

Compressibility of finite languages by grammars

Stefan Hetzl

Institute of Discrete Mathematics and Geometry
Vienna University of Technology

joint work with Sebastian Eberhard

Descriptive Complexity of Formal Systems (DCFS) 2015
Waterloo, Ontario, Canada

June 26, 2015

- ▶ Grammar based compression
- ▶ Smallest grammar problem
(compression of a single word by a CFG)
- ▶ This talk: compression of a finite language by a grammar
incompressible sequence of finite languages
- ▶ Motivation: application in proof theory

- ▶ The smallest grammar problem(s)
- ▶ Incompressible languages
- ▶ Trees and proofs

The smallest grammar problem

- ▶ **Problem.**

Given $w \in \Sigma^*$, find minimal CFG G with $L(G) = \{w\}$
here: minimal w.r.t. sum of lengths of rhs of production rules

- ▶ **Decision Problem.**

Given $w \in \Sigma^*$ and $k \in \mathbb{N}$, is there a CFG G with $L(G) = \{w\}$
and $\text{size}(G) \leq k$?

- ▶ Decision problem **NP**-complete [Storer, Szymanski '78]

- ▶ Approximation: linear-time algorithms with logarithmic approximation ratio
[Charikar et al. '02], [Rytter '03], [Sakamoto '05],
[Charikar et al. '05], [Jež '13], [Jež '14]

- ▶ Practically efficient approximation algorithms

- ▶ Sequitur [Nevill-Manning, Witten '97]
- ▶ Re-Pair [Larsson, Moffat '99]

Our variant of the smallest grammar problem

- ▶ A grammar $G = (N, \Sigma, P, S)$ is called *right-linear* if all productions are of the form $A \rightarrow wB$ or $A \rightarrow w$ for $w \in \Sigma^*$.
- ▶ **Definition.** $A <_G^1 B$ if there is $A \rightarrow u \in P$ s.t. B occurs in u . Define $<_G$ as transitive closure of $<_G^1$.
- ▶ **Definition.** RLAG: right-linear acyclic grammar
- ▶ **Problem.**
Given finite $L \subseteq \Sigma^*$, find minimal RLAG G with $L(G) \supseteq L$.
here: minimal w.r.t. number of production rules
- ▶ $|G|$ is number of production rules

Many smallest grammar problems

- ▶ **Problem (traditional).**

Given $w \in \Sigma^*$, find minimal CFG G with $L(G) = \{w\}$

here: minimal w.r.t. sum of lengths of rhs of production rules

- ▶ **Problem (this talk).**

Given finite $L \subseteq \Sigma^*$, find minimal RLAG G with $L(G) \supseteq L$.

here: minimal w.r.t. number of production rules

Many smallest grammar problems

- ▶ **Problem (traditional).**

Given $w \in \Sigma^*$, find minimal CFG G with $L(G) = \{w\}$
here: minimal w.r.t. sum of lengths of rhs of production rules

- ▶ **Problem (this talk).**

Given finite $L \subseteq \Sigma^*$, find minimal RLAG G with $L(G) \supseteq L$.
here: minimal w.r.t. number of production rules

- ▶ Many smallest grammar problems:

- ▶ RLAG / ACFG / TRATG / ...
- ▶ Size / number of production rules / ...
- ▶ $L(G) \supseteq L, L(G) = L$

- ▶ Compression of a *finite language*

- ▶ Emphasis on formalism for compression

- ▶ Operations on compressed representation

- ✓ The smallest grammar problem(s)
- ▶ Incompressible languages
- ▶ Trees and proofs

- ▶ **Definition.** Finite L is called incompressible if every RLAG G with $L(G) \supseteq L$ satisfies $|G| \geq |L|$.
- ▶ **Definition.** A sequence $(L_n)_{n \geq 1}$ is called incompressible if there is an $M \in \mathbb{N}$ s.t. for all $n \geq M$ the language L_n is incompressible.

- ▶ **Definition.** Finite L is called incompressible if every RLAG G with $L(G) \supseteq L$ satisfies $|G| \geq |L|$.
- ▶ **Definition.** A sequence $(L_n)_{n \geq 1}$ is called incompressible if there is an $M \in \mathbb{N}$ s.t. for all $n \geq M$ the language L_n is incompressible.
- ▶ $L_n = \{a\}$ is incompressible.
- ▶ $L_n = \{a_1, \dots, a_n\}$ is incompressible.
- ▶ Is there incompressible $(L_n)_{n \geq 1}$ s.t.
 - ▶ alphabet is finite and
 - ▶ $|L_n|$ is unbounded ?

Incompressible languages

- ▶ $\Sigma = \{\mathbf{0}, \mathbf{1}, \mathbf{s}\}$
- ▶ Write $\mathbf{b}_l(i) \in \{\mathbf{0}, \mathbf{1}\}^l$ for l -bit binary representation of i .
- ▶ For $n \geq 1$ define

$$l(n) = \lceil \log_2(n) \rceil$$

$$k(n) = \lceil \frac{9n}{l(n) + 1} \rceil$$

$$L_n = \{(\mathbf{s}\mathbf{b}_{l(n)}(i))^{k(n)} \mid 0 \leq i \leq n - 1\}$$

- ▶ $|L_n| = n$
- ▶ Length of all $w \in L_n$ is $O(n)$

Incompressible languages – Example

For $n = 10$ we have $l(n) = 4$ and $k(n) = 18$ and $L_n =$

s0000s0000	...	s0000
s0001s0001	...	s0001
⋮	⋮	⋮
s1001s1001	...	s1001

Definition. Building block, segment.

Theorem. $(L_n)_{n \geq 1}$ is incompressible.

Proof Sketch.

1. W.r.t. compressibility: reduced RLAGs enough
2. Reduced RLAG that covers L_n has only short productions
3. Short productions cannot cover many segments
4. Compressing grammar must cover many segments per production

3 and 4 contradict.

- ▶ **Corollary.** There is no sequence $(G_n)_{n \geq 1}$ of RLAGs and $M \in \mathbb{N}$ s.t. $L(G_n) = L_n$ and $|G_n| < |L_n|$ for all $n \geq M$.
- ▶ **Theorem.** There is a sequence $(G_n)_{n \geq 1}$ of acyclic CFGs which compresses $(L_n)_{n \geq 1}$.

Proof. Let P_n be

$$\begin{aligned} S &\rightarrow (\mathbf{s}A_1)^{k(n)}, \\ A_1 &\rightarrow \mathbf{0}A_2 \mid \mathbf{1}A_2, \\ &\vdots \\ A_{l(n)} &\rightarrow \mathbf{0} \mid \mathbf{1}. \end{aligned}$$

Then $|P_n| = 2\lceil \log(n) \rceil + 1 < n = |L_n|$.

- ✓ The smallest grammar problem(s)
- ✓ Incompressible languages
- ▶ Trees and proofs

- ▶ Rigid tree languages [Jacquemard, Clay, Vacher '09]
- ▶ **Definition.** A *regular tree grammar* is a tuple (N, Σ, P, S) s.t. all productions are of the form $A \rightarrow t$ with $t \in T(\Sigma \cup N)$.
- ▶ **Definition.** $<_G$ on N as for word grammars.
- ▶ **Definition.** A derivation $S \Longrightarrow_G^* t$ satisfies *rigidity condition* if it uses at most one A -production for every nonterminal A .
- ▶ **Definition.** A totally rigid acyclic tree (TRAT) grammar is an acyclic regular tree grammar $G = (N, \Sigma, P, S)$. Define $L(G) = \{t \in T(\Sigma) \mid S \Longrightarrow_G^* t \text{ satisfying rigidity condition}\}$.
- ▶ *Example.* $S \rightarrow f(A, B), A \rightarrow g(B), B \rightarrow c \mid d$
as regular tree grammar:

$$L = \{f(g(c), c), f(g(c), d), f(g(d), c), f(g(d), d)\}$$

as TRATG:

$$L = \{f(g(c), c), f(g(d), d)\}$$

From word languages to tree languages

- ▶ For alphabet Σ define $\Sigma^T = \{f_x \mid x \in \Sigma\} \cup \{e\}$
- ▶ Map words to trees, e.g.: $(abaac)^T = f_a(f_b(f_a(f_a(f_c(e))))))$
- ▶ \cdot^T maps RLAG to TRATG
- ▶ **Lemma.** If L is RLA-incompressible, then L^T is TRAT-incompressible.
- ▶ **Corollary.** $(L_n^T)_{n \geq 1}$ is TRAT-incompressible.

A corollary in proof theory

- ▶ Inference rule “cut”: use of a lemma in a proof
- ▶ **Theorem** [H '12].
cut-free proof ... trivial tree grammar: tree language
proof with Π_1 -cuts ... (non-trivial) TRAT grammar
- ▶ Cut-elimination gives trivial bounds on compressibility
 \Rightarrow Π_1 -compression: exponential
- ▶ We construct formulas ψ_n in first-order predicate logic s.t.
 - ▶ cut-free complexity $O((2^n)^2)$
 - ▶ Π_1 -cut complexity 2^n \Rightarrow only quadratic

- ▶ Sequence of incompressible languages
- ▶ Compressing finite languages is interesting

Open Questions / Future Work

- ▶ Complexity of smallest grammar problem(s) for finite languages

We know: Decision problem for $\text{TRATG}(2)$, number of production rules, $L(G) \supseteq L$ is **NP**-complete.

- ▶ Approximation ratios?
- ▶ Practically efficient algorithms?