

Unambiguous Parikh automata and holonomic series

Florent Koechlin
LIPN Back-to-school day, Villetaneuse

October 10th, 2023

Curriculum



Arnaud Carayol
Cyril Nicaud
Pablo Rotondo
Alin Bostan (Inria Saclay)



Mathilde Bouvel
Valentin Féray (IECL)
Xavier Goac

My research areas



Combinatorics

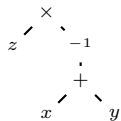
$$\sum_{n \geq 0} a_n x^n$$



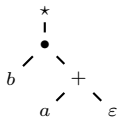
My research areas



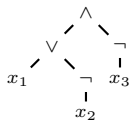
Random
expression trees



$$z \times (x + y)^{-1}$$



$$(b \cdot (a + \varepsilon))^*$$



$$(x_1 \vee \neg x_2) \wedge \neg x_3$$



Thesis



Combinatorics

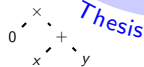
$$\sum_{n \geq 0} a_n x^n$$



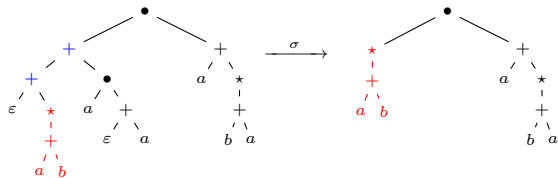
My research areas



Random
expression trees



Reduction by absorbing pattern:



Combinatorics

$$\sum_{n \geq 0} a_n x^n$$

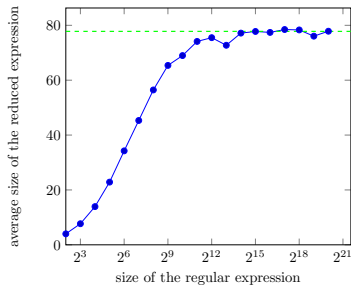


My research areas



Uniform trees

Random expression trees



Combinatorics

$$\sum_{n \geq 0} a_n x^n$$



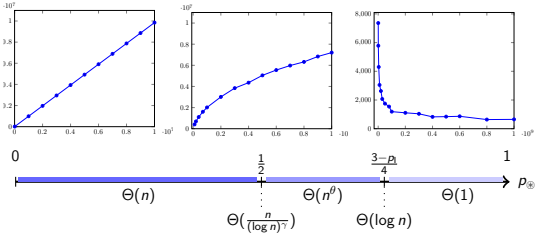
My research areas



Uniform trees
BST distribution

Random expression trees

Thesis



Combinatorics

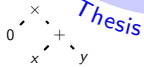
$\sum_{n \geq 0} a_n x^n$

My research areas



Uniform trees
BST distribution

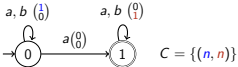
Random expression trees



Thesis

Combinatorics

$\sum_{n \geq 0} a_n x^n$



Unambiguous PA
Bounded CFL

Automata theory

Thesis

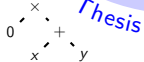
Computer Algebra

My research areas

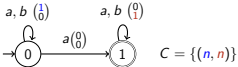


Uniform trees
BST distribution

Random expression trees



Thesis



Unambiguous PA
Bounded CFL

Automata theory

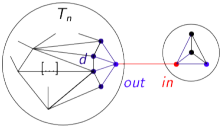
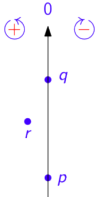
Thesis

Combinatorics

$$\sum_{n \geq 0} a_n x^n$$

Computer Algebra

Postdoc
Combinatorial geometry

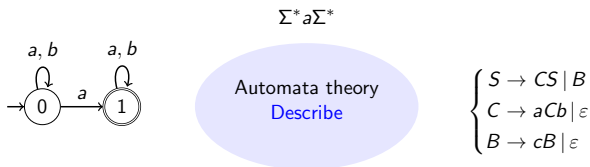


Automata theory

Language: abstract set of words built over a finite alphabet Σ .

Example

The set of words over $\{a, b\}$ containing at least one letter a :
a, ba, ab, aa, baa, aba, abb, ...



Link between describing and counting

Generating series

Let L be a language, ℓ_n the number of words in L of length n :

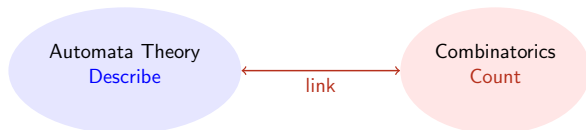
$$L(x) = \sum_{n=0}^{+\infty} \ell_n x^n$$

Example

$$b^* a (a + b)^* \rightarrow L(x) = \sum_{n \geq 0} (2^n - 1) x^n = \frac{x}{(1-x)(1-2x)}$$

Example

$$\text{Well bracketed words} \rightarrow L(x) = \sum_{n \geq 0} \frac{1}{n+1} \binom{2n}{n} x^{2n} = \frac{1 - \sqrt{1-4x^2}}{2x^2}$$

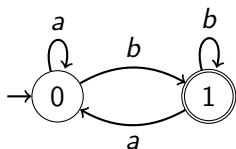


Link between automata and generating series

regular
languages

\subseteq

unambiguous
context-free languages



$$\begin{cases} S \rightarrow aSB \mid \varepsilon \\ B \rightarrow cB \mid bS \end{cases}$$

$$\begin{cases} q_0(x) = xq_0(x) + xq_1(x) \\ q_1(x) = 1 + xq_1(x) + xq_0(x) \end{cases}$$

$$q_0(x) = \frac{x}{1-2x}$$

$$\begin{cases} S(x) = xS(x)B(x) + 1 \\ B(x) = xB(x) + xS(x) \end{cases}$$

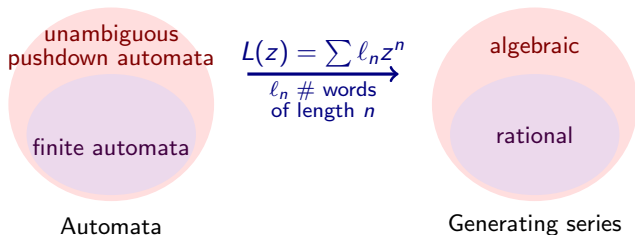
$$x^2S(x)^2 - (1-x)S(x) + 1 - x = 0$$

rational series
 $L(x) = \frac{P(x)}{Q(x)}$

\subseteq

algebraic series
 $P(L(x), x) = 0$

Link between two hierarchies



Two remarkable applications :

- analytic proofs of **inherent ambiguity** [Flajolet 87]

Analytic criteria for inherent ambiguity

Theorem [Flajolet '87]

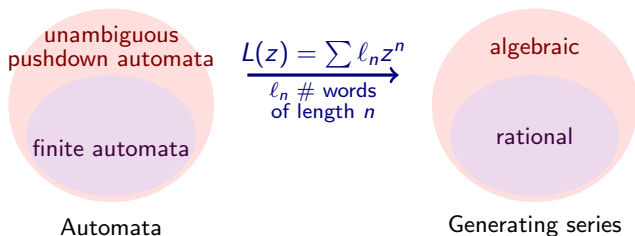
$\Omega_3 = \{w \in \{a, b, c\}^* : |w|_a \neq |w|_b \text{ or } |w|_b \neq |w|_c\}$ is inherently ambiguous.

Analytic proof:

- Suppose that $\Omega_3(x)$ is algebraic
- Let $I = (a + b + c)^* \setminus \Omega_3$
- Then $I(x) = \frac{1}{1-3x} - \Omega_3(x)$ would be algebraic by closure properties
- But $I = \{w \in \{a, b, c\}^* : |w|_a = |w|_b = |w|_c\}$

$$[x^{3n}]I(x) = \binom{3n}{n, n, n} = \frac{(3n)!}{(n!)^3} \sim_{n \rightarrow \infty} 3^{3n} \frac{\sqrt{3}}{2\pi n}$$

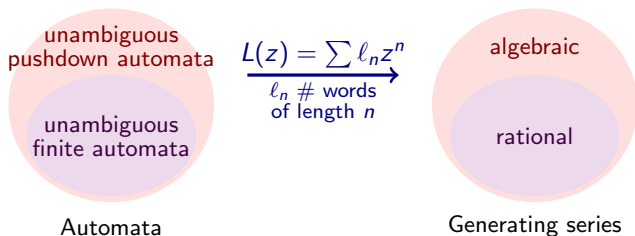
Link between two hierarchies



Two remarkable applications :

- analytic proofs of **inherent ambiguity** [Flajolet 87]
- **polynomial algorithm** for the inclusion problem for unambiguous NFA's [Stearns & Hunt 85]

Link between two hierarchies



Two remarkable applications :

- analytic proofs of **inherent ambiguity** [Flajolet 87]
- **polynomial algorithm** for the inclusion problem for unambiguous NFA's [Stearns & Hunt 85]

Inclusion problem for unambiguous automata

Problem : Given \mathcal{A} and \mathcal{B} two unambiguous NFA, $L(\mathcal{A}) \subseteq L(\mathcal{B})$?

Proposition: If $L(\mathcal{A}) \subsetneq L(\mathcal{B})$, there exists a **small witness** $w \in L(\mathcal{B}) \setminus L(\mathcal{A})$ of size at most $|Q_{\mathcal{A}}| + |Q_{\mathcal{B}}|$

- $C(x) := \sum c_n x^n = B(x) - A(x)$ is **rational**
- The coefficients of $C(x)$ satisfy a linear recurrence:

$$\forall n \geq r, c_n = \alpha_1 c_{n-1} + \dots + \alpha_r c_{n-r}$$

- the **order** r is at most $|Q_{\mathcal{A}}| + |Q_{\mathcal{B}}|$



Inclusion problem for unambiguous automata

Problem : Given \mathcal{A} and \mathcal{B} two unambiguous NFA, $L(\mathcal{A}) \subseteq L(\mathcal{B})$?

Proposition: If $L(\mathcal{A}) \subsetneq L(\mathcal{B})$, there exists a **small witness** $w \in L(\mathcal{B}) \setminus L(\mathcal{A})$ of size at most $|Q_{\mathcal{A}}| + |Q_{\mathcal{B}}|$

- $C(x) := \sum c_n x^n = B(x) - A(x)$ is **rational**
- The coefficients of $C(x)$ satisfy a linear recurrence:

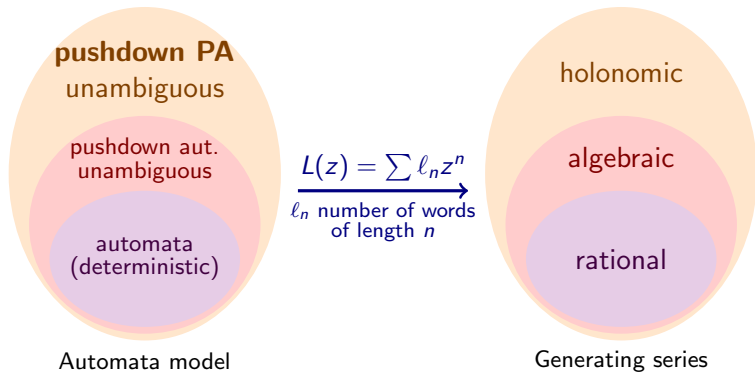
$$\forall n \geq r, c_n = \alpha_1 c_{n-1} + \dots + \alpha_r c_{n-r}$$

- the **order** r is at most $|Q_{\mathcal{A}}| + |Q_{\mathcal{B}}|$

Théorème [Stearns and Hunt 85] : The inclusion problem for unambiguous NFA is polynomial.

- $L(\mathcal{A}) \not\subseteq L(\mathcal{B}) \Leftrightarrow L(\mathcal{A}) \cap L(\mathcal{B}) \subsetneq L(\mathcal{A})$
- Compute coefficients up to $|Q_{\mathcal{A}}||Q_{\mathcal{B}}| + |Q_{\mathcal{A}}|$ (**dynamic prog.**)

Extension of hierarchy



- During my phd: find a suitable model **corresponding** to holonomic series, that is **relevant** from a modeling / automata point of view

Holonomic series [Stanley 80]

Rational: $P(x)f(x) = Q(x)$

Algebraic: $P(x, f(x)) = 0$.

Definition: A series $f(x) = \sum_n a_n x^n$ is **holonomic** (or **D-finite**) if it satisfies a linear differential equation:

$$P_k(x)f^{(k)}(x) + \dots + P_0(x)f(x) = 0 \quad \text{avec } P_i(x) \in \mathbb{Q}[x]$$

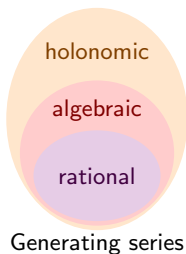
Alternative definition: the coefficients a_n satisfy a **linear recurrence** $p_r(n)a_{n+r} + \dots + p_0(n)a_n = 0$

Example: $F(x) = e^x := \sum \frac{x^n}{n!}$ is **holonomic** but is **not algebraic**

- differential equation: $F' - F = 0$
- recurrence relation: $(n+1)f_{n+1} - f_n = 0$

Generalizable in several variables

Closed by sum, product, composition with algebraic series, Hadamard product...



Generating series

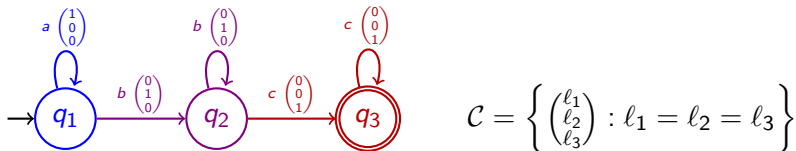
Parikh automata (PA) [Klaedtke, Rueß '03]

Motivation: $\{a^n b^n c^n\}$ is simple but not context-free

Idea: add **constrained** counters

- transitions labelled by vectors in \mathbb{N}^d
- test on the final value **at the end of the run**

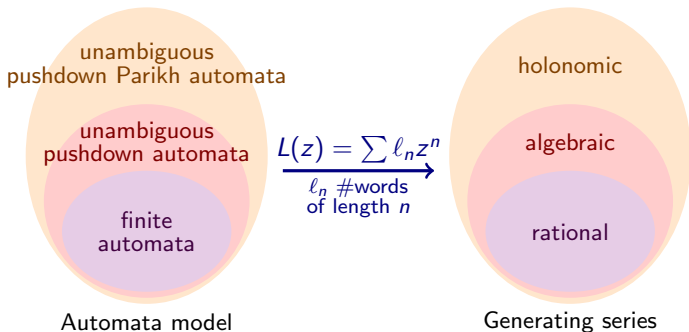
Parikh automaton: finite automaton on $\Sigma \times \mathbb{N}^d$ with semilinear constraints



$$q_1 \xrightarrow{\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} / a} q_1 \xrightarrow{\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} / a} q_1 \xrightarrow{\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} / b} q_2 \xrightarrow{\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} / b} q_2 \xrightarrow{\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} / b} q_3 \xrightarrow{\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} / c} q_3, w = aabbcc, v = \begin{pmatrix} 2 \\ 2 \\ 2 \end{pmatrix} \in \mathcal{C}$$

Can be extended with a stack.

Extension



Two remarkable applications :

- analytic proofs of **inherent ambiguity** for PA
- doubly exponential bound for an **algorithm** for the inclusion problem for uPA

Algorithmic application: inclusion problem for uPA



Pose $L_C := L_B \setminus L_A$

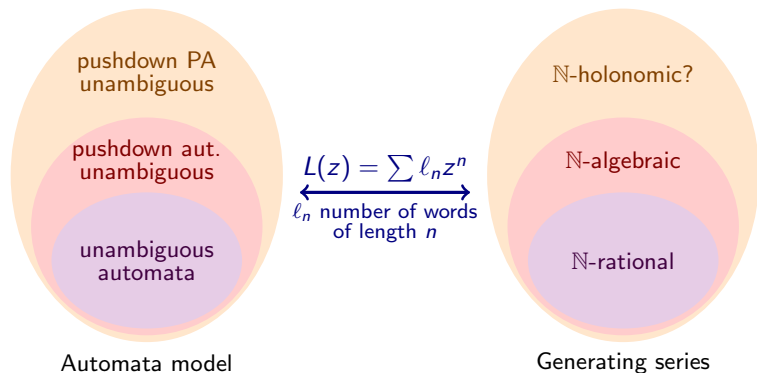
idea: replace $L_C \stackrel{?}{=} \emptyset$ by $C(x) \stackrel{?}{=} 0$

- "Compute" $A(x)$ and $B(x)$ from \mathcal{A} and \mathcal{B}
→ possible by **unambiguity**
- Differential equation satisfied by $C(x) = B(x) - A(x)$
- Linear recurrence satisfied by $c_n = b_n - a_n$

$$-p_r(n)c_{n+r} = p_{r-1}(n)c_{n+r-1} + \dots + p_0(n)c_n$$

- Bound B such that $c_n = 0$ for $n \leq B$ implies $\forall n, c_n = 0$
 $C(x) = x^{100}$ satisfies $x C'(x) - 100 C(x) = 0$ and $(n - 100)c_n = 0$

Garrabrant & Pak's conjecture



Subclass of holonomic series:

- natural from combinatorics
- rich enough to have nice closure and algorithmic properties

Garrabrant & Pak's conjecture

If $f(x, y) = \sum_{n,m} a_{n,m} x^n y^m$, the **diagonal** of f is defined by

$$\Delta f(x) = \sum_{n \in \mathbb{N}} a_{n,n} x^n.$$

$$C(x) = \frac{1 - \sqrt{1 - 4x}}{2x} = \Delta \frac{1 - x/y}{1 - x - y}$$

A series is the diagonal of a \mathbb{N} -rational series if and only if it the generating series of an unambiguous PA.

Conjecture [Garrabrant & Pak '14]:

The series of the Catalan numbers is not the **diagonal** of any \mathbb{N} -rational series.

Garrabrant & Pak's conjecture

If $f(x, y) = \sum_{n,m} a_{n,m} x^n y^m$, the **diagonal** of f is defined by

$$\Delta f(x) = \sum_{n \in \mathbb{N}} a_{n,n} x^n.$$

$$C(x) = \frac{1 - \sqrt{1 - 4x}}{2x} = \Delta \frac{1 - x/y}{1 - x - y}$$

A series is the diagonal of a \mathbb{N} -rational series if and only if it the generating series of an unambiguous PA.

Conjecture [Garrabrant & Pak '14]:

The series of the Catalan numbers is not the **diagonal** of any \mathbb{N} -rational series.

THANK YOU !

Sommaire annexes

Publications

Schéma projet recherche

Automates Parikh et séries holonomes

Séries rationnelles

Types d'ordre

Semilinéaire

Weakly unambiguous, autres modèles

Éléments absorbants

Focus regexp

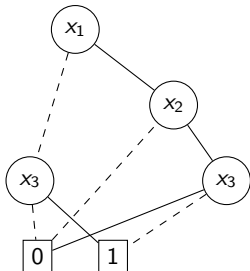
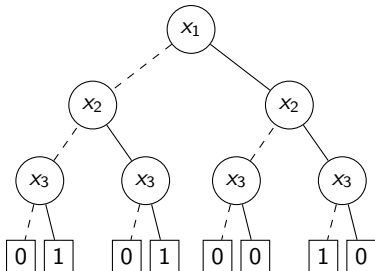
Distribution ABR

Illustration problème de l'inclusion

Pistes de recherche sur les fonctions booléennes

Structures compactes: circuits et BDD

$$f = (x_1 \wedge x_2 \wedge \neg x_3) \vee (\neg x_1 \wedge x_3)$$



Pistes de recherche:

- les raisons qui font que l'élément absorbant / neutre n'agit plus aussi fortement que dans les structures d'arbres ?
- généraliser ces structures compactes à d'autres classes d'expressions?

Conférences

- Arnaud Carayol, Philippe Duchon, Florent Koechlin and Cyril Nicaud. One Drop of Non-Determinism in a Random Deterministic Automaton. STACS'23.
- Florent Koechlin. New analytic techniques for proving the inherent ambiguity of context-free languages. FSTTCS'22
- Florent Koechlin and Pablo Rotondo. Analysis of an efficient reduction algorithm for random regular expressions based on universality detection. CSR'21
- Florent Koechlin and Pablo Rotondo. Absorbing patterns in BST-like expression-trees. STACS'21
- Alin Bostan, Arnaud Carayol, Florent Koechlin and Cyril Nicaud. Weakly-unambiguous Parikh automata and their link to holonomic series. ICALP'20
- Florent Koechlin, Cyril Nicaud and Pablo Rotondo. On the Degeneracy of Random Expressions Specified by Systems of Combinatorial Equations. DLT'20
- Florent Koechlin, Cyril Nicaud and Pablo Rotondo. Uniform Random Expressions Lack Expressivity. MFCS'19

Versions longues

- Florent Koechlin, Cyril Nicaud and Pablo Rotondo. Simplifications of Uniform Expressions Specified by Systems. IJFCS'21, Volume No. 32, Issue No. 06, pp. 733 - 760.

Soumissions en cours

- Florent Koechlin et Pablo Rotondo. Analysis of an efficient reduction algorithm for random regular expressions based on universality detection. Soumis à Theory of Computing Systems (TCS), en tant que special issue de CSR'21.
- Florent Koechlin et Pablo Rotondo, The effects of semantic simplifications on random BST-like expression-trees, Soumis à Discrete Mathematics

Projet de recherche : Interaction entre combinatoire et théorie des automates

retour

Langages formels

Questions pertinentes

Problèmes combinatoires associés

algébriques bornés

décidabilité ambiguïté

stratification des semilinéaires

forme irréductible de séries rationnelles

algébriques

finie ambiguïté

comportement asymptotique
de séries \mathbb{N} -algébriques

de Parikh (à pile)

critères ambiguïté

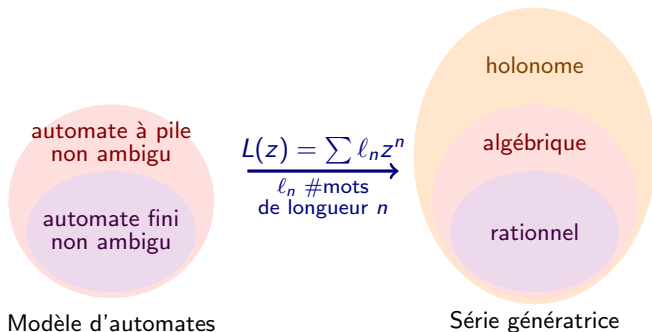
diagonales de séries \mathbb{N} -rationnelles

complexité inclusion,
universalité

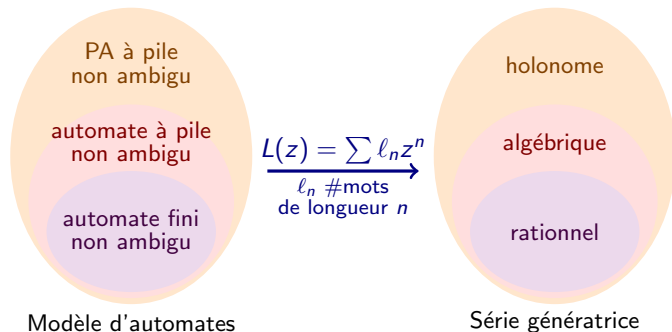
algorithmes de calcul de diagonales,
télescopage créatif

VASS, registres,
poids, ...

suites récurrentes,
méthode du noyau



- Dans ma thèse: trouver un modèle **correspondant** aux séries holonomes, **pertinent** en modélisation et théorie des automates



- Dans ma thèse: trouver un modèle **correspondant** aux séries holonomes, **pertinent** en modélisation et théorie des automates

Rationnelle: $P(x)f(x) = Q(x)$

Algébrique: $P(x, f(x)) = 0$.

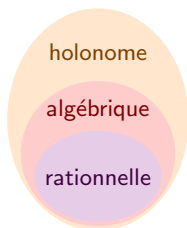
Holonyme: Équation différentielle linéaire

$$P_k(x)f^{(k)}(x) + \dots + P_0(x)f(x) = 0 \text{ avec } P_i(x) \in \mathbb{Q}[x]$$

Définition alternative: les coefficients a_n satisfont une **réurrence linéaire** $p_r(n)a_{n+r} + \dots + p_0(n)a_n = 0$

Exemple: $F(x) = e^x$ holonyme mais pas algébrique

- équation différentielle: $F' - F = 0$
- récurrence: $(n+1)f_{n+1} - f_n = 0$



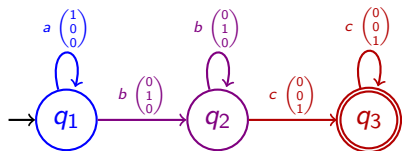
Série génératrice

Généralisable à plusieurs variables: systèmes aux dérivées partielles

Automates de Parikh (PA) [Klaedtke, Rueß '03]

[retour](#)

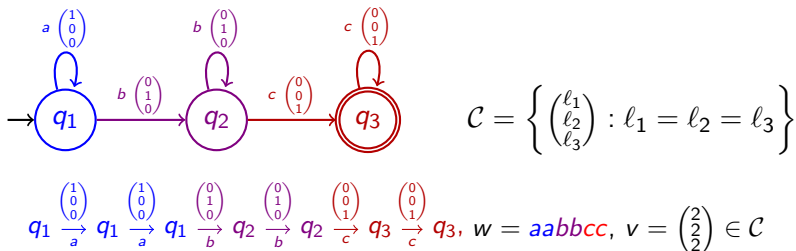
Automate de Parikh: automate fini sur $\Sigma \times \mathbb{N}^d$ avec des contraintes semilinéaires sur des vecteurs d'entiers



$$\mathcal{C} = \left\{ \begin{pmatrix} \ell_1 \\ \ell_2 \\ \ell_3 \end{pmatrix} : \ell_1 = \ell_2 = \ell_3 \right\}$$

$$q_1 \xrightarrow{a \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}} q_1 \xrightarrow{a \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}} q_1 \xrightarrow{b \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}} q_2 \xrightarrow{b \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}} q_2 \xrightarrow{c \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}} q_3 \xrightarrow{c \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}} q_3, w = aabbcc, v = \begin{pmatrix} 2 \\ 2 \\ 2 \end{pmatrix} \in \mathcal{C}$$

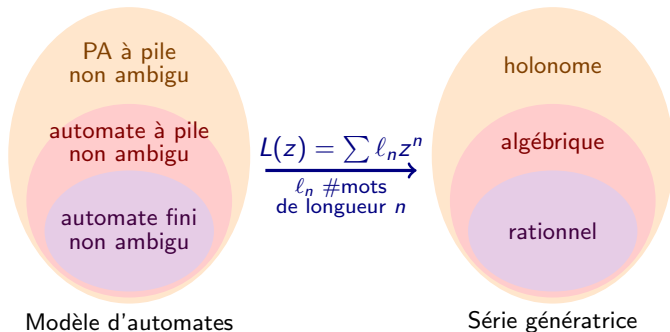
Automate de Parikh: automate fini sur $\Sigma \times \mathbb{N}^d$ avec des contraintes semilinéaires sur des vecteurs d'entiers



La série d'un langage de Parikh (à pile) non ambigu est **holonome**.

Systèmes associés holonomes:

$$\begin{cases} (a - 27 a^2 bc) \frac{\partial^2}{\partial a^2} F + (1 - 54 abc) \frac{\partial}{\partial a} F + 6 bc F = 0 \\ (b - 27 b^2 ca) \frac{\partial^2}{\partial b^2} F + (1 - 54 bca) \frac{\partial}{\partial b} F + 6 ca F = 0 \\ (c - 27 c^2 ab) \frac{\partial^2}{\partial c^2} F + (1 - 54 cab) \frac{\partial}{\partial c} F + 6 ab F = 0 \end{cases}$$

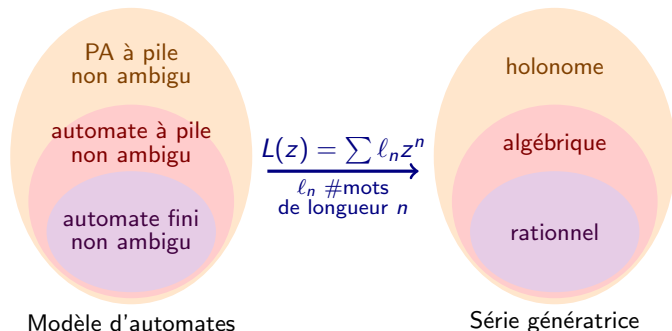


Contributions [ICALP '20]:

- existence de langages **intrinsèquement ambigus**, même avec pile et compteurs + **techniques** pour en trouver
- l'**inclusion** des PA non ambigus est décidable en 2-EXPTIME
→ analyse fine des tailles des séries par des outils de **calcul formel**

Focus: séries provenant de la théorie des langages

retour



Systèmes associés holonomes:

$$\begin{cases} (a - 27 a^2 bc) \frac{\partial^2}{\partial a^2} F + (1 - 54 abc) \frac{\partial}{\partial a} F + 6 bc F = 0 \\ (b - 27 b^2 ca) \frac{\partial^2}{\partial b^2} F + (1 - 54 bca) \frac{\partial}{\partial b} F + 6 ca F = 0 \\ (c - 27 c^2 ab) \frac{\partial^2}{\partial c^2} F + (1 - 54 cab) \frac{\partial}{\partial c} F + 6 ab F = 0 \end{cases}$$

Problème: Certains langages ambigus ont une série rationnelle.

→ $a^n b^m c^p$ avec $n = m$ ou $m = p$ et $(ab)^* c + a^* (bc)^*$

Séries génératrices associées sont *rationnelles*.

Exemple: $L = \{a^n b^m c^p \text{ avec } n \neq m \text{ ou } m \neq p\}$:

$$S(a, b, c) = \frac{a+b+c-ab-ac-bc}{(1-a)(1-b)(1-c)(1-abc)} \quad [\text{Makarov 21, Koechlin 22}]$$

Bornés: $L \subseteq w_1^* \dots w_d^*$

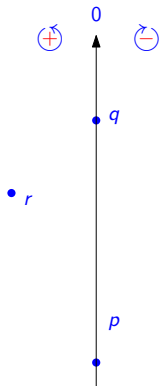
Conjecture [Ginsburg & Ullian '66] L'intrinsèque ambiguïté des langages algébriques bornés est décidable.

L'approche par les séries ouvre de nouvelles approches, qui s'avèrent plus efficaces que la plupart des critères existants.

Types d'ordre

retour

Orientation, type d'ordre

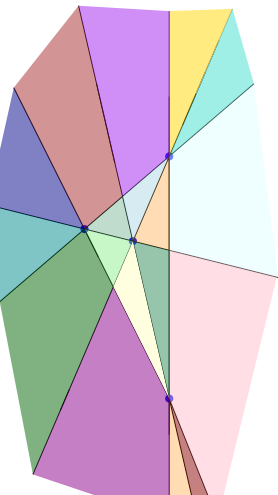


Orientation de trois points

$$\circ \chi(p, q, r) = \begin{cases} +1 & \text{si } r \text{ est à gauche de } (pq) \\ 0 & \text{si } r \text{ est sur } (pq) \\ -1 & \text{si } r \text{ est à droite de } (pq) \end{cases}$$

$$\circ \chi(p, q, r) = \text{sign} \begin{vmatrix} x_p & x_q & x_r \\ y_p & y_q & y_r \\ 1 & 1 & 1 \end{vmatrix}$$

Orientation, type d'ordre



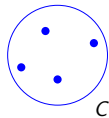
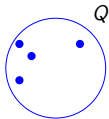
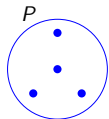
Orientation de trois points

$$\circ \chi(p, q, r) = \begin{cases} +1 & \text{si } r \text{ est à gauche de } (pq) \\ 0 & \text{si } r \text{ est sur } (pq) \\ -1 & \text{si } r \text{ est à droite de } (pq) \end{cases}$$

$$\circ \chi(p, q, r) = \text{sign} \begin{vmatrix} x_p & x_q & x_r \\ y_p & y_q & y_r \\ 1 & 1 & 1 \end{vmatrix}$$

$$\circ f(z) = (\chi(p, q, z))_{p, q \in P} \text{ constant dans les cellules}$$

Orientation, type d'ordre

**Type d'ordre**

- Deux ensembles P et Q de points de \mathbb{R}^2 *ont le même type d'ordre* s'il y a une bijection $f : P \rightarrow Q$ qui préserve les orientations :

$$\forall p, q, r \in P, \quad \chi(p, q, r) = \chi(f(p), f(q), f(r)).$$

- Type d'ordre = *classe d'équivalence* pour cette relation

Motivation

Dans la suite, les types d'ordre considérés sont *simples* : $\chi(p, q, r) = \pm 1$.

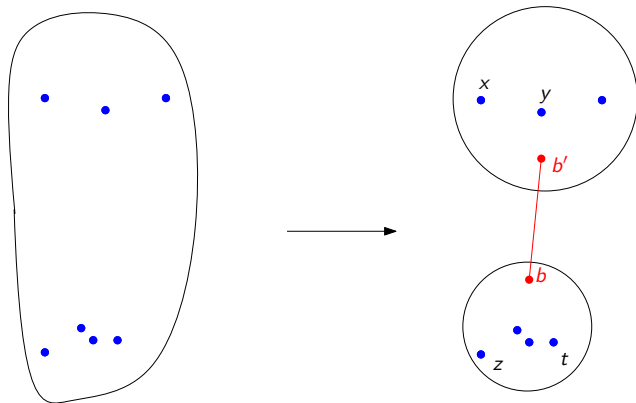
Motivation :

- Algorithmes exacts :
 - $[p, q] \in \text{Conv}(P) \Leftrightarrow \forall x \in P \setminus \{p, q\}, \chi(p, q, x) = \text{cst}$
 - $x \in \Delta(p, q, r) \Leftrightarrow \chi(p, q, x) = \chi(q, r, x) = \chi(r, p, x)$
 - (p, q) coupe $[a, b] \Leftrightarrow \chi(p, q, a) = -\chi(q, r, b)$
 - $[p, q]$ et $[a, b]$ se coupent $\Leftrightarrow \dots$
- Benchmarks

Difficultés

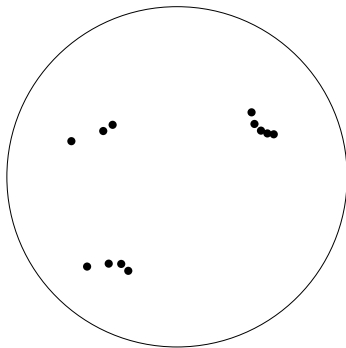
- Difficiles à décrire exhaustivement ($n = 11$ [Aichholzer et al.])
- Difficiles à compter $t_n = n^{3n+o(n)}$ [Alon, Warren]

Split selon un module

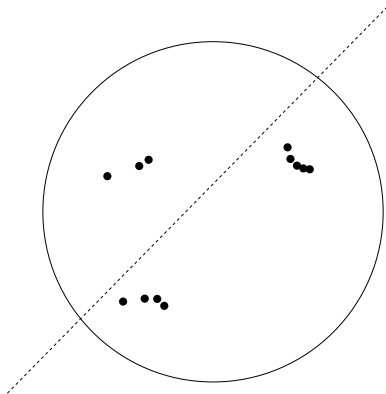


Règle : $\chi(x, y, z) = \chi(x, y, b')$ et $\chi(x, z, t) = \chi(b, z, t)$

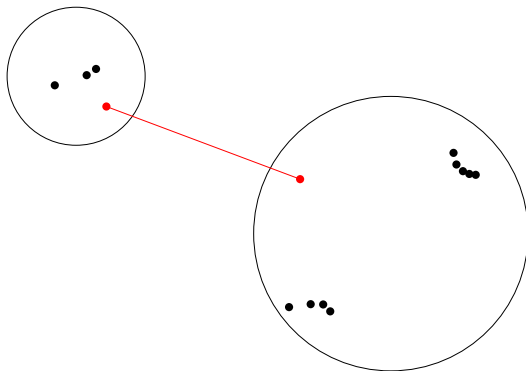
Décomposition modulaire



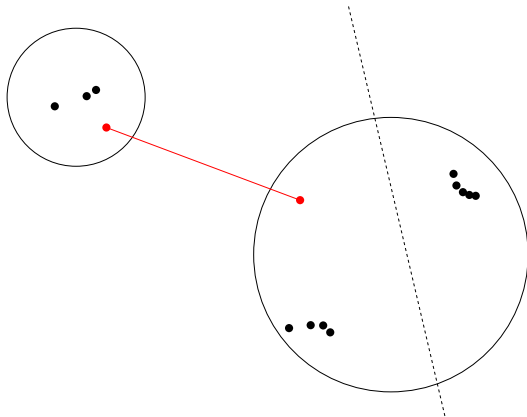
Décomposition modulaire



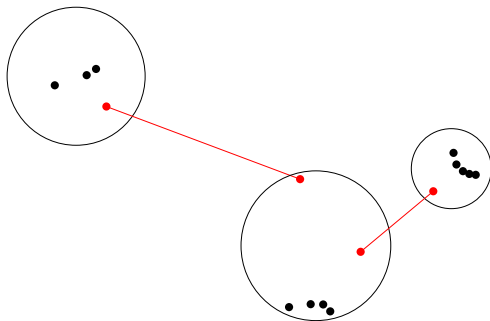
Décomposition modulaire



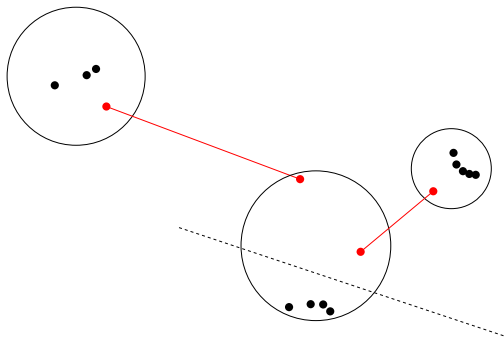
Décomposition modulaire



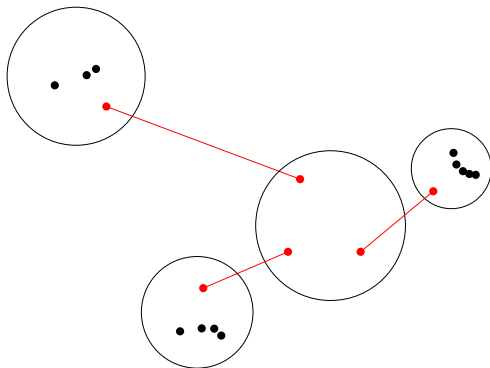
Décomposition modulaire



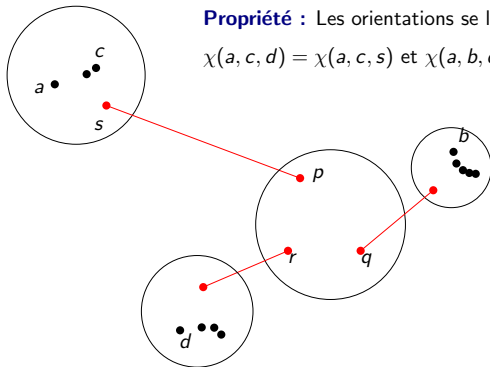
Décomposition modulaire



Décomposition modulaire



Décomposition modulaire



Propriété : Les orientations se lisent en suivant les proxies !

$$\chi(a, c, d) = \chi(a, c, s) \text{ et } \chi(a, b, d) = \chi(p, q, r)$$

Décomposition modulaire

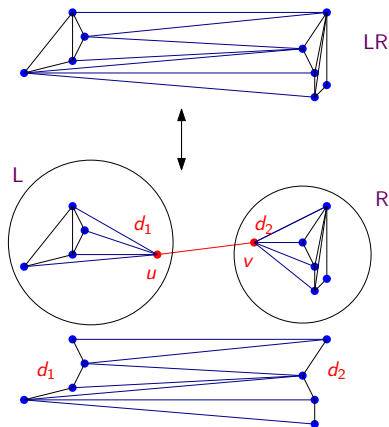
Décomposition modulaire

- On part d'un noeud seul contenant un type d'ordre représenté par un ensemble de points
- On applique $\Rightarrow := \rightarrow_{split \text{ non convexe}} \cup \rightarrow_{fusion \text{ convexe}}$

Résultat

\Rightarrow est *confluente* : à isomorphisme de graphe prêt, un type d'ordre a une *unique décomposition modulaire normale*

Calcul de triangulations



Notations

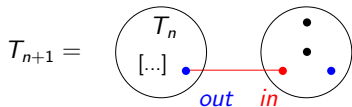
- $\tau(LR)$: nombre de triangulations de LR .
- $\tau_{d_1}(L)$: nombre de triangulations de L telles que u est d'arité d_1 .
- $\tau_{d_2}(R)$: nombre de triangulations de R telles que v est d'arité d_2 .
- $\binom{n}{k} := \binom{k+n-1}{k}$

Bijection :

$$\tau(LR) = \sum_{d_1, d_2} \binom{d_2}{d_1 - 1} \tau_{d_1}(L) \tau_{d_2}(R)$$

Un exemple analysable : une chaîne

Relation de récurrence



Polynôme des triangulations

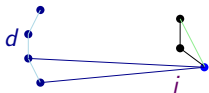
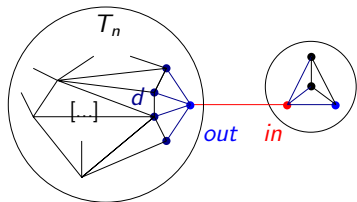
- $F_n(u) = \sum_{\Delta \in \mathcal{T}_n} u^{\deg_{\Delta}(out)}$
- $F_n(1) = \tau(T_n)$
- $F_1(u) = u^3$

Bijection :

$$\tau(T_{n+1}) = \sum_{d_1, d_2} \binom{d_2}{d_1 - 1} \tau_{d_1}(T_n) \tau_{d_2}(T_1)$$

Un exemple analysable : une chaîne

Relation de récurrence



Polynôme des triangulations

- $F_n(u) = \sum_{\Delta \in \mathcal{T}_n} u^{\deg_{\Delta}(out)}$
- $F_n(1) = \tau(\mathcal{T}_n)$
- $F_1(u) = u^3$

Récurrence

$$F_{n+1}(u) = u^3 \sum_{d \geq 0} [u^d] F_n(u) \sum_{i=0}^{d-1} u^i (d-i)$$

Méthode du noyau

Réurrence :
$$F_{n+1}(u) = u^3 \sum_{d \geq 0} [u^d] F_n(u) \sum_{i=0}^{d-1} u^i (d-i)$$

Série multivariée

- $F(z, u) = \sum_{n \geq 0} F_n(u) z^n$
- $F(z, 1) = \sum_{n \geq 0} \tau(T_n) z^n$

Relation vérifiée par F :

$$F(z, u) \left(1 - \frac{u^4 z}{(u-1)^2} \right) = u^3 z \left(1 - \frac{\partial_u F(z, 1)}{u-1} - \frac{F(z, 1)}{(u-1)^2} \right)$$

Méthode du noyau

- Le terme de gauche est *linéaire* en $\partial_u F(z, 1), F(z, 1)$
- Chaque racine $u(z)$ de $1 - \frac{u^4 z}{(u-1)^2}$ fournit une équation linéaire. Si on en trouve assez on peut résoudre et trouver $F(z, 1)$

Parikh automata (PA) [Klaedtke, Rueß '03]

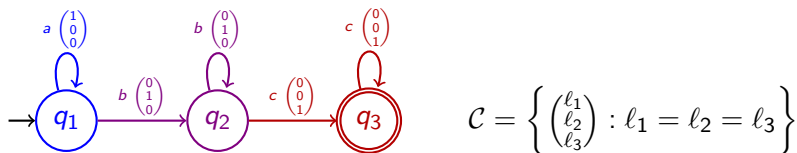
retour

Motivation: $\{a^n b^n c^n\}$ is simple but not context-free

Idea: add **constrained** counters

- transitions labelled by vectors in \mathbb{N}^d
- test on the final value **at the end of the run**

Parikh automaton: finite automaton on $\Sigma \times \mathbb{N}^d$ with semilinear constraints



$$q_1 \xrightarrow{\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} / a} q_1 \xrightarrow{\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} / a} q_1 \xrightarrow{\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} / b} q_2 \xrightarrow{\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} / b} q_2 \xrightarrow{\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} / b} q_3 \xrightarrow{\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} / c} q_3, w = aabbcc, v = \begin{pmatrix} 2 \\ 2 \\ 2 \end{pmatrix} \in C$$

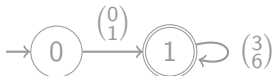
Can be extended with a stack.

$\bigwedge \bigvee$ of linear inequalities or equalities modulo constants

$$\{(3n, 6n + 1) : n \in \mathbb{N}\} = \{(x_1, x_2) : x_1 \equiv 0[3] \wedge x_2 = 2x_1 + 1\}$$

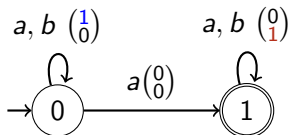
Equivalent definitions

- Finite union of linear sets $\vec{c} + P^*$ where $P = \{p_1, \dots, p_r\}$
 $(0, 1) + \{(3, 6)\}^*$
- Presburger arithmetic
 $\Phi(x_1, x_2) := \exists x, x_1 - 3x = 0 \wedge 1 + 2x_1 - x_2 = 0$
- (Unambiguous) rationnel subsets of $(\mathbb{N}^d, +)$



Universality: undecidable for non deterministic Parikh automata

Weak unambiguity: at most one accepting run (final state + semilinear constraint)



$$C = \{(n, n) : n \in \mathbb{N}\}$$

$$L(\mathcal{A}) = \{\text{words with an } a \text{ in the middle}\} = \{\dots, abba**a**bab, \dots\}$$

\neq unambiguous Parikh automata [Cadilhac, Finkel, McKenzie 13]

Theorem [Bostan, Carayol, K., Nicaud '20] : The class of weakly unambiguous Parikh languages coincide with :

- RCM of [Castiglione, Massazza '17]
- **unambiguous two-way RBCM** [Ibarra '78]
⇒ stronger version of [Castiglione, Massazza '17]'s conjecture

Théorème [ICALP '20] : La classe des automates de Parikh à pile non ambigus of weakly unambiguous pushdown Parikh languages coincide with :

- LCL adapted from [Massazza '93]
- **unambiguous one-way RBCM** with a stack [Ibarra '78]

Description of expression trees

retour

1 equation $E = a + b + \varepsilon + \overset{\star}{\underset{|}{E}} + \overset{+}{\underset{\wedge}{E E}} + \overset{\bullet}{\underset{\wedge}{E E}}$

Systems of equations
$$\begin{cases} \mathcal{L}_R = \overset{\star}{\underset{|}{S}} + \mathcal{S}, \\ \mathcal{S} = a + b + \varepsilon + \overset{+}{\underset{\wedge}{\mathcal{L}_R \mathcal{L}_R}} + \overset{\bullet}{\underset{\wedge}{\mathcal{L}_R \mathcal{L}_R}}. \end{cases}$$

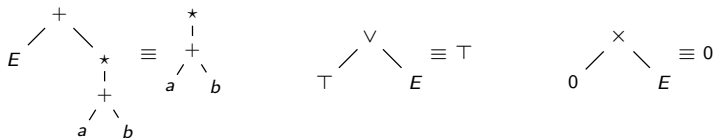
useful for

- average case analysis (uniform distribution)
- counting of syntactic trees, bounds on represented objects [Lee, Shallit '05]



Semantic rule: absorbing pattern

retour



Absorbing operator: \otimes operator of arity $a \geq 2$

Absorbing pattern: constant tree \mathcal{P} of size p

Semantic simplification: $c_1 \overset{\otimes}{\dots} c_a \equiv \mathcal{P}$, whenever $C_i = \mathcal{P}$ of one $i \in \{1, \dots, a\}$

Example

- $\mathcal{P} = (a + b)^*$ is absorbing for $+$
- 0 is absorbing for \times
- \top is absorbing for \vee , \perp is absorbing for \wedge

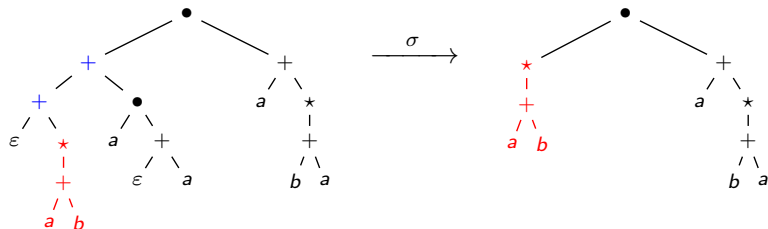
Simplification with absorbing pattern

retour

Bottom-up simplification:

$$\begin{array}{c} \circledast \\ / \quad \backslash \\ C_1 \cdots C_a \end{array} \rightsquigarrow \mathcal{P}, \text{ if } C_i = \mathcal{P} \text{ for some } i \in \{1, \dots, a\}.$$

Example ($\mathcal{P} =$  and $\circledast = +$)



Question : What is the average size of $\sigma(T)$?

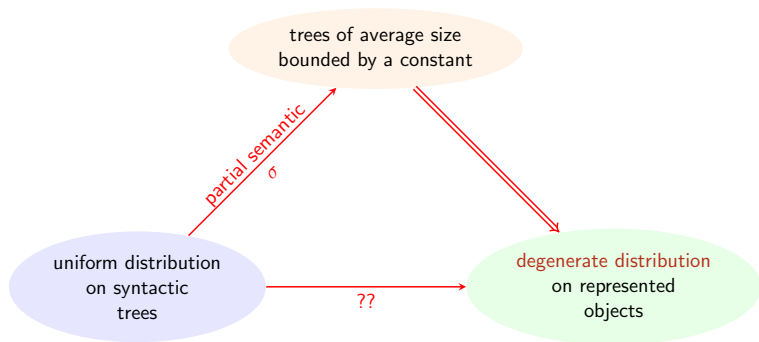
Uniform distribution is often degenerate

retour

Theorem : (K., Nicaud, Rotondo 20)

Let \mathcal{F} a set of expressions described by a well founded system, with an absorbing pattern.

Then, **uniform** expressions have, after reduction, an average size **bounded by a constant C** independent of n .



- Systèmes naturels pour décrire des expressions

$$\left\{ \begin{array}{l} \mathcal{B} = \overline{\mathcal{I}}_S + \mathcal{S}, \\ \mathcal{S} = \perp + \top + x + \bigwedge_B \bigwedge_B + \bigvee_B \bigvee_B \end{array} \right. \quad \left\{ \begin{array}{l} \mathcal{L} = \mathcal{L}_+ + \mathcal{L}_\bullet + \mathcal{L}_2, \\ \mathcal{L}_+ = \mathcal{L} \begin{array}{c} \diagup \quad \diagdown \\ \quad \mathcal{L}_\bullet + \mathcal{L}_2 \end{array}, \\ \mathcal{L}_\bullet = \mathcal{L} \begin{array}{c} \diagup \quad \diagdown \\ \quad \mathcal{L}_+ + \mathcal{L}_2 \end{array}, \\ \mathcal{L}_2 = a + b + \varepsilon + \begin{array}{c} \star \\ \downarrow \\ \mathcal{L} \end{array} \end{array} \right. \quad [\text{Lee, Shallit '05}]$$

- Éléments absorbants difficiles à éviter : $((a + b)^*)^*$, $(a^*b^*)^*$...

$$\left[\begin{array}{l} \text{Broda, Machiavelo,} \\ \text{Moreira, Reis'21} \end{array} \right] \left\{ \begin{array}{l} \alpha = \varepsilon + a + b + \begin{array}{c} \bullet \\ \alpha \quad \alpha \end{array} + \begin{array}{c} \star \\ \alpha \end{array} + \begin{array}{c} \dagger \\ \alpha_P \quad \alpha_P \end{array} \\ \alpha_P = \varepsilon + a + b + \begin{array}{c} \bullet \\ \alpha \quad \alpha \end{array} + \begin{array}{c} \star \\ \alpha_\Sigma \end{array} + \begin{array}{c} \dagger \\ \alpha_P \quad \alpha_P \end{array} \\ \alpha_\Sigma = \varepsilon + a + b + \begin{array}{c} \bullet \\ \alpha \quad \alpha \end{array} + \begin{array}{c} \star \\ \alpha \end{array} + \gamma \\ \gamma = \begin{array}{c} \dagger \\ \alpha_{ab} \quad \alpha_{ab} \end{array} + \begin{array}{c} \dagger \\ \alpha_{ab} \quad a \end{array} + \begin{array}{c} \dagger \\ \alpha_{ab} \quad b \end{array} + \begin{array}{c} \dagger \\ a \quad \alpha_{ab} \end{array} + \begin{array}{c} \dagger \\ b \quad \alpha_{ab} \end{array} + \begin{array}{c} \dagger \\ a \quad a \end{array} + \begin{array}{c} \dagger \\ b \quad b \end{array} \\ \alpha_{ab} = \varepsilon + \begin{array}{c} \bullet \\ \alpha \quad \alpha \end{array} + \begin{array}{c} \star \\ \alpha_\Sigma \end{array} + \begin{array}{c} \dagger \\ \alpha_P \quad \alpha_P \end{array} \end{array} \right.$$

→ évite $(a + b)^*$ mais pas $((a + b)^*)^*$!

Natural conditions ?

retour

- Natural conditions ensuring the system is describing expression trees correctly
- The system must not prevent the simplification from happening
- Some technical restrictions to keep the system manageable without having to know anything else about the semantics of the objects represented by the trees

Natural conditions ?

retour

- (H_1) The graph of unit rules is acyclic ($\mathcal{L}_R = \mathcal{T} + \dots$)
- (H_2) The system is *non-ambiguous*: each complete expression can be built in at most one way.
- (H_3) The system is *aperiodic*.
- (H_4) The dependency graph of the system is strongly connected.
- (H_5) [Reduction is **non-trivial**] There is a rule of the form

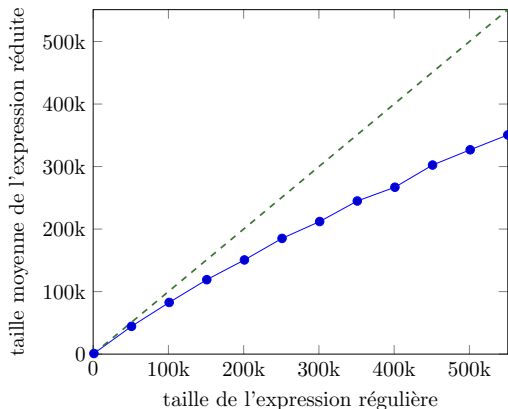
$$\begin{array}{c} \textcircled{*} \\ / \quad \backslash \\ T_1 \cdots T_a \end{array} \quad \text{with } \mathcal{P} \in \sigma(\mathcal{L}(T_i)) \text{ and } a(T_j) \geq 1, i \neq j.$$

- (H_6) The system is *not linear*: there is a rule of arity at least 2.

Theorem (K., Nicaud, Rotondo, 2020)

Under the hypotheses H_ , all moments associated to the size after reduction of a uniform tree of size n are bounded by a constant.*

$$\mathcal{L}_R = a + b + \varepsilon + \overset{*}{\mathcal{L}_R} + \overset{\bullet}{\mathcal{L}_R \mathcal{L}_R} + \overset{+}{\mathcal{L}_R \mathcal{L}_R} \quad \text{avec } \mathcal{P} = \begin{array}{c} \overset{*}{+} \\ \swarrow \quad \searrow \\ a \quad b \end{array}$$



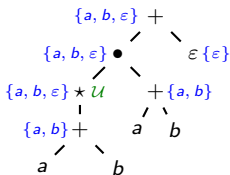
Proposition (K, Nicaud, Rotondo '19):

La taille réduite moyenne tend vers **3 624 217**.

Focus sur $\mathcal{L}_R = a + b + \varepsilon + \overset{*}{\uparrow} \mathcal{L}_R + \overset{\bullet}{\wedge} \mathcal{L}_R \mathcal{L}_R + \overset{+}{\wedge} \mathcal{L}_R \mathcal{L}_R$

retour

- $\mathcal{P} = \overset{*}{\uparrow} \begin{matrix} a \\ + \\ b \end{matrix}$, mais marche aussi avec tout $\overset{*}{\uparrow} \begin{matrix} L_a \\ + \\ L_b \end{matrix}$ ou $\overset{*}{\uparrow} L_{a,b}$
- Savoir si E reconnaît a ou b est facile.
- On a aussi besoin d' ε pour la concaténation



$$(a + b)^* \cdot (a + b) + \varepsilon$$

Règles:

$$\overset{*}{\uparrow} E \equiv \mathcal{U} \text{ si } \{a, b\} \subseteq \mathcal{L}(E)$$

$$\overset{+}{\wedge} \begin{matrix} U \\ E \end{matrix} \equiv \mathcal{U}$$

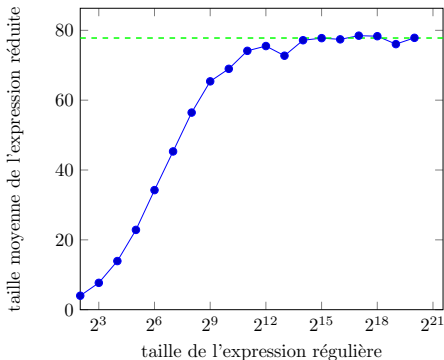
$$\overset{\bullet}{\wedge} \begin{matrix} U \\ E \end{matrix} \equiv \mathcal{U} \text{ si } \varepsilon \in \mathcal{L}(E)$$

Focus sur $\mathcal{L}_R = a + b + \varepsilon + \overset{*}{\mathcal{L}_R} + \overset{\bullet}{\mathcal{L}_R} \overset{\wedge}{\mathcal{L}_R} + \overset{+}{\mathcal{L}_R} \overset{\wedge}{\mathcal{L}_R}$

retour

- $\mathcal{P} = \overset{*}{+} \begin{matrix} a \\ b \end{matrix}$, mais marche aussi avec tout $\overset{*}{+} \begin{matrix} L_a \\ L_b \end{matrix}$ ou $\overset{*}{L}_{a,b}$
- Savoir si E reconnaît a ou b est facile.
- On a aussi besoin d' ε pour la concaténation

Théorème (K., Rotondo '21):
La taille moyenne tend vers une constante $C \approx 77,8$.



$$\mathcal{L}_R = a + b + \varepsilon + \overset{*}{\mathcal{L}_R} + \overset{\bullet}{\mathcal{L}_R \wedge \mathcal{L}_R} + \overset{+}{\mathcal{L}_R \wedge \mathcal{L}_R}$$

L'analyse introduit deux classes intéressantes :

\mathcal{U} : les entièrement réductibles et $\mathcal{T}_{\Sigma, \varepsilon}$: reconnaissent Σ et ε .

$$\mathcal{U} \subsetneq \text{expressions universelles} \subsetneq \mathcal{T}_{\Sigma, \varepsilon}$$

Théorème (K., Rotondo '21):

La proportion d'expressions universelles sur un alphabet à deux lettres est asymptotiquement comprise entre **31%** et **46%**.

```
function RandomFormula( $n$ ):  
  if  $n = 1$  then  
    |  $p :=$  random symbol in  $AP \cup \{\top, \perp\}$ ;  
    | return  $p$ ;  
  else if  $n = 2$  then  
    |  $op :=$  random operator in  $\{\neg, \mathbf{X}, \square, \diamond\}$ ;  
    |  $f :=$  RandomFormula(1);  
    | return  $op f$ ;  
  else  
    |  $op :=$  random operator in  $\{\neg, \mathbf{X}, \square, \diamond, \wedge, \vee, \rightarrow, \leftrightarrow, \mathbf{U}, \mathbf{R}\}$ ;  
    | if  $op$  in  $\{\neg, \mathbf{X}, \square, \diamond\}$  then  
    |   |  $f :=$  RandomFormula( $n - 1$ );  
    |   | return  $op f$ ;  
    | else  
    |   |  $x :=$  uniform integer in  $[1, n - 2]$ ;  
    |   |  $f_1 :=$  RandomFormula( $x$ );  
    |   |  $f_2 :=$  RandomFormula( $n - x - 1$ );  
    |   | return  $(f_1 op f_2)$ ;
```

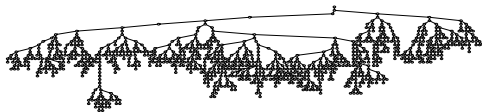
Algorithme 1: Pseudo-code utilisé 1btt [Tauriainen '00]

Distribution ABR: différente de l'uniforme



uniforme
 $h \sim c\sqrt{n}$

[Flajolet, Odlyzko '82]



ABR

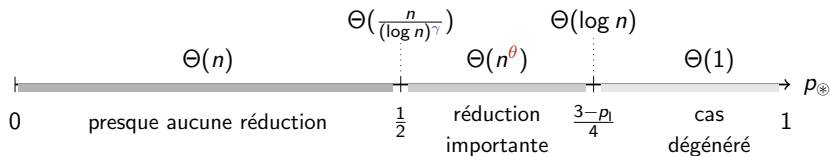
$h \sim d \log n$

[Devroye '86]

La distribution ABR est plus complexe

[retour](#)

Théorème [K. Rotondo '21]: cinq régimes pour la taille moyenne réduite, dépendant de la probabilité p_{\otimes} de l'opérateur absorbant



avec $\gamma = \frac{2}{1-p_1}$ et $\theta = 1 - \frac{4p_{\otimes}-2}{1-p_1}$.

- **Réurrences** sur e_n l'espérance de la taille réduite et sur γ_n la probabilité des entièrement réductibles
- **Équations différentielles**
sur les séries génératrices $E(z)$ et $A(z) = \sum_n \gamma_n z^n$
 $E(z)$: **équation linéaire** non homogène du premier ordre dépendant de $A(z)$
 $A(z)$: **équation de Riccati** $A'(z) = a(z) + b(z)A(z) - p_{\otimes} A(z)^2$
- Principe de la combinatoire analytique : étude de A et E vues comme des fonctions **analytiques** $\mathbb{C} \rightarrow \mathbb{C}$, autour de leur singularité $z = 1$
- Différentes techniques : intégrations singulières pour $E(z)$, linéarisation et méthode de Frobenius pour $A(z)$

Inclusion problem for unambiguous automata

Problem : Given \mathcal{A} and \mathcal{B} two unambiguous NFA, $L(\mathcal{A}) \subseteq L(\mathcal{B})$?

Proposition: If $L(\mathcal{A}) \subsetneq L(\mathcal{B})$, there exists a **small witness** $w \in L(\mathcal{B}) \setminus L(\mathcal{A})$ of size at most $|Q_{\mathcal{A}}| + |Q_{\mathcal{B}}|$

- $C(x) := \sum c_n x^n = B(x) - A(x)$ is **rational**
- The coefficients of $C(x)$ satisfy a linear recurrence:

$$\forall n \geq r, c_n = \alpha_1 c_{n-1} + \dots + \alpha_r c_{n-r}$$

- the **order** r is at most $|Q_{\mathcal{A}}| + |Q_{\mathcal{B}}|$



Problem : Given \mathcal{A} and \mathcal{B} two unambiguous NFA, $L(\mathcal{A}) \subseteq L(\mathcal{B})$?

Proposition: If $L(\mathcal{A}) \subsetneq L(\mathcal{B})$, there exists a **small witness** $w \in L(\mathcal{B}) \setminus L(\mathcal{A})$ of size at most $|Q_{\mathcal{A}}| + |Q_{\mathcal{B}}|$

- $C(x) := \sum c_n x^n = B(x) - A(x)$ is **rational**
- The coefficients of $C(x)$ satisfy a linear recurrence:

$$\forall n \geq r, c_n = \alpha_1 c_{n-1} + \dots + \alpha_r c_{n-r}$$

- the **order** r is at most $|Q_{\mathcal{A}}| + |Q_{\mathcal{B}}|$

Théorème [Stearns and Hunt 85] : The inclusion problem for unambiguous NFA is polynomial.

- $L(\mathcal{A}) \not\subseteq L(\mathcal{B}) \Leftrightarrow L(\mathcal{A}) \cap L(\mathcal{B}) \subsetneq L(\mathcal{A})$
- Compute coefficients up to $|Q_{\mathcal{A}}||Q_{\mathcal{B}}| + |Q_{\mathcal{A}}|$ (**dynamic prog.**)

Algorithmic application: inclusion problem for wuPA

retour



Pose $L_C := L_B \setminus L_A$

idea: replace $L_C \stackrel{?}{=} \emptyset$ by $C(x) \stackrel{?}{=} 0$

- "Compute" $A(x)$ and $B(x)$ from \mathcal{A} and \mathcal{B}
→ possible by **weak-unambiguity**
- Differential equation satisfied by $C(x) = B(x) - A(x)$
- Linear recurrence satisfied by $c_n = b_n - a_n$

$$-p_r(n)c_{n+r} = p_{r-1}(n)c_{n+r-1} + \dots + p_0(n)c_n$$

- Bound B such that $c_n = 0$ for $n \leq B$ implies $\forall n, c_n = 0$
 $C(x) = x^{100}$ satisfies $xC'(x) - 100C(x) = 0$ and $(n - 100)c_n = 0$

\mathcal{A} weakly-unambiguous PA with a semilinear constraint C of dimension d (of the vectors)

$$A(x) = \underbrace{\bar{A}(x, y_1, \dots, y_d)}_{\substack{\text{rat. series} \\ \text{obtained from } \mathcal{A}}} \odot \underbrace{C(x, y_1, \dots, y_d)}_{\substack{\text{rat. series} \\ \text{obtained from } C}} \quad [y_1 \rightarrow 1, \dots, y_d \rightarrow 1]$$

At each step, we bound the size of the system of differential equation satisfied by the series (order + coefficient)

- **Main difficulty:** the Hadamard product \odot
- Detailed analysis of the proof of closure from [\[Lipshitz 89\]](#)