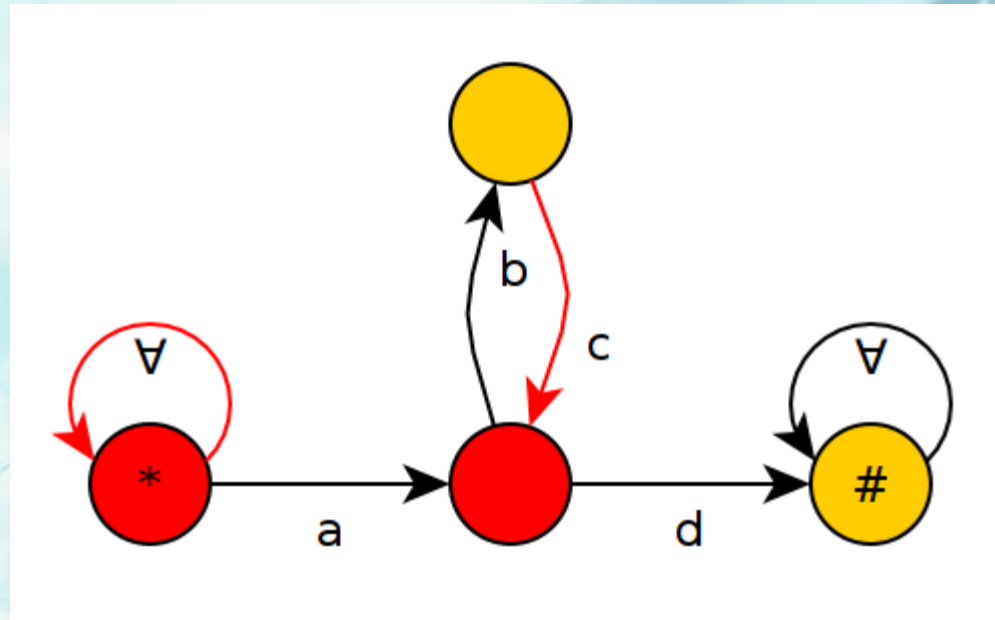
Example

xyza**b**cbcdxyz



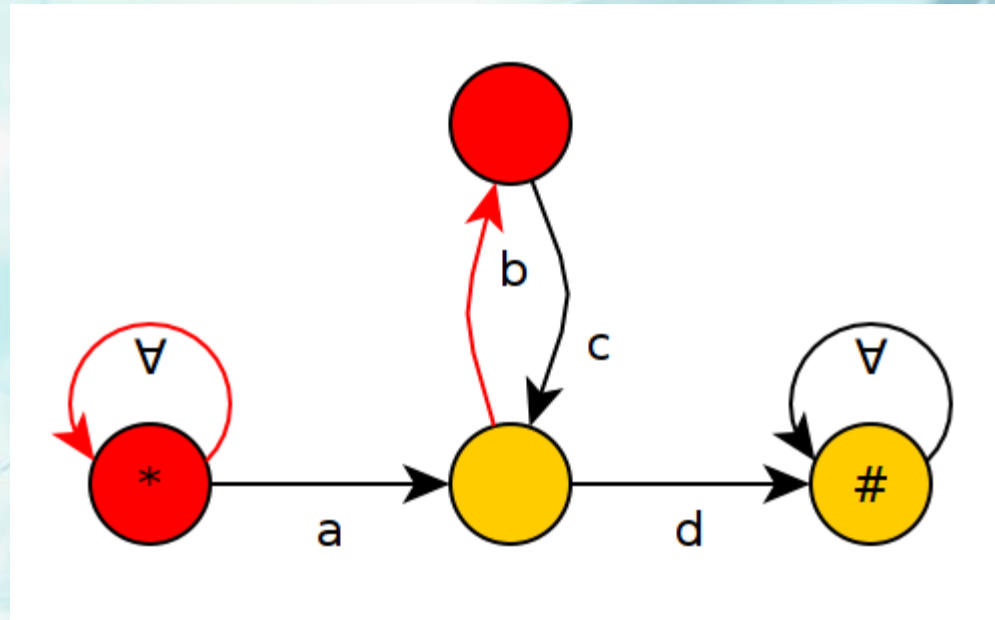
Example

xyzab**c**bcdxyz



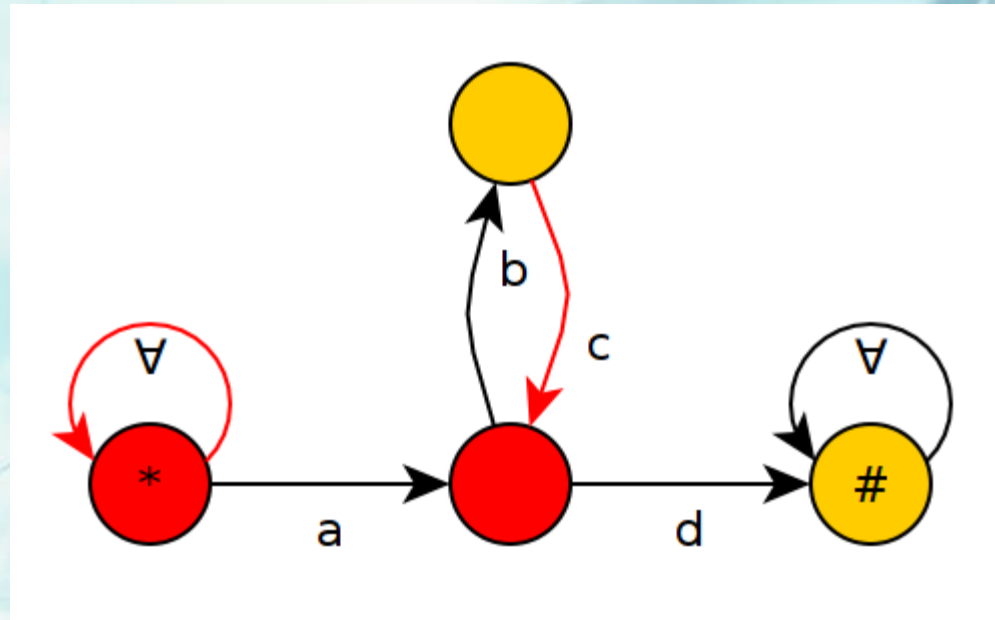
Example

xyzabc**b**cdxyz



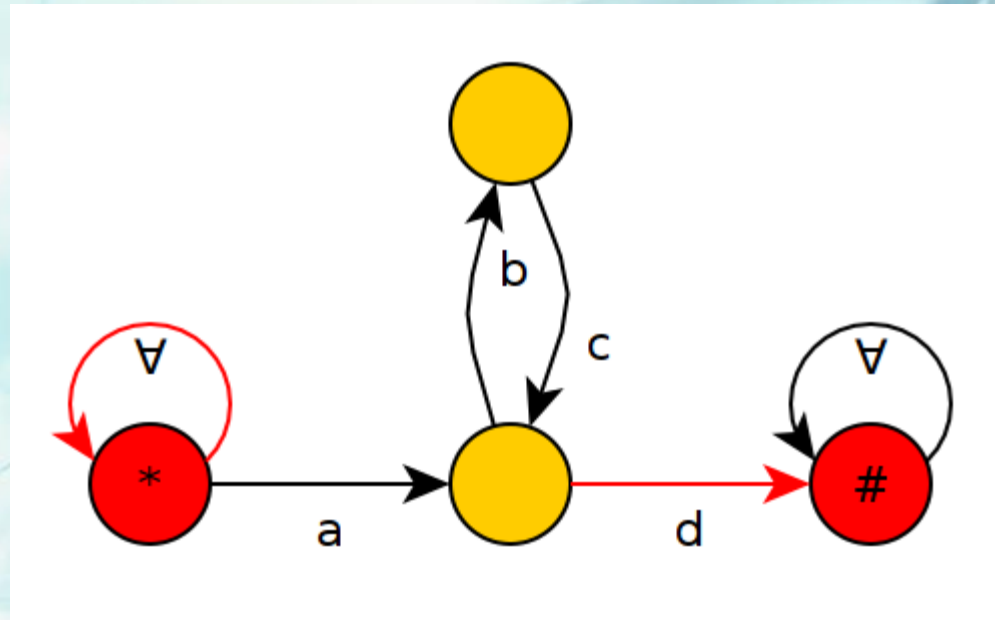
Example

xyzabcbcdxyz



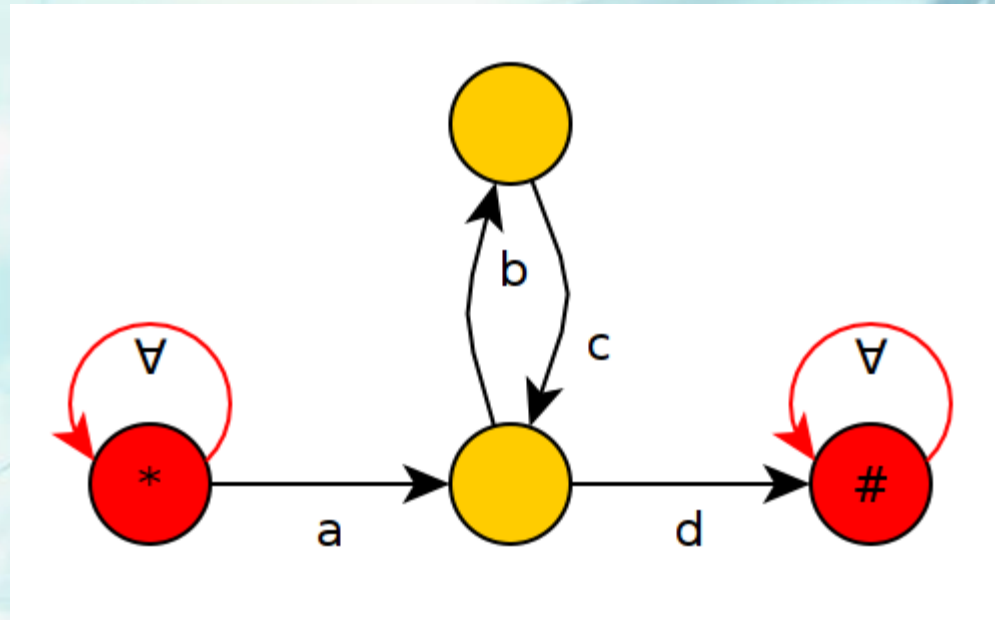
Example

xyzabc**cd**xyz



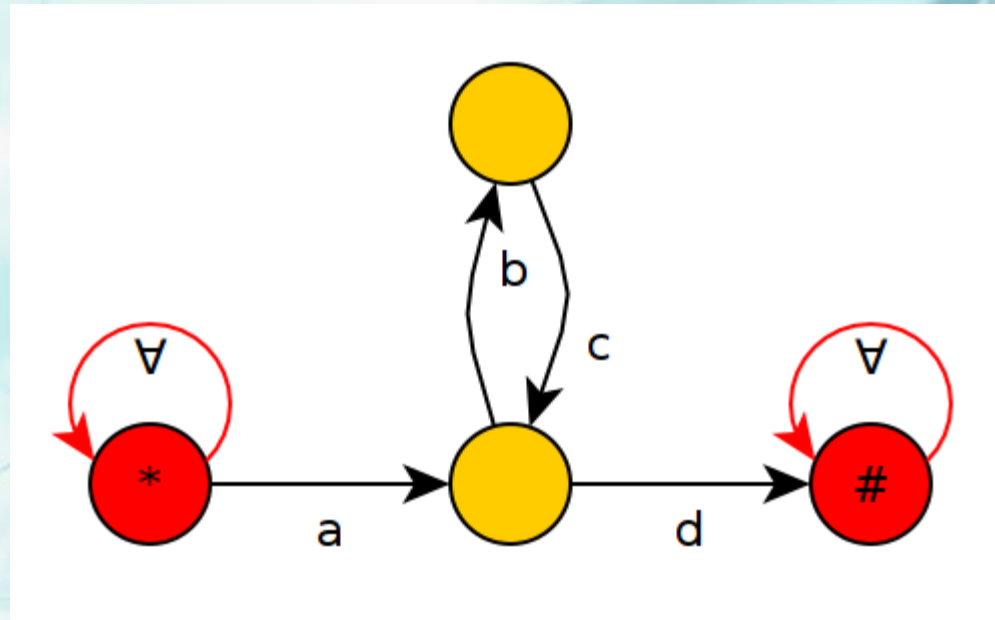
Example

xyzabc**bc**dxyz



Example

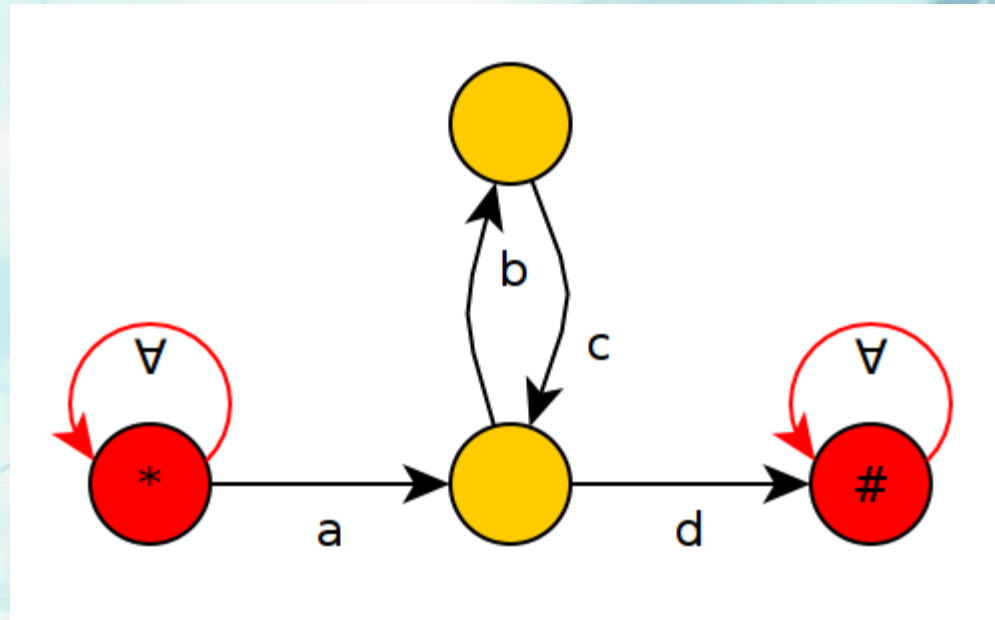
xyzabc**bc**xyz



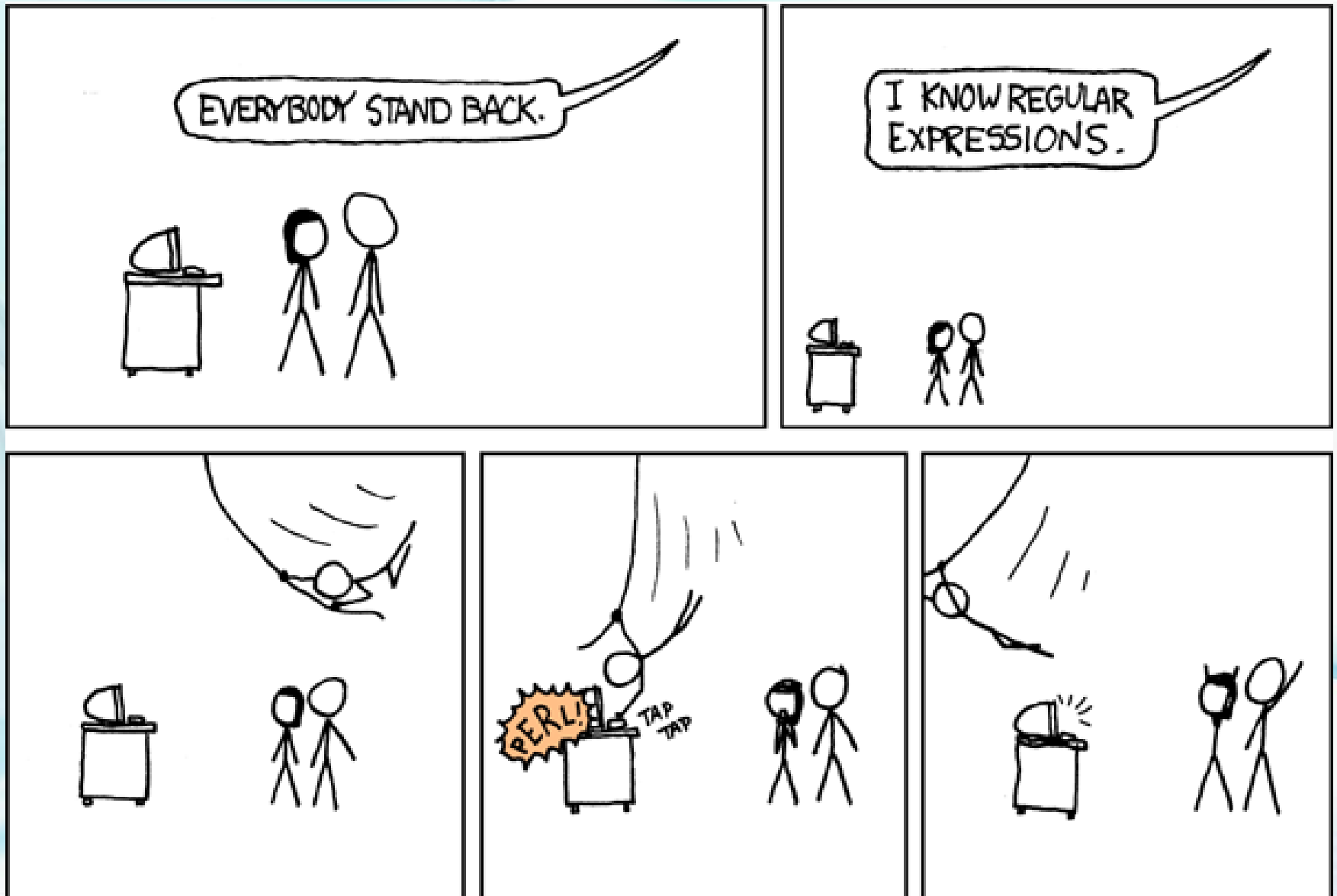
Example

xyzabcbcdxyz

Finish! Match!



Okay, now all of us know...



Regex based search

- PostgreSQL can regex based search :)
- It's only a sequential search for a while :(



Inverted indexes on q-grams

Q-grams

- Q-gram is substring of length q which can be used as a signature of original string
- Widely used in various string processing tasks

Inverted index on q-grams

- Maintain association between q-gram and all the strings where it mentioned.
- `pg_trgm` has an implementation for $q = 3$

pg_trgm

1. "regular expressions",
2. "expressive speech",
3. "regular speaker"

=>

' e': {1,2}	'egu': {1,3}	'pre': {1,2}
' r': {1,3}	'er ': {3}	'reg': {1,3}
' s': {2,3}	'ess': {1,2}	'res': {1,2}
' ex': {1,2}	'exp': {1,2}	'sio': {1}
' re': {1,3}	'gul': {1,3}	'siv': {2}
' sp': {2,3}	'ion': {1}	'spe': {2,3}
'ach': {2}	'ive': {2}	'ssi': {1,2}
'ake': {3}	'ker': {3}	'ula': {1,3}
'ar ': {1,3}	'lar': {1,3}	've ': {2}
'ch ': {2}	'ns ': {1}	'xpr': {1,2}
'eac': {2}	'ons': {1}	
'eak': {3}	'pea': {2,3}	

V-grams or multigrams

- Each q-gram can have specific q
- Selectivity of all q-grams are similar (and low enough)
- More effective index search!

V-grams or multigrams

Problems:

- Hard to maintain online
- ...

Patent trololo!





How to use it for regex search?

General idea

$/[ab]cde/ \Rightarrow (acd \text{ OR } bcd) \text{ AND } cde$

acd: {1,4,5}, bcd: {2,3,4}, cde:{2,4,6}

	acd	bcd	cde	(acd OR bcd) AND cde	recheck
1	t	f	f	f	
2	f	t	t	t	→ f
3	f	t	f	f	
4	t	t	t	t	→ t
5	t	f	f	f	
6	f	f	t	f	

General idea

$/[ab]cde/ \Rightarrow (acd \text{ OR } bcd) \text{ AND } cde$

How to do this in general case?



Existing approaches for q-gram extraction

Scholar paper

Junghoo Ch and Sridhar Rajagopalan,
A fast regular expression indexing engine, Proceedings 18th
International Conference on Data
Engineering, 2002

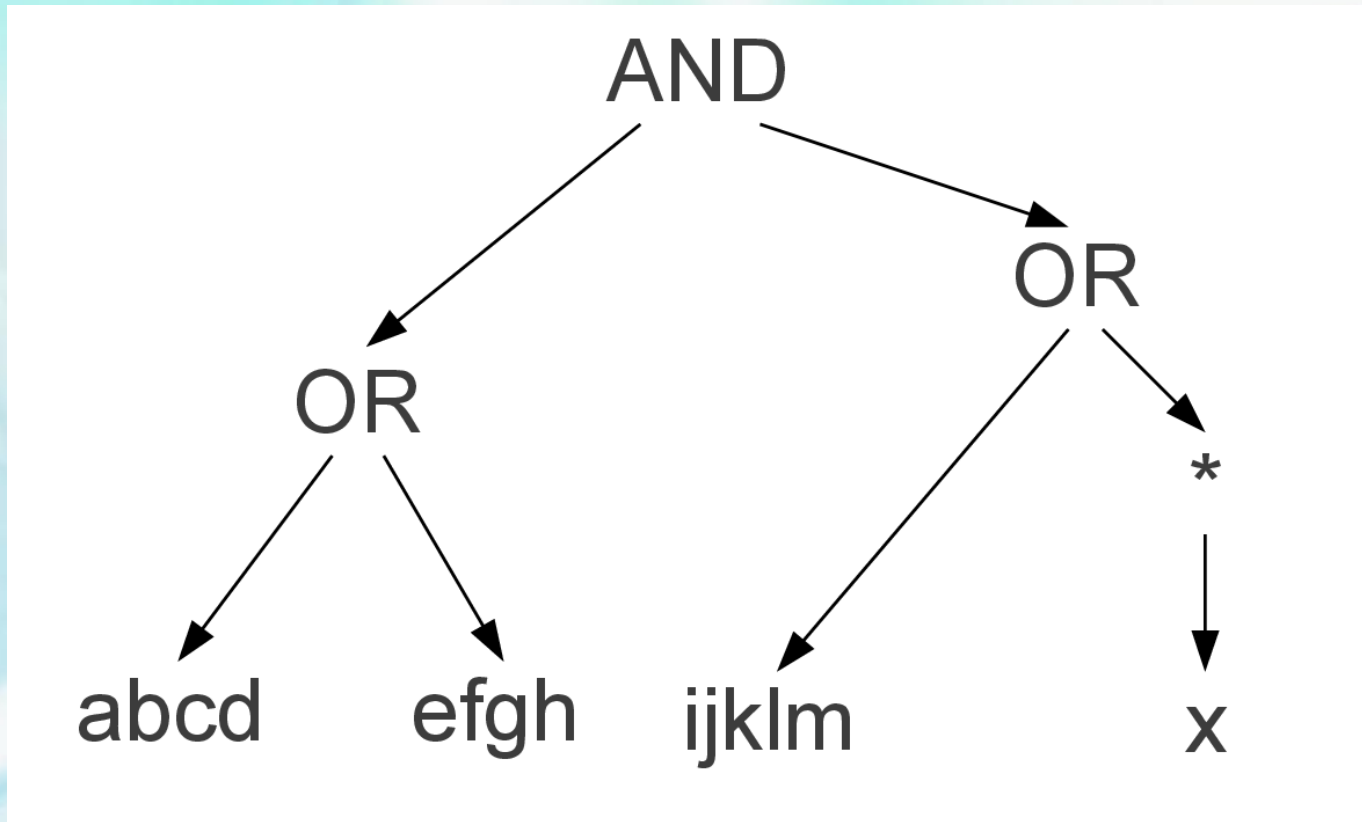
Still widely referenced as state of art
work about indexing for regular
expressions.

FREE method

- Extract tree of continuous string fraction from regex.
- Transform those continuous fractions to multigrams (q-grams with variable q).
- Use inverted index on multigrams for query evaluation

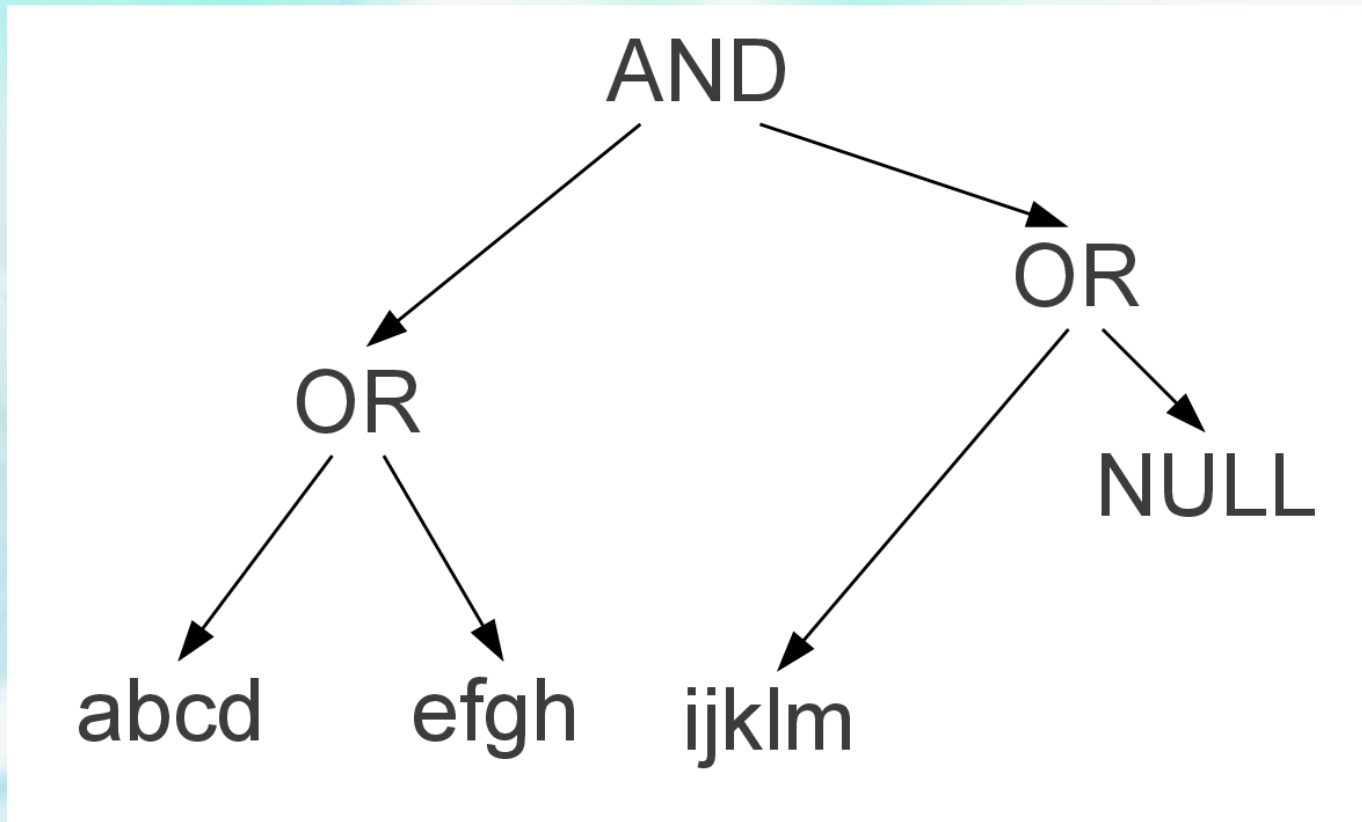
FREE method: example

Tree for `/(abcd|efgh)(ijklm|x*)/`



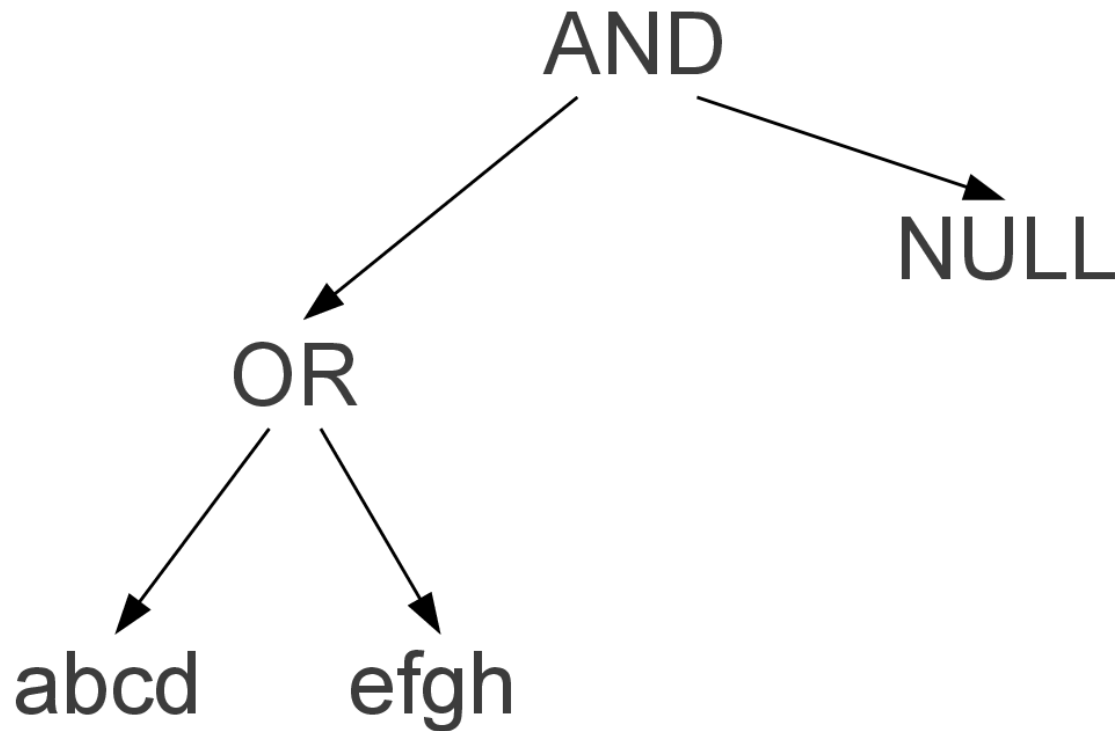
Scholar paper: example

Replace "*" nodes with NULL



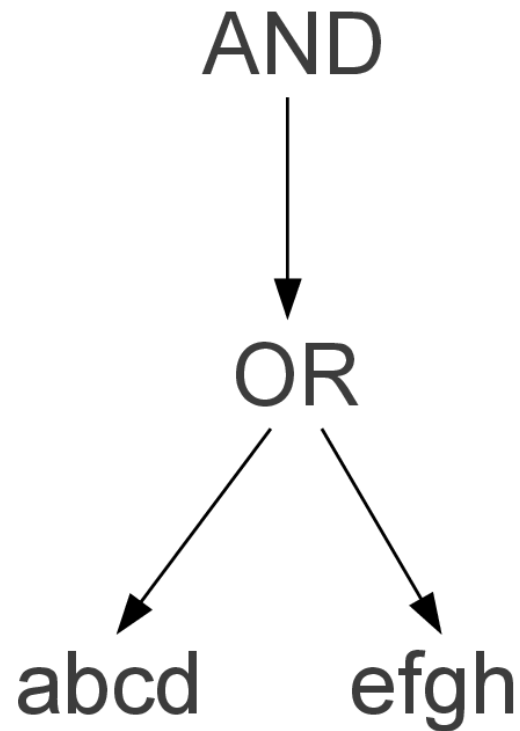
Scholar paper: example

NULL "eats" parent OR node



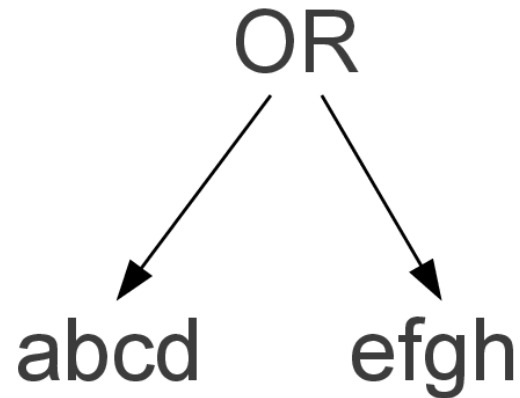
Scholar paper: example

AND node "eats" child NULL



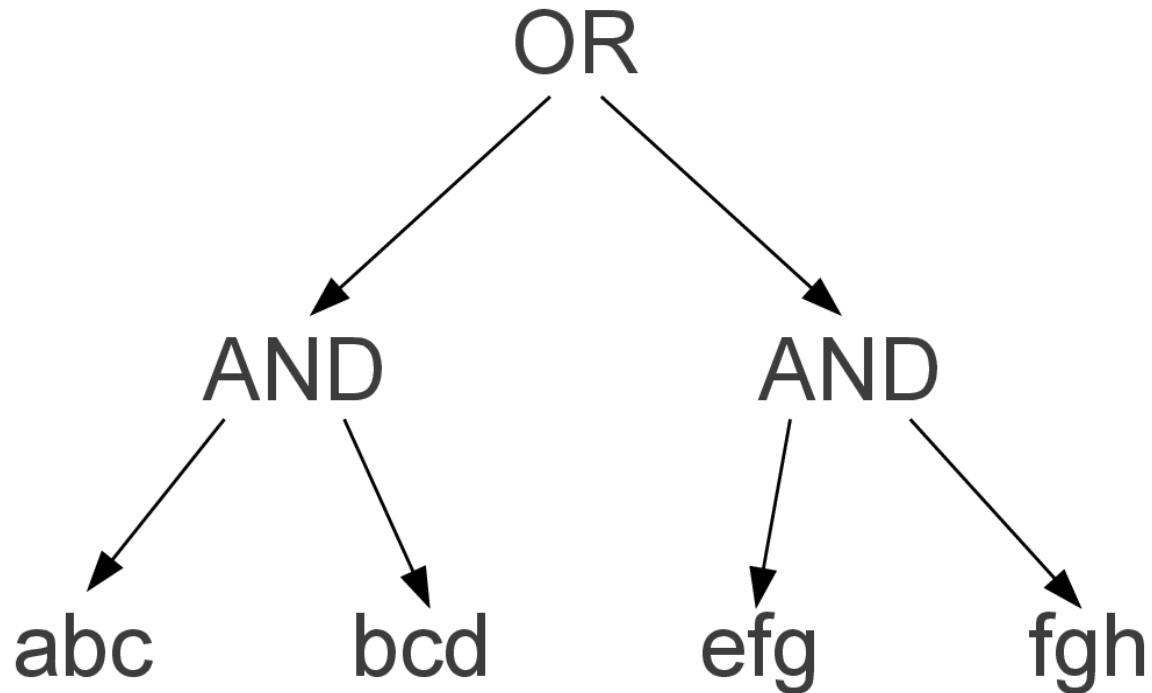
Scholar paper: example

Simplify a bit



Scholar paper: example

Expand continuous string fractions into trigrams



Google code search

- Was launched in 2006.
- Supports regex search.
- Google guys are smart. It can't be a sequential scan.
- It also seemed to be something better than previous technique.
- We don't know what... :(

We **didn't** know what until...

- Google code search was closed in 2011 :(
- Russ Cox has published description of indexing technique in January 2012
<http://swtch.com/~rsc/regexp/regexp4.html>
- More than 5 years of intrigue!

Google code search method

- Get 5 characteristics about each part of regex: emptyable, exact, prefix, suffix, match.
- Recursively union them (with possible simplification)
- Use inverted index of trigrams for query evaluation (similar to pg_trgm)

Google code search method

Original regex: `/a(bc)+d/`

`a`: {exact: a}

`bc`: {exact: bc}

`d`: {exact: d}

`(bc)+`: {prefix:bc, suffix: bc}

`a(bc)+`: {prefix:abc, suffix:bc}

`a(bc)+d`: {prefix:abc, suffix:bcd}

Google code search method

`/a(bc)+d/`



`{prefix:abc, suffix:bcd}`



`abc AND bcd`



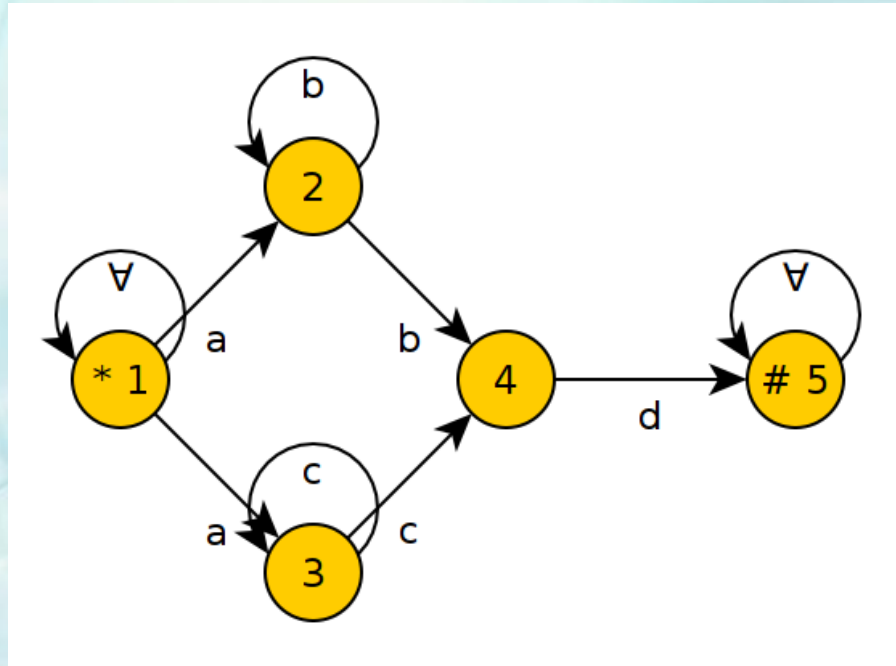
Proposed method

Proposed method

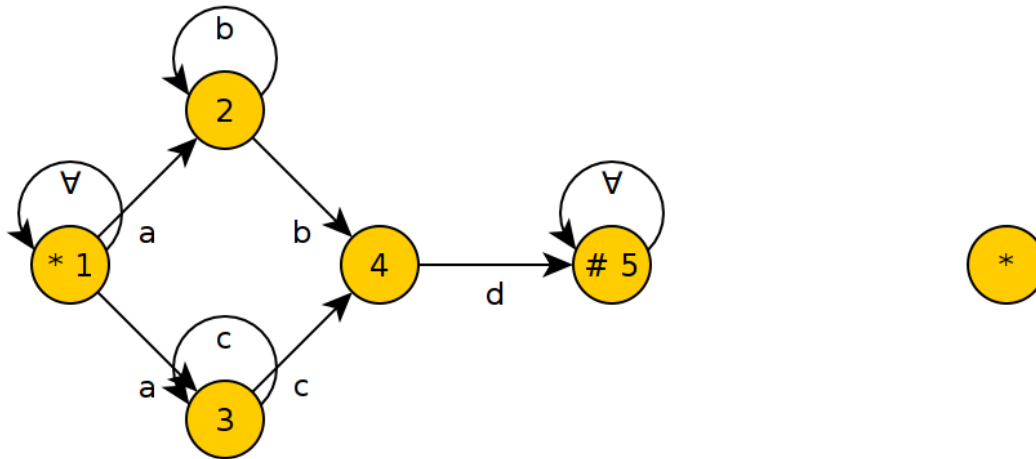
- Transform automaton into automaton like graph on trigrams
- Simplify that graph
- Use `pg_trgm` indexes

Transformation example

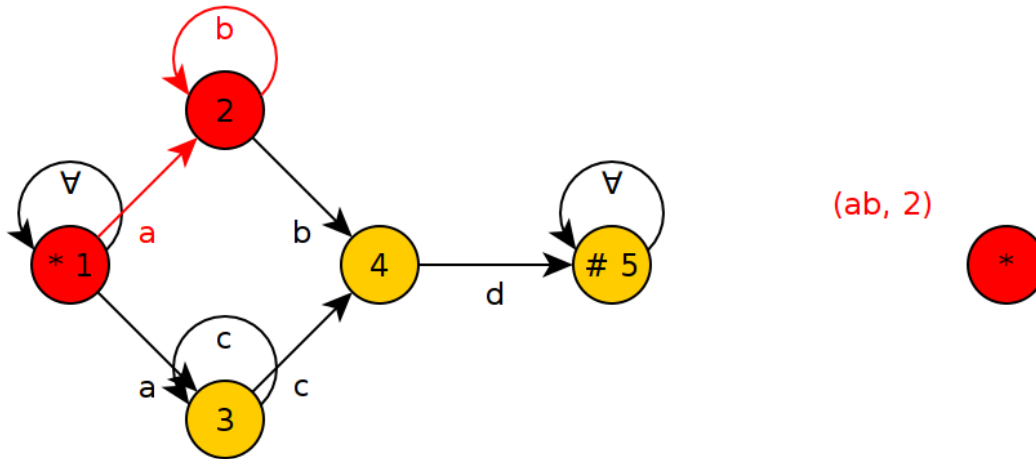
/a(b+|c+)d/



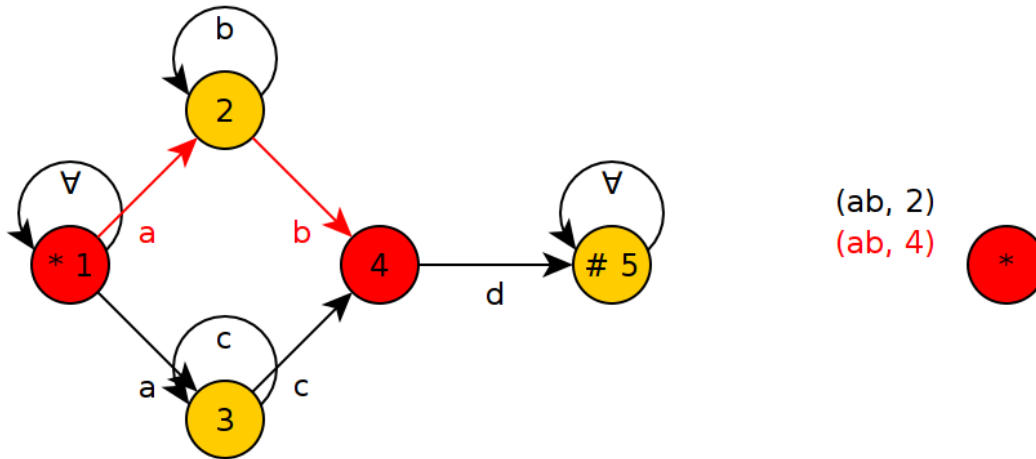
Transformation example



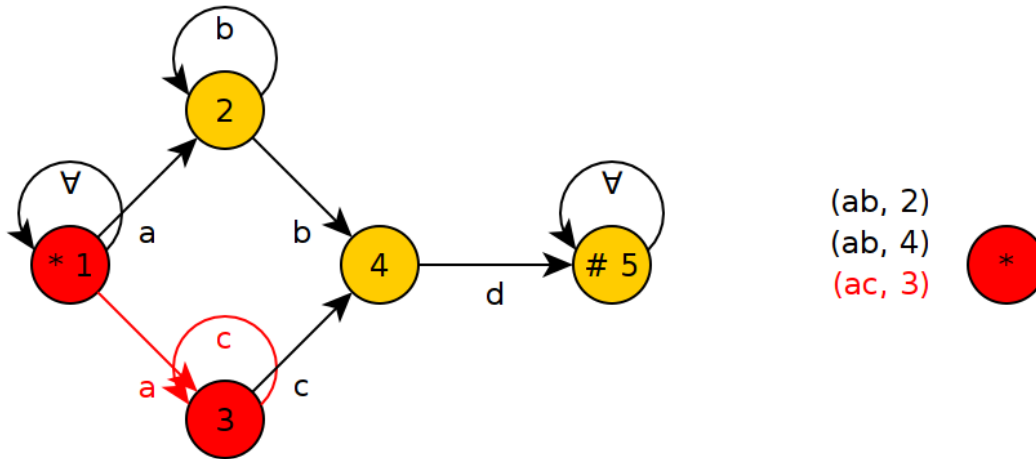
Transformation example



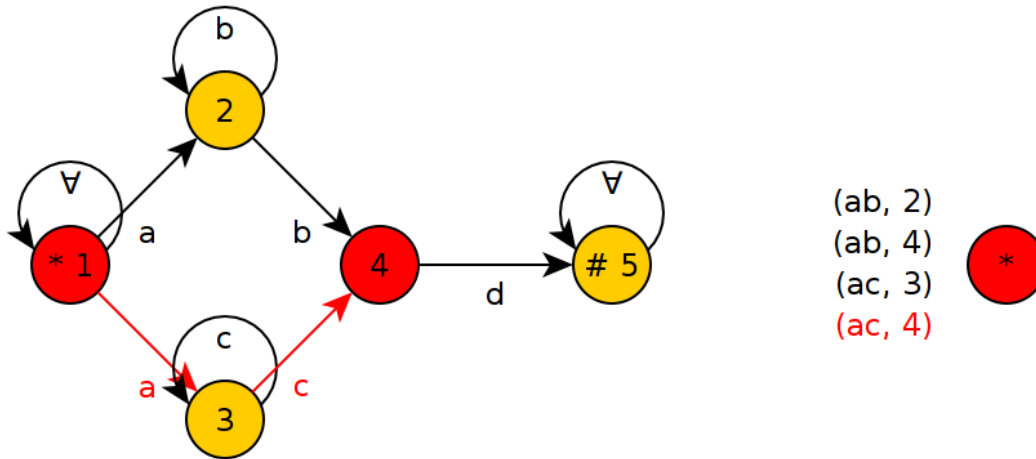
Transformation example



Transformation example



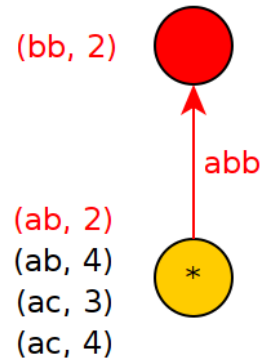
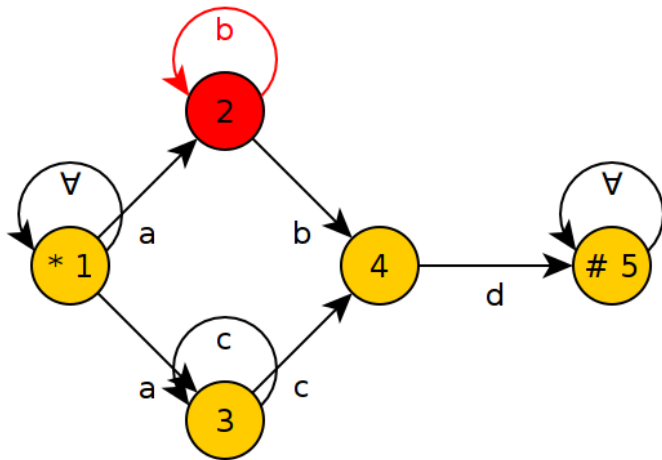
Transformation example



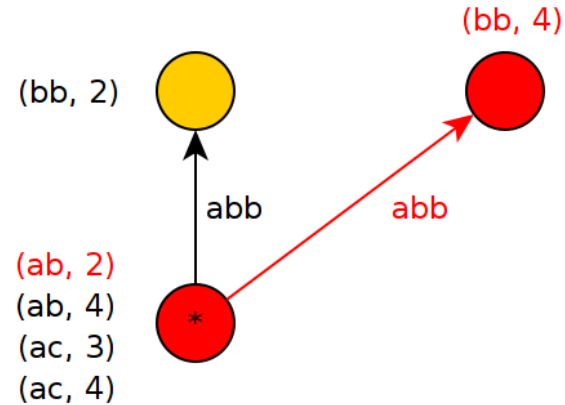
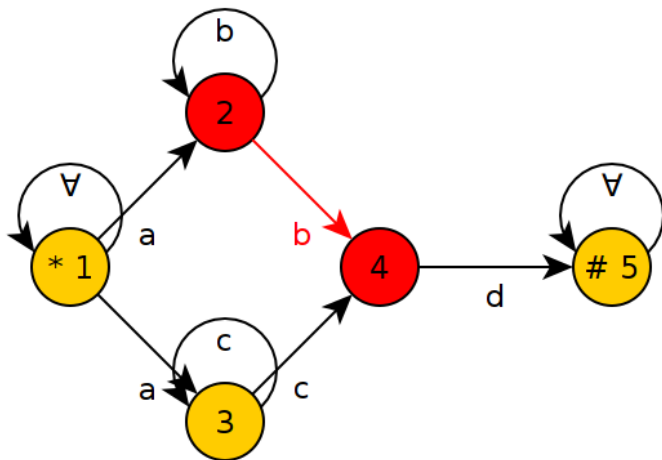
(ab, 2)
(ab, 4)
(ac, 3)
(ac, 4)



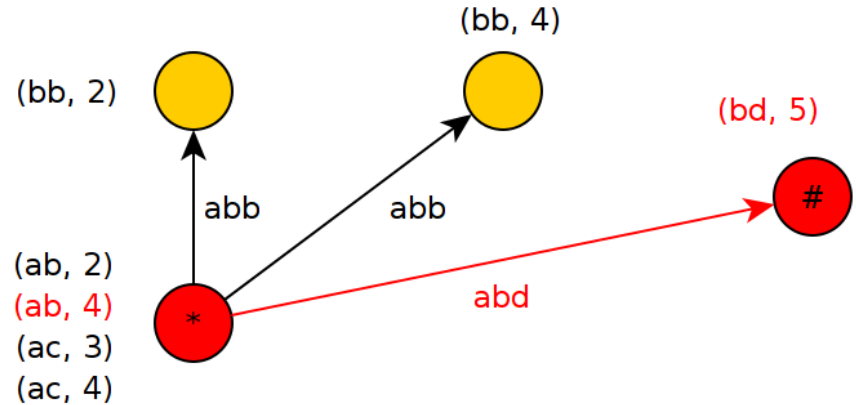
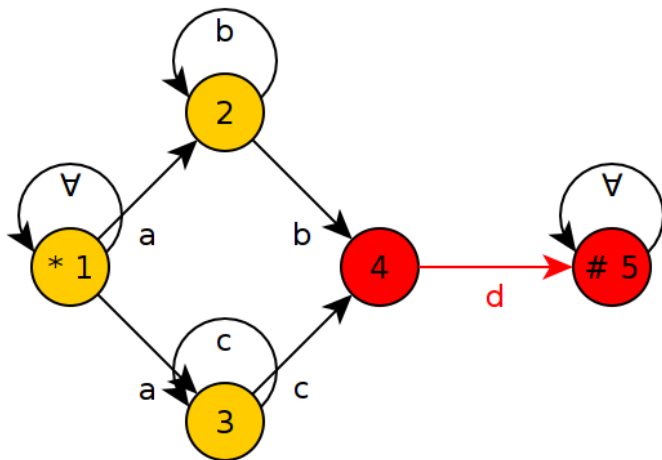
Transformation example



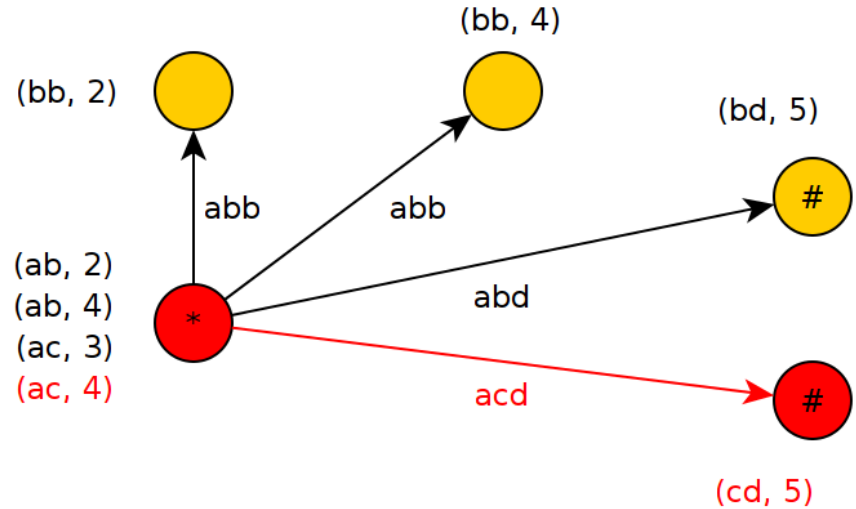
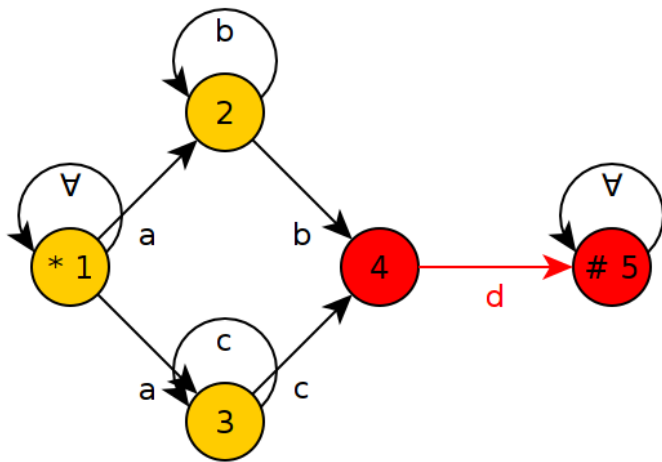
Transformation example



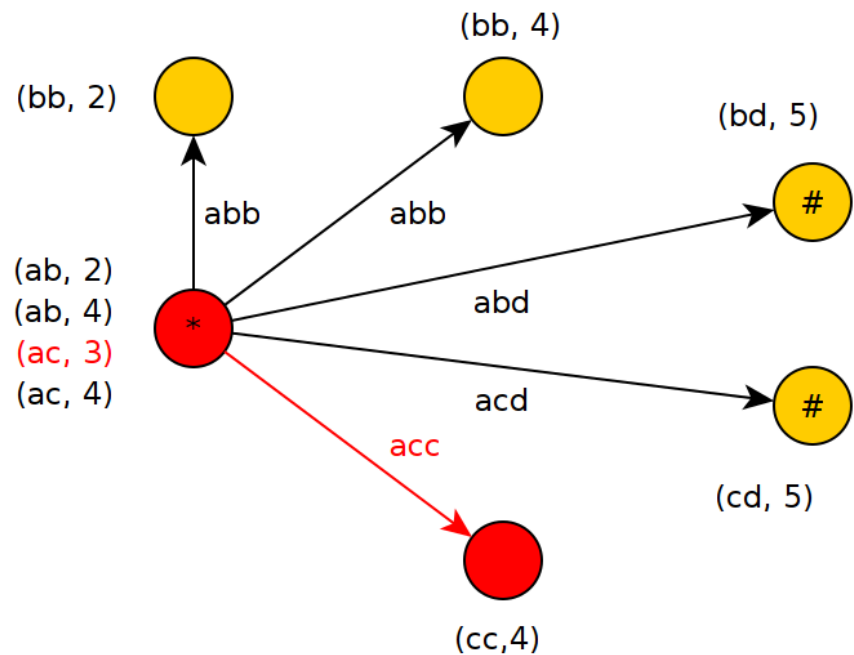
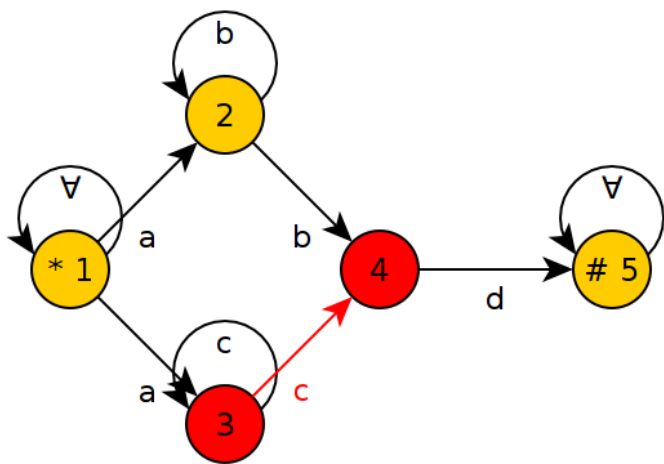
Transformation example



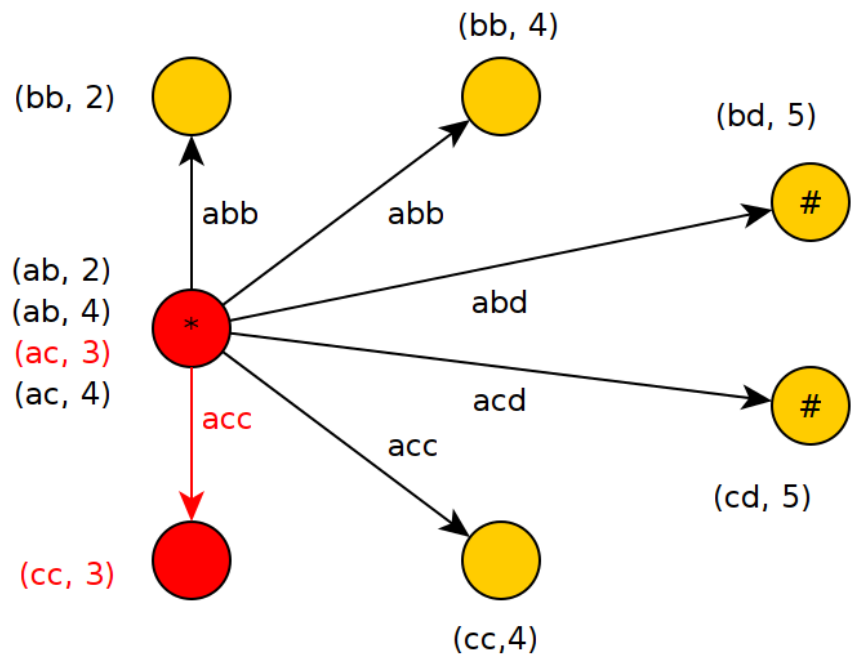
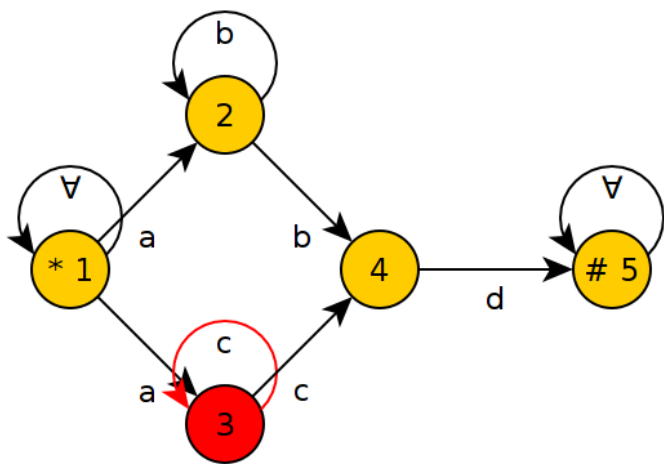
Transformation example



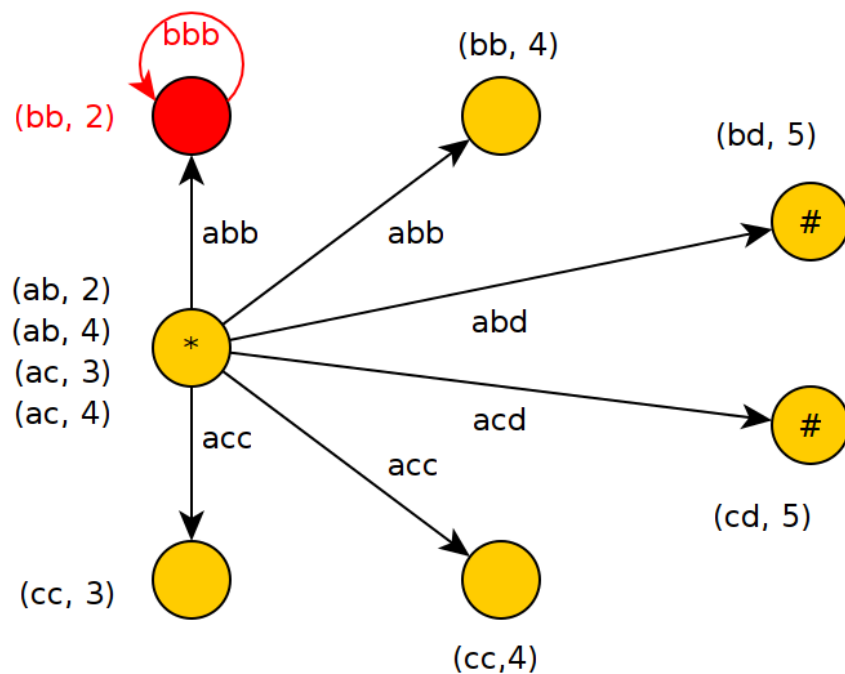
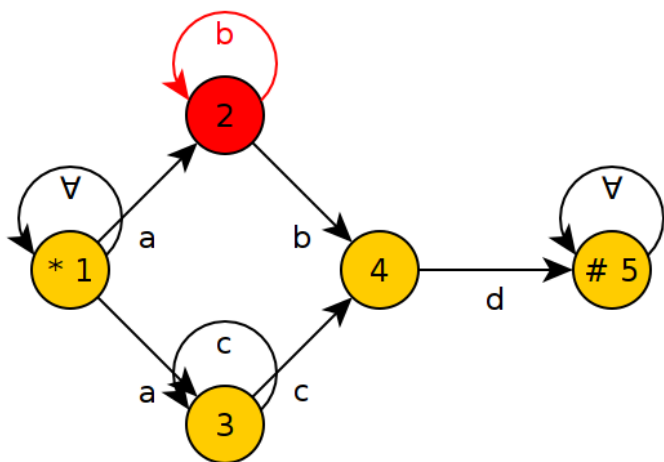
Transformation example



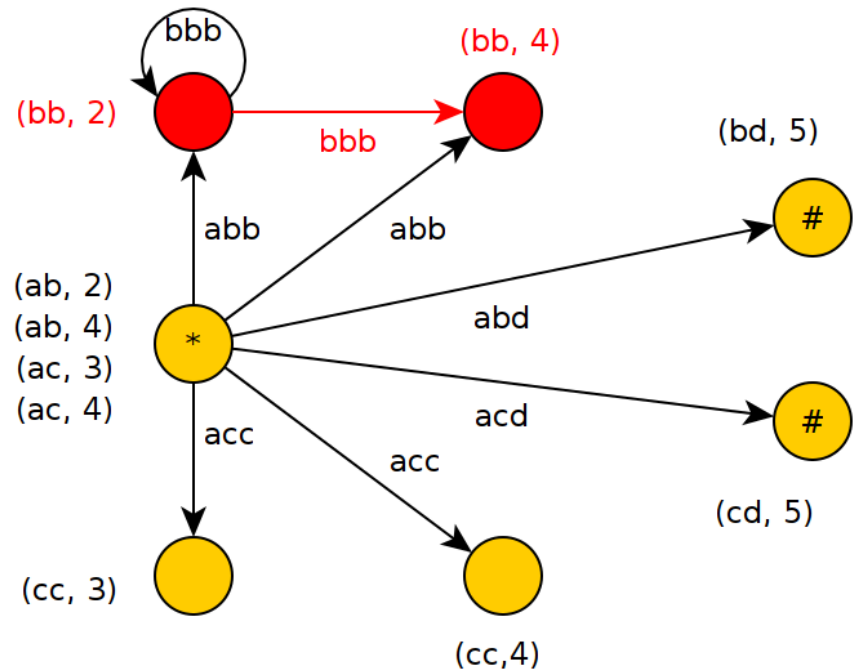
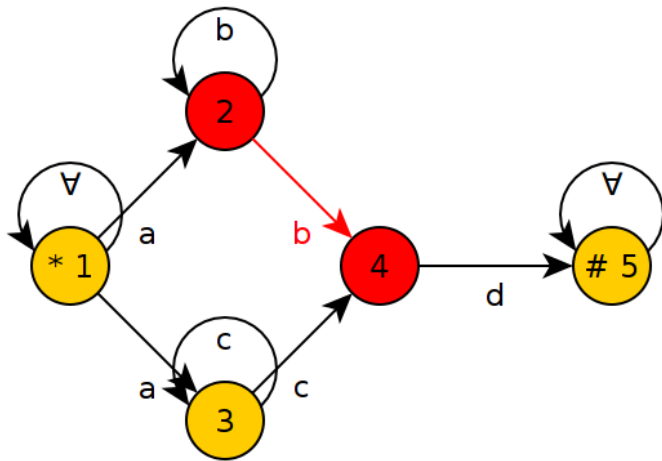
Transformation example



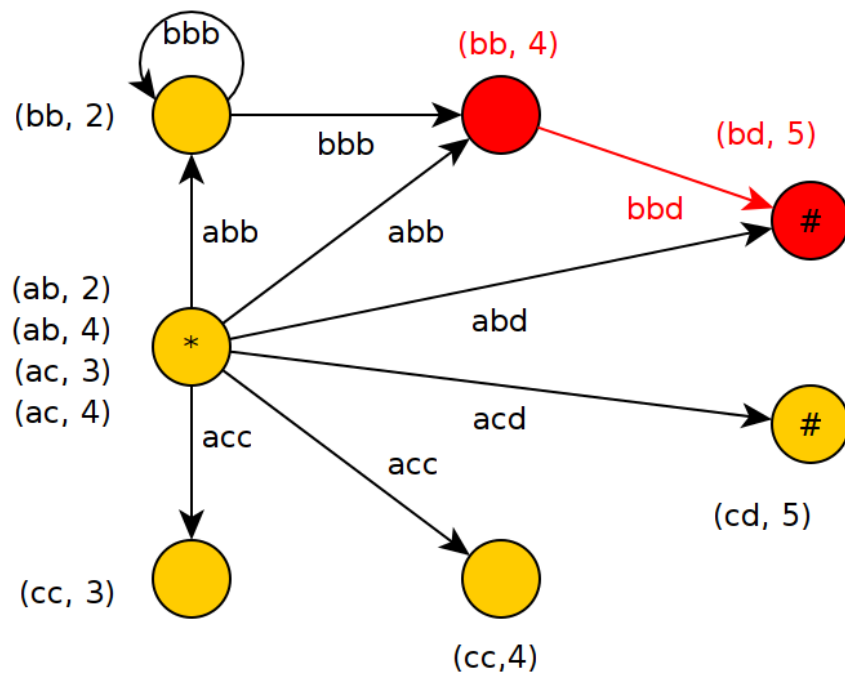
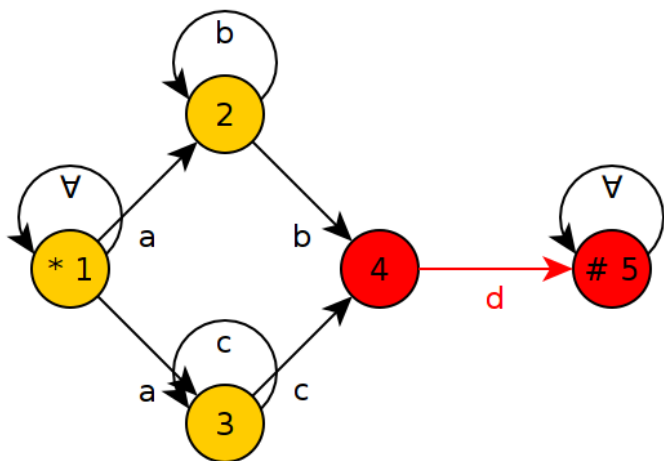
Transformation example



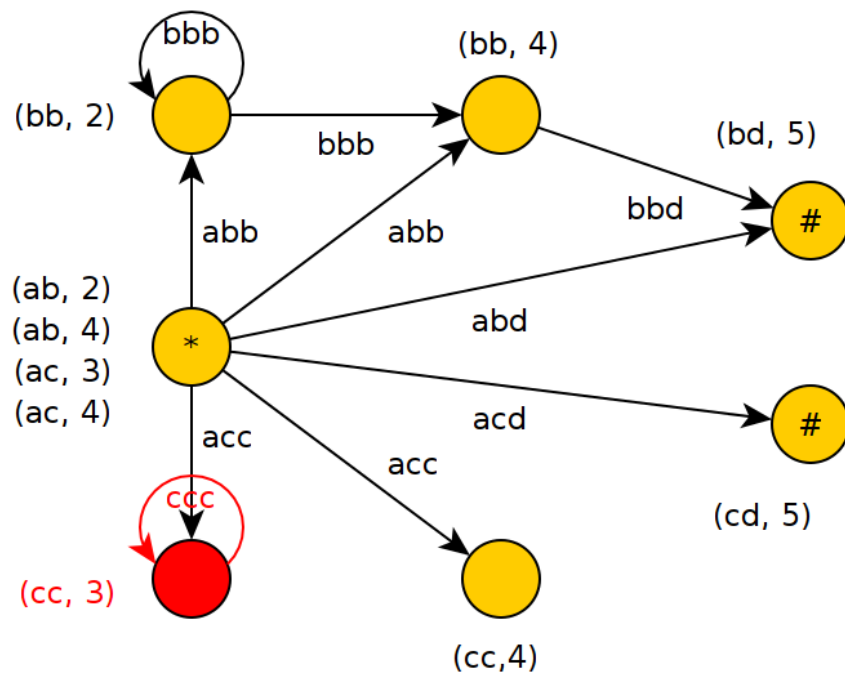
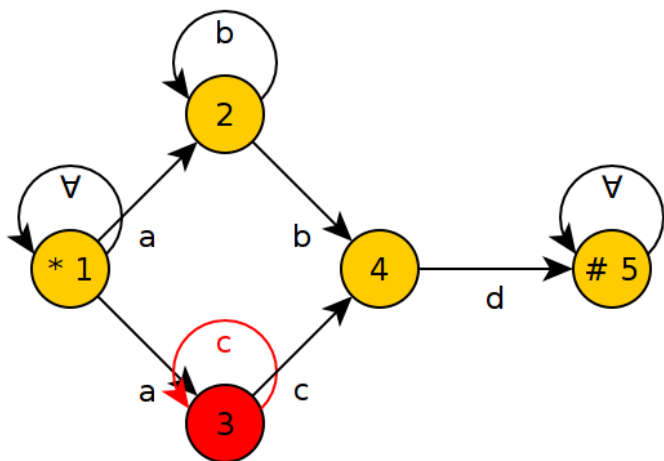
Transformation example



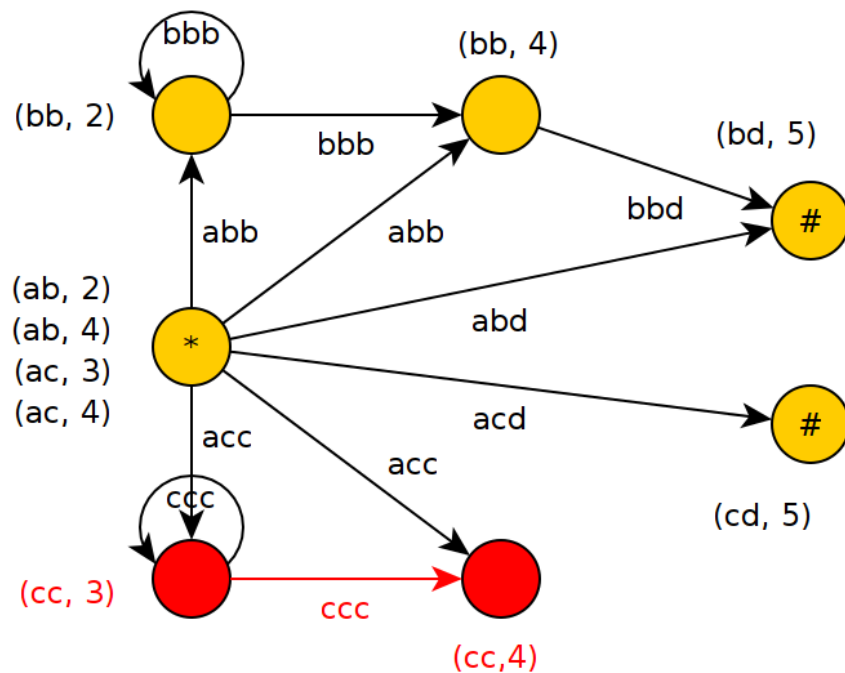
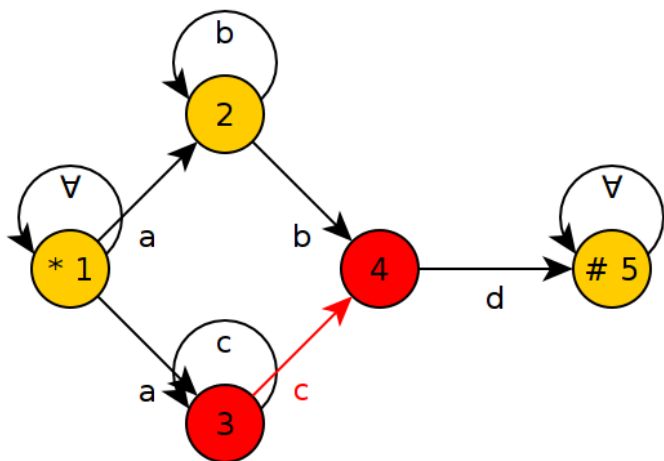
Transformation example



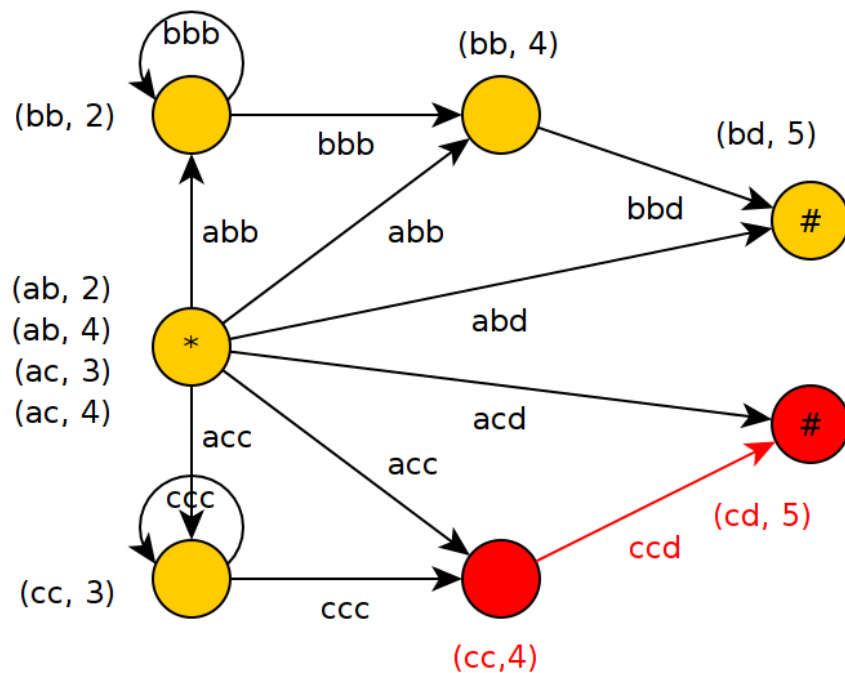
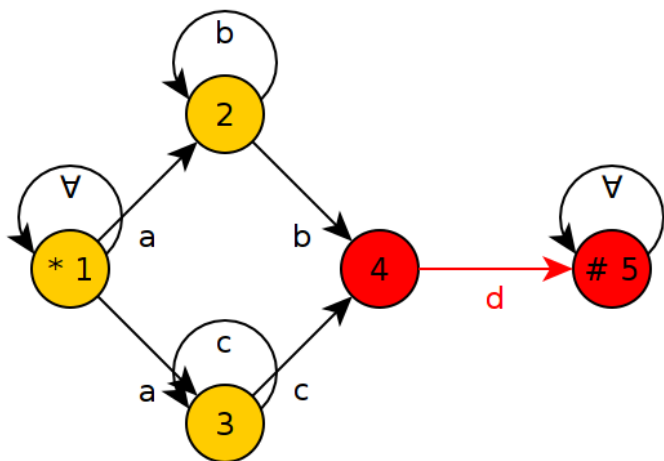
Transformation example



Transformation example

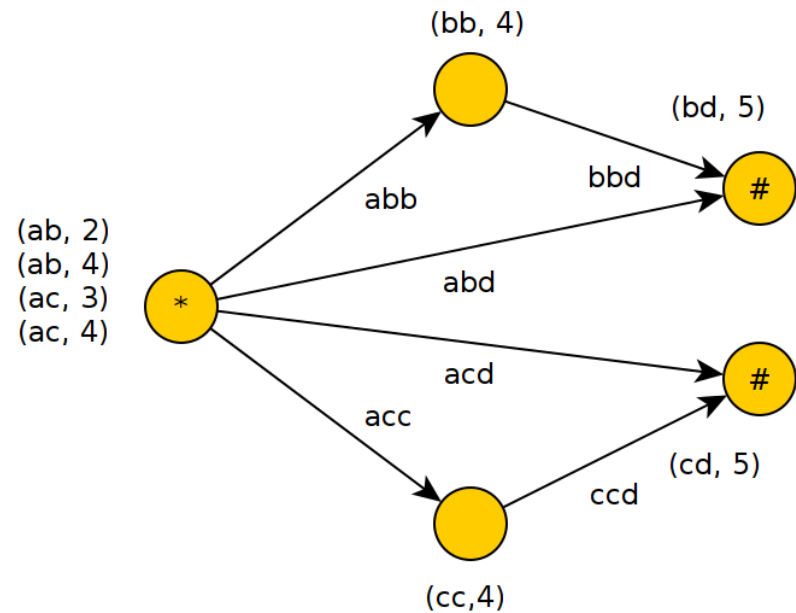
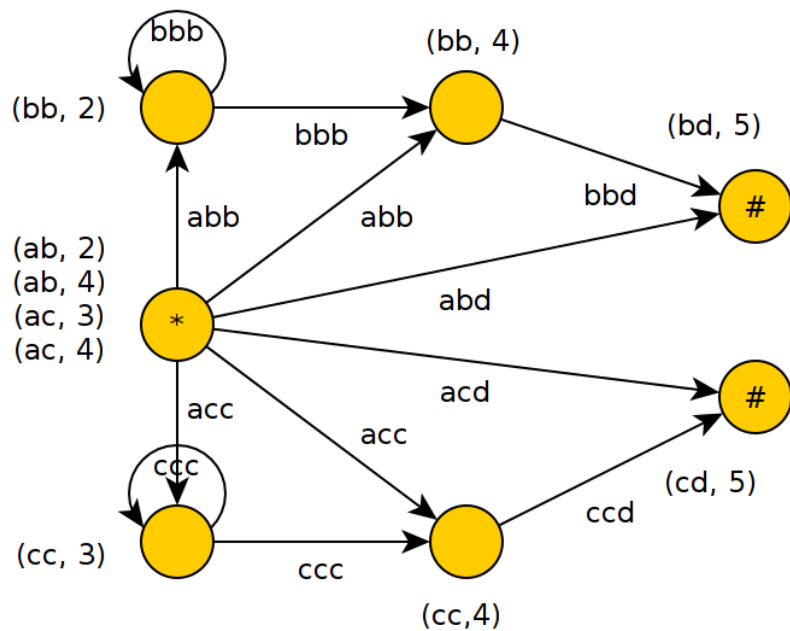


Transformation example



Transformation example

Result could be simplified



Transformation example

Implemented simplification technique: collect following matrix.

abd	abb	bbd	acd	acc	ccd
1	0	0	0	0	0
0	1	1	0	0	0
0	0	0	1	0	0
0	0	0	0	1	1

Means: $abd \text{ OR } (abb \text{ AND } bbd) \text{ OR } acd \text{ OR } (acc \text{ AND } ccd)$

Comparison on examples

Regex: `/(abc|cba)def/`

FREE: `(abc OR cba) AND def`

GSC:

`def AND ((abc AND bcd AND cde)
OR (ade AND bad AND cba))`

My:

`(abc AND bcd AND cde AND def) OR
(ade AND bad AND cba AND def)`

Comparison on examples

Regex: /abc+de/

FREE: nothing

GSC: abc AND cde

My:

(abc AND cde AND bcd) OR

(abc AND cde AND bcc AND ccd)

Comparison on examples

Regex: `/(abc*)+de/`

FREE: nothing

GSC: nothing

My:

(abd AND bde) OR

(abc AND bcd AND cde) OR

(abc AND bcc AND ccd AND cde)

Comparison on examples

Regex: `/ab(cd)*ef/`

FREE: nothing

GSC: nothing

My:

(abe AND bef) OR

(abc AND bde AND cde AND def)

Performance results

2.5 M DBLP paper titles of 47 avg. length

Regex	Index scan	Seq scan
/database.*(sql query)/	773 ms	18653 ms
/postgres(ql)?/	268 ms	17574 ms
/plan+er/	253 ms	12885 ms
/((nucl anino).*acid/	200 ms	20085 ms
/[aei](bc)+a/	2 ms	13195 ms

WIP patch for pg_trgm

WIP patch was posted to mailing list:

<http://archives.postgresql.org/pgsql-hackers/2011-11/msg01297.php>

Any feedback is welcome.

Problems

- Possible large transformed graph
- Possible large simplified presentation and it's high computational complexity
- Usage of trigrams rather than v-grams or multigrams

What did we miss?

- Difference between deterministic and non-deterministic automata
- Grouping characters into colors
- Handling of start/end of string/line

Help needed

Regular expressions and string collections from real-life tasks for proving effectiveness of proposed method.

The background features a complex, abstract pattern of overlapping, flowing lines in various shades of teal, light blue, and pale green. The lines are thin and have a soft, ethereal quality, creating a sense of movement and depth. A solid black horizontal band is positioned across the middle of the image, providing a high-contrast background for the white text.

Thank you for attention!