# *Choosing Word Occurrences for the Smallest Grammar Problem*

Rafael Carrascosa[1], <u>Matthias Gallé</u>[2],
François Coste[2], Gabriel Infante-Lopez[1]

[1]NLP Group
U. N. de Córdoba
Argentina

[2]Symbiose Project
IRISA / INRIA
France

UMR **IRISA**



LATA
May, 25[th] 2010

# Smallest Grammar Problem

### Problem Definition
Given a sequence $s$, find a context-free grammar $G(s)$ of minimal size that generates exactly this and only this sequence.

# Smallest Grammar Problem

### Problem Definition
Given a sequence $s$, find a context-free grammar $G(s)$ of minimal size that generates exactly this and only this sequence.

### Example
$s =$ "how much wood would a woodchuck chuck if a woodchuck could chuck wood?", a possible $G(s)$ (not necessarily minimal) is

| | | |
|---|---|---|
| $S$ | $\rightarrow$ | how much $N_2$ w$N_3$ $N_4$ $N_1$ if $N_4$ c$N_3$ $N_1$ $N_2$ ? |
| $N_1$ | $\rightarrow$ | chuck |
| $N_2$ | $\rightarrow$ | wood |
| $N_3$ | $\rightarrow$ | ould |
| $N_4$ | $\rightarrow$ | a $N_2 N_1$ |

# Smallest Grammar Problem

### Problem Definition
Given a sequence $s$, find a context-free grammar $G(s)$ of minimal size that generates exactly this and only this sequence.

### Applications

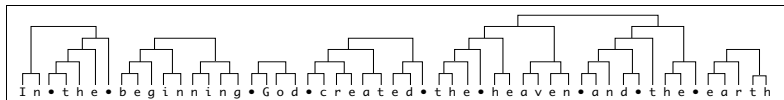- Data Compression
- Sequence Complexity
- Structure Discovery

# Smallest Grammar Problem

### Problem Definition
Given a sequence $s$, find a context-free grammar $G(s)$ of minimal size that generates exactly this and only this sequence.

### Applications

- Data Compression
- Sequence Complexity
- Structure Discovery



In•the•beginning•God•created•the•heaven•and•the•earth

©Nevill-Manning 1997

# Smallest Grammar Problem

### Problem Definition
Given a sequence $s$, find a context-free grammar $G(s)$ of minimal size that generates exactly this and only this sequence.

### Remark
Not only $S$, but any non-terminal of the grammar generates only one sequence of terminal symbols: $cons :: \mathcal{N} \to \Sigma^*$

# Smallest Grammar Problem

## Problem Definition

Given a sequence $s$, find a context-free grammar $G(s)$ of minimal size that generates exactly this and only this sequence.

## Remark

Not only $S$, but any non-terminal of the grammar generates only one sequence of terminal symbols: $cons :: \mathcal{N} \to \Sigma^*$

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $S$ | $\to$ | how much $N_2$ w$N_3$ $N_4$ $N_1$ if $N_4$ c$N_3$ $N_1$ $N_2$ ? | | | $cons(S)$ | $=$ | $s$ |
| $N_1$ | $\to$ | chuck | | | $cons(N_1)$ | $=$ | chuck |
| $N_2$ | $\to$ | wood | $\Rightarrow$ | | $cons(N_2)$ | $=$ | wood |
| $N_3$ | $\to$ | ould | | | $cons(N_3)$ | $=$ | ould |
| $N_4$ | $\to$ | a $N_2 N_1$ | | | $cons(N_4)$ | $=$ | a woodchuck |

# Smallest Grammar Problem

### Problem Definition
Given a sequence $s$, find a context-free grammar $G(s)$ of minimal
size that generates exactly this and only this sequence.

### Size of a Grammar

$$|G| = \sum_{N \to \omega \in \mathcal{P}} (|\omega| + 1)$$

# Smallest Grammar Problem

### Problem Definition

Given a sequence $s$, find a context-free grammar $G(s)$ of minimal size that generates exactly this and only this sequence.

### Size of a Grammar

$$|G| = \sum_{N \rightarrow \omega \in \mathcal{P}} (|\omega| + 1)$$

| | | |
|---|---|---|
| $S$ | $\rightarrow$ | how much $N_2$ w$N_3$ $N_4$ $N_1$ if $N_4$ c$N_3$ $N_1$ $N_2$ ? |
| $N_1$ | $\rightarrow$ | chuck |
| $N_2$ | $\rightarrow$ | wood |
| $N_3$ | $\rightarrow$ | ould |
| $N_4$ | $\rightarrow$ | a $N_2 N_1$ |

$$\Downarrow$$

how much $N_2$ w$N_3$ $N_4$ $N_1$ if $N_4$ c$N_3$ $N_1$ $N_2$ | chuck | wood | ould | a $N_2$ $N_1$ |

# Previous Approaches

1. Practical algorithms: Sequitur (and offline friends). 1996

   "Compression and Explanation Using Hierarchical Grammars". Nevill-Manning & Witten. The Computer

   Journal. 1997

2. Compression theoretical framework: Grammar Based Code. 2000

   "Grammar-based codes: a new class of universal lossless source codes". Kieffer & Yang. IEEE T on

   Information Theory. 2000

3. Approximation ratio to the size of a Smallest Grammar in the worst case. 2002

   "The Smallest Grammar Problem", Charikar et.al. IEEE T on Information Theory. 2005

# Previous Approaches

1. **Practical algorithms: Sequitur (and offline friends). 1996**

   "Compression and Explanation Using Hierarchical Grammars". Nevill-Manning & Witten. The Computer Journal. 1997

2. **Compression theoretical framework: Grammar Based Code. 2000**

   "Grammar-based codes: a new class of universal lossless source codes". Kieffer & Yang. IEEE T on Information Theory. 2000

3. **Approximation ratio to the size of a Smallest Grammar in the worst case. 2002**

   "The Smallest Grammar Problem", Charikar et.al. IEEE T on Information Theory. 2005

- Maximal Length (ML): take longest repeat, replace all occurrences with new symbol, iterate.

# Offline algorithms

- Maximal Length (ML): take longest repeat, replace all occurrences with new symbol, iterate.

$S \rightarrow$  how_much_wood_would_a_woodchuck_chuck_ if_a_woodchuck_could_chuck_wood?

- Maximal Length (ML): take longest repeat, replace all occurrences with new symbol, iterate.

$S \rightarrow$ how_much_wood_would_a_woodchuck_chuck_
if_a_woodchuck_could_chuck_wood?

# Offline algorithms

- Maximal Length (ML): take longest repeat, replace all occurrences with new symbol, iterate.

$S \rightarrow$ how_much_wood_would_a_woodchuck_chuck_
    if_a_woodchuck_could_chuck_wood?

$\Downarrow$

$S \rightarrow$ how_much_wood_would$N_1$huck_if$N_1$ould_chuck_wood?

$N_1 \rightarrow$ _a_woodchuck_c

# Offline algorithms

- Maximal Length (ML): take longest repeat, replace all occurrences with new symbol, iterate.

$S \quad \rightarrow \quad$ how_much_wood_would_a_woodchuck_chuck_
if_a_woodchuck_could_chuck_wood?

$\Downarrow$

$S \quad \rightarrow \quad$ how_much_wood_would$N_1$huck_if$N_1$ould_chuck_wood?

$N_1 \quad \rightarrow \quad$ _a_woodchuck_c

4

# Offline algorithms

- Maximal Length (ML): take longest repeat, replace all occurrences with new symbol, iterate.

$S \quad \rightarrow \quad$ how_much_wood_would_a_woodchuck_chuck_ if_a_woodchuck_could_chuck_wood?

$\qquad \Downarrow$

$S \quad \rightarrow \quad$ how_much_wood_would$N_1$huck_if$N_1$ould_chuck_wood?

$N_1 \quad \rightarrow \quad$ _a_woodchuck_c

$\qquad \Downarrow$

$S \quad \rightarrow \quad$ how_much_wood_would$N_1$huck_if_$N_1$ould_$N_2$wood?

$N_1 \quad \rightarrow \quad$ _a_wood$N_2$c

$N_2 \quad \rightarrow \quad$ chuck_

# Offline algorithms

- Maximal Length (ML): take longest repeat, replace all occurrences with new symbol, iterate.

  Bentley & McIlroy "Data compression using long common strings". DCC. 1999.

  Nakamura, et.al. "Linear-Time Text Compression by Longest-First Substitution". MDPI Algorithms. 1999

- Most Frequent (MF): take most frequent repeat, replace all occurrences with new symbol, iterate

  Larsson & Moffat. "Offline Dictionary-Based Compression". DCC. 1999

- Most Compressive (MC): take repeat that compress the best, replace with new symbol, iterate

  Apostolico & Lonardi. "Off-line compression by greedy textual substitution" Proceedings of IEEE. 2000
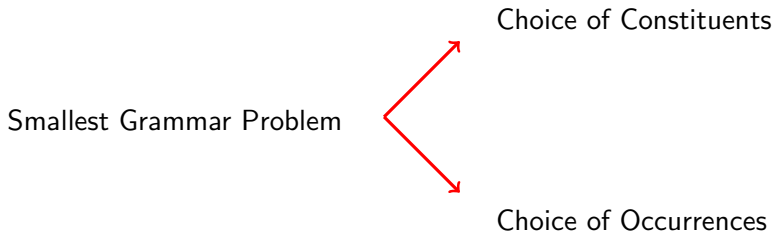
# A General Framework: IRR

IRR (Iterative Repeat Replacement) framework
Input: a sequence $s$, a score function $f$

1. Initialize Grammar by $S \rightarrow s$
2. take repeat $\omega$ that maximizes $f$ over $G$
3. **if** replacing $\omega$ would yield a bigger grammar than $G$
   **then**
   > 3.1 **return** $G$

   **else**
   > 3.1 replace all (non-overlapping) occurrences of $\omega$ in $G$ by new symbol $N$
   > 3.2 add rule $N \rightarrow \omega$ to $G$
   > 3.3 goto 2

Complexity: $\mathcal{O}(n^3)$

# Results on Canterbury Corpus

| sequence | Sequitur | IRR-ML | IRR-MF | IRR-MC |
|---|---|---|---|---|
| alice29.txt | 19.9% | 37.1% | 8.9% | 41000 |
| asyoulik.txt | 17.7% | 37.8% | 8.0% | 37474 |
| cp.html | 22.2% | 21.6% | 10.4% | 8048 |
| fields.c | 20.3% | 18.6% | 16.1% | 3416 |
| grammar.lsp | 20.2% | 20.7% | 15.1% | 1473 |
| kennedy.xls | 4.6% | 7.7% | 0.3% | 166924 |
| lcet10.txt | 24.5% | 45.0% | 8.0% | 90099 |
| plrabn12.txt | 14.9% | 45.2% | 5.8% | 124198 |
| ptt5 | 23.4% | 26.1% | 6.4% | 45135 |
| sum | 25.6% | 15.6% | 11.9% | 12207 |
| xargs.1 | 16.1% | 16.2% | 11.8% | 2006 |
| *average* | *19.0%* | *26.5%* | *9.3%* | |

Extends and confirms results of Nevill-Manning & Witten "On-Line and Off-Line Heuristics for Inferring Hierarchies of Repetitions in Sequences". Proc. of the IEEE. vol 80 no 11. November 2000

# A General Framework: IRR

IRR (Iterative Repeat Replacement) framework
Input: a sequence $s$, a score function $f$

1. Initialize Grammar by $S \rightarrow s$
2. take repeat $\omega$ that maximizes $f$ over $G$
3. **if** replacing $\omega$ would yield a bigger grammar than $G$
   **then**
      3.1 **return** $G$
   **else**
      3.1 replace all (non-overlapping) occurrences of $\omega$ in $G$ by new symbol $N$
      3.2 add rule $N \rightarrow \omega$ to $G$
      3.3 goto 2

Complexity: $\mathcal{O}(n^3)$

# Split the Problem

Smallest Grammar Problem

Choice of Constituents

Choice of Occurrences

# A General Framework: IRR

IRR (Iterative Repeat Replacement) framework
Input: a sequence $s$, a score function $f$

1. Initialize Grammar by $S \rightarrow s$
2. take repeat $\omega$ that maximizes $f$ over $G$
3. **if** replacing $\omega$ would yield a bigger grammar than $G$
   **then**
      3.1 **return** $G$
   **else**
      3.1 replace all (non-overlapping) occurrences of $\omega$ in $G$ by new symbol $N$
      3.2 add rule $N \rightarrow \omega$ to $G$
      3.3 goto 2

Complexity: $\mathcal{O}(n^3)$

# A General Framework: IRRCOO

IRRCOO (Iterative Repeat Replacement with Choice of Occurrence Optimization) framework
Input: a sequence $s$, a score function $f$

1. Initialize Grammar by $S \rightarrow s$

2. take repeat $\omega$ that maximizes $f$ over $G$

3. **if** replacing $\omega$ would yield a bigger grammar than $G$
   **then**
      3.1 **return** $G$
   **else**
      3.1 $G \leftarrow mgp(cons(G) \cup cons(\omega))$
      3.2 goto 2

# Choice of Occurrences

## Minimal Grammar Parsing (MGP) Problem

Given sequences $\Omega = \{s = w_0, w_1, \ldots, w_m\}$, find a context-free grammar of minimal size that has non-terminals $\{S = N_0, N_1, \ldots N_m\}$ such that $\text{cons}(N_i) = w_i$.
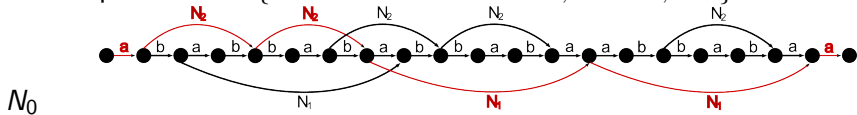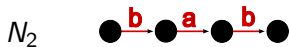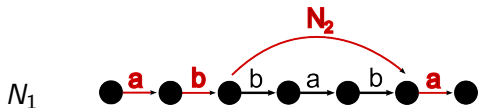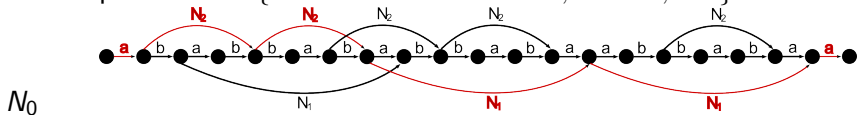
# Choice of Occurrences: an Example

Given sequences $\Omega = \{ababbababbabaabbabaa, abbaba, bab\}$

Given sequences $\Omega = \{abababababbababaabbabaa, abbaba, bab\}$



$N_0$

$N_1$

$N_2$

# Choice of Occurrences: an Example

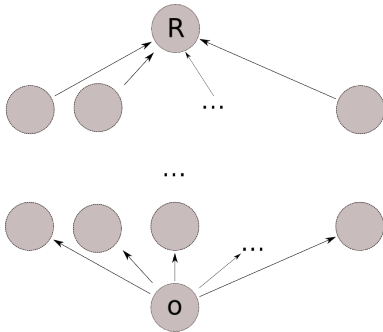Given sequences $\Omega = \{ababbababbabaabbabaa, abbaba, bab\}$

Given sequences $\Omega = \{ababbababbabaabbabaa, abbaba, bab\}$



$N_0$

$N_1$

$N_2$

A minimal grammar for $\Omega$ is

$$N_0 \rightarrow aN_2N_2N_1N_1a$$
$$N_1 \rightarrow abN_2a$$
$$N_2 \rightarrow bab$$

# Choice of Occurrences: an Example

Given sequences $\Omega = \{ababbababbbabaabbabaa, abbaba, bab\}$



$N_0$

$N_1$

$N_2$

$mgp$ can be computed in $\mathcal{O}(n^3)$

Smallest Grammar Problem

Choice of Constituents

Choice of Occurrences

# A Search Space for the SGP

Given $s$, take the lattice $\langle \mathcal{R}(s), \subseteq \rangle$ and associate a score to each node $\eta$: the size of the grammar $mgp(\eta \cup \{s\})$. A smallest grammar will have associated a node with minimal score.

# A Search Space for the SGP

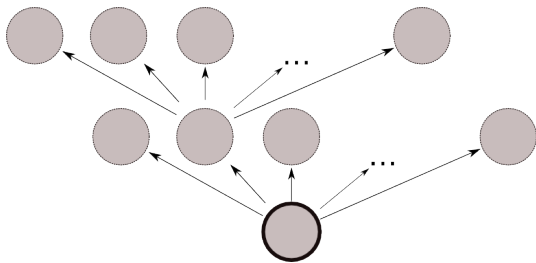### Lattice is a good search space

For every sequence $s$, there is a node $\eta$ in $\langle \mathcal{R}(s), \subseteq \rangle$ such that $mgp(\eta \cup \{s\})$ is a smallest grammar.

### Not the case for IRR search space

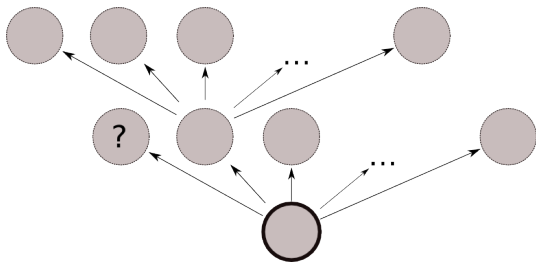But, there exists a sequence $s$ such that for any score function $f$, $IRR(s, f)$ does not return a smallest grammar $\boxed{\blacktriangleright \text{Proof}}$
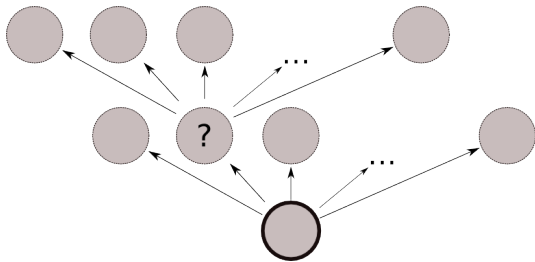
# The ZZ Algorithm

bottom-up phase: given node $\eta$, compute scores of nodes $\eta \cup \{w_i\}$ and take node with smallest score.

# The ZZ Algorithm

bottom-up phase: given node $\eta$, compute scores of nodes $\eta \cup \{w_i\}$ and take node with smallest score.

bottom-up phase: given node $\eta$, compute scores of nodes $\eta \cup \{w_i\}$ and take node with smallest score.
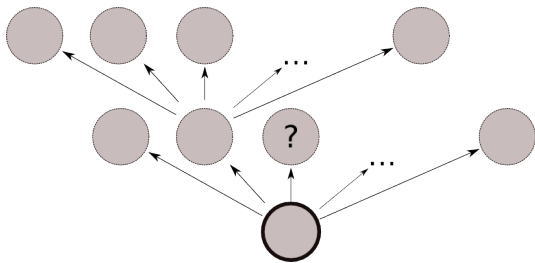
# The ZZ Algorithm

bottom-up phase: given node $\eta$, compute scores of nodes $\eta \cup \{w_i\}$ and take node with smallest score.

# The ZZ Algorithm

bottom-up phase: given node $\eta$, compute scores of nodes $\eta \cup \{w_i\}$ and take node with smallest score.

bottom-up phase: given node $\eta$, compute scores of nodes $\eta \cup \{w_i\}$ and take node with smallest score.
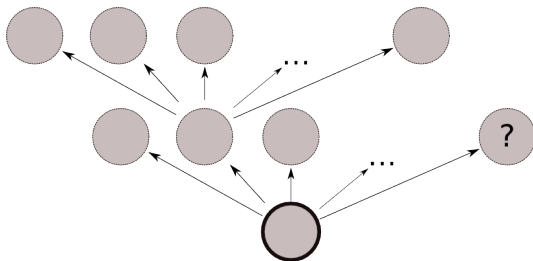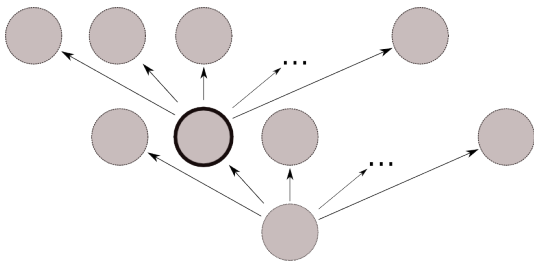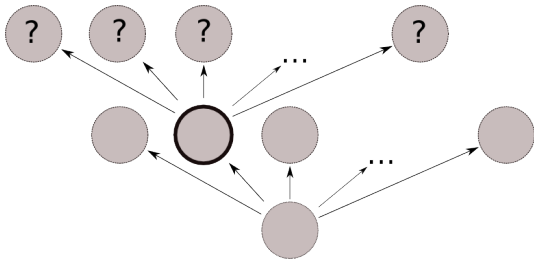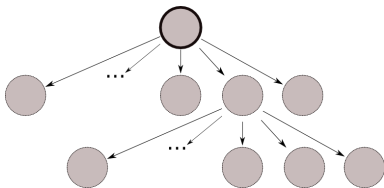
# The ZZ Algorithm

bottom-up phase: given node $\eta$, compute scores of nodes $\eta \cup \{w_i\}$ and take node with smallest score.

# The ZZ Algorithm

bottom-up phase: given node $\eta$, compute scores of nodes $\eta \cup \{w_i\}$ and take node with smallest score.

top-down phase: given node $\eta$, compute scores of nodes $\eta \setminus \{w_i\}$ and take node with smallest score.

# The ZZ Algorithm

bottom-up phase: given node $\eta$, compute scores of nodes $\eta \cup \{w_i\}$ and take node with smallest score.
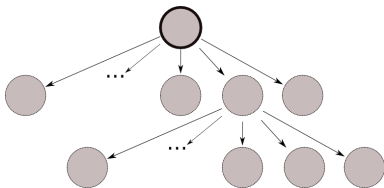
top-down phase: given node $\eta$, compute scores of nodes $\eta \setminus \{w_i\}$ and take node with smallest score.



ZZ: succession of both phases. Is in $\mathcal{O}(n^7)$

## Results on Canterbury Corpus

| sequence | IRRCOO-MC | ZZ | IRR-MC |
|---|---|---|---|
| alice29.txt | -4.3% | -8.0% | 41000 |
| asyoulik.txt | -2.9% | -6.6% | 37474 |
| cp.html | -1.3% | -3.5% | 8048 |
| fields.c | -1.3% | -3.1% | 3416 |
| grammar.lsp | -0.1% | -0.5% | 1473 |
| kennedy.xls | -0.1% | -0.1% | 166924 |
| lcet10.txt | -1.7% | – | 90099 |
| plrabn12.txt | -5.5% | – | 124198 |
| ptt5 | -2.6% | – | 45135 |
| sum | -0.8% | -1.5% | 12207 |
| xargs.1 | -0.8% | -1.7% | 2006 |
| *average* | *-2.0%* | *-3.1%* | |

## New Results

| Classification | sequence name | length | IRRMGP* | size improvement |
|---|---|---|---|---|
| Virus | P. lambda | 48 Knt | 13061 | -4.25% |
| Bacterium | E. coli | 4.6 Mnt | 741435 | -8.82% |
| Protist | T. pseudonana chrI | 3 Mnt | 509203 | -8.15% |
| Fungus | S. cerevisiae | 12.1 Mnt | 1742489 | -9.68% |
| Alga | O. tauri | 12.5 Mnt | 1801936 | -8.78% |

# Back to Structure

How similar are the structures returned by the different algorithms?

# Back to Structure

How similar are the structures returned by the different algorithms?
Standard measure to compare parse trees:

- Unlabeled Precision and Recall (F-measure)
- Unlabeled Non Crossing Precision and Recall (F-measure)

Dan Klein. "The Unsupervised Learning of Natural Language Structure". Phd Thesis. U Stanford. 2005

# Similarity of Structure

| sequence | algorithm vs IRR-MC | size gain | $U_F$ | $UNC_F$ |
|---|---|---|---|---|
| fields.c | ZZ | 3.1 % | 77.8 | 85.3 |
| | IRRCOO-MC | 1.3 % | 84.1 | 88.7 |
| cp.html | ZZ | 3.5 % | 66.3 | 75.0 |
| | IRRCOO-MC | 1.3 % | 81.4 | 84.8 |
| alice.txt | ZZ | 8.0 % | 36.6 | 38.6 |
| | IRRCOO-MC | 4.3 % | 63.9 | 66.0 |
| asyoulike.txt | ZZ | 6.6 % | 34.6 | 35.8 |
| | IRRCOO-MC | 2.9 % | 55.1 | 56.9 |

# Conclusions and Perspectices

- ⋆ Split SGP into two complementary problems: choice of constituents and choice of occurrences
- ⋆ Definition of a search space that contains a solution....
- ⋆ ... and to define algorithms which find smaller grammars than state-of-the-art.

# Conclusions and Perspectices

- ⋆ Split SGP into two complementary problems: choice of constituents and choice of occurrences
- ⋆ Definition of a search space that contains a solution....
- ⋆ ... and to define algorithms which find smaller grammars than state-of-the-art.
- • Promising results on DNA sequences (whole genomes)

## Conclusions and Perspectices

- ⋆ Split SGP into two complementary problems: choice of constituents and choice of occurrences
- ⋆ Definition of a search space that contains a solution....
- ⋆ ... and to define algorithms which find smaller grammars than state-of-the-art.
- • Promising results on DNA sequences (whole genomes)
- • Focus on the structure. Meaning of (dis)similarity.

# The End

$S \rightarrow$ thDkAforBr_attenC._DoAhave_Dy_quesCs?
$A \rightarrow$ B_
$B \rightarrow$ _you
$C \rightarrow$ tion
$D \rightarrow$ an

# Parse Tree Similarity Measures

$$UNC_P(P_1, P_2) = \frac{|\{b \in \text{brackets}(P_1) : b \text{ does not cross brackets}(P_2)\}|}{|\text{brackets}(P_1)|}$$

$$UNC_R(P_1, P_2) = \frac{|\{b \in \text{brackets}(P_2) : b \text{ does not cross brackets}(P_1)\}|}{|\text{brackets}(P_2)|}$$

$$UNC_F(P_1, P_2) = \frac{2}{UNC_P(P_1, P_2)^{-1} + UNC_R(P_1, P_2)^{-1}}$$

$$U_P(P_1, P_2) = \frac{|\text{brackets}(P_1) \cap \text{brackets}(P_2)|}{|\text{brackets}(P_1)|}$$

$$U_R(P_1, P_2) = \frac{|\text{brackets}(P_1) \cap \text{brackets}(P_2)|}{|\text{brackets}(P_2)|}$$

$$U_F(P_1, P_2) = \frac{2}{U_P(P_1, P_2)^{-1} + U_R(P_1, P_2)^{-1}}$$

Example

$xaxbxcx|_1xbxcxax|_2xcxaxbx|_3xaxcxbx|_4xbxaxcx|_5xcxbxax|_6xax|_7xbx|_8xcx$

### Example

$xaxbxcx|_1xbxcxax|_2xcxaxbx|_3xaxcxbx|_4xbxaxcx|_5xcxbxax|_6xax|_7xbx|_8xcx$

A smallest grammar is:

$$S \;\rightarrow\; AbC|_1BcA|_2CaB|_3AcB|_4BaC|_5CbA|_6A|_7B|_8C$$
$$A \;\rightarrow\; xax$$
$$B \;\rightarrow\; xbx$$
$$C \;\rightarrow\; xcx$$

### Example

$xaxbxcx|_1xbxcxax|_2xcxaxbx|_3xaxcxbx|_4xbxaxcx|_5xcxbxax|_6xax|_7xbx|_8xcx$

But what IRR can do is like:

$$
\begin{aligned}
S &\rightarrow Abxcx|_1xbxcA|_2xcAbx|_3Acxbx|_4xbAcx|_5xcxbA|_6A|_7xbx|_8xcx \\
A &\rightarrow xax \\
&\Downarrow \\
S &\rightarrow Abxcx|_1BcA|_2xcAbx|_3AcB|_4xbAcx|_5xcxbA|_6A|_7B|_8xcx \\
A &\rightarrow xax \\
B &\rightarrow xbx \\
&\Downarrow \\
S &\rightarrow AbC|_1BcA|_2xcAbx|_3AcB|_4xbAcx|_5CbA|_6A|_7B|_8C \\
A &\rightarrow xax \\
B &\rightarrow xbx \\
C &\rightarrow xcx
\end{aligned}
$$