

## The Shortest Superstring Problem: The Overlap Lemma

Last updated Nov 30, 2007

Overlap Lemma: Let  $c$  and  $c'$  be cycles in  $\mathcal{C}$  and  $r, r'$  their respective representative strings. Then

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

where  $wt(c)$  is the cost of cycle  $c$ .

Cycle Lemma:

If every string in  $S' \subseteq S$  is a substring of  $t^\infty$  for some string  $t$

$\Rightarrow$  there is a cycle of weight at most  $|t|$  in the prefix graph covering all of the vertices (corresponding to strings) in  $S$ .

Proof:

Let  $t = t_1 t_2 \dots t_k$ .

Suppose  $S' = \{s_1, \dots, s_m\}$ .

Now let  $j_i$  be the starting point of the first occurrence of  $s_i$  in  $t^\infty$ . This must be in the first copy of  $t$  (why).

Note that all of the  $j_i$  are different from each other since no substring in  $S$  is a substring of any other.

Now sort the strings by their starting points and consider the cycle  $c$  that visits the vertices corresponding to the strings in the sorted order. This cycle has length at most  $t$  so we are done.

GCD Lemma: Let  $X$  be a prefix of both  $\alpha^\infty$  and  $\beta^\infty$  with  $|X| \geq |\alpha| + |\beta|$ . Then

1. If  $|\alpha| = |\beta|$  then  $\alpha = \beta$ .
2. If  $|\alpha| > |\beta|$  then  $X$  is a prefix of  $\gamma^\infty$  where  $\gamma = X[1]X[2] \dots X[|\alpha| - |\beta|]$ .

Proof: (i) is obvious. To prove (ii) set  $p = |\alpha|$ ,  $q = |\beta|$ . By definition,  $\forall, 0 < i \leq q$  **and**  $0 < j \leq p$ ,

$$X[i + p] = X[i] \quad \text{and} \quad X[j + q] = X[j]$$

We now show that  $\forall i, 0 < i \leq |X| - (p - q)$ ,  
 $X[i + (p - q)] = X[i]$ .

First assume that  $0 < i \leq q$ . Then

$$\begin{aligned} X[i + (p - q)] &= X[i + (p - q) + q] \\ &= X[i + p] = X[i] \end{aligned}$$

Now assume that  $q < i \leq |X| - (p - q)$ . Then

$$\begin{aligned} X[i + (p - q)] &= X[i + (p - q) - p] \\ &= X[i - q] = X[i] \end{aligned}$$

Corollary: Let  $X$  be a prefix of both  $\alpha^\infty$  and  $\beta^\infty$  with  $|X| \geq |\alpha| + |\beta|$ . Then  $X$  is a prefix of  $\gamma^\infty$  where  $\gamma = X[1]X[2] \dots X[\gcd(|\alpha|, |\beta|)]$ . Thus

$$\gamma^\infty = \alpha^\infty = \beta^\infty.$$

Overlap Lemma: Let  $c$  and  $c'$  be cycles in  $\mathcal{C}$  and  $r, r'$  their respective representative strings. Then

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

where  $\text{wt}(c)$  is the cost of cycle  $c$ .

Proof: Assume the contrary, that

$$\text{overlap}(r, r') \geq \text{wt}(c) + \text{wt}(c')$$

Let  $\alpha$  be the prefix of length  $\text{wt}(c)$  of  $\text{overlap}(r, r')$  and  $\alpha'$  the prefix of length  $\text{wt}(c')$  of  $\text{overlap}(r, r')$ . Notice that

- 1. Every string “in”  $c$  is a substring of  $\alpha^\infty$ .**
- 2. Every string “in”  $c'$  is a substring of  $(\alpha')^\infty$ .**
- 3.  $\text{overlap}(r, r')$  is a prefix of both  $\alpha^\infty$  and  $(\alpha')^\infty$ .**

From the GCD Lemma we know that the string  $\gamma$  containing the first  $\text{gcd}(\text{wt}(c), \text{wt}(c'))$  characters of  $\text{overlap}(r, r')$  satisfies

$$\gamma^\infty = \alpha^\infty = (\alpha')^\infty.$$

We just saw that

$$\gamma^\infty = \alpha^\infty = (\alpha')^\infty$$

so  $\gamma^\infty$  contains every string in  $c$  and every string in  $c'$ .  
Furthermore, by construction,

$$|\gamma| = \gcd(\text{wt}(c), \text{wt}(c'))$$

so, from the Cycle Lemma, we therefore have that there is a cycle of weight at most  $\gcd(\text{wt}(c), \text{wt}(c'))$  covering all strings in  $c$  and  $c'$ .

This contradicts the minimality of  $\mathcal{C}$ . Thus

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