Analysis of an efficient reduction algorithm for random regular expressions based on universality detection

#### Pablo Rotondo

LIGM, Université Gustave Eiffel

Joint work with Florent Koechlin

**CSR 2021**, 30 June, 2021.

## Plan of the talk

1. Introduction: regular expression trees, uniform distribution

- 2. Semantic reductions: absorbing patterns, universality
- 3. Main results: expected size, proportion of universals
- 4. Techniques for the proof
- 5. Conclusions and further work

## Introduction: context

### Problem

Automatically test a program taking regular expressions as input

$$(a+b) \cdot b^{\star}$$
,  $(b \cdot (a+\varepsilon))^{\star}$ ,  $(a \cdot a^{\star}) + (b+a)^{\star}$ .

## Introduction: context

### Problem

Automatically test a program taking regular expressions as input

$$(a+b) \cdot b^{\star}, \qquad (b \cdot (a+\varepsilon))^{\star}, \qquad (a \cdot a^{\star}) + (b+a)^{\star}.$$

Example: automata constructions



h

b2

Introduction: random regular expressions



### Introduction: random regular expressions



#### **Generate a random expression tree**

- Realistic distribution
- Simple implementation, possibility of theoretical analysis.

#### **Expression trees:**

trees defined inductively,

$$\mathcal{L} = a_1 + \ldots + a_k + \varepsilon + \overset{\star}{\underset{\mathcal{L}}{\vdash}} + \overset{\bullet}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}},$$

• size |T| = number of nodes.

#### **Expression trees:**

trees defined inductively,

$$\mathcal{L} = a_1 + \ldots + a_k + \varepsilon + \overset{\star}{\underset{\mathcal{L}}{\vdash}} + \overset{\bullet}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}},$$

• size 
$$|T| =$$
 number of nodes.

**Idea:** Fix target size n,

#### **Expression trees:**

trees defined inductively,

$$\mathcal{L} = a_1 + \ldots + a_k + \varepsilon + \overset{\star}{\underset{\mathcal{L}}{\vdash}} + \overset{\bullet}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}},$$

• size |T| = number of nodes.

Idea: Fix target size n, pick tree T of size |T| = n uniformly

#### **Expression trees:**

trees defined inductively,

$$\mathcal{L} = a_1 + \ldots + a_k + \varepsilon + \overset{\star}{\underset{\mathcal{L}}{\vdash}} + \overset{\bullet}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}},$$

• size 
$$|T| =$$
 number of nodes.

**Idea:** Fix target size n, pick tree T of size |T| = n uniformly

natural a priori choice,

#### **Expression trees:**

trees defined inductively,

$$\mathcal{L} = a_1 + \ldots + a_k + \varepsilon + \overset{\star}{\underset{\mathcal{L}}{\vdash}} + \overset{\bullet}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}},$$

• size 
$$|T| =$$
 number of nodes.

**Idea:** Fix target size n, pick tree T of size |T| = n uniformly

natural a priori choice,

 efficient sampling (Boltzmann, Recursive, Devroye's constrainted GW),

#### **Expression trees:**

trees defined inductively,

$$\mathcal{L} = a_1 + \ldots + a_k + \varepsilon + \overset{\star}{\underset{\mathcal{L}}{\vdash}} + \overset{\bullet}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}},$$

• size 
$$|T| =$$
 number of nodes.

**Idea:** Fix target size n, pick tree T of size |T| = n uniformly

- natural a priori choice,
- efficient sampling

(Boltzmann, Recursive, Devroye's constrainted GW),

amenable to theoretical study (Analytic Combinatorics).

#### **Expression trees:**

trees defined inductively,

$$\mathcal{L} = a_1 + \ldots + a_k + \varepsilon + \overset{\star}{\underset{\mathcal{L}}{\vdash}} + \overset{\bullet}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}} + \overset{+}{\underset{\mathcal{L}}{\wedge}},$$

• size 
$$|T| =$$
 number of nodes.

**Idea:** Fix target size n, pick tree T of size |T| = n uniformly

- natural a priori choice,
- efficient sampling

(Boltzmann, Recursive, Devroye's constrainted GW),

amenable to theoretical study (Analytic Combinatorics).

 $\Longrightarrow$  Model used in numerous practical and theoretical works

Uniform expression trees [Koechlin,Nicaud,R. 2020]

Expected size after (linear) reduction is bounded O(1).

Uniform expression trees [Koechlin,Nicaud,R. 2020] Expected size after (linear) reduction is bounded O(1).

Universal result: not only regular expressions,

Uniform expression trees [Koechlin,Nicaud,R. 2020] Expected size after (linear) reduction is bounded O(1).

- Universal result: not only regular expressions,
- Absorbing patterns: only semantic hypothesis, absorbing pattern  $\mathcal{P}$ ,

Uniform expression trees [Koechlin,Nicaud,R. 2020] Expected size after (linear) reduction is bounded O(1).

- Universal result: not only regular expressions,
- Absorbing patterns: only semantic hypothesis, absorbing pattern  $\mathcal{P}$ ,

simplest case, false  $\land$  (...)  $\equiv$  false.

Uniform expression trees [Koechlin,Nicaud,R. 2020] Expected size after (linear) reduction is bounded O(1).

- Universal result: not only regular expressions,
- Absorbing patterns: only semantic hypothesis, absorbing pattern  $\mathcal{P}$ ,

$$\overset{\circledast}{\underset{\mathcal{P}}{\overset{}}_{T}} \rightsquigarrow \mathcal{P} \qquad \overset{\circledast}{\underset{T}{\overset{}}_{\mathcal{P}}} \rightsquigarrow \mathcal{P}$$

simplest case,  $false \land (...) \equiv false$ .

Wide variety of examples:



Uniform expression trees [Koechlin,Nicaud,R. 2020] Expected size after (linear) reduction is bounded O(1).

- Universal result: not only regular expressions,
- Absorbing patterns: only semantic hypothesis, absorbing pattern  $\mathcal{P}$ ,

$$\overset{\circledast}{\underset{\mathcal{P}}{\overset{}}_{T}} \rightsquigarrow \mathcal{P} \qquad \overset{\circledast}{\underset{T}{\overset{}}_{\mathcal{P}}} \rightsquigarrow \mathcal{P}$$

simplest case,  $false \land (...) \equiv false$ .

Wide variety of examples:



What does this say about regular expressions? O(1) ?

### Regular expressions: reduction by absorbing pattern

**Hidden constant** O(1): for regular expressions on two letters, the limit size after reduction is 3 624 217.



### Regular expressions: reduction by absorbing pattern

**Hidden constant** O(1): for regular expressions on two letters, the limit size after reduction is 3 624 217.



Question. Are uniform regular expressions useful nonetheless?





▶ We consider more specific algorithm based on *universality detection* 



expression is universal  $\Leftrightarrow$  equivalent to  $\Sigma^{\star}$ ,



▶ We consider more specific algorithm based on *universality detection* expression is universal  $\Leftrightarrow$  equivalent to  $\Sigma^{\star}$ ,

 $\Rightarrow$  substitute universal subtrees by smallest universal tree  ${\cal U}$  .

Idea: substitute universal subtrees by smallest universal tree  $\ensuremath{\mathcal{U}}$  .

Idea: substitute universal subtrees by smallest universal tree  ${\mathcal U}$  .

▶ We define bottom-up propagation rules

$$\overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{\to}{\underset{\mathcal{L}}{\overset{}}} \overset{U}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{\to}{\underset{\mathcal{U}}{\overset{}}} \overset{U}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{}}} \overset{\to}{\underset{\mathcal{U}}{\overset{}}} \overset{U}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\overset{\bullet}}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\overset{\bullet}}{\underset{\mathcal{U}}} \overset{\overset}{\overset}{\overset}{\overset}{\overset}{\overset}{\overset}}{\overset}{\overset}{\overset}{\overset}{\overset$$

Idea: substitute universal subtrees by smallest universal tree  ${\mathcal U}$  .

▶ We define bottom-up propagation rules

$$\overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{\to}{\underset{\mathcal{L}}{\overset{}}} \overset{U}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\overset{\bullet}}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\overset{\bullet}}{\underset{\mathcal{U}}} \overset{\bullet}{\overset{\bullet}}} \overset{\bullet}{\overset{\bullet}}{\overset{\bullet}}{\overset{\bullet}} \overset{\bullet}{\overset{$$

• Examples for  $\Sigma = \{a, b\}$ ,

$$\{a, b, \varepsilon\} \bullet \mathcal{U}$$

$$\{a, b, \varepsilon\} \star \mathcal{U} \qquad \mathbf{1}$$

$$\{a, b, \varepsilon\} + \mathcal{U} \qquad \mathbf{1}$$

$$\{a, b\} + \qquad a$$

$$a \qquad b$$

$$(\mathbf{I}) : (a + b)^{\star} \cdot a^{\star}$$

Idea: substitute universal subtrees by smallest universal tree  ${\mathcal U}$  .

▶ We define bottom-up propagation rules

$$\overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{L}}{\overset{}}} \overset{}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\overset}{\overset}}} \overset{\bullet}{\overset}{\overset{\bullet}}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\overset}}{\underset{\mathcal{U}}} \overset{\bullet}{\overset}}{\overset{\bullet}}{\overset{\bullet}}} \overset{\bullet}{\overset}{\overset}}$$

▶ Examples for  $\Sigma = \{a, b\}$ ,

Idea: substitute universal subtrees by smallest universal tree  ${\mathcal U}$  .

▶ We define bottom-up propagation rules

$$\overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{L}}{\overset{}}} \overset{}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\overset}{\overset{\bullet}}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}}{\overset{\overset}{\overset}{\overset}}{\overset}{\overset}}} \overset{\overset}{\overset}{\overset}{\overset$$

► Examples for  $\Sigma = \{a, b\}$ ,

$$\begin{array}{c} \{a,b,\varepsilon\} \bullet \mathcal{U} \\ \{a,b,\varepsilon\} \star \mathcal{U} \\ \{a,b,\varepsilon\} \star \mathcal{U} \\ a,b,\varepsilon\} \star \mathcal{U} \\ b,\varepsilon\} \star \mathcal{U} \\ c,c,\varepsilon\} \star \mathcal{U} \\ c$$

▶ Detection is only partial: example  $\Sigma \cdot \Sigma^{\star} + \varepsilon$ 

Idea: substitute universal subtrees by smallest universal tree  $\ensuremath{\mathcal{U}}$  .

▶ We define bottom-up propagation rules

$$\overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{\to}{\underset{\mathcal{L}}{\overset{}}} \overset{U}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{\to}{\underset{\mathcal{U}}{\overset{}}} \overset{U}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{}}} \overset{\to}{\underset{\mathcal{U}}{\overset{}}} \overset{U}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}{\overset{\bullet}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\overset{\bullet}}{\underset{\mathcal{U}}} \overset{\bullet}{\underset{\mathcal{U}}} \overset{\bullet}{\overset{\bullet}}{\underset{\mathcal{U}}} \overset{\overset}{\overset}{\overset}{\overset}{\overset}{\overset}{\overset}}{\overset}{\overset}{\overset}{\overset}{\overset$$

► Examples for  $\Sigma = \{a, b\}$ ,

$$\begin{array}{c} \{a,b,\varepsilon\} \bullet \mathcal{U} \\ \{a,b,\varepsilon\} \star \mathcal{U} \\ \{a,b,\varepsilon\} \star \mathcal{U} \\ a,b,\varepsilon\} \star \mathcal{U} \\ b,c\} \star \mathcal{U} \\ c,c\} \\ c,$$

► Detection is only partial: example Σ · Σ\* + ε ⇒ universality problem is PSPACE-complete !

Rewriting rules:

$$\overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{L}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}} \overset{+}{\underset{\mathcal{U}}} \overset{+}}{\overset{}}} \overset{}}{\overset{}}{\overset{}}{\overset{}}} \overset{}}{\overset{}}}{\overset{}}} \overset{}}{\overset{}}{\overset{$$

Rewriting rules:

$$\overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{L}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\overset{}}{\overset{}}}{\overset{}}} \overset{}}{\overset{}}} \overset{}}{\overset{}}{\overset{}}{\overset{}}}{\overset{}}}{\overset{}}}{\overset{}}}{\overset{}}{\overset{}}{\overset{}}{\overset{}$$

#### Main result

Consider the regular expression trees over  $\Sigma = \{a_1, \ldots, a_k\}$ . Take the bottom-up simplification  $\sigma$  induced by our rewriting rules.

Rewriting rules:

$$\overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{L}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}} \overset{}}{\overset{}} \overset{}}{\overset{}}}{\overset{}}} \overset{}}{\overset{}}} \overset{}}{\overset{}}{\overset{}}{\overset{}}$$

#### Main result

Consider the regular expression trees over  $\Sigma = \{a_1, \ldots, a_k\}$ . Take the bottom-up simplification  $\sigma$  induced by our rewriting rules.

Then the expected size of the simplification of a random uniform tree tends to a **constant** as the size n tends to infinity.

Moreover, the constant can be computed efficiently

$ \Sigma $	2	3	4	5
$\lim \mathbb{E}_n[ \sigma(T) ]$	77.79724	495.59151	$2518.20513\ldots$	$11694.43727\ldots$

Rewriting rules:

$$\overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{L}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{-}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}}} \overset{+}{\underset{\mathcal{U}}{\overset{}}} \overset{+}{\underset{\mathcal{U}}} \overset{+}{\underset{\mathcal{U}}} \overset{+}}{\overset{}}} \overset{}}{\overset{}}{\overset{}}{\overset{}}} \overset{}}{\overset{}}}{\overset{}}} \overset{}}{\overset{}}{\overset{$$

#### Main result

Consider the regular expression trees over  $\Sigma = \{a_1, \ldots, a_k\}$ . Take the bottom-up simplification  $\sigma$  induced by our rewriting rules.

Then the expected size of the simplification of a random uniform tree tends to a **constant** as the size n tends to infinity.

Moreover, the constant can be computed efficiently

$ \Sigma $	2	3	4	5
$\lim \mathbb{E}_n[ \sigma(T) ]$	77.79724	495.59151	$2518.20513\ldots$	$11694.43727\ldots$

Note. Compare  $\sim 77.8$  against previous  $\sim 3.6 \times 10^6$  for two letters.

## Results: plots



## Results: plots



### Proposition

For n large enough, the proportion  $Pr_n(univ.)$  of universal expressions trees belongs to the intervals:

$ \Sigma $	2	3	4	5
interval	(0.31, 0.46)	(0.13, 0.27)	(0.062, 0.15)	(0.028, 0.077)

### Proposition

For n large enough, the proportion  $Pr_n(univ.)$  of universal expressions trees belongs to the intervals:

$ \Sigma $	2	3	4	5
interval	(0.31, 0.46)	(0.13, 0.27)	(0.062, 0.15)	(0.028, 0.077)

- Preponderance of universal expression trees: between 31% and 46% for two letters {a, b}
- Uniform model not adapted to sampling regular languages

We employ Analytic Combinatorics to study the expectation,

### Bivariate generating function

$$L(z,u) := \sum_{T \in \mathcal{L}} u^{|\sigma(T)|} z^{|T|} \implies \mathbb{E}_n[|\sigma(T)|] = \frac{[z^n]\partial_u L(z,u)|_{u=1}}{[z^n]L(z,u)|_{u=1}},$$

encodes input and output sizes.

We employ Analytic Combinatorics to study the expectation,

Bivariate generating function

$$L(z,u) := \sum_{T \in \mathcal{L}} u^{|\sigma(T)|} z^{|T|} \implies \mathbb{E}_n[|\sigma(T)|] = \frac{[z^n]\partial_u L(z,u)|_{u=1}}{[z^n]L(z,u)|_{u=1}},$$

encodes input and output sizes.

Symbolic Step. We find a formal equation describing L(z, u).
 Here this is done from a *combinatorial specificiation*

$$\boldsymbol{y}(z, u) = \boldsymbol{F}(z, u; \boldsymbol{y}(z, u)).$$

We employ Analytic Combinatorics to study the expectation,

Bivariate generating function

$$L(z,u) := \sum_{T \in \mathcal{L}} u^{|\sigma(T)|} z^{|T|} \implies \mathbb{E}_n[|\sigma(T)|] = \frac{[z^n]\partial_u L(z,u)|_{u=1}}{[z^n]L(z,u)|_{u=1}},$$

encodes input and output sizes.

Symbolic Step. We find a formal equation describing L(z, u).
 Here this is done from a *combinatorial specificiation*

$$\boldsymbol{y}(z, u) = \boldsymbol{F}(z, u; \boldsymbol{y}(z, u)).$$

▶ Analytic Step. A *Transfer Theorem* links the behaviour at dominant singularities  $\rho \in \mathbb{C}$  to asymptotics of coefficients

$$L(z) \sim_{z \to \rho} \lambda (1 - z/\rho)^{-\alpha} \Longrightarrow [z^n] L(z) \sim \lambda n^{\alpha - 1} / \Gamma(\alpha) \rho^{-n}.$$

We employ Analytic Combinatorics to study the expectation,

Bivariate generating function

$$L(z,u) := \sum_{T \in \mathcal{L}} u^{|\sigma(T)|} z^{|T|} \implies \mathbb{E}_n[|\sigma(T)|] = \frac{[z^n]\partial_u L(z,u)|_{u=1}}{[z^n]L(z,u)|_{u=1}},$$

encodes input and output sizes.

Symbolic Step. We find a formal equation describing L(z, u).
 Here this is done from a *combinatorial specificiation*

$$\boldsymbol{y}(z, u) = \boldsymbol{F}(z, u; \boldsymbol{y}(z, u)).$$

▶ Analytic Step. A *Transfer Theorem* links the behaviour at dominant singularities  $\rho \in \mathbb{C}$  to asymptotics of coefficients

$$L(z) \sim_{z \to \rho} \lambda (1 - z/\rho)^{-\alpha} \Longrightarrow [z^n] L(z) \sim \lambda n^{\alpha - 1} / \Gamma(\alpha) \rho^{-n}.$$

 $\Rightarrow$  Study asymptotics over  $z\in\mathbb{C}$ 

# Combinatorial specification: two letters $\Sigma = \{a, b\}$

For every  $X \subseteq \{a, b\}$  introduce:

- *T*<sub>X,ε</sub> the set of trees recognizing every letter in X and ε, and no letter not in X
- *T*<sub>X,ε̄</sub> the set of trees recognizing every letter in X, and no letter not in X, nor ε

## Combinatorial specification: two letters $\Sigma = \{a, b\}$

For every  $X \subseteq \{a, b\}$  introduce:

- *T*<sub>X,ε</sub> the set of trees recognizing every letter in X and ε, and no letter not in X
- *T*<sub>X,ε̄</sub> the set of trees recognizing every letter in X, and no letter not in X, nor ε

$$\begin{split} \mathcal{T}_{X,\varepsilon} &= \varepsilon \mathbf{1}_{X=\emptyset} + \overset{\star}{\tau_{X,\varepsilon}} + \overset{\star}{\tau_{X,\overline{\varepsilon}}} + \sum_{(S,S'):S\cup S'=X} \tau_{S,\varepsilon} \overset{\wedge}{\tau_{S',\varepsilon}} \\ &+ \sum_{(S,S'):S\cup S'=X} \tau_{S,\varepsilon} \overset{+}{\tau_{S',\varepsilon}} + \sum_{(S,S'):S\cup S'=X} \tau_{S,\varepsilon} \overset{+}{\tau_{S',\varepsilon}} + \sum_{(S,S'):S\cup S'=X} \tau_{S,\overline{\varepsilon}} \overset{+}{\tau_{S',\varepsilon}} \\ \mathcal{T}_{X,\overline{\varepsilon}} &= X \mathbf{1}_{|X|=1} + \sum_{S\subseteq\Sigma} \tau_{X,\overline{\varepsilon}} \overset{\wedge}{\tau_{S,\varepsilon}} + \sum_{S\subseteq\Sigma} \tau_{S,\varepsilon} \overset{\wedge}{\tau_{X,\overline{\varepsilon}}} + \mathbf{1}_{X=\emptyset} \sum_{S,S'\subseteq\Sigma} \tau_{S,\overline{\varepsilon}} \overset{\wedge}{\tau_{S',\overline{\varepsilon}}} \\ &+ \sum_{(S,S'):S\cup S'=X} \tau_{S,\overline{\varepsilon}} \overset{+}{\tau_{S',\overline{\varepsilon}}}, \end{split}$$

### Example: combinatorial specification

Trees recognizing the letter a and no other letter, and not recognizing arepsilon



### Definition (Fully reducible expressions)

A regular expression tree T is fully reducible when  $\sigma(T) = \mathcal{U}$ . In other words, it is recognized as universal by our algorithm.

### Definition (Fully reducible expressions)

A regular expression tree T is fully reducible when  $\sigma(T) = U$ . In other words, it is recognized as universal by our algorithm.

**Dictate** the reduction process: leaves of reduced expression.

### Definition (Fully reducible expressions)

A regular expression tree T is fully reducible when  $\sigma(T) = \mathcal{U}$ . In other words, it is recognized as universal by our algorithm.

Dictate the reduction process: leaves of reduced expression.

Size after reduction  $p = |\mathcal{U}|$  for  $T \in \mathcal{R}$ .

### Definition (Fully reducible expressions)

A regular expression tree T is fully reducible when  $\sigma(T) = U$ . In other words, it is recognized as universal by our algorithm.

- Dictate the reduction process: leaves of reduced expression.
- Size after reduction  $p = |\mathcal{U}|$  for  $T \in \mathcal{R}$ .
- $\blacktriangleright$  The class of fully reducible trees  ${\cal R}$  satisfies the equation:

$$\mathcal{R} = \overset{\star}{\mathcal{T}_{\Sigma,\overline{\varepsilon}}} + \overset{\star}{\mathcal{T}_{\Sigma,\varepsilon}} + \overset{+}{\mathcal{N}_{\mathcal{L}}} + \overset{+}{\mathcal{N}_{\mathcal{L}}} + \overset{+}{\mathcal{N}_{\mathcal{L}}} + \overset{\bullet}{\mathcal{N}_{\mathcal{L}}} + \overset{\bullet}{\mathcal{N}_{\mathcal{E}}} + \overset{\bullet}{\mathcal{T}_{\varepsilon} \setminus \mathcal{R}} \mathcal{R}.$$

 $\implies$  completes the combinatorial specification of L(z, u).

## Solving efficiently: auxiliary classes

• every tree : 
$$\mathcal{L} = \bigcup_X \mathcal{T}_{X,\varepsilon} \cup \mathcal{T}_{X,\overline{\varepsilon}}$$
  
 $\mathcal{L} = a + b + \varepsilon + \overset{*}{\mathcal{L}} + \overset{\bullet}{\mathcal{L}}_{\mathcal{L}} + \overset{+}{\mathcal{L}}_{\mathcal{L}}^{\dagger}_{\mathcal{L}}$   
 $L(z) = 3z + zL(z) + 2z(L(z))^2$   
• trees recognizing  $\varepsilon$  :  $\mathcal{T}_{\varepsilon} = \bigcup_X \mathcal{T}_{X,\varepsilon}$   
 $\mathcal{T}_{\varepsilon} = \varepsilon + \overset{*}{\mathcal{L}} + \overset{\bullet}{\mathcal{T}_{\varepsilon}} + \overset{+}{\mathcal{T}_{\varepsilon}} \overset{+}{\mathcal{L}} + \overset{+}{\mathcal{L}} \overset{+}{\mathcal{T}_{\varepsilon}} \tau_{\varepsilon}$   
 $T_{\varepsilon}(z) = \frac{z + zL(z)}{1 - 2zL(z)}$   
• trees not recognizing  $\varepsilon$  :  $T_{\overline{\varepsilon}}(z) = L(z) - T_{\varepsilon}(z)$ 

### The system becomes triangular

2

$$\begin{split} T_{\emptyset,\overline{\varepsilon}}(z) &= function(T_{\emptyset,\overline{\varepsilon}}(z)) \\ T_{\{a\},\overline{\varepsilon}}(z) &= function(T_{\{a\},\overline{\varepsilon}}(z),T_{\emptyset,\overline{\varepsilon}}(z)) \\ T_{\{b\},\overline{\varepsilon}}(z) &= function(T_{\{b\},\overline{\varepsilon}}(z),T_{\emptyset,\overline{\varepsilon}}(z)) \\ T_{\{a,b\},\overline{\varepsilon}}(z) &= function(T_{\{a,b\},\overline{\varepsilon}}(z),T_{\{a\},\overline{\varepsilon}}(z),T_{\{b\},\overline{\varepsilon}}(z),T_{\emptyset,\overline{\varepsilon}}(z)) \\ T_{\emptyset,\varepsilon}(z) &= function(T_{\emptyset,\varepsilon}(z),T_{\emptyset,\overline{\varepsilon}}(z)) \end{split}$$

 $T_{\{a,b\},\varepsilon}(z) = function(T_{\{a,b\},\varepsilon}(z), \text{and everyone above})$ 

**>** Each equation is of degree  $2 \Rightarrow$  exactly solvable

$$\begin{split} T_{\{a,b\},\overline{z}}(z) &= \tfrac{1}{4z} \Big( -\sqrt{\Delta(z)} + 2\sqrt{(2z+2)}\,\sqrt{\Delta(z)} - 6z^2 + 2} - \sqrt{(2z+2)}\,\sqrt{\Delta(z)} + 10z^2 + 2 - z - 1 \Big)\,, \\ \text{where } \Delta(z) \text{ is the determinant of the equation for } L(z). \end{split}$$

The expression

$$T_{\{a,b\},\overline{z}}(z) = \frac{1}{4z} \left( -\sqrt{\Delta(z)} + 2\sqrt{(2z+2)}\sqrt{\Delta(z)} - 6z^2 + 2 - \sqrt{(2z+2)}\sqrt{\Delta(z)} + 10z^2 + 2 - z - 1 \right),$$

implies a square-root behaviour

$$T_{\{a,b\},\overline{\varepsilon}}(z) \sim A - B\sqrt{1-z/\rho}$$

for z close to dominant singularity  $\rho.$ 

The expression

 $T_{\{a,b\},\overline{\varepsilon}}(z) = \frac{1}{4z} \left( -\sqrt{\Delta(z)} + 2\sqrt{(2z+2)}\sqrt{\Delta(z)} - 6z^2 + 2 - \sqrt{(2z+2)}\sqrt{\Delta(z)} + 10z^2 + 2 - z - 1 \right),$ 

implies a square-root behaviour

$$T_{\{a,b\},\overline{\varepsilon}}(z) \sim A - B\sqrt{1 - z/\rho}$$

for z close to dominant singularity  $\rho$ .

More generally

▶ square-root behaviour generalizes to  $T_{X,\varepsilon}$  and  $T_{X,\overline{\varepsilon}}$ ,

The expression

 $T_{\{a,b\},\overline{\varepsilon}}(z) = \frac{1}{4z} \left( -\sqrt{\Delta(z)} + 2\sqrt{(2z+2)}\sqrt{\Delta(z)} - 6z^2 + 2 - \sqrt{(2z+2)}\sqrt{\Delta(z)} + 10z^2 + 2 - z - 1 \right),$ 

implies a square-root behaviour

$$T_{\{a,b\},\overline{\varepsilon}}(z) \sim A - B\sqrt{1 - z/\rho}$$

for z close to dominant singularity  $\rho$ .

More generally

Square-root behaviour generalizes to  $T_{X,\varepsilon}$  and  $T_{X,\overline{\varepsilon}}$ , and for every  $k = |\Sigma| \Rightarrow$  use Drmota's Theorem.

The expression

 $T_{\{a,b\},\overline{z}}(z) = \frac{1}{4z} \left( -\sqrt{\Delta(z)} + 2\sqrt{(2z+2)}\sqrt{\Delta(z)} - 6z^2 + 2 - \sqrt{(2z+2)}\sqrt{\Delta(z)} + 10z^2 + 2 - z - 1 \right),$ 

implies a square-root behaviour

$$T_{\{a,b\},\overline{\varepsilon}}(z) \sim A - B\sqrt{1 - z/\rho}$$

for z close to dominant singularity  $\rho$ .

More generally

square-root behaviour generalizes to T<sub>X,ε</sub> and T<sub>X,ε</sub>, and for every k = |Σ| ⇒ use Drmota's Theorem.
 then to ∂<sub>u</sub>L(z, u)|<sub>u=1</sub>, numerator of the expectation. [Closure]

The expression

 $T_{\{a,b\},\overline{z}}(z) = \frac{1}{4z} \left( -\sqrt{\Delta(z)} + 2\sqrt{(2z+2)}\sqrt{\Delta(z)} - 6z^2 + 2 - \sqrt{(2z+2)}\sqrt{\Delta(z)} + 10z^2 + 2 - z - 1 \right),$ 

implies a square-root behaviour

$$T_{\{a,b\},\overline{\varepsilon}}(z) \sim A - B\sqrt{1 - z/\rho}$$

for z close to dominant singularity  $\rho$ .

More generally

square-root behaviour generalizes to T<sub>X,ε</sub> and T<sub>X,ε</sub>, and for every k = |Σ| ⇒ use Drmota's Theorem.
 then to ∂<sub>u</sub>L(z, u)|<sub>u=1</sub>, numerator of the expectation. [Closure]

► Coefficients A and B determine asymptotics [Transfer Theorem]

The expression

 $T_{\{a,b\},\overline{\varepsilon}}(z) = \frac{1}{4z} \Big( -\sqrt{\Delta(z)} + 2\sqrt{(2z+2)}\sqrt{\Delta(z)} - 6z^2 + 2 - \sqrt{(2z+2)}\sqrt{\Delta(z)} + 10z^2 + 2 - z - 1 \Big) \,,$ 

implies a square-root behaviour

$$T_{\{a,b\},\overline{\varepsilon}}(z) \sim A - B\sqrt{1 - z/\rho}$$

for z close to dominant singularity  $\rho$ .

More generally

square-root behaviour generalizes to T<sub>X,ε</sub> and T<sub>X,ε</sub>, and for every k = |Σ| ⇒ use Drmota's Theorem.
 then to ∂<sub>u</sub>L(z, u)|<sub>u=1</sub>, numerator of the expectation. [Closure]

Coefficients A and B determine asymptotics [Transfer Theorem] we show how to compute these efficiently.

## Conclusion and further work

- We have shown a simple linear algorithm, reducing uniform regular expressions to small constant size.
- Therefore, uniform random regular expression trees tend to describe very limited languages.

## Conclusion and further work

- We have shown a simple linear algorithm, reducing uniform regular expressions to small constant size.
- Therefore, uniform random regular expression trees tend to describe very limited languages.

Future work

- Other distributions seem more appropriate (BST, ...)
- Algorithm (partially) detects universality, improvements ?

## Conclusion and further work

- We have shown a simple linear algorithm, reducing uniform regular expressions to small constant size.
- Therefore, uniform random regular expression trees tend to describe very limited languages.

Future work

- Other distributions seem more appropriate (BST, ...)
- Algorithm (partially) detects universality, improvements ?

Thank you!