

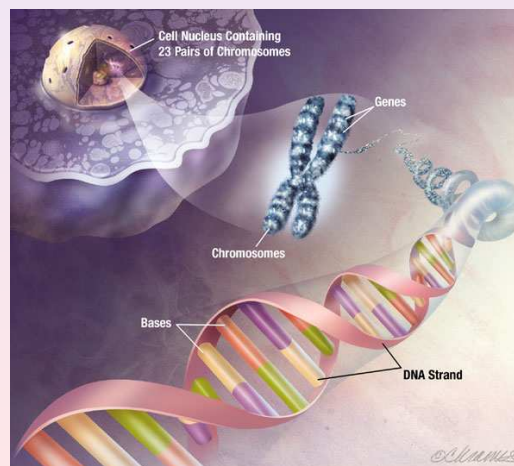
Approximating the Shortest Superstring Problem

Martin Paluszewski
University of Copenhagen

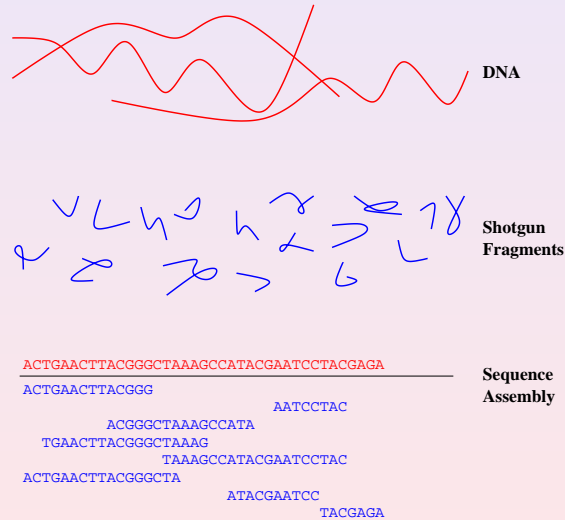
12/2-2008

Sequencing of DNA

- DNA, a string of 4 letters (A, G, C, T)
- Virus DNA 5×10^4 letters (base pairs)
- Human DNA 3×10^9 letters (base pairs)



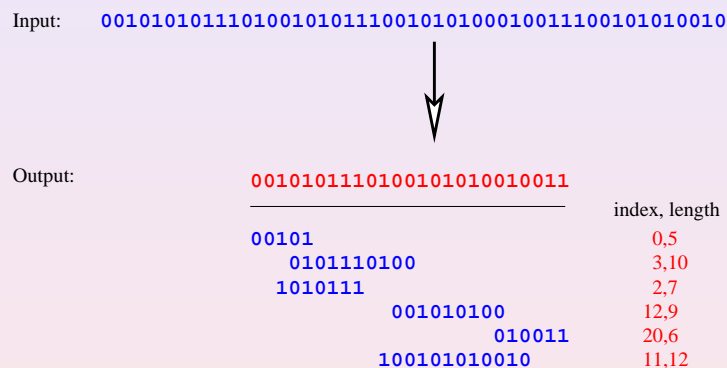
Shotgun Method



Problem: How can we automate this?



Data Compression



Output of compression:

- A shortest superstring
- An ordered list of start index and fragment lengths.



Shortest Superstring Problem (SSP)

- Given a finite alphabet Σ , and set of n strings, $S = \{s_1, \dots, s_n\} \subseteq \Sigma^+$.
- Find a shortest string s that contains each s_i as a substring.
- Without loss of generality, we may assume that no string s_i is a substring of another string s_j , $j \neq i$.

Example

Input (S)

CATGC
CTAAGT
GCTA
TTCA
ATGCATC

$\Sigma = \{A, G, C, T\}$

Output

$s = \text{GCTAAGTTCATGCATC}$
.....CATGC...
.CTAAGT.....
GCTA.....
.....TTCA.....
.....ATGCATC

NP-hard!

Solution Methods

- Exact Algorithms
- Metaheuristics
- Approximation Algorithms

Outline

- 1 Greedy, conjecture: $ALG \leq 2 \cdot OPT$
 - 2 Set cover, $ALG \leq 2 \cdot H_n \cdot OPT$
 $H_n = 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n}$
 - 3 Cycle cover, $ALG \leq 4 \cdot OPT$
 - 4 Cycle cover (greedy), $ALG \leq 3 \cdot OPT$
- Best known: $ALG \leq 2.5 \cdot OPT$ (Sweedyk)

Greedy Algorithm

- 1: **Greedy Shortest Superstring**
- 2: **input:** A set of strings S .
- 3: **output:** A short superstring of S .
- 4: $T \leftarrow S$
- 5: **while** $|T| > 1$ **do**
- 6: Let a and b be the most overlapping strings of T
- 7: Replace a and b with the string obtained by overlapping a and b
- 8: **end while**
- 9: T contains a superstring of S

Example

- $S = T = \{\text{CATGC}, \text{CTAAGT}, \text{GCTA}, \text{TTCA}, \text{ATGCATC}\}$
- $T = \{\text{CATGCATC}, \text{CTAAGT}, \text{GCTA}, \text{TTCA}\}$
- $T = \{\text{CATGCATC}, \text{GCTAAGT}, \text{TTCA}\}$
- $T = \{\text{TTCATGCATC}, \text{GCTAAGT}\}$
- $T = \{\text{GCTAAGTTTCATGCATC}\}$

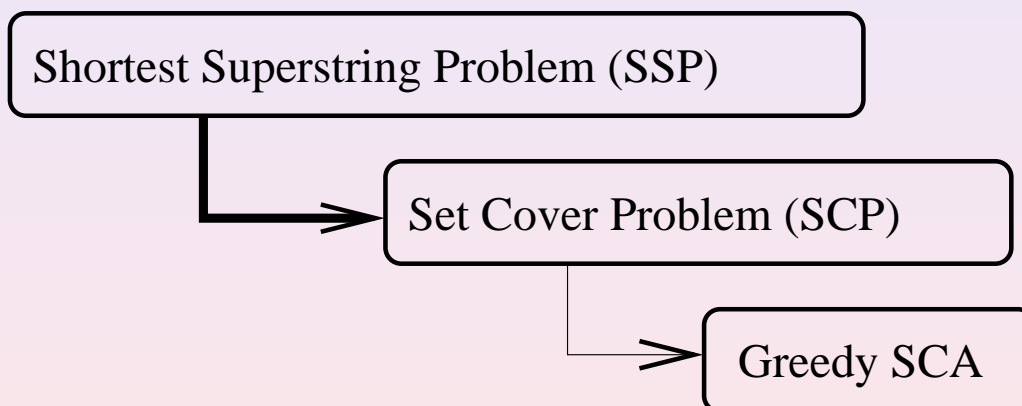
Approximation guarantee

- $ALG \leq 4 \cdot OPT$ (proved by Blum et. al.)
- $ALG \leq 2 \cdot OPT$ (conjectured)

Conjectured worst case

$$S = \{ab^k, b^k c, b^{k+1}\}$$

Approximating SSP Using Set Cover



- **Set Cover Problem:**
Choose some sets that cover all elements with least cost
- **Elements**
- **Subsets**
- **Cost of a subset**

- **Set Cover Problem:**
Choose some sets that cover all elements with least cost
- **Elements**
 - The input strings
- **Subsets**
 - σ_{ijk} = string obtained by overlapping input strings s_i and s_j , k letters.
 - $\beta = S \cup \sigma_{ijk}$, all i,j,k
 - Let $\pi \in \beta$
 - $\text{set}(\pi) = \{s \in S \mid s \text{ is a substr. of } \pi\}$
- **Cost of a subset**
 - $\text{set}(\pi)$ is $|\pi|$
- A solution to SSP is the concatenation of π obtained from SCP

Example

$S = \{CATGC, CTAAGT, GCTA, TTCA, ATGCATC\}$

π	Set	Cost
CATGC..... ...CTAAGT CATGCTAAGT	CATGC, CTAAGT, GCTA	10
CATGC.. ...GCTA CATGCTA	CATGC, GCTA	7
.....CATGC ATGCATC.... ATGCATCATGC	CATGC, ATGCATC	11
CTAAGT..TTCA CTAAGTTCA	CTAAGT, TTCA	9
ATGCATC.....CTAAGT ATGCATCTAAGT	CTAAGT, ATGCATC	12

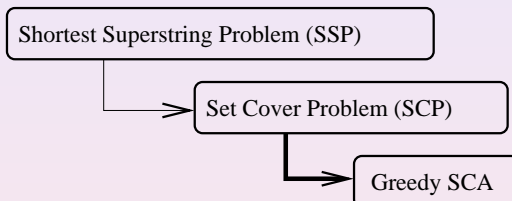


Example

GCTA..... ...ATGCATC GCTATGCATC	GCTA, ATGCATC	10
TTCA..... ...ATGCATC TTCATGCATC	TTCA, ATGCATC, CATGC	10
GCTA.. .CTAAGT GCTAAGT	GCTA, CTAAGT	7
TTCA.. ..CATGC TTCATGC	CATGC, TTCA	7
CATGC.. .ATGCATC CATGCATC	CATGC, ATGCATC	8
CATGC	CATGC	5
CTAAGT	CTAAGT	6
GCTA	GCTA	4
TTCA	TTCA	4
ATGCATC	ATGCATC	7



Approximation



Lemma

$$\text{OPT}_{SSP} \leq \text{OPT}_{SCP} \leq 2 \cdot \text{OPT}_{SSP}$$

Approximation

Proof

First inequality: $\text{OPT}_{SSP} \leq \text{OPT}_{SCP}$

Let s be the string obtained from an optimal solution to the set cover problem. Then,

$$\text{OPT}_{SCP} = |s|$$

Since s covers all strings, s is a valid solution to the SSP and:

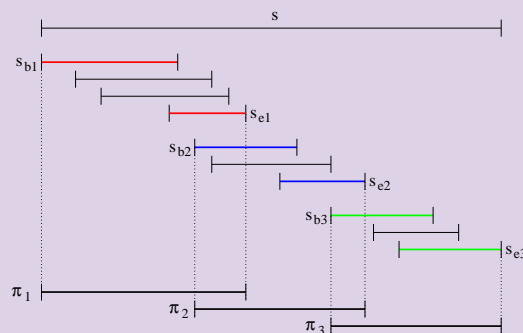
$$|s| \geq \text{OPT}_{SSP}$$

Approximation

Proof, continued

Second inequality: $\text{OPT}_{SCP} \leq 2 \cdot \text{OPT}_{SSP}$

Show that some set cover always can be constructed such that the inequality holds.



Approximation

Proof, continued

- $\text{set}(\pi_1), \text{set}(\pi_2), \dots, \text{set}(\pi_t)$ is a solution to SCP (not necessarily optimal).
- Each input string is covered by at most two π strings. (π_i cannot overlap with π_{i+2}).
- $\sum_i |\pi_i| \leq 2 \cdot \text{OPT}_{SSP}$
- $\text{OPT}_{SCP} \leq 2 \cdot \text{OPT}_{SSP}$ \square

Approximation

Proof, continued

- $\text{GREEDY}_{SCP} \leq H_n \text{OPT}_{SCP}$
 $H_n = 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n}$
 (Vijay V. Vazirani)

$$\text{GREEDY}_{SCP} \leq 2H_n \text{OPT}_{SSP}$$

$$H_{10} \simeq 2.9, H_{100} \simeq 5.2, H_{1000} \simeq 7.5$$

Factor 4 Algorithm (cycle cover)

Some definitions

- **overlap**(s_i, s_j): The maximum overlap between s_i and s_j ,

i	ACGGCTAT....
jCTATTAGC
overlapCTAT....
- **prefix**(s_i, s_j): First letters of s_i where $\text{overlap}(s_i, s_j)$ is removed

i	ACGGCTAT....
jCTATTAGC
prefix	ACGG.....

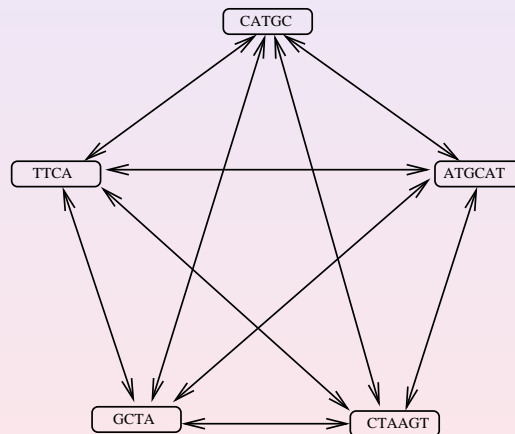
Important property of a shortest superstring s

$$s = \text{prefix}(s_1, s_2) \circ \text{prefix}(s_2, s_3) \circ \dots \circ \text{prefix}(s_n, s_1) \circ \text{overlap}(s_n, s_1)$$

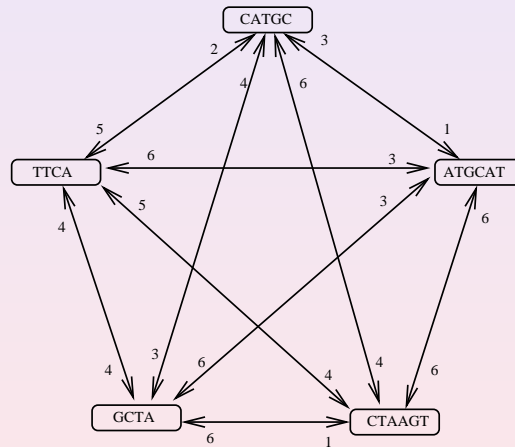
Example

s : TCTAAGTTCATGCATC
 1 TCTA.....
 2 .CTAAGT.....
 3TTCA.....
 4CATGC...
 5ATGCATC
 1TCTA

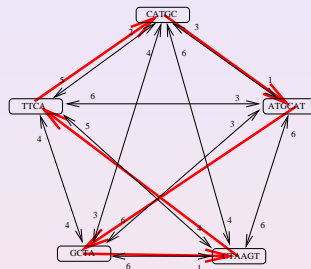
Prefix Graph



Prefix Graph



TSP in Prefix Graph



$$\text{OPT}_{\text{TSP}} = |\text{prefix}(s_a, s_b)| + |\text{prefix}(s_b, s_c)| + \dots + |\text{prefix}(s_n, s_a)|$$

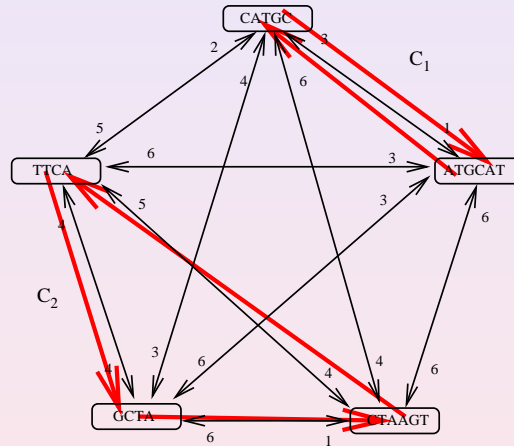
$$\text{OPT}_{\text{SSP}} = |\text{prefix}(s_1, s_2)| + |\text{prefix}(s_2, s_3)| + \dots + |\text{prefix}(s_n, s_1)| + |\text{overlap}(s_n, s_1)|$$

Lower bound:

$$\text{OPT}_{\text{TSP}} \leq \text{OPT}_{\text{SSP}}$$



TSP is NP-hard \rightarrow cycle cover problem

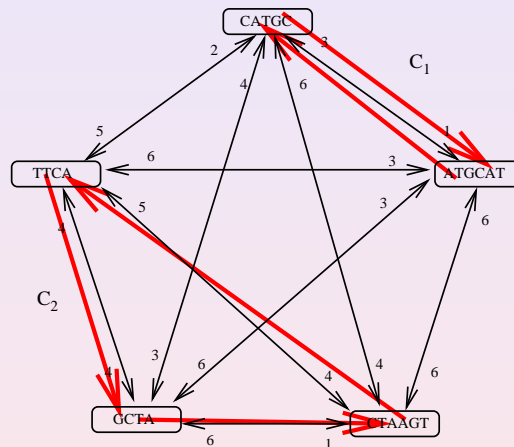


$$\text{OPT}_{\text{CCP}} \leq \text{OPT}_{\text{TSP}} \leq \text{OPT}_{\text{SSP}}$$

Algorithm (factor 4)

- 1 Construct the prefix graph corresponding to strings in S
- 2 Find a minimum weight cycle cover of the prefix graph, $C = \{c_1, \dots, c_k\}$
- 3 For each cycle $c_i \in C$, arbitrarily choose a representative string s_{i_1}
- 4 For cycle $c_i \in C$, construct:
 $\alpha(c_i) = \text{prefix}(s_{i_1}, s_{i_2}) \circ \dots \circ \text{prefix}(s_{i_{l-1}}, s_{i_l}) \circ \text{prefix}(s_{i_l}, s_{i_1})$
- 5 For cycle $c_i \in C$, construct:
 $\sigma(c_i) = \alpha(c_i) \circ s_{i_1}$
- 6 Output $\sigma(c_1) \circ \dots \circ \sigma(c_k)$

Example



- $\sigma(c_1) = C \text{ ATG CATGC}$
- $\sigma(c_2) = G \text{ CTAAG TTCA GCTA}$
- $s = C \text{ ATG CATGC G CTAAG TTCA GCTA}$

$$\text{ALG}_{CC} \leq 4 \cdot \text{OPT}_{SSP}$$

Proof (sketch):

$$s = \sigma(c_1) \circ \dots \circ \sigma(c_k) \quad (1)$$

$$s = \alpha(c_1) \circ s_{1_1} \circ \dots \circ \alpha(c_k) \circ s_{k_1} \quad (2)$$

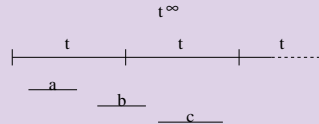
$$|s| = |\alpha(c_1) \circ \dots \circ \alpha(c_k)| + |s_{1_1} \circ \dots \circ s_{k_1}| \quad (3)$$

- Since $\text{OPT}_{CCP} = |\alpha(c)|$ the α strings are bounded by OPT_{SSP} .
- $\text{OPT}_{CCP} \leq \text{OPT}_{TSP} \leq \text{OPT}_{SSP}$
- We need to show that concatenation of representative strings is bounded by $3 \cdot \text{OPT}_{SSP}$

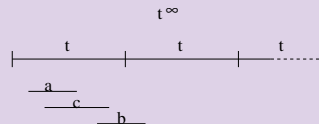
Lemma 1

If each string in $S' \subseteq S$ is a substring of t^∞ for a string t , then there is a cycle of weight at most $|t|$ in the prefix graph covering all vertices corresponding to strings in S'

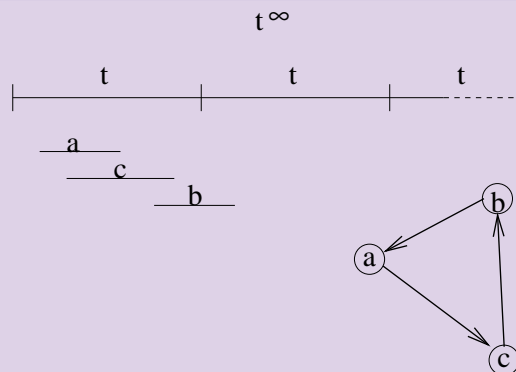
Proof (sketch)



- Sort strings



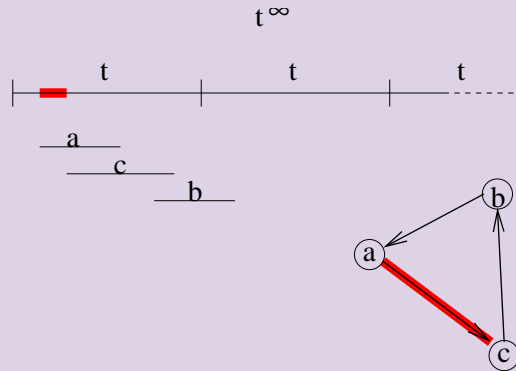
Proof (sketch), continued



$$\text{Cost}(a \rightarrow c \rightarrow b \rightarrow a) \leq |t|$$

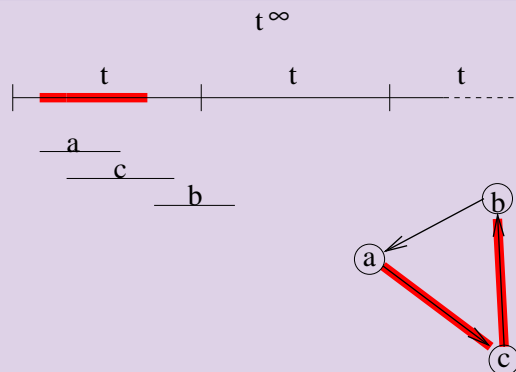


Proof (sketch), continued



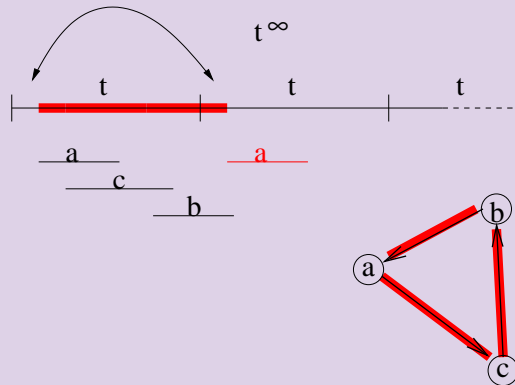
$$\text{Cost}(a \rightarrow c \rightarrow b \rightarrow a) \leq |t|$$

Proof (sketch), continued



$$\text{Cost}(a \rightarrow c \rightarrow b \rightarrow a) \leq |t|$$

Proof (sketch), continued

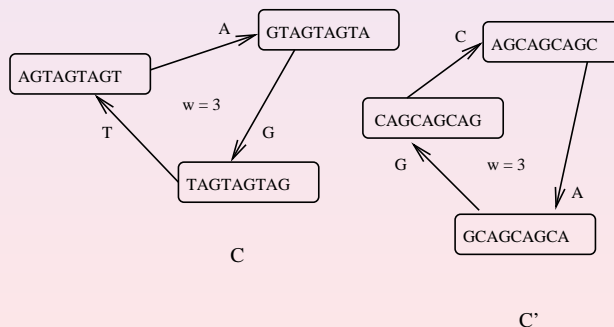


$$\text{Cost}(a \rightarrow c \rightarrow b \rightarrow a) \leq |t|$$

Lemma 2

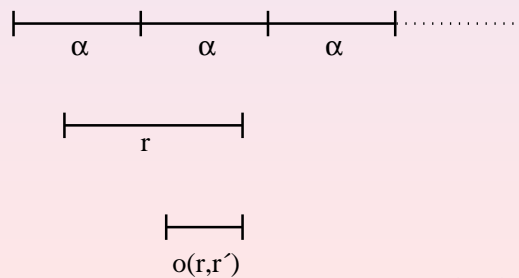
Let C be an optimal cycle cover. Let c and c' be two cycles in C , and let r, r' be strings from each cycle. Then

$$|\text{overlap}(r, r')| < w(c) + w(c')$$



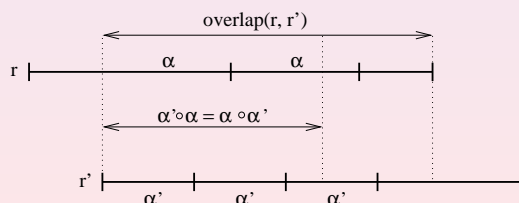
Proof

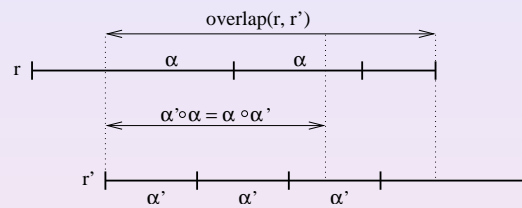
- Suppose $|\text{overlap}(r, r')| \geq w(c) + w(c')$
 $\alpha(c_i) = \text{prefix}(r, s_{i_2}) \circ \dots \circ \text{prefix}(s_{i_{l-1}}, s_{i_l}) \circ \text{prefix}(s_{i_l}, r)$
- r is a substring of $\alpha(c)^\infty$, so $\text{overlap}(r, r')$ is also a substring of $\alpha(c)^\infty$.
- $|\text{overlap}(r, r')| \geq w(c)$, so the overlap consists of repeating strings $\alpha(c)$.



Proof

- Suppose $|\text{overlap}(r, r')| \geq w(c) + w(c')$
 $\alpha(c_i) = \text{prefix}(r, s_{i_2}) \circ \dots \circ \text{prefix}(s_{i_{l-1}}, s_{i_l}) \circ \text{prefix}(s_{i_l}, r)$
- r is a substring of $\alpha(c)^\infty$, so $\text{overlap}(r, r')$ is also a substring of $\alpha(c)^\infty$.
- $|\text{overlap}(r, r')| \geq w(c)$, so the overlap consists of repeating strings $\alpha(c)$.





Proof, continued

α and α' commute, so

$$\alpha^k \circ \alpha'^k = \alpha'^k \circ \alpha^k$$

for all $k > 0$, so

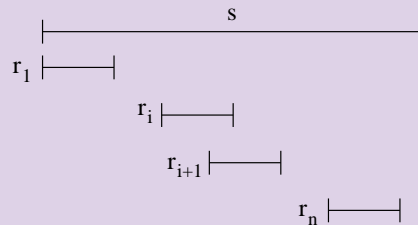
$$\alpha^\infty = \alpha'^\infty$$

Proof, continued

- All strings from c and c' are substring of α^∞ .
- There is a cycle covering all strings in c and c' with weight at most $|\alpha|$ (lemma 1)
- Contradiction.

Proof of factor 4 algorithm

- r_i are sorted representative strings



$$\text{OPT}_{SSP} \geq \sum_{i=1}^k |r_i| - \sum_{i=1}^{k-1} |\text{overlap}(r_i, r_{i+1})|$$

Proof of factor 4 algorithm, continued.

$$\text{OPT}_{SSP} \geq \sum_{i=1}^k |r_i| - \sum_{i=1}^{k-1} |\text{overlap}(r_i, r_{i+1})|$$

Lemma 2:

$$|\text{overlap}(r, r')| < w(c) + w(c')$$

$$\text{OPT}_{SSP} \geq \sum_{i=1}^k |r_i| - 2 \sum_{i=1}^k w(c_i)$$

$$\sum_{i=1}^k |r_i| \leq \text{OPT}_{SSP} + 2 \sum_{i=1}^k w(c_i) \leq 3 \text{OPT}_{SSP}$$

$$\text{ALG} = \sum_{i=1}^k |\sigma(c_i)| = w(C) + \sum_{i=1}^k |r_i| \leq 4 \cdot \text{OPT}_{SSP}$$

Factor 3 Algorithm

- Construct the prefix graph corresponding to strings in S
- Find a minimum weight cycle cover of the prefix graph,
 $C = \{c_1, \dots, c_k\}$
- Run the greedy superstring algorithm on $\{\sigma(c_1), \dots, \sigma(c_k)\}$
and output the resulting string τ

Factor 3 Algorithm

Compression

- $\|S\| = \sum_{x \in S} |x|$
- Compression = $\|S\| - |s|$
- Max. compression \leftrightarrow min. superstring
- Lemma: Greedy superstring algorithm achieves at least half the optimal compression:

$$\|S\| - \text{GREEDY} \geq \frac{1}{2}(\|S\| - \text{OPT})$$

Lemma 1

$$|\tau| \leq \text{OPT}_\sigma + w(C)$$

Proof

- Assume $\sigma(c_1), \dots, \sigma(c_k)$ appear in this order in the superstring of S_σ
- Maximum compression is:

$$\|S_\sigma\| - \text{OPT}_\sigma = \sum_{i=1}^{k-1} |\text{overlap}(\sigma(c_i), \sigma(c_{i+1}))|$$

Proof, continued

$$\|S_\sigma\| - \text{OPT}_\sigma = \sum_{i=1}^{k-1} |\text{overlap}(\sigma(c_i), \sigma(c_{i+1}))|$$

- $|\text{overlap}(r, r')| < w(c) + w(c')$ (lemma 2)
- So maximum compression is:
 $\|S_\sigma\| - \text{OPT}_\sigma \leq 2w(C)$
- Greedy algorithm gives:
 $\|S_\sigma\| - |\tau| \geq \frac{1}{2}(\|S_\sigma\| - \text{OPT}_\sigma)$
- $2(|\tau| - \text{OPT}_\sigma) \leq \|S_\sigma\| - \text{OPT}_\sigma \leq 2w(C)$

Lemma 2

$$\text{OPT}_\sigma \leq \text{OPT}_{SSP} + w(C)$$

Proof

- Let $S_r = \{r_1, \dots, r_k\}$ (repr. strings)
- Each $\sigma(c_i)$ begins and ends with r_i , so
 $\|S_\sigma\| - \text{OPT}_\sigma \geq \|S_r\| - \text{OPT}_r$
- $\|S_\sigma\| = \|S_r\| + w(C)$
- $\text{OPT}_\sigma \leq \text{OPT}_r + w(C)$

Lemma 1

$$|\tau| \leq \text{OPT}_\sigma + w(C)$$

Lemma 2

$$\text{OPT}_\sigma \leq \text{OPT}_{SSP} + w(C)$$

Combine lemma 1 and lemma 2 to get:

$$|\tau| \leq 3 \cdot \text{OPT}_{SSP}$$