

# REDoS Detection in “Domino” Regular Expressions by Ambiguity Analysis

Nepeivoda A., Belikova Yu., Shevchenko K., Teriykha M.,  
Knyazikhin D., Delman A., Terentyeva A.  
PSI + BMSTU

SYRCoSE 2023

# REGULAR EXPRESSIONS

*(...known to be easier to write than to read)*

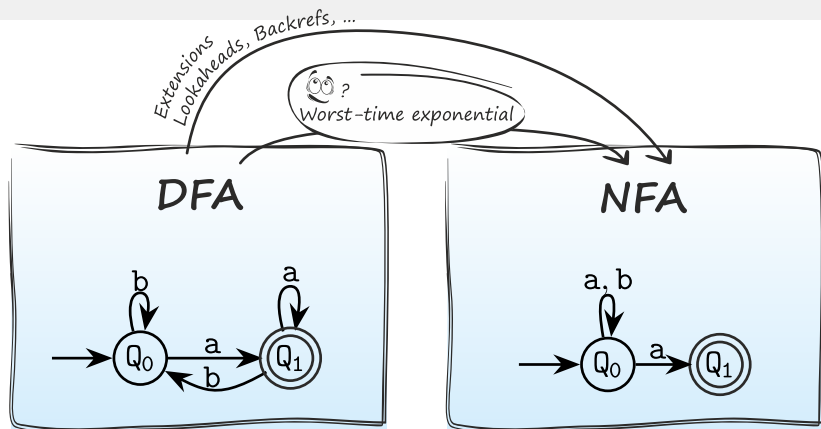
Academic Regular Expressions ( $\mathcal{R}$ ) is a set of expressions over alphabet  $\Sigma$  closed *w.r.t.*:

- concatenation,
- union,
- iteration (Kleene star).

Syntactic sugar:

- positive iteration,
- option,
- intervals,
- negation

# REGULAR EXPRESSIONS

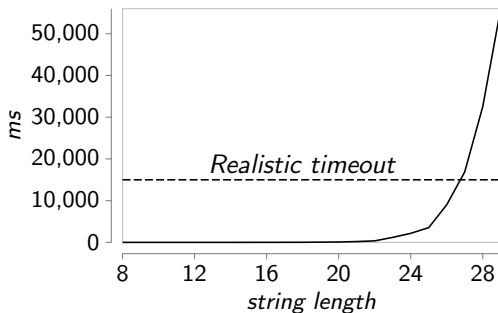


- Go
- RUST
- RE2

- PYTHON
- JAVASCRIPT
- JAVA
- PCRE

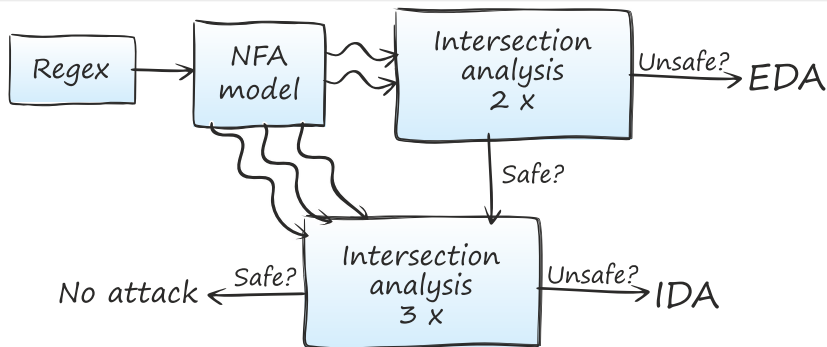
# CATASTROPHIC BACKTRACKING

- PYTHON regex engine
- Regular expression:  $(a^*b^*)^*$
- String:  $a^nc$



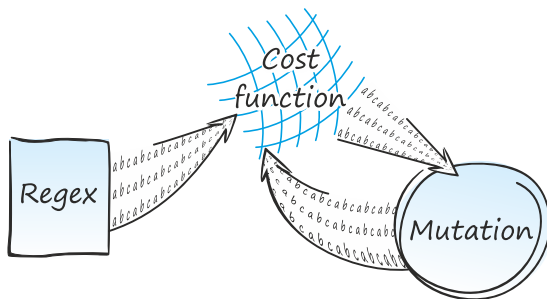
# REDoS DETECTION

- Static analysis
- Fuzz
- Combined + pattern-based methods



# REDoS DETECTION

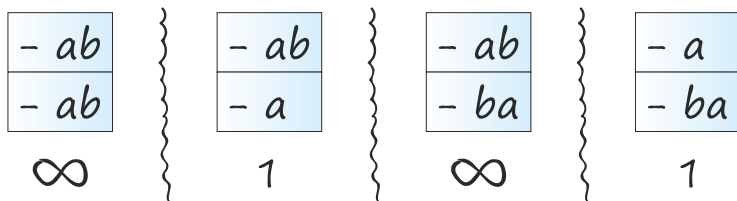
- Static analysis
- **Fuzz**
- Combined + pattern-based methods



- [illegible]

- 4 / 15

## COMPLEX AMBIGUITIES

$$(ab)^*a(ba)^*$$


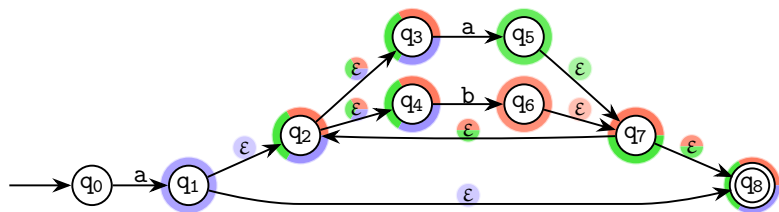
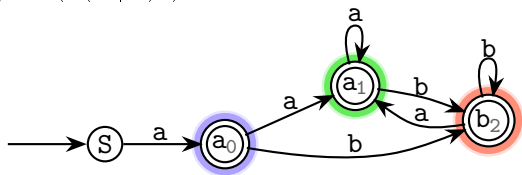
$$ababa$$




# RESEARCH QUESTIONS

- Are the “domino” regexes really dangerous?
  - what sorts of ambiguities can they contain?
  - what is the complexity of their matching in real regex engines?
- If yes, how to deal with them?
  - use another sort of static analysis?
  - or use heuristic-based approaches?

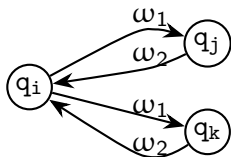
## NFA MODELS

(a) Thompson( $a(a \mid b)^*$ ) with colored  $\epsilon$ -closures(b) Glushkov( $a(a \mid b)^*$ )

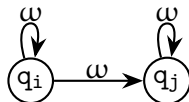
# DEGREE OF AMBIGUITY

## Definition

A degree of ambiguity of an NFA  $\mathcal{A}$  is a worst-case bound on the number of paths recognizing an input string (in a length of the string).



(a) Exponential ambiguity  
( $q_j$  or  $q_k$  may coincide with  $q_i$ )

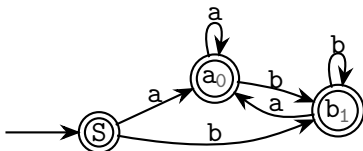


(b) Polynomial ambiguity

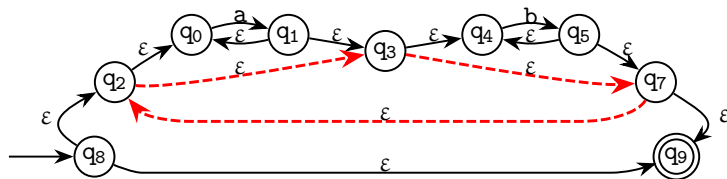
## SSNF PROPERTY

Regex  $(a^*b^*)^*$ .

Glushkov construction is deterministic:



Thompson construction is highly ambiguous and leads to critical timeouts:



# SSNF PROPERTY

*Weideman in his PhD (2015) uses the similar example to prove that the Glushkov construction is inadequate for REDoS detection. But still, Glushkov automaton can be used to detect ambiguities, preceded with a simple linear-time test!*

A regex  $r$  is in a (strong) star-normal form (SSNF) if it does not contain a subexpression  $(r_0)^*$  s.t.  $\varepsilon \in \mathcal{L}(r_0)$

If  $r$  is in the SSNF, then

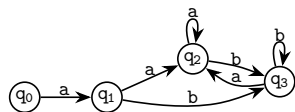
$$\text{Ambiguity}(\text{Thompson}(r)) = \text{Ambiguity}(\text{Glushkov}(r)).$$

# TRANSFORMATION MONOID

- Usual approach to detect ambiguities in NFA: iterative self-intersection.
  - Considering the transition function as the state acting over the letter set:  $F(q_i, \alpha) = F_{q_i}(\alpha) = \{q_{j_1}, \dots, q_{j_k}\}$ .
  - Costly, if the regex has a plenty of overlaps.
- Consider the transition function as the letter acting over the state set:  $F(q_i, \alpha) = F_\alpha(q_i) = \{q_{j_1}, \dots, q_{j_k}\}$ , and then  $F_\alpha(q_i) \circ F_\beta(q_i) = F_{q_i}(\beta\alpha) = F_{\beta\alpha}(q_i)$ .
  - State set is finite  $\Rightarrow$  set of the equivalence classes of the actions is finite.
  - The ambiguity criterion is then easily formulated, given the actions imposed by the members of the equivalence classes.

# TRANSFORMATION MONOID

Given an  $\varepsilon$ -free automaton  $\mathcal{A}$  over the alphabet  $\Sigma$ , its transformation monoid  $\mathcal{TM} = \text{Transmonoid}(\mathcal{A})$  is the monoid of transformations imposed by elements of  $\Sigma^*$  on the states of  $\mathcal{A}$ .



(a) Labelled transition system of  $\mathcal{A}$

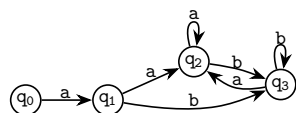
	$q_0$	$q_1$	$q_2$	$q_3$
a	$\{q_1\}$	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$
b	$\{\}$	$\{q_3\}$	$\{q_3\}$	$\{q_3\}$
aa	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$
ab	$\{q_3\}$	$\{q_3\}$	$\{q_3\}$	$\{q_3\}$
ba	$\{\}$	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$

(b) Equivalence classes

$bb \mapsto b$        $aaa \mapsto aa$   
 $aab \mapsto ab$      $aba \mapsto aa$   
 $baa \mapsto ba$      $bab \mapsto bb$

(c) Rewriting rules

# TRANSFORMATION MONOID



(a) Labelled transition system of  $\mathcal{A}$

	$q_0$	$q_1$	$q_2$	$q_3$
a	$\{q_1\}$	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$
b	$\{\}$	$\{q_3\}$	$\{q_3\}$	$\{q_3\}$
aa	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$
ab	$\{q_3\}$	$\{q_3\}$	$\{q_3\}$	$\{q_3\}$
ba	$\{\}$	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$

(b) Equivalence classes

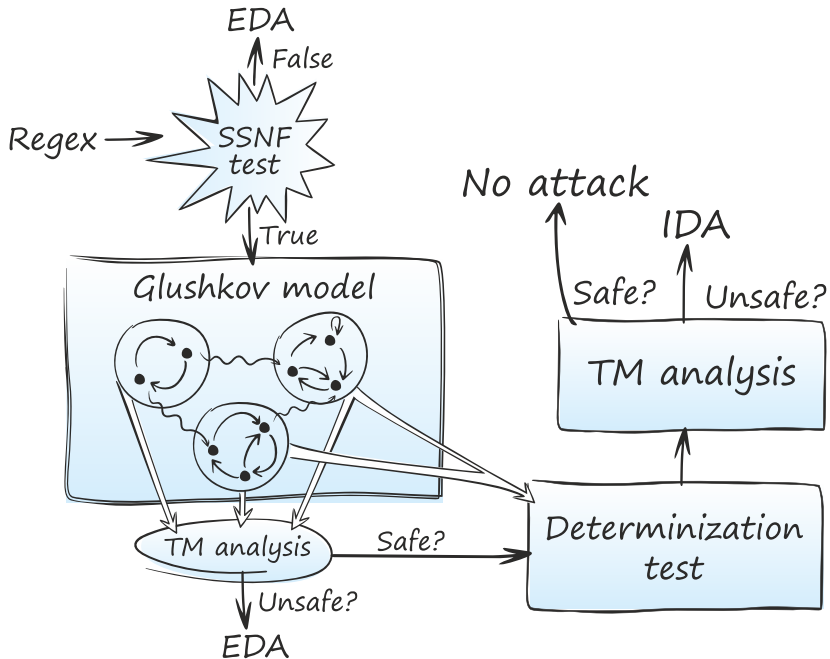
$bb \mapsto b$        $aaa \mapsto aa$   
 $aab \mapsto ab$      $aba \mapsto aa$   
 $baa \mapsto ba$      $bab \mapsto bb$

(c) Rewriting rules

An  $\varepsilon$ -free  $\mathcal{A}$  satisfies IDA  $\Leftrightarrow$  its transformation monoid contains an equivalence class  $\omega$  s.t. for some states  $q_i, q_j$ ,  $q_i \in F_\omega(q_i)$ ,  $q_j \in F_\omega(q_j)$ , and  $q_j \in F_\omega(q_i)$ .

- Costly, if the regex has few overlaps;
- Less costly, if applied to the strongly connected components (orbits) of the NFA;
- Since IDA (not EDA) can occur in the pair of the orbits, a preliminary test helps to accelerate the monoid criterion.





# OTHER TOOLS

*(open-source tools, described in research papers, and with adequate specification)*

- REGEXSTATICANALYZER, intersection-based analysis of academic regexes;
- RESCUE, fuzzing-based analysis of extended regexes;
- REVEALER, combined analysis of extended regexes.

*(did not succeed to find a working implementation)*

- REGULATOR, fuzzing based on regex structure;
- REDOSHUNTER, pattern-based approach, the patterns given in the paper do not detect dominoes;
- REDOSCALPEL, the same problem with the patterns.

# TEST SET

- length  $< 50$  terms, alphabet  $< 6$  letters (to construct non-trivial dominoes);
- alternation under an iteration, or several iterations (to have possible ambiguities).

## Ambiguity Classes

- Purely polynomial — critical on long input strings
- Non-SSNF — critical on short strings
- Polynomial under iteration — another critical case, non-trivial to detect

## RESULTS

Tool	Exponential		Polynomial		Safe		Unsafe		Timeouts
	$\mu$ (s)	$\sigma$ (s)	$\mu$ (s)	$\sigma$ (s)	$\mu$ (s)	$\sigma$ (s)	$\mu$ (s)	$\sigma$ (s)	
RSA	1.895	2.614	3.480	3.748	0.836	0.341	2.578	3.221	13
ReScue	–	–	–	–	0.940	1.724	8.803	6.263	43
Revealer	0.410	0.035	0.402	0.021	0.320	0.065	0.409	0.033	0
Our method	0.846	1.059	1.178	1.259	0.484	0.400	1.014	1.169	0

Tool	F <sub>1</sub> -score	Total error rate	Vuln. error rate
RSA	0.90	0.13	0.00
ReScue	0.39	–	–
Revealer	0.55	0.45	0.04
Our method	1.00	0.00	0.00

# CONCLUSION

Despite the success on the domino test set, extension to the whole class of regexes requires further developments:

- to apply alphabet factorization,
- to make use of the given approach for more precise analysis of extended regexes with complex overlaps,
- to combine the monoid-based and intersection-based ambiguity detection algorithms, based on heuristics detecting possible overlaps.