

# On String Compression by Context-free Grammars

**Rüdiger Reischuk**

joint work with Jan Arpe

presented at Data Compression Conference, Utah 2006



Institut für Theoretische Informatik  
Universität zu Lübeck

2. Lübeck-Tartu Workshop  
September 26, 2006



- 1 Overview on Research at ITCS
- 2 Compression of Strings
  - The Problem Setting
  - Grammar-based Compression
  - Optimization Problem MGC
- 3 Transformation between Alphabets
  - Blockcodes
  - Upper and Lower Bounds
  - Special Blockcodes
- 4 Approximability
- 5 Conclusion

# Research at ITCS

- algorithmic complexity

# Research at ITCS

- algorithmic complexity
  - design efficient algorithms: randomized, parallel, online

# Research at ITCS

- algorithmic complexity
  - design efficient algorithms: randomized, parallel, online
  - algorithmic learning
    - discrete optimization, molecular biology

# Research at ITCS

- algorithmic complexity
  - design efficient algorithms: randomized, parallel, online
  - algorithmic learning
    - discrete optimization, molecular biology
  - complexity classes
    - lower bounds on problem complexity

# Research at ITCS

- algorithmic complexity
  - design efficient algorithms: randomized, parallel, online
  - algorithmic learning
    - discrete optimization, molecular biology
  - complexity classes
    - lower bounds on problem complexity
- formal methods and analysis: distributed systems

# Research at ITCS

- algorithmic complexity
  - design efficient algorithms: randomized, parallel, online
  - algorithmic learning
    - discrete optimization, molecular biology
  - complexity classes
    - lower bounds on problem complexity
- formal methods and analysis: distributed systems
- security, cryptography, compression of digital data



# Research at ITCS

- algorithmic complexity
  - design efficient algorithms: randomized, parallel, online
  - algorithmic learning
    - discrete optimization, molecular biology
  - complexity classes
    - lower bounds on problem complexity
- formal methods and analysis: distributed systems
- security, cryptography, compression of digital data
  - information hiding, steganography, watermarking

# Research at ITCS

- algorithmic complexity
  - design efficient algorithms: randomized, parallel, online
  - algorithmic learning
    - discrete optimization, molecular biology
  - complexity classes
    - lower bounds on problem complexity
- formal methods and analysis: distributed systems
- security, cryptography, compression of digital data
  - information hiding, steganography, watermarking
  - authentication and privacy issues in networks
    - string compression, signatures for digital data

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - 01010101  $\rightsquigarrow$  (0 1)<sup>4</sup>

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (0\ 1)^4$
  - $11110100001000101111 \rightsquigarrow$

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (0\ 1)^4$
  - $11110100001000101111 \rightsquigarrow$   
*largest prime less than 1M in binary*

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (0\ 1)^4$
  - $11110100001000101111 \rightsquigarrow$   
*largest prime less than 1M in binary*
  - ...

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (0\ 1)^4$
  - $11110100001000101111 \rightsquigarrow$   
*largest prime less than 1M in binary*
  - ...
- ① universal schemes



# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (0\ 1)^4$
  - $11110100001000101111 \rightsquigarrow$   
*largest prime less than 1M in binary*
  - ...
- ① universal schemes
  - $|\text{compress}(x)| = \text{Kolmogorov complexity } K(x)$   
not computable!

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (0\ 1)^4$
  - $11110100001000101111 \rightsquigarrow$   
*largest prime less than 1M in binary*
  - ...
- ① universal schemes
  - $|\text{compress}(x)| = \text{Kolmogorov complexity } K(x)$   
not computable!
- ② resource bounded versions

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (01)^4$
  - $11110100001000101111 \rightsquigarrow$   
*largest prime less than 1M in binary*
  - ...
- ① universal schemes
  - $|\text{compress}(x)| = \text{Kolmogorov complexity } K(x)$   
not computable!
- ② resource bounded versions
  - $K^T(x)$  for some time bound  $T$   
still very expensive to estimate!

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (01)^4$
  - $11110100001000101111 \rightsquigarrow$   
*largest prime less than 1M in binary*
  - ...
- ① universal schemes
  - $|\text{compress}(x)| = \text{Kolmogorov complexity } K(x)$   
not computable!
- ② resource bounded versions
  - $K^T(x)$  for some time bound  $T$   
still very expensive to estimate!
- ③ more efficient, practical methods:

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (01)^4$
  - $11110100001000101111 \rightsquigarrow$   
*largest prime less than 1M in binary*
  - ...
- ① universal schemes
  - $|\text{compress}(x)| = \text{Kolmogorov complexity } K(x)$   
not computable!
- ② resource bounded versions
  - $K^T(x)$  for some time bound  $T$   
still very expensive to estimate!
- ③ more efficient, practical methods:
  - Lempel-Ziv methodology: dictionaries with pointers

# Compression of Strings

- given a string  $x$  find a string  $y = \text{compress}(x)$  that specifies  $x$  in a more compact way
  - $01010101 \rightsquigarrow (01)^4$
  - $11110100001000101111 \rightsquigarrow$   
*largest prime less than 1M in binary*
  - ...
- ① universal schemes
  - $|\text{compress}(x)| = \text{Kolmogorov complexity } K(x)$   
not computable!
- ② resource bounded versions
  - $K^T(x)$  for some time bound  $T$   
still very expensive to estimate!
- ③ more efficient, practical methods:
  - Lempel-Ziv methodology: dictionaries with pointers
  - simple computational model: formal grammars

# Grammar-based Compression

given a string  $x \in \Sigma^*$

# Grammar-based Compression

given a string  $x \in \Sigma^*$

- 1 find a **context-free grammar**  $G_x$  such that



# Grammar-based Compression

given a string  $x \in \Sigma^*$

- 1 find a **context-free grammar**  $G_x$  such that  $L(G_x) = \{x\}$

# Grammar-based Compression

given a string  $x \in \Sigma^*$

- 1 find a context-free grammar  $G_x$  such that  $L(G_x) = \{x\}$
- 2 design a suitable encoding of  $G_x$ :  $\text{code}(G_x)$

# Grammar-based Compression

given a string  $x \in \Sigma^*$

- 1 find a context-free grammar  $G_x$  such that  $L(G_x) = \{x\}$
- 2 design a suitable encoding of  $G_x$ :  $\text{code}(G_x)$
- 3  **$\text{compress}(x) := \text{code}(G_x)$**

# Grammar-based Compression

given a string  $x \in \Sigma^*$

- 1 find a context-free grammar  $G_x$  such that  $L(G_x) = \{x\}$
- 2 design a suitable encoding of  $G_x$ :  $\text{code}(G_x)$
- 3  **$\text{compress}(x) := \text{code}(G_x)$**

further requirements for  $G_x$ :

# Grammar-based Compression

given a string  $x \in \Sigma^*$

- 1 find a context-free grammar  $G_x$  such that  $L(G_x) = \{x\}$
- 2 design a suitable encoding of  $G_x$ :  $\text{code}(G_x)$
- 3  **$\text{compress}(x) := \text{code}(G_x)$**

further requirements for  $G_x$ :

- deterministic and acyclic

# Grammar-based Compression

given a string  $x \in \Sigma^*$

- 1 find a context-free grammar  $G_x$  such that  $L(G_x) = \{x\}$
- 2 design a suitable encoding of  $G_x$ :  $\text{code}(G_x)$
- 3 **compress(x) := code(G<sub>x</sub>)**

further requirements for  $G_x$ :

- deterministic and acyclic

in short: „ $G_x$  is a grammar for  $x$ .“

# Grammar-based Compression

given a string  $x \in \Sigma^*$

- 1 find a context-free grammar  $G_x$  such that  $L(G_x) = \{x\}$
- 2 design a suitable encoding of  $G_x$ :  $\text{code}(G_x)$
- 3 **compress(x) := code( $G_x$ )**

further requirements for  $G_x$ :

- deterministic and acyclic

in short: „ $G_x$  is a grammar for  $x$ .“

main concern in the following: find a **small grammar  $G_x$  for  $x$ !**

# Grammar-based Compression

given a string  $x \in \Sigma^*$

- 1 find a context-free grammar  $G_x$  such that  $L(G_x) = \{x\}$
- 2 design a suitable encoding of  $G_x$ :  $\text{code}(G_x)$
- 3 **compress(x) := code(G<sub>x</sub>)**

further requirements for  $G_x$ :

- deterministic and acyclic

in short: „ $G_x$  is a grammar for  $x$ .“

main concern in the following: find a **small grammar  $G_x$  for  $x$ !**  
coding of  $G_x$  can be done by standard methods



# Example 1

## greedy strategy:

iteratively replace longest subwords that appear more than once by a nonterminal

# Example 1

Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

# Example 1

## Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

## Grammar $G_x$ for $x$

$S \rightarrow \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

# Example 1

## Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

## Grammar $G_x$ for $x$

$S \rightarrow \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

# Example 1

## Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

## Grammar $G_x$ for $x$

$S \rightarrow \text{Ars} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{A}$

$A \rightarrow \text{fische}$

# Example 1

## Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

## Grammar $G_x$ for $x$

$S \rightarrow \text{Ars} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup A$

$A \rightarrow \text{fische}$

# Example 1

## Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

## Grammar $G_x$ for $x$

$S \rightarrow \text{Ars} \sqcup \text{fritz} \sqcup \text{Bt} \sqcup \text{frische} \sqcup \text{A}$

$A \rightarrow \text{Be}$

$B \rightarrow \text{fisch}$

# Example 1

## Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

## Grammar $G_x$ for $x$

$S \rightarrow \text{Ars} \sqcup \text{fritz} \sqcup \text{Bt} \sqcup \text{frische} \sqcup \text{A}$

$A \rightarrow \text{Be}$

$B \rightarrow \text{fisch}$



# Example 1

## Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

## Grammar $G_x$ for $x$

$S \rightarrow \text{ArsCtz} \sqcup \text{BtCsche} \sqcup A$

$A \rightarrow \text{Be}$

$B \rightarrow \text{fisch}$

$C \rightarrow \sqcup \text{fri}$

# Example 1

## Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

## Grammar $G_x$ for $x$

$S \rightarrow \text{ArsCtz} \sqcup \text{BtCsche} \sqcup A$   
 $A \rightarrow \text{Be}$   
 $B \rightarrow \text{fisch}$   
 $C \rightarrow \sqcup \text{fri}$

# Example 1

## Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

## Grammar $G_x$ for $x$

$S \rightarrow \text{ArsCtz} \sqcup \text{BtCDe} \sqcup A$

$A \rightarrow \text{Be}$

$B \rightarrow \text{fiD}$

$C \rightarrow \sqcup \text{fri}$

$D \rightarrow \text{sch}$

## Example 1

### Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

### Grammar $G_x$ for $x$

$S \rightarrow \text{ArsCtz} \sqcup \text{BtCDe} \sqcup A$

$A \rightarrow \text{Be}$

$B \rightarrow \text{fiD}$

$C \rightarrow \sqcup \text{fri}$

$D \rightarrow \text{sch}$

### Compression Rate

$$|x| = 36$$

## Example 1

### Input

$x = \text{fischers} \sqcup \text{fritz} \sqcup \text{fischt} \sqcup \text{frische} \sqcup \text{fische}$

### Grammar $G_x$ for $x$

$S \rightarrow \text{ArsCtz} \sqcup \text{BtCDe} \sqcup A$

$A \rightarrow \text{Be}$

$B \rightarrow \text{fiD}$

$C \rightarrow \sqcup \text{fri}$

$D \rightarrow \text{sch}$

### Compression Rate

$|x| = 36 \quad \rightsquigarrow \quad |G_x| = 26$

# Optimization Problem

## Minimum Grammar Compression

**MGC**

# Optimization Problem

## Minimum Grammar Compression

### MGC

- input:  $x \in \Sigma^*$

# Optimization Problem

## Minimum Grammar Compression

### MGC

- input:  $x \in \Sigma^*$
- solution: grammar  $G_x$  for  $x$



# Optimization Problem

## Minimum Grammar Compression

### MGC

- input:  $x \in \Sigma^*$
- solution: grammar  $G_x$  for  $x$
- goal: minimize

$$|G_x| := \sum_{A \rightarrow \alpha \text{ in } G_x} |\alpha|$$

# Optimization Problem

## Minimum Grammar Compression

### MGC

- input:  $x \in \Sigma^*$
- solution: grammar  $G_x$  for  $x$
- goal: minimize

$$|G_x| := \sum_{A \rightarrow \alpha \text{ in } G_x} |\alpha|$$

- notation:  $m^*(x) := \min |G_x|$

# Optimization Problem

## Minimum Grammar Compression

### **k-MGC** bounded alphabet size

- input:  $x \in \Sigma^*$ ,  $|\Sigma| \leq k$
- solution: grammar  $G_x$  for  $x$
- goal: minimize

$$|G_x| := \sum_{A \rightarrow \alpha \text{ in } G_x} |\alpha|$$

- notation:  $m^*(x) := \min |G_x|$

# Optimization Problem

## Minimum Grammar Compression

### 2-MGC binary alphabet

- input:  $x \in \Sigma^*$ ,  $|\Sigma| = 2$
- solution: grammar  $G_x$  for  $x$
- goal: minimize

$$|G_x| := \sum_{A \rightarrow \alpha \text{ in } G_x} |\alpha|$$

- notation:  $m^*(x) := \min |G_x|$

# Known Results about MGC

- $m^*(x) \geq \Omega(\log |x|)$

# Known Results about MGC

- $m^*(x) \geq \Omega(\log |x|)$
- $G_x$  is a grammar for  $x$  of size  $m \implies \forall \ell \in [1, \dots, |x|]$   
number of diff. substrings of length  $\ell$  in  $x$  is at most  $m \cdot \ell$   
[Lehman 2002]

# Known Results about MGC

- $m^*(x) \geq \Omega(\log |x|)$
- $G_x$  is a grammar for  $x$  of size  $m \implies \forall \ell \in [1, \dots, |x|]$   
number of diff. substrings of length  $\ell$  in  $x$  is at most  $m \cdot \ell$   
[Lehman 2002]
- MGC is  $\mathcal{NP}$ -complete [Storer 1977]

# Known Results about MGC

- $m^*(x) \geq \Omega(\log |x|)$
- $G_x$  is a grammar for  $x$  of size  $m \implies \forall \ell \in [1, \dots, |x|]$   
number of diff. substrings of length  $\ell$  in  $x$  is at most  $m \cdot \ell$   
[Lehman 2002]
- MGC is  $\mathcal{NP}$ -complete [Storer 1977]
- $k$ -MGC is  $\mathcal{NP}$ -complete for  $k \geq 3$



# Known Results about MGC

- $m^*(x) \geq \Omega(\log |x|)$
- $G_x$  is a grammar for  $x$  of size  $m \implies \forall \ell \in [1, \dots, |x|]$   
number of diff. substrings of length  $\ell$  in  $x$  is at most  $m \cdot \ell$   
[Lehman 2002]
- MGC is  $\mathcal{NP}$ -complete [Storer 1977]
- $k$ -MGC is  $\mathcal{NP}$ -complete for  $k \geq 3$
  
- Open question: is 2-MGC  $\mathcal{NP}$ -complete as well?

# Known Results about MGC

- $m^*(x) \geq \Omega(\log |x|)$
- $G_x$  is a grammar for  $x$  of size  $m \implies \forall \ell \in [1, \dots, |x|]$   
number of diff. substrings of length  $\ell$  in  $x$  is at most  $m \cdot \ell$   
[Lehman 2002]
- MGC is  $\mathcal{NP}$ -complete [Storer 1977]
- $k$ -MGC is  $\mathcal{NP}$ -complete for  $k \geq 3$
  
- Open question: is 2-MGC  $\mathcal{NP}$ -complete as well?
- $\implies$  transformation between alphabets

# Known Results: Approximability

- several greedy strategies have been proposed  
performance can be very bad for specific examples

## Know Results: Approximability

- several greedy strategies have been proposed  
performance can be very bad for specific examples
- MGC is approximable with a factor  $O(\log |x|)$   
[Charikar, Lehman, Liu, Panigraphy, Prabhakaran, Rasala,  
Sahai, Shelat: STOC'02] + [Rytter: CPM'02]

## Know Results: Approximability

- several greedy strategies have been proposed  
performance can be very bad for specific examples
- MGC is approximable with a factor  $O(\log |x|)$   
[Charikar, Lehman, Liu, Panigraphy, Prabhakaran, Rasala,  
Sahai, Shelat: STOC'02] + [Rytter: CPM'02]
- Open question:
  - better upper bounds – even for simpler problems?
  - improvements for  $k$ -MGC?

## Know Results: Approximability

- several greedy strategies have been proposed  
performance can be very bad for specific examples
- MGC is approximable with a factor  $O(\log |x|)$   
[Charikar, Lehman, Liu, Panigrahy, Prabhakaran, Rasala,  
Sahai, Shelat: STOC'02] + [Rytter: CPM'02]
- Open question:
  - better upper bounds – even for simpler problems?
  - improvements for  $k$ -MGC?
- MGC is not  $1 + \frac{1}{8568}$ -approximable (unless  $\mathcal{P} = \mathcal{NP}$ )  
[Lehman, Shelat: SODA'02]

## Know Results: Approximability

- several greedy strategies have been proposed  
performance can be very bad for specific examples
- MGC is approximable with a factor  $O(\log |x|)$   
[Charikar, Lehman, Liu, Panigraphy, Prabhakaran, Rasala, Sahai, Shelat: STOC'02] + [Rytter: CPM'02]
- Open question:
  - better upper bounds – even for simpler problems?
  - improvements for  $k$ -MGC?
- MGC is not  $1 + \frac{1}{8568}$ -approximable (unless  $\mathcal{P} = \mathcal{NP}$ )  
[Lehman, Shelat: SODA'02]
- Open question:
  - lower bound for  $k$ -MGC?
  - nonconstant lower bound for MGC?

# Transformation between Alphabets

How do compression properties transfer between alphabets?

- $\Gamma$  arbitrary alphabet,  $\Sigma$  alphabet of size  $k$



# Transformation between Alphabets

How do compression properties transfer between alphabets?

- $\Gamma$  arbitrary alphabet,  $\Sigma$  alphabet of size  $k$
- $\varphi : \Gamma^* \rightarrow \Sigma^*$  is an  $l$ -blockcode

# Transformation between Alphabets

How do compression properties transfer between alphabets?

- $\Gamma$  arbitrary alphabet,  $\Sigma$  alphabet of size  $k$
- $\varphi : \Gamma^* \rightarrow \Sigma^*$  is an  $l$ -blockcode
  - iff homomorphism + injective

# Transformation between Alphabets

How do compression properties transfer between alphabets?

- $\Gamma$  arbitrary alphabet,  $\Sigma$  alphabet of size  $k$
- $\varphi : \Gamma^* \rightarrow \Sigma^*$  is an  $l$ -blockcode
  - iff homomorphism + injective
  - $|\varphi(\gamma)| = l$  for all  $\gamma \in \Gamma$

# Transformation between Alphabets

How do compression properties transfer between alphabets?

- $\Gamma$  arbitrary alphabet,  $\Sigma$  alphabet of size  $k$
- $\varphi : \Gamma^* \rightarrow \Sigma^*$  is an  $l$ -blockcode
  - iff homomorphism + injective
  - $|\varphi(\gamma)| = l$  for all  $\gamma \in \Gamma$
- Question:

# Transformation between Alphabets

How do compression properties transfer between alphabets?

- $\Gamma$  arbitrary alphabet,  $\Sigma$  alphabet of size  $k$
- $\varphi : \Gamma^* \rightarrow \Sigma^*$  is an  $l$ -blockcode
  - iff homomorphism + injective
  - $|\varphi(\gamma)| = l$  for all  $\gamma \in \Gamma$
- Question:
  - relations between  $G_x$  and  $G_{\varphi(x)}$ ?

# Transformation between Alphabets

How do compression properties transfer between alphabets?

- $\Gamma$  arbitrary alphabet,  $\Sigma$  alphabet of size  $k$
- $\varphi : \Gamma^* \rightarrow \Sigma^*$  is an  $\ell$ -blockcode
  - iff homomorphism + injective
  - $|\varphi(\gamma)| = \ell$  for all  $\gamma \in \Gamma$
- Question:
  - relations between  $G_x$  and  $G_{\varphi(x)}$ ?
  - in particular,  $m^*(x) \leftrightarrow m^*(\varphi(x))$ ?

# From $G_x$ to $G_{\varphi(x)}$

## Example 2

*string*

$x = \text{abcaabc}$

*transformation*

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

---

# From $G_x$ to $G_{\varphi(x)}$

## Example 2

*string*

$x = \text{abcaabc}$

*transformation*

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

---

$G_x :$

S	→	PaP
P	→	abc



# From $G_x$ to $G_{\varphi(x)}$

## Example 2

*string*

$x = \text{abcaabc}$

*transformation*

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

$G_x$			$G_{\varphi(x)}$	
S	→	PaP	$S_a$	→ 1T1
P	→	abc	$S_b$	→ T10
			$S_c$	→ 01T
			T	→ 00

# From $G_x$ to $G_{\varphi(x)}$

## Example 2

<i>string</i>	<i>transformation</i>
$x = \text{abcaabc}$	$\varphi(a) = 1001, \varphi(b) = 0010$ $\varphi(c) = 0100$
$G_x :$	
S → PaP	$S_a \rightarrow 1T1$
P → abc	$S_b \rightarrow T10$
	$S_c \rightarrow 01T$
	$T \rightarrow 00$
$G_{\varphi(x)} :$	
S → PS <sub>a</sub> P	$S_a \rightarrow 1T1$
P → S <sub>a</sub> S <sub>b</sub> S <sub>c</sub>	$S_b \rightarrow T10$
	$S_c \rightarrow 01T$
	$T \rightarrow 00$

# From $G_x$ to $G_{\varphi(x)}$

## Example 2

<i>string</i>		<i>transformation</i>	
$x = abcaabc$		$\varphi(a) = 1001, \varphi(b) = 0010$ $\varphi(c) = 0100$	
$G_x :$	$S \rightarrow PaP$ $P \rightarrow abc$		$S_a \rightarrow 1T1$ $S_b \rightarrow T10$ $S_c \rightarrow 01T$ $T \rightarrow 00$
$G_{\varphi(x)} :$	$S \rightarrow PS_aP$ $P \rightarrow S_aS_bS_c$		$S_a \rightarrow 1T1$ $S_b \rightarrow T10$ $S_c \rightarrow 01T$ $T \rightarrow 00$

$$|G_{\varphi(x)}| = |G_x| + |G_{\varphi}| = 6 + 11 = 17$$

## From $G_x$ to $G_{\varphi(x)}$

### Theorem

$G_x$  grammar for  $x$ ,  $G_\varphi$  grammar for  $\varphi(\Gamma)$

## From $G_x$ to $G_{\varphi(x)}$

### Theorem

$G_x$  grammar for  $x$ ,  $G_\varphi$  grammar for  $\varphi(\Gamma)$

$\implies G_{\varphi(x)} = G_x \cup G_\varphi$  is a grammar for  $\varphi(x)$ , that means

$$|G_{\varphi(x)}| \leq |G_x| + |G_\varphi|$$

## From $G_x$ to $G_{\varphi(x)}$

### Theorem

$G_x$  grammar for  $x$ ,  $G_\varphi$  grammar for  $\varphi(\Gamma)$

$\implies G_{\varphi(x)} = G_x \cup G_\varphi$  is a grammar for  $\varphi(x)$ , that means

$$|G_{\varphi(x)}| \leq |G_x| + |G_\varphi|$$

$$m^*(\varphi(x)) \leq m^*(x) + m^*(\varphi(\Gamma))$$

# From $G_{\varphi(x)}$ to $G_x$

## Example 3

$x = abcaabc$

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

---

$G_{\varphi(x)}$ : S  $\rightarrow$  PAP  
 P  $\rightarrow$  ABC

---

A  $\rightarrow$  1T1  
 B  $\rightarrow$  T10  
 C  $\rightarrow$  01T  
 T  $\rightarrow$  00

# From $G_{\varphi(x)}$ to $G_x$

## Example 3

$x = abcaabc$

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

---

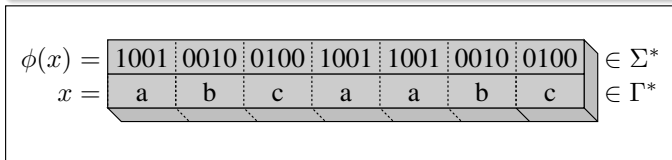
$G_{\varphi(x)}$ : S  $\rightarrow$  PAP  
 P  $\rightarrow$  ABC

A  $\rightarrow$  1T1

B  $\rightarrow$  T10

C  $\rightarrow$  01T

T  $\rightarrow$  00





# From $G_{\varphi(x)}$ to $G_x$

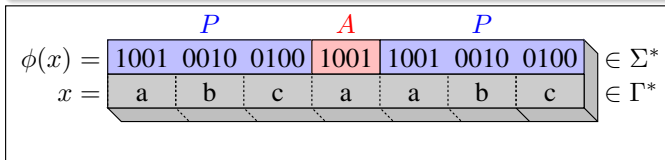
## Example 3

$x = abcaabc$

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

$G_{\varphi(x)}$ :	S	→	PAP	A	→	1T1
	P	→	ABC	B	→	T10
				C	→	01T
				T	→	00



# From $G_{\varphi(x)}$ to $G_x$

## Example 3

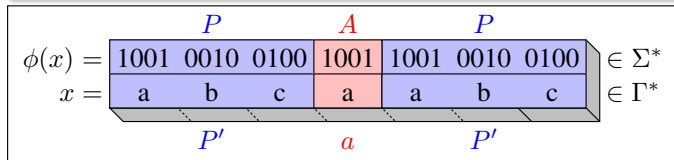
$x = abcaabc$

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

---

$G_x: S' \rightarrow P'aP'$   
 $P' \rightarrow abc$



# From $G_{\varphi(x)}$ to $G_x$

## Example 3

$x = abcaabc$

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

---

$G_x: S' \rightarrow P'aP'$   
 $P' \rightarrow abc$

$$|G_x| = |G_{\varphi(x)}| - |G_{\varphi}| = 17 - 11 = 6$$

## From $G_{\varphi(x)}$ to $G_x$

### Example 3

$x = \text{abcaabc}$

$\varphi(\text{a}) = 1001, \varphi(\text{b}) = 0010$

$\varphi(\text{c}) = 0100$

---

$G_x: S' \rightarrow P'aP'$

$P' \rightarrow \text{abc}$

$$|G_x| = |G_{\varphi(x)}| - |G_{\varphi}| = 17 - 11 = 6$$

? always true ?

$$|G_x| \leq |G_{\varphi(x)}| - m^*(\varphi(\Gamma)) \quad ?$$

# From $G_{\varphi(x)}$ to $G_x$

## Example 4

$x = abcaabc$

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

---

$G_{\varphi(x)}: S \rightarrow UUV1UU$

$U \rightarrow VV$

$V \rightarrow 100$

# From $G_{\varphi(x)}$ to $G_x$

## Example 4

$x = \text{abcaabc}$

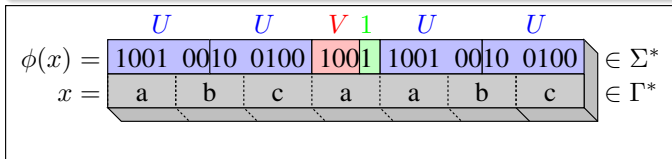
$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

$G_{\varphi(x)}: S \rightarrow \text{UUV1UU}$

$U \rightarrow \text{VV}$

$V \rightarrow 100$



# From $G_{\varphi(x)}$ to $G_x$

## Example 4

$x = \text{abcaabc}$

$\varphi(\text{a}) = 1001, \varphi(\text{b}) = 0010$

$\varphi(\text{c}) = 0100$

---

$G_{\varphi(x)}: S \rightarrow \text{UUV1UU}$

$U \rightarrow \text{VV}$

$V \rightarrow 100$

$$|G_{\varphi(x)}| - m^*(\varphi(\Gamma)) = 11 - 11 = 0$$

# From $G_{\varphi(x)}$ to $G_x$

## Example 4

$x = \text{abcaabc}$

$\varphi(\text{a}) = 1001, \varphi(\text{b}) = 0010$

$\varphi(\text{c}) = 0100$

---

$G_{\varphi(x)}: S \rightarrow \text{UUV1UU}$

$U \rightarrow \text{VV}$

$V \rightarrow 100$

$$|G_{\varphi(x)}| - m^*(\varphi(\Gamma)) = 11 - 11 = 0$$

$$|G_x| = 6 \quad ???$$



# From $G_{\varphi(x)}$ to $G_x$

## Example 4

$x = \text{abcaabc}$

$\varphi(a) = 1001, \varphi(b) = 0010$

$\varphi(c) = 0100$

---

$G_{\varphi(x)}: S \rightarrow \text{UUV1UU}$

$U \rightarrow \text{VV}$

$V \rightarrow 100$

$$|G_{\varphi(x)}| - m^*(\varphi(\Gamma)) = 11 - 11 = 0$$

$$|G_x| = 6 \quad ???$$

$$|G_x| \leq |G_{\varphi(x)}| \quad ???$$

## From $G_{\varphi(x)}$ to $G_x$ : lower bound

Small alphabets may compress much better

### Theorem

For every  $\ell$  there exists an alphabet  $\Gamma$  of size  $2^\ell$ ,  
a binary blockcode  $\varphi : \Gamma^* \rightarrow \{0,1\}^*$  and  
a string  $x \in \Gamma^*$  of size about  $2^{2\ell}$  such that

$$m^*(x) \geq \Omega(\ell \cdot m^*(\varphi(x)))$$

## From $G_{\varphi(x)}$ to $G_x$ : lower bound

Small alphabets may compress much better

### Theorem

For every  $\ell$  there exists an alphabet  $\Gamma$  of size  $2^\ell$ ,  
a binary blockcode  $\varphi : \Gamma^* \rightarrow \{0,1\}^*$  and  
a string  $x \in \Gamma^*$  of size about  $2^{2\ell}$  such that

$$m^*(x) \geq \Omega(\ell \cdot m^*(\varphi(x)))$$

$$m^*(\varphi(x)) \leq O(2^\ell), \quad m^*(\varphi(\Gamma)) \leq O(2^\ell), \quad m^*(x) \geq \Omega(\ell \cdot 2^\ell)$$

## From $G_{\varphi(x)}$ to $G_x$ : upper bound

This loss of compression is the worst for  $l$ -blockcodes

### Theorem

If  $\varphi$  is an  $l$ -blockcode and  $G_{\varphi(x)}$  a grammar for a coded string  $\varphi(x)$

$\implies \exists$  grammar  $G_x$  for  $x$  with  $|G_x| \leq 2 \cdot l \cdot |G_{\varphi(x)}|$

## From $G_{\varphi(x)}$ to $G_x$ : upper bound

This loss of compression is the worst for  $l$ -blockcodes

### Theorem

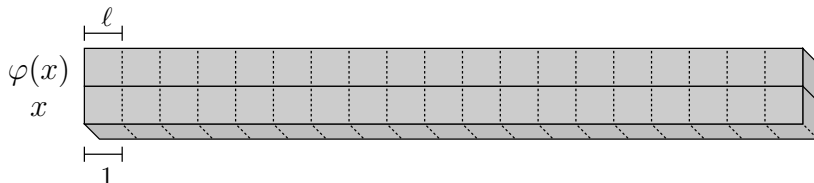
If  $\varphi$  is an  $l$ -blockcode and  $G_{\varphi(x)}$  a grammar for a coded string  $\varphi(x)$

$\implies \exists$  grammar  $G_x$  for  $x$  with  $|G_x| \leq 2 \cdot l \cdot |G_{\varphi(x)}|$

$$m^*(x) \leq 2 \cdot l \cdot m^*(\varphi(x)) \quad \forall x$$

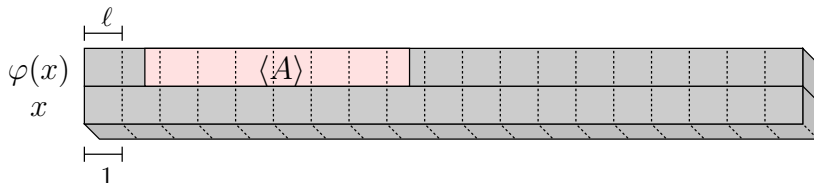
# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha$



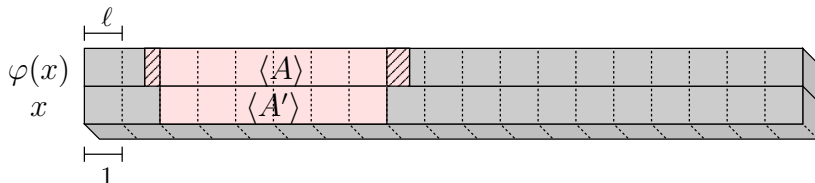
# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha$



# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha$

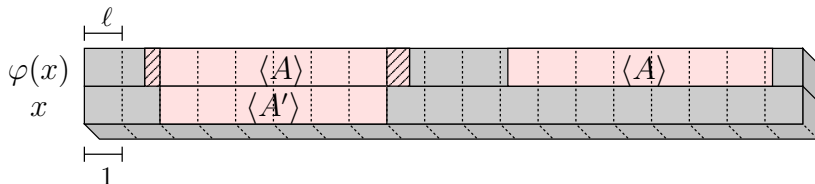


in  $G_x$ :  $A' \rightarrow \alpha'$



# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

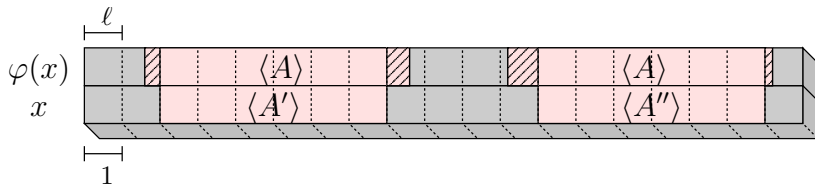
in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha$



in  $G_x$ :  $A' \rightarrow \alpha'$

# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

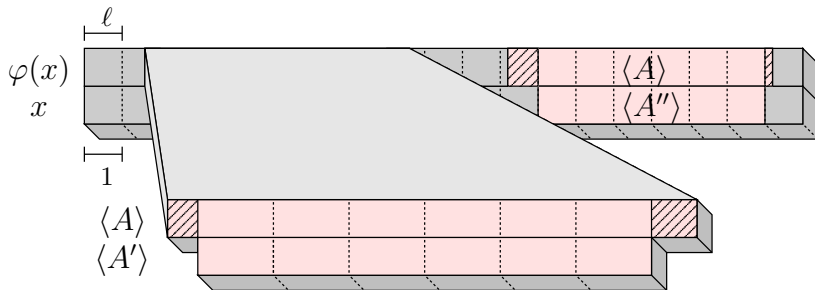
in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha$



in  $G_x$ :  $A' \rightarrow \alpha'$   $A'' \rightarrow \alpha'' \dots$  max.  $l$  copies

# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

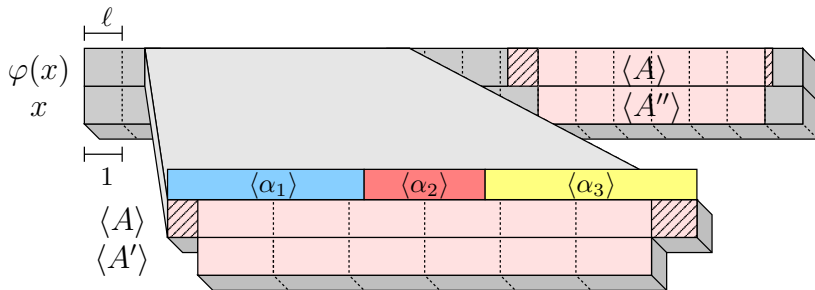
in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha$



in  $G_x$ :  $A' \rightarrow \alpha'$

# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

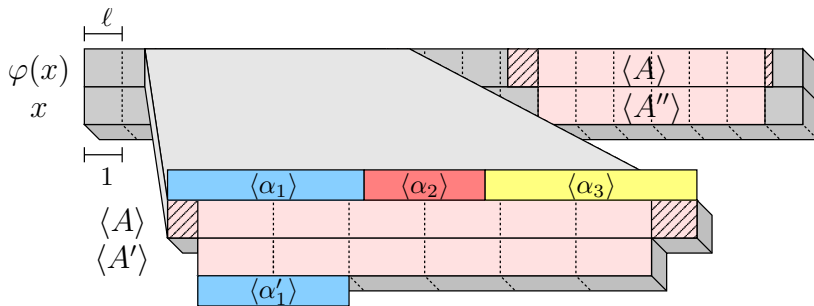
in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha = \alpha_1\alpha_2\alpha_3$



in  $G_x$ :  $A' \rightarrow \alpha'$

# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

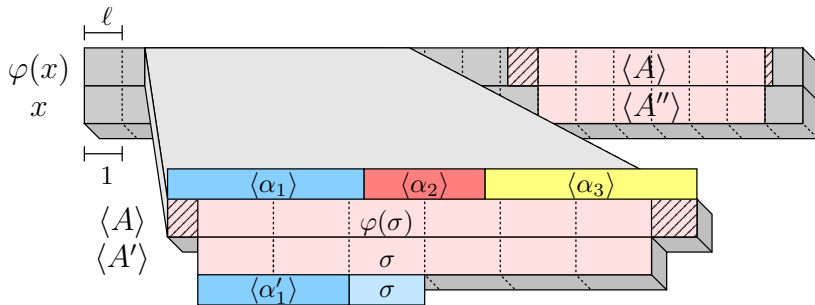
in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha = \alpha_1\alpha_2\alpha_3$



in  $G_x$ :  $A' \rightarrow \alpha' = \alpha'_1$

# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

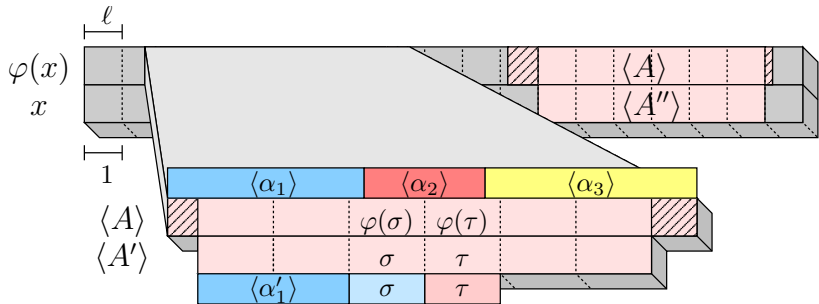
in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha = \alpha_1\alpha_2\alpha_3$



in  $G_x$ :  $A' \rightarrow \alpha' = \alpha'_1\sigma$

# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

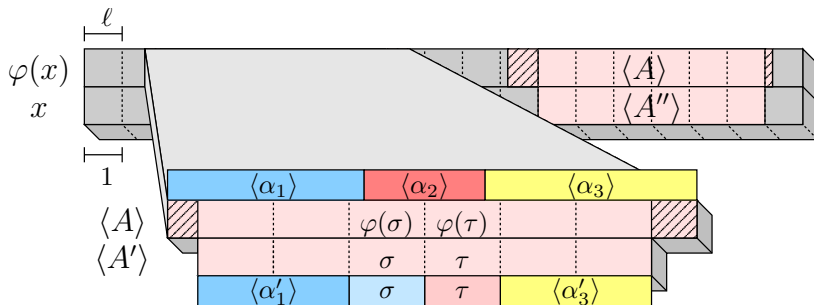
in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha = \alpha_1\alpha_2\alpha_3$



in  $G_x$ :  $A' \rightarrow \alpha' = \alpha'_1\sigma\tau$

# From $G_{\varphi(x)}$ to $G_x$ : upper bound - idea of proof

in  $G_{\varphi(x)}$ :  $A \rightarrow \alpha = \alpha_1\alpha_2\alpha_3$



in  $G_x$ :  $A' \rightarrow \alpha' = \alpha'_1\sigma\tau\alpha'_3$



# Better Bounds for Special Blockcodes

## Overlap-free Blockcodes

$\Gamma = \{\gamma_0, \dots, \gamma_k\}$ ,  $k < 2^m$ ,  $b_m(i) :=$  binary repr. of  $i$  of length  $m$

# Better Bounds for Special Blockcodes

## Overlap-free Blockcodes

$\Gamma = \{\gamma_0, \dots, \gamma_k\}$ ,  $k < 2^m$ ,  $b_m(i) :=$  binary repr. of  $i$  of length  $m$

①  $\varphi(\gamma_i) := 0^m 1 b_m(i) 1$

# Better Bounds for Special Blockcodes

## Overlap-free Blockcodes

$\Gamma = \{\gamma_0, \dots, \gamma_k\}$ ,  $k < 2^m$ ,  $b_m(i) :=$  binary repr. of  $i$  of length  $m$

①  $\varphi(\gamma_i) := 0^m 1 b_m(i) 1$  **overlap-free  $\ell = (2m + 2)$ -blockcode**

# Better Bounds for Special Blockcodes

## Overlap-free Blockcodes

$\Gamma = \{\gamma_0, \dots, \gamma_k\}$ ,  $k < 2^m$ ,  $b_m(i) :=$  binary repr. of  $i$  of length  $m$

- 1  $\varphi(\gamma_i) := 0^m 1 b_m(i) 1$  **overlap-free  $\ell = (2m + 2)$ -blockcode**
- 2  $G_x$  grammar for  $x \in \Gamma^*$ ,

# Better Bounds for Special Blockcodes

## Overlap-free Blockcodes

$\Gamma = \{\gamma_0, \dots, \gamma_k\}$ ,  $k < 2^m$ ,  $b_m(i) :=$  binary repr. of  $i$  of length  $m$

①  $\varphi(\gamma_i) := 0^m 1 b_m(i) 1$  **overlap-free  $\ell = (2m + 2)$ -blockcode**

②  $G_x$  grammar for  $x \in \Gamma^*$ ,  
 $\implies \exists$  grammar  $G_{\varphi(x)}$  for  $\varphi(x)$  with

$$|G_{\varphi(x)}| \leq (12 + o(1)) \cdot |G_x|$$

# Better Bounds for Special Blockcodes

## Overlap-free Blockcodes

$\Gamma = \{\gamma_0, \dots, \gamma_k\}$ ,  $k < 2^m$ ,  $b_m(i) :=$  binary repr. of  $i$  of length  $m$

①  $\varphi(\gamma_i) := 0^m 1 b_m(i) 1$  **overlap-free  $\ell = (2m + 2)$ -blockcode**

②  $G_x$  grammar for  $x \in \Gamma^*$ ,  
 $\implies \exists$  grammar  $G_{\varphi(x)}$  for  $\varphi(x)$  with

$$|G_{\varphi(x)}| \leq (12 + o(1)) \cdot |G_x|$$

③  $G_{\varphi(x)}$  grammar for  $\varphi(x)$

# Better Bounds for Special Blockcodes

## Overlap-free Blockcodes

$\Gamma = \{\gamma_0, \dots, \gamma_k\}$ ,  $k < 2^m$ ,  $b_m(i) :=$  binary repr. of  $i$  of length  $m$

①  $\varphi(\gamma_i) := 0^m 1 b_m(i) 1$  **overlap-free  $\ell = (2m + 2)$ -blockcode**

②  $G_x$  grammar for  $x \in \Gamma^*$ ,  
 $\implies \exists$  grammar  $G_{\varphi(x)}$  for  $\varphi(x)$  with

$$|G_{\varphi(x)}| \leq (12 + o(1)) \cdot |G_x|$$

③  $G_{\varphi(x)}$  grammar for  $\varphi(x)$   
 $\implies \exists$  grammar  $G_x$  for  $x$  with

$$|G_x| \leq 2 \cdot |G_{\varphi(x)}|$$

# Approximability

## Corollary

2-MGC  $c$ -approximable



# Approximability

## Corollary

2-MGC  $c$ -approximable

$\Rightarrow$  MGC  $(24c + \varepsilon)$ -approximable

# Approximability

## Corollary

2-MGC  $c$ -approximable

$\Rightarrow$  MGC  $(24c + \varepsilon)$ -approximable

## Corollary

nonconstant lower bound for MGC

# Approximability

## Corollary

2-MGC  $c$ -approximable

$\Rightarrow$  MGC  $(24c + \varepsilon)$ -approximable

## Corollary

nonconstant lower bound for MGC

$\Rightarrow$  nonconstant lower bound for 2-MGC

# Approximability

## Corollary

2-MGC  $c$ -approximable

$\Rightarrow$  MGC  $(24c + \varepsilon)$ -approximable

## Corollary

nonconstant lower bound for MGC

$\Rightarrow$  nonconstant lower bound for 2-MGC

$\Rightarrow$  2-MGC  $\notin$  P (if  $\mathcal{P} \neq \mathcal{NP}$ )

# Conclusion

## Bounds on the Compression Ratio

- optimal compression of a string and its  $\ell$ -blockcoding:

$$\frac{m^*(x)}{2\ell} \leq m^*(\varphi(x)) \leq m^*(x) + m^*(\varphi(\Gamma))$$

# Conclusion

## Bounds on the Compression Ratio

- optimal compression of a string and its  $\ell$ -blockcoding:

$$\frac{m^*(x)}{2\ell} \leq m^*(\varphi(x)) \leq m^*(x) + m^*(\varphi(\Gamma))$$

- bounds are asymptotically sharp

# Conclusion

## Bounds on the Compression Ratio

- optimal compression of a string and its  $\ell$ -blockcoding:

$$\frac{m^*(x)}{2\ell} \leq m^*(\varphi(x)) \leq m^*(x) + m^*(\varphi(\Gamma))$$

- bounds are asymptotically sharp
- improvements for overlap-free blockcodes:

$$m^*(x) \in \Theta(m^*(\varphi(x)))$$

# Conclusion

## Bounds on the Compression Ratio

- optimal compression of a string and its  $\ell$ -blockcoding:

$$\frac{m^*(x)}{2\ell} \leq m^*(\varphi(x)) \leq m^*(x) + m^*(\varphi(\Gamma))$$

- bounds are asymptotically sharp
- improvements for overlap-free blockcodes:

$$m^*(x) \in \Theta(m^*(\varphi(x)))$$

- implications for the approximability



# Conclusion

## Open Problems

- nonconstant lower bound for MGC

# Conclusion

## Open Problems

- nonconstant lower bound for MGC
- $\mathcal{NP}$ -completeness of 2-MGC

# Conclusion

## Open Problems

- nonconstant lower bound for MGC
- $\mathcal{NP}$ -completeness of 2-MGC
- optimal compression in the neighborhood of a string

# Conclusion

## Open Problems

- nonconstant lower bound for MGC
- $\mathcal{NP}$ -completeness of 2-MGC
- optimal compression in the neighborhood of a string
  - given a string  $x$  and a distance  $d$   
with respect to a metric – for example, the edit distance –  
find a string  $z$  in the  $d$ -neighborhood of  $x$  with maximal  
grammar compression

# Conclusion

## Open Problems

- nonconstant lower bound for MGC
- $\mathcal{NP}$ -completeness of 2-MGC
- optimal compression in the neighborhood of a string
  - given a string  $x$  and a distance  $d$  with respect to a metric – for example, the edit distance – find a string  $z$  in the  $d$ -neighborhood of  $x$  with maximal grammar compression
  - bounds on the relation between distances of strings and their difference in optimal grammar compression

# Conclusion

## Open Problems

- nonconstant lower bound for MGC
- $\mathcal{NP}$ -completeness of 2-MGC
- optimal compression in the neighborhood of a string
  - given a string  $x$  and a distance  $d$  with respect to a metric – for example, the edit distance – find a string  $z$  in the  $d$ -neighborhood of  $x$  with maximal grammar compression
  - bounds on the relation between distances of strings and their difference in optimal grammar compression
  - edit grammars [CLL+02]