

String Matching with Variable Length Gaps

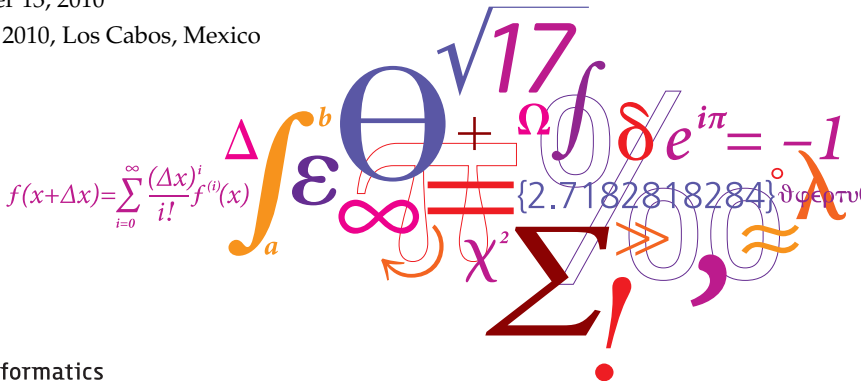


By Philip Bille, Inge Li Gørtz, Hjalte Wedel Vildhøj and David Kofoed Wind

Presented by Hjalte Wedel Vildhøj

October 13, 2010

SPIRE 2010, Los Cabos, Mexico



DTU Informatics

Department of Informatics and Mathematical Modelling

The Variable Length Gap Problem

Given some string $T \in \Sigma^+$ and a *variable length gap pattern*

$$P = P_1 \cdot g\{a_1, b_1\} \cdot P_2 \cdot g\{a_2, b_2\} \cdots g\{a_{k-1}, b_{k-1}\} \cdot P_k .$$

Find the *end positions* for all occurrences of P in T .

The Variable Length Gap Problem

Given some string $T \in \Sigma^+$ and a *variable length gap pattern*

$$P = P_1 \cdot \underbrace{g\{a_1, b_1\}} \cdot P_2 \cdot g\{a_2, b_2\} \cdots g\{a_{k-1}, b_{k-1}\} \cdot P_k .$$

Some $x \in \Sigma^*$ s.t. $a_1 \leq |x| \leq b_1$

Find the *end positions* for all occurrences of P in T .

The Variable Length Gap Problem

Given some string $T \in \Sigma^+$ and a *variable length gap pattern*

$$P = P_1 \cdot \underbrace{g\{a_1, b_1\}} \cdot P_2 \cdot g\{a_2, b_2\} \cdots g\{a_{k-1}, b_{k-1}\} \cdot P_k .$$

Some $x \in \Sigma^*$ s.t. $a_1 \leq |x| \leq b_1$

Find the *end positions* for all occurrences of P in T .

Example: $P = A \cdot g\{6, 7\} \cdot CC \cdot g\{2, 6\} \cdot GT$

$T = \text{ATCGGCTCCAGACCAGTACCCGTTCCGTGGT}$

Solution: $\{\}$

The Variable Length Gap Problem

Given some string $T \in \Sigma^+$ and a *variable length gap pattern*

$$P = P_1 \cdot \underbrace{g\{a_1, b_1\}} \cdot P_2 \cdot g\{a_2, b_2\} \cdots g\{a_{k-1}, b_{k-1}\} \cdot P_k .$$

Some $x \in \Sigma^*$ s.t. $a_1 \leq |x| \leq b_1$

Find the *end positions* for all occurrences of P in T .

Example: $P = A \cdot g\{6, 7\} \cdot CC \cdot g\{2, 6\} \cdot GT$

$T =$ A TCGGCT CC AGACCA GT ACCCGTTCCGTGGT

Solution: $\{17\}$

The Variable Length Gap Problem

Given some string $T \in \Sigma^+$ and a *variable length gap pattern*

$$P = P_1 \cdot \underbrace{g\{a_1, b_1\}} \cdot P_2 \cdot g\{a_2, b_2\} \cdots g\{a_{k-1}, b_{k-1}\} \cdot P_k .$$

Some $x \in \Sigma^*$ s.t. $a_1 \leq |x| \leq b_1$

Find the *end positions* for all occurrences of P in T .

Example: $P = A \cdot g\{6, 7\} \cdot CC \cdot g\{2, 6\} \cdot GT$

$T =$ ATCGGCTCCAGACCAGTACC CGTTCCGTGGT

Solution: $\{17\}$

Not a valid match!

The Variable Length Gap Problem

Given some string $T \in \Sigma^+$ and a *variable length gap pattern*

$$P = P_1 \cdot \underbrace{g\{a_1, b_1\}} \cdot P_2 \cdot g\{a_2, b_2\} \cdots g\{a_{k-1}, b_{k-1}\} \cdot P_k .$$

Some $x \in \Sigma^*$ s.t. $a_1 \leq |x| \leq b_1$

Find the *end positions* for all occurrences of P in T .

Example: $P = A \cdot g\{6, 7\} \cdot CC \cdot g\{2, 6\} \cdot GT$

$T =$ ATCGGCTCCAG **A** $\overbrace{\text{CCAGTA}}^6$ **CC** $\overbrace{\text{CGTTCC}}^6$ **GT**GGT

Solution: $\{17, 28\}$

end pos in T

The Variable Length Gap Problem

Given some string $T \in \Sigma^+$ and a *variable length gap pattern*

$$P = P_1 \cdot \underbrace{g\{a_1, b_1\}} \cdot P_2 \cdot g\{a_2, b_2\} \cdots g\{a_{k-1}, b_{k-1}\} \cdot P_k .$$

Some $x \in \Sigma^*$ s.t. $a_1 \leq |x| \leq b_1$

Find the *end positions* for all occurrences of P in T .

Example: $P = A \cdot g\{6, 7\} \cdot CC \cdot g\{2, 6\} \cdot GT$

$T =$ ATCGGCTCCAG **A** $\overbrace{\text{CCAGTAC}}^7$ **CC** $\overbrace{\text{GTTCC}}^5$ **GT**GGT

Solution: $\{17, 28\}$

end pos in T

The Variable Length Gap Problem

Given some string $T \in \Sigma^+$ and a *variable length gap pattern*

$$P = P_1 \cdot \underbrace{g\{a_1, b_1\}} \cdot P_2 \cdot g\{a_2, b_2\} \cdots g\{a_{k-1}, b_{k-1}\} \cdot P_k .$$

Some $x \in \Sigma^*$ s.t. $a_1 \leq |x| \leq b_1$

Find the *end positions* for all occurrences of P in T .

Example: $P = A \cdot g\{6, 7\} \cdot CC \cdot g\{2, 6\} \cdot GT$

$T = \text{ATCGGCTCCAGACCAGT} \text{A} \overbrace{\text{CCCGTT}}^7 \text{CC} \overbrace{\text{GTG}}^3 \text{GT}$

Solution: $\{17, 28, 31\}$

end pos in T

The Variable Length Gap Problem

Given some string $T \in \Sigma^+$ and a *variable length gap pattern*

$$P = P_1 \cdot \underbrace{g\{a_1, b_1\}} \cdot P_2 \cdot g\{a_2, b_2\} \cdots g\{a_{k-1}, b_{k-1}\} \cdot P_k .$$

Some $x \in \Sigma^*$ s.t. $a_1 \leq |x| \leq b_1$

Find the *end positions* for all occurrences of P in T .

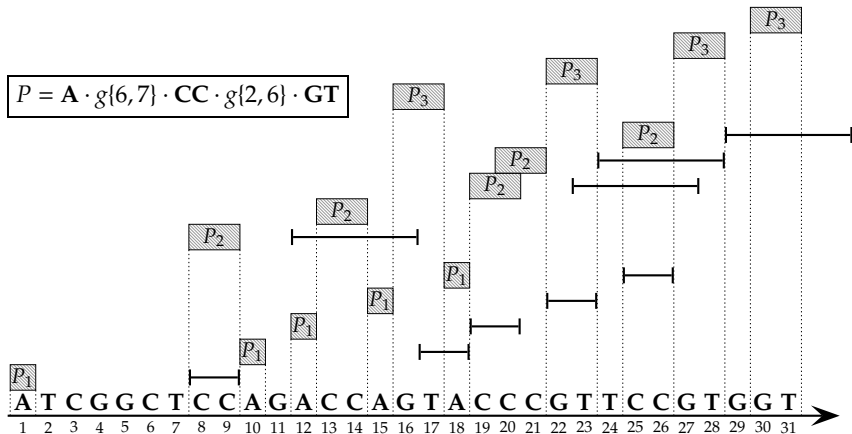
Example: $P = A \cdot g\{6, 7\} \cdot CC \cdot g\{2, 6\} \cdot GT$

$T = \text{ATCGGCTCCAGACCAGTACCCGTTCCGTGGT}$

Solution: $\{17, 28, 31\}$

A Closer Look At The Problem

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



A Closer Look At The Problem

Parameters

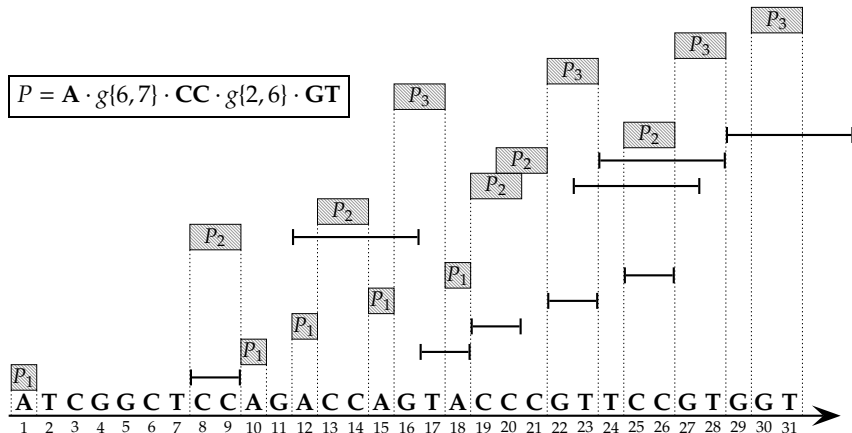
$$n = |T|$$

$\alpha = \# \text{ occ. of } P_1, P_2, \dots, P_k \text{ in } T$

$$m = \sum_{i=1}^k |P_i|$$

$$A = \sum_{i=1}^k a_i$$

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



A Closer Look At The Problem

Parameters

$$n = |T|$$

$$\alpha = \# \text{ occ. of } P_1, P_2, \dots, P_k \text{ in } T$$

$$m = \sum_{i=1}^k |P_i|$$

$$A = \sum_{i=1}^k a_i$$

Known Upper Bounds

<i>By</i>	<i>Time</i>	<i>Space</i>
Bille & Thorup ¹	$O\left(n\left(k\frac{\log w}{w} + \log k\right) + m \log m + A\right)$	$O(m + A)$
Morgante et al. ²	$O((n + m) \log k + \alpha)$	$O(m + \alpha)$

¹P. Bille and M. Thorup. Regular expression matching with multi-strings and intervals. In *Proc. 21st SODA, 2010*

²M. Morgante, A. Policriti, N. Vitacolonna, and A. Zuccolo. Structured motifs search. *J. Comput. Bio.*, 12(8):1065-1082, 2005

A Closer Look At The Problem

Parameters

$$n = |T|$$

$$\alpha = \# \text{ occ. of } P_1, P_2, \dots, P_k \text{ in } T$$

$$m = \sum_{i=1}^k |P_i|$$

$$A = \sum_{i=1}^k a_i$$

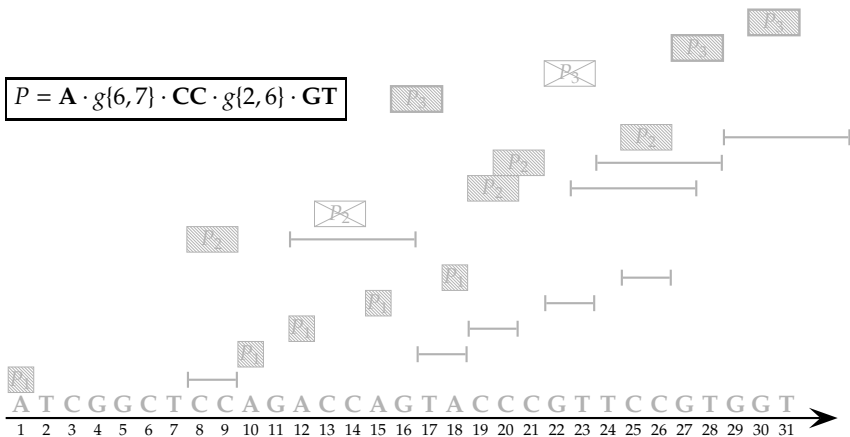
Known Upper Bounds

<i>By</i>	<i>Time</i>	<i>Space</i>
Bille & Thorup	$O\left(n\left(k\frac{\log w}{w} + \log k\right) + m \log m + A\right)$	$O(m + A)$
Morgante et al.	$O((n + m) \log k + \alpha)$	$O(m + \alpha)$

Can you get the best of both?

Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$

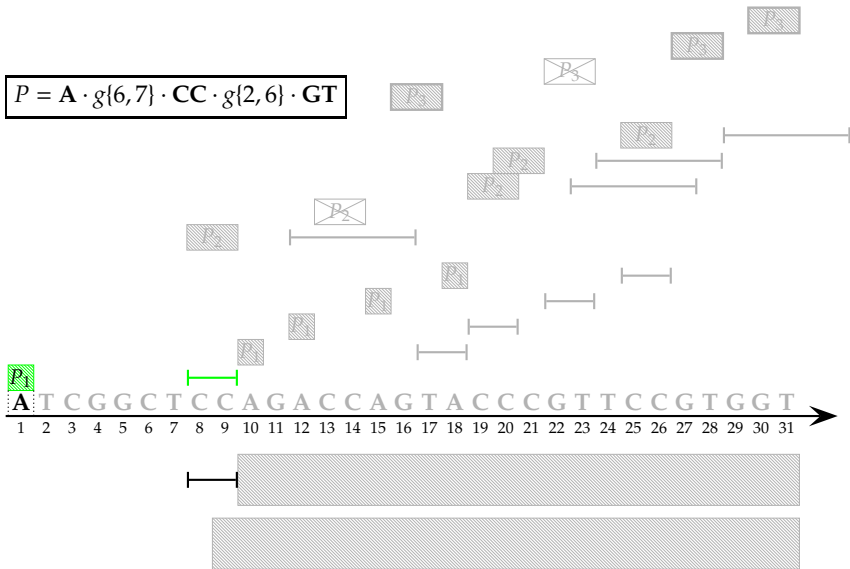


L_2

L_3

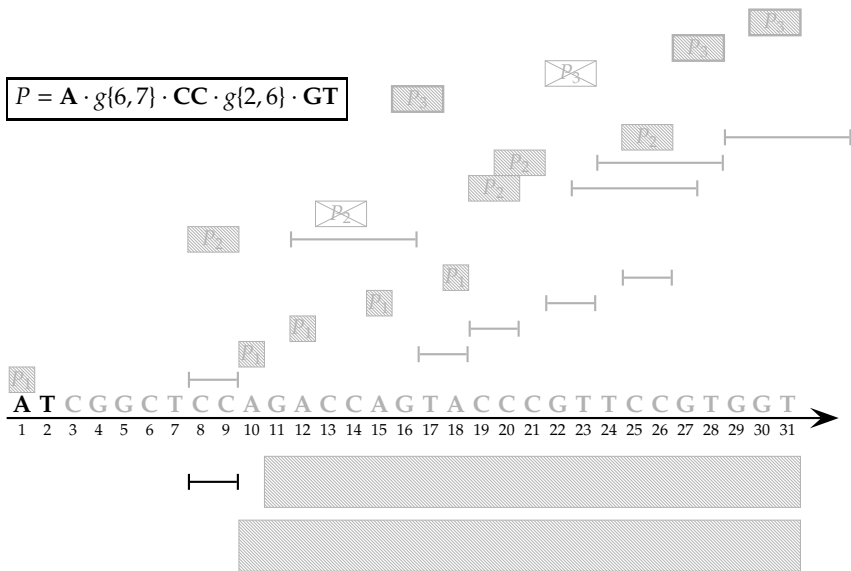
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



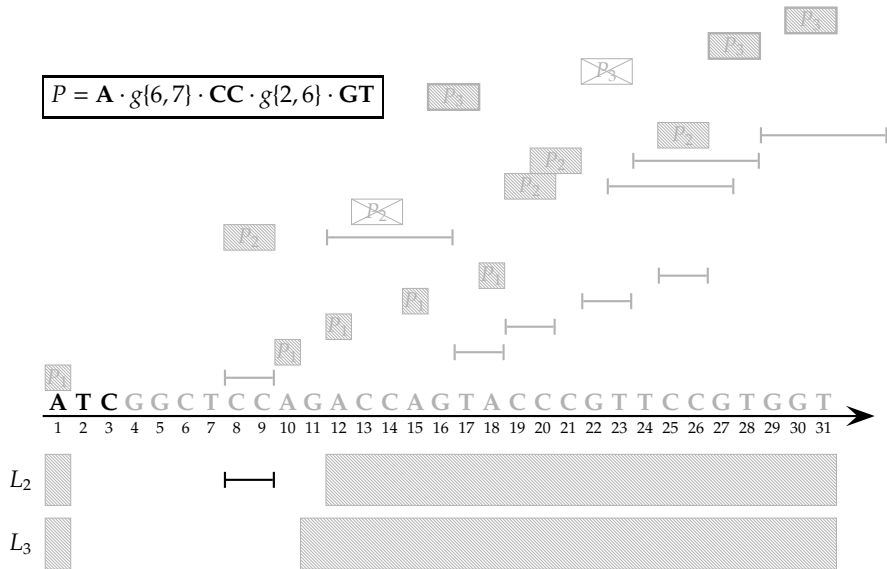
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



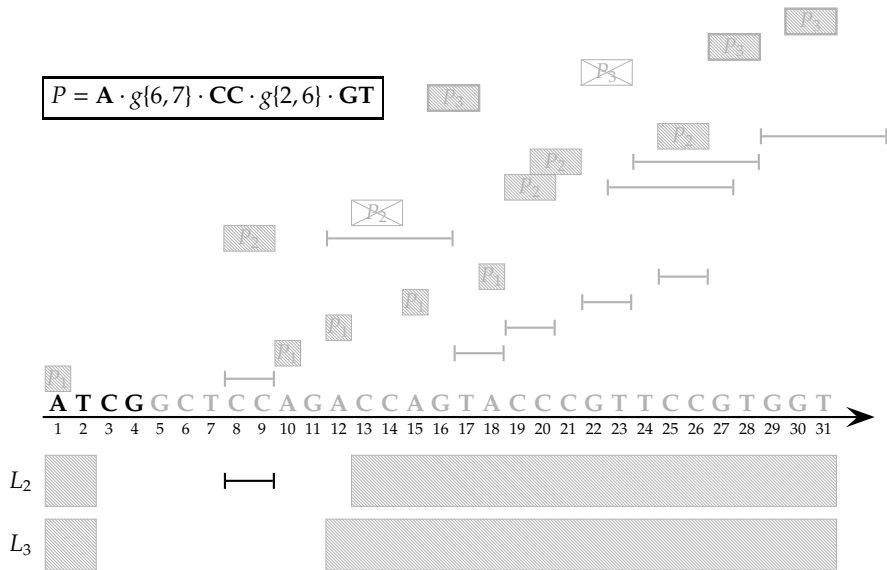
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



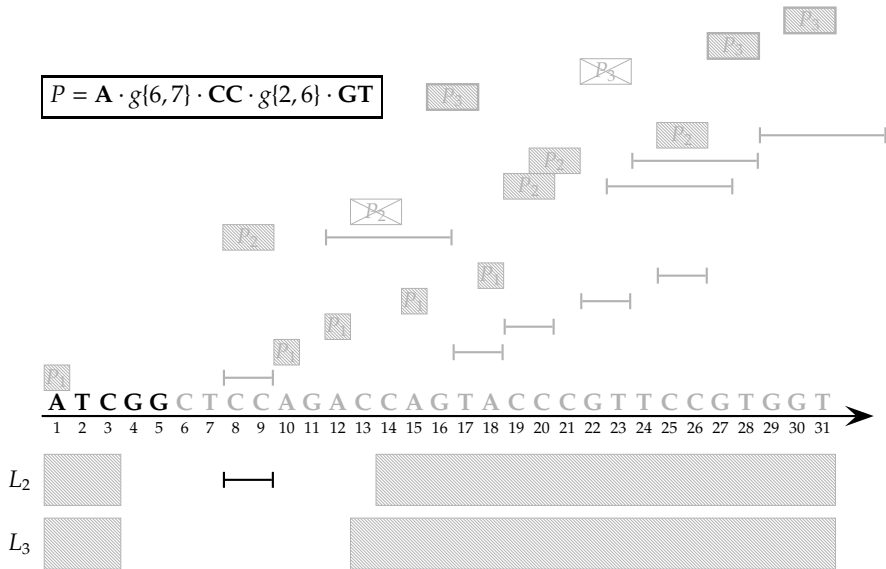
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



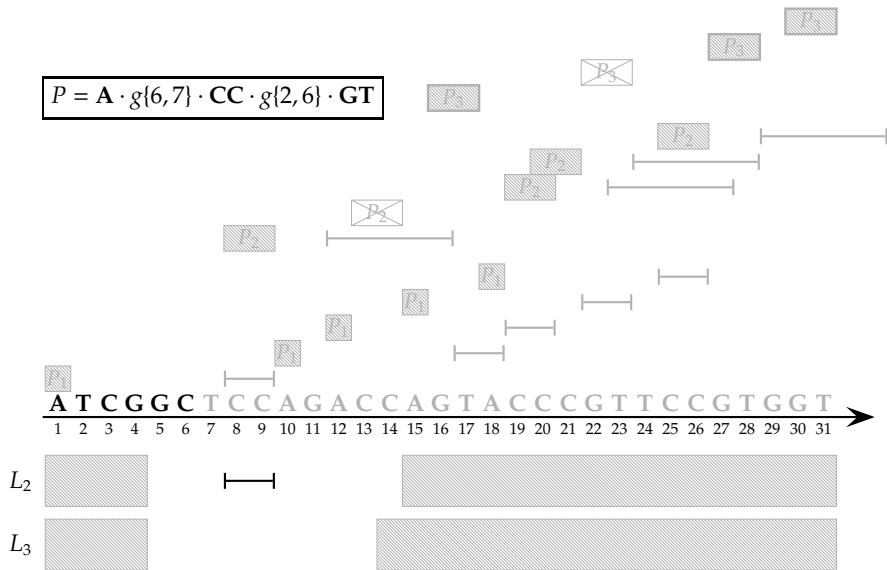
Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



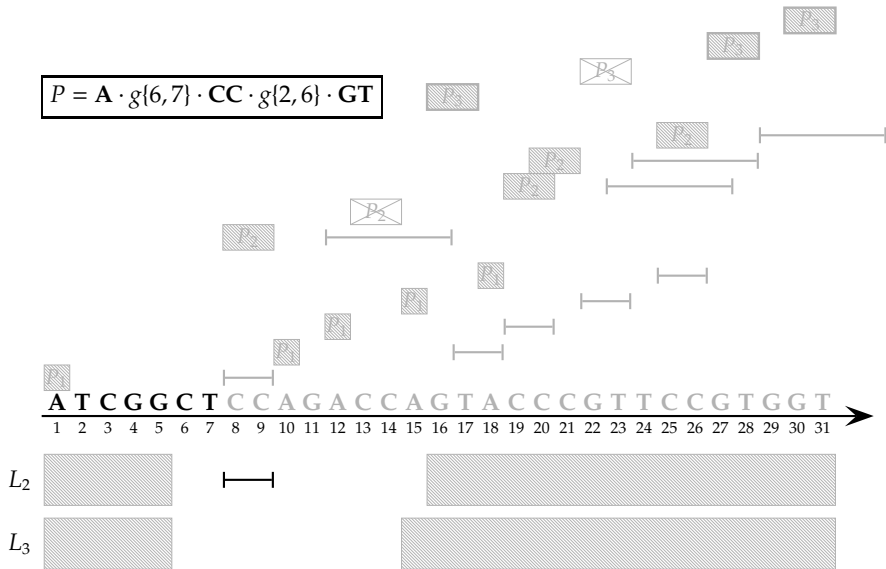
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



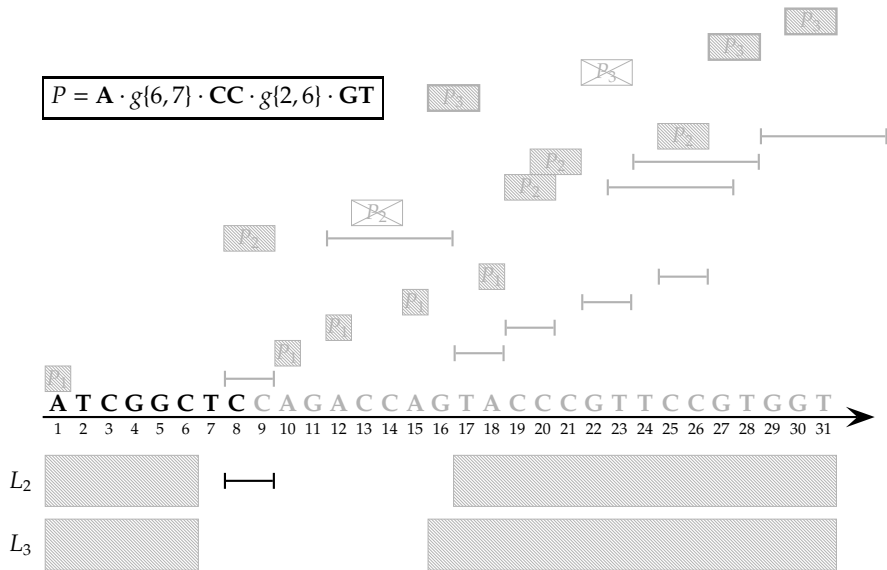
Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



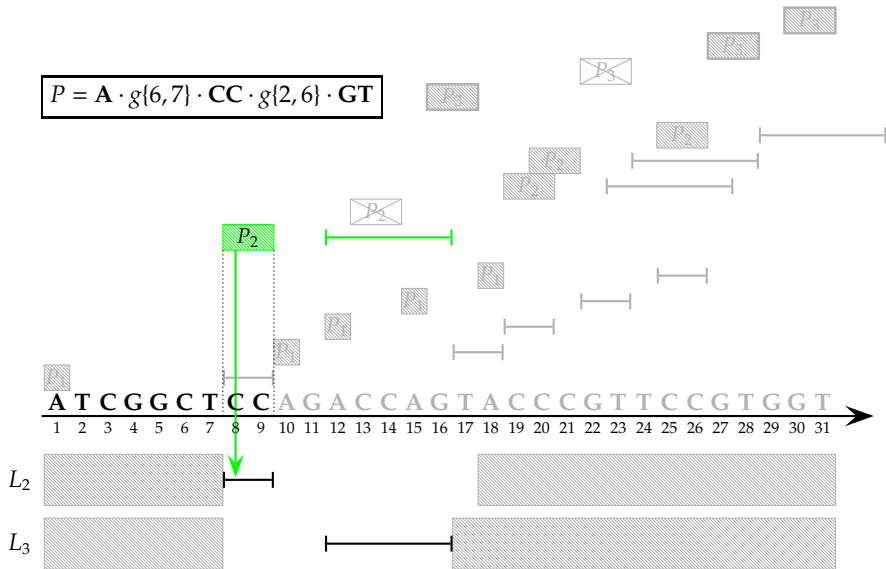
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



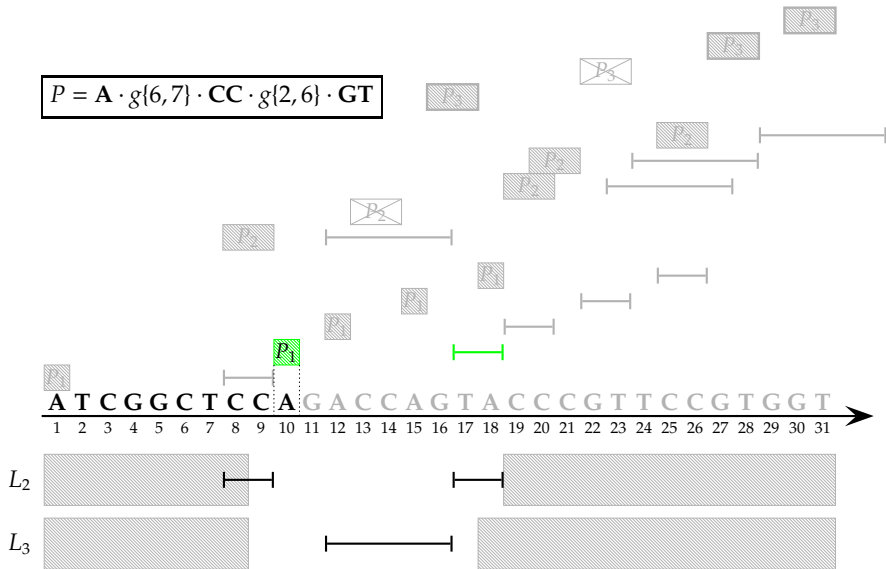
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



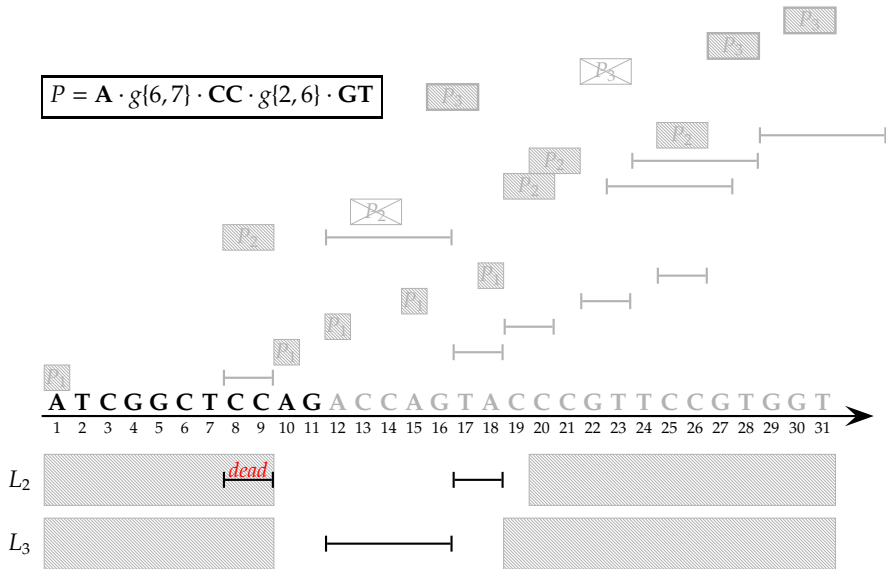
Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



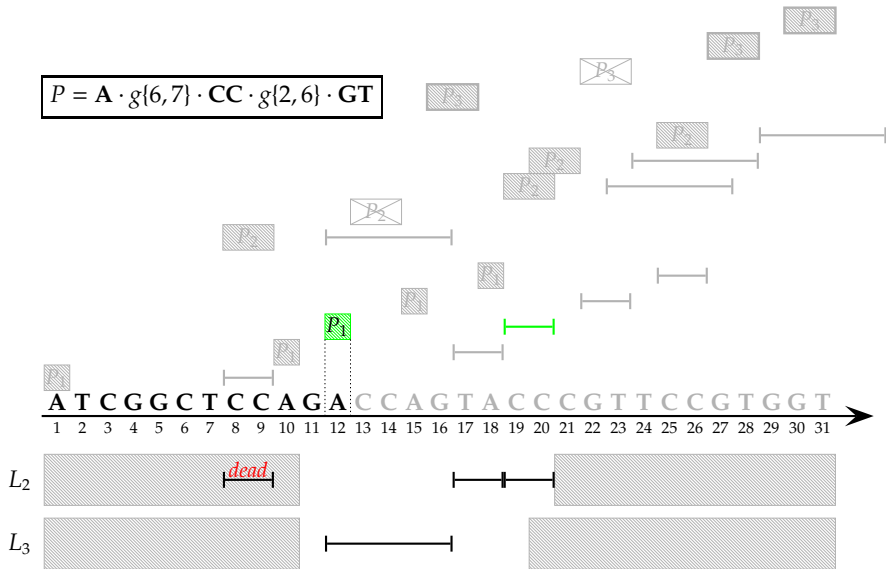
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



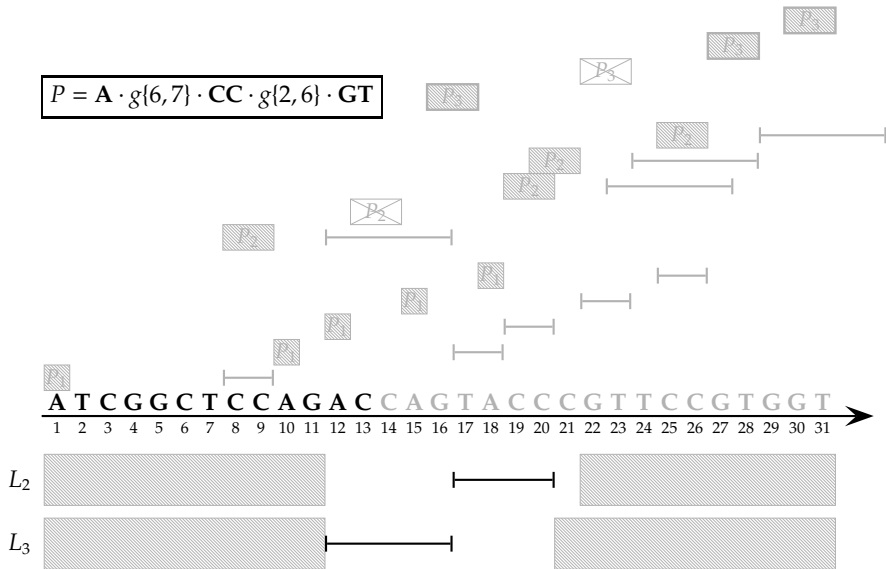
Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



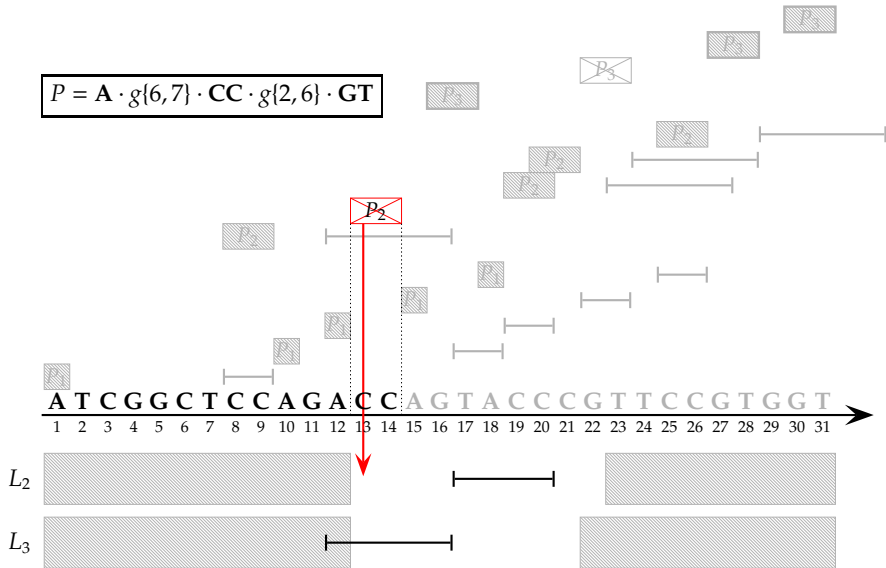
Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



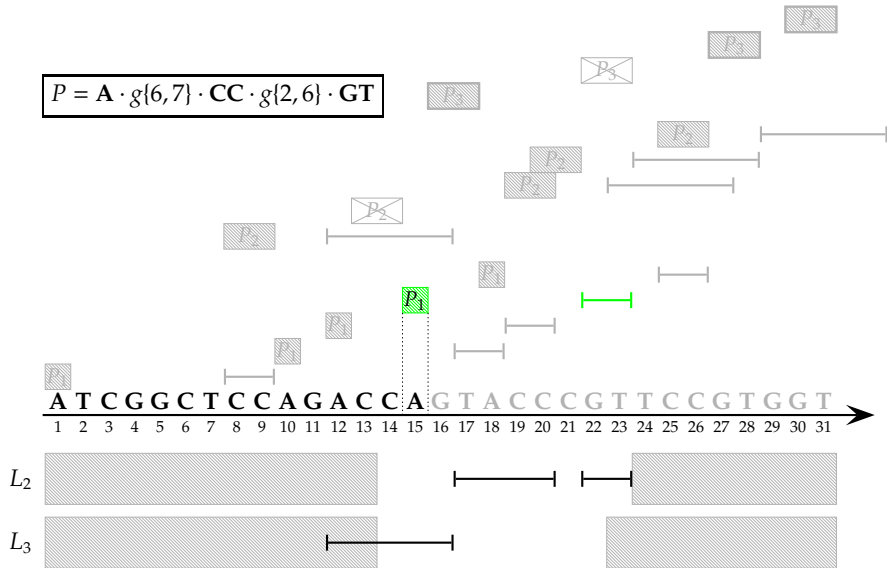
Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



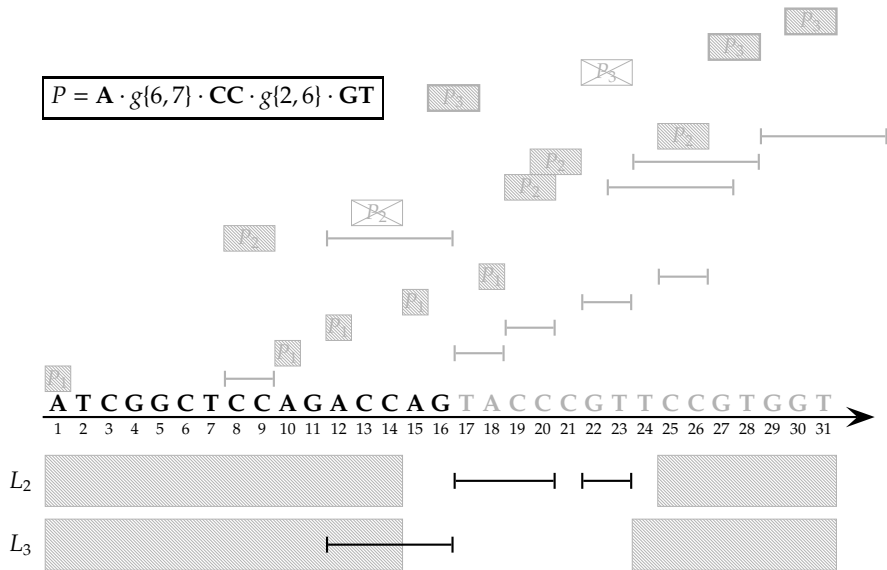
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



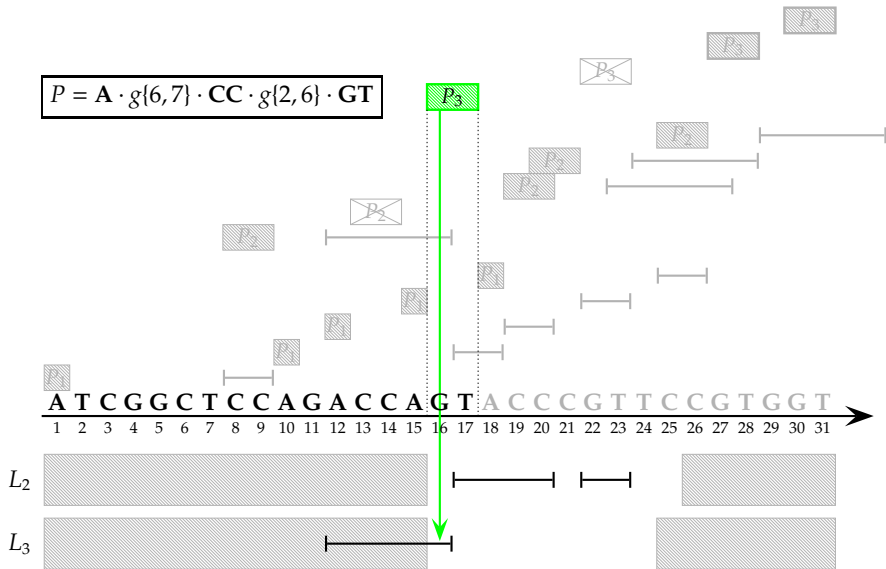
Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$



Illustrating the Algorithm

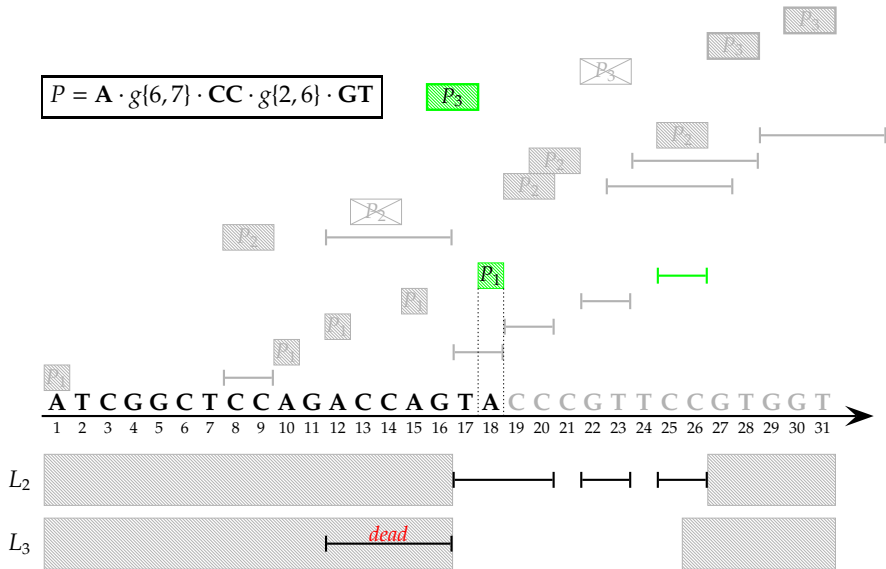
$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

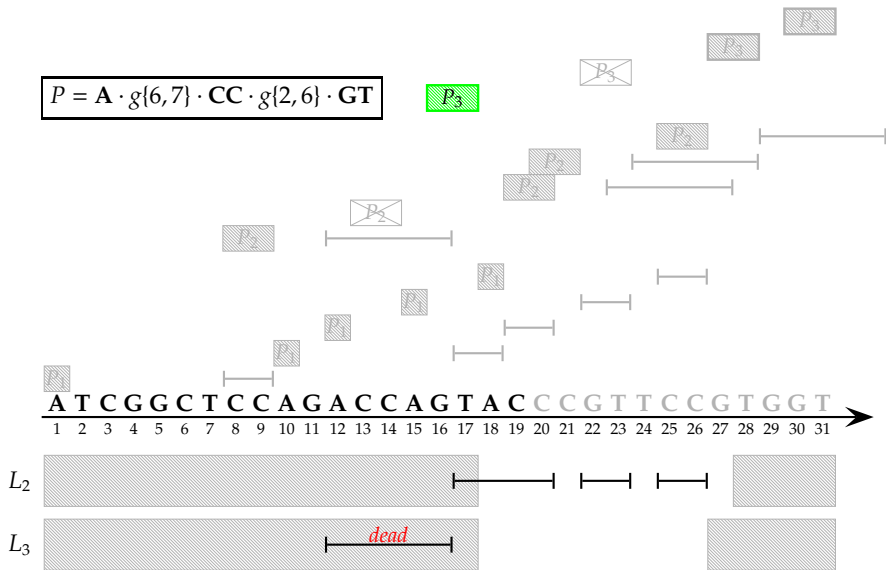
P_3



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

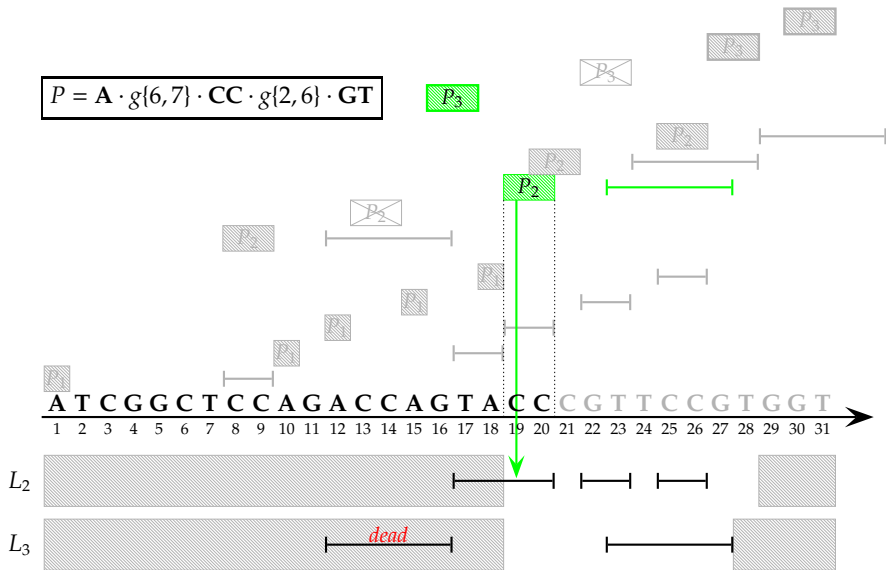
P_3



Illustrating the Algorithm

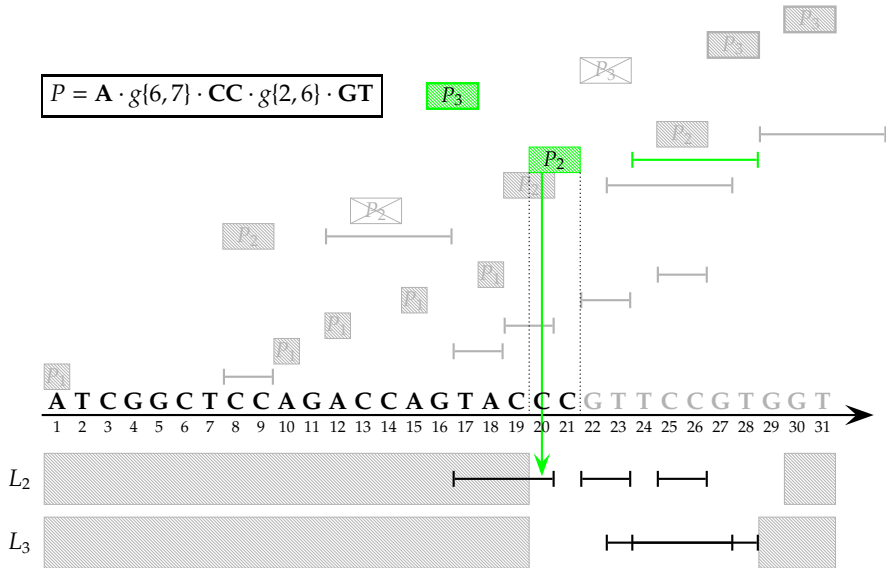
$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

P_3



Illustrating the Algorithm

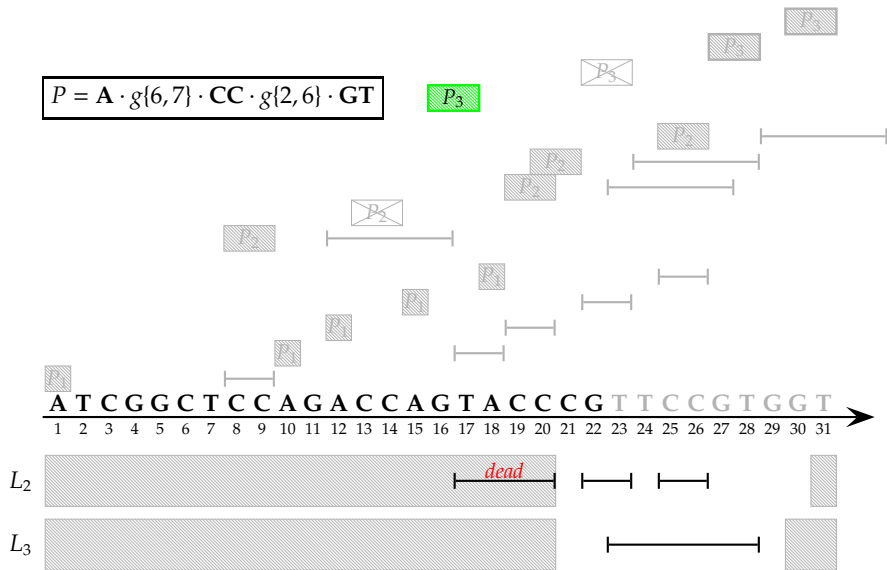
$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$

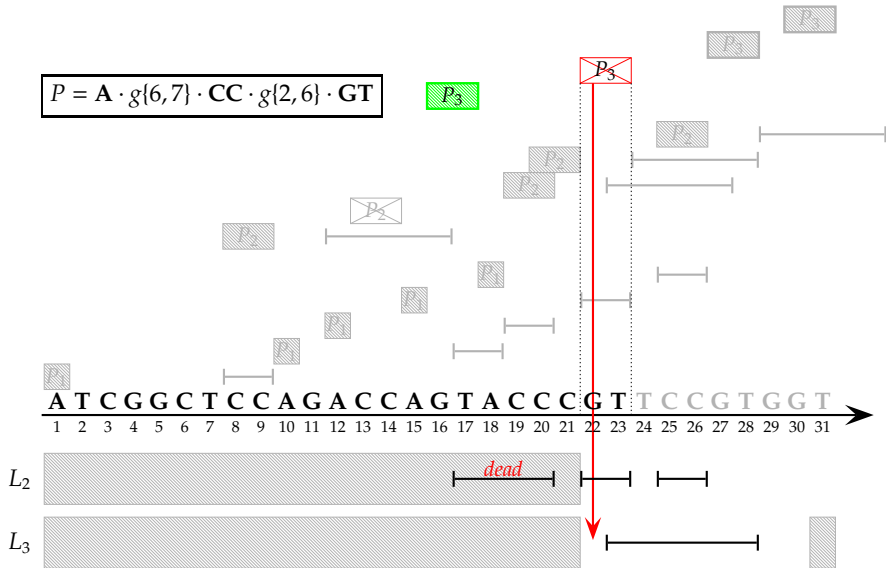
P_3



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

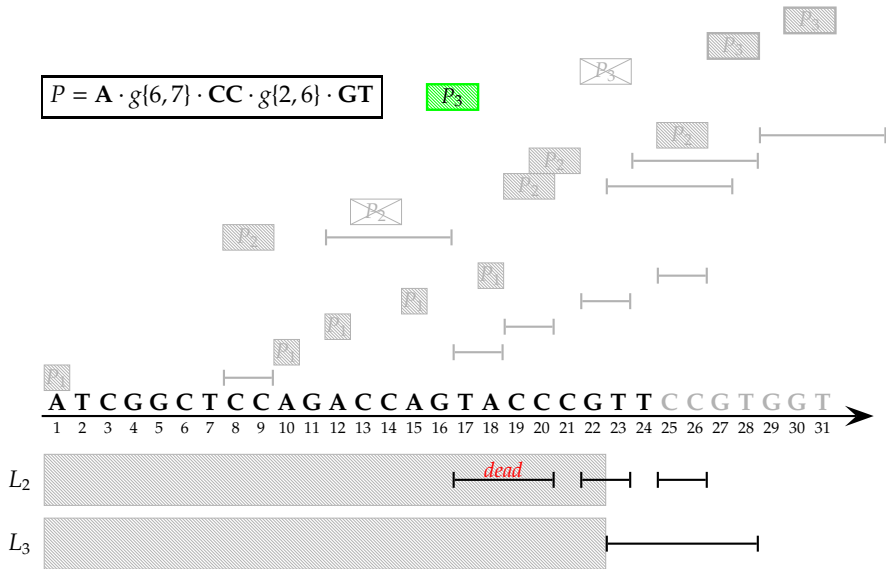
P_3



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

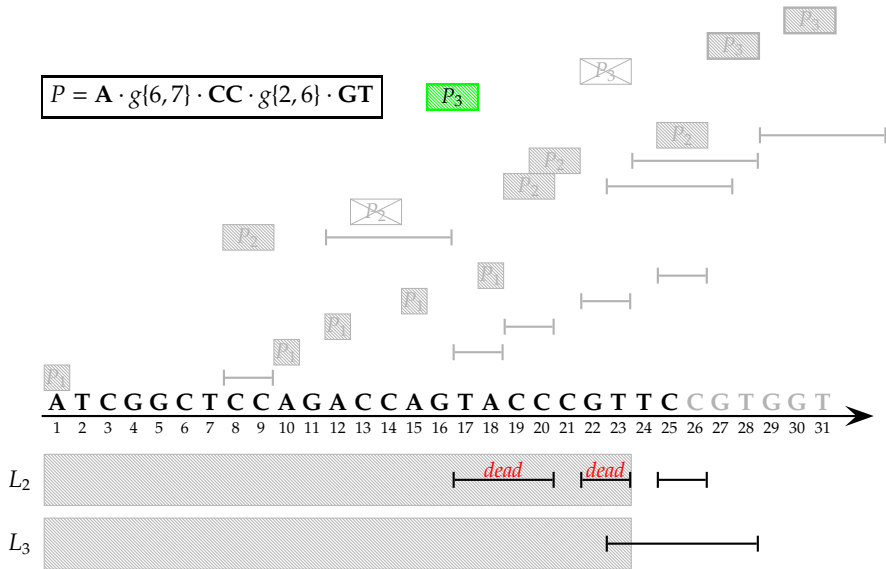
P_3



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

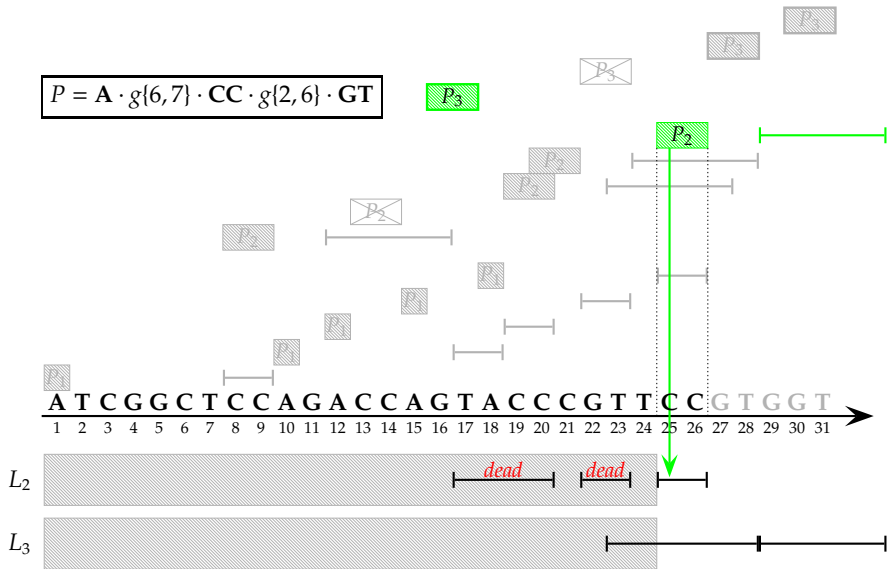
P_3



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

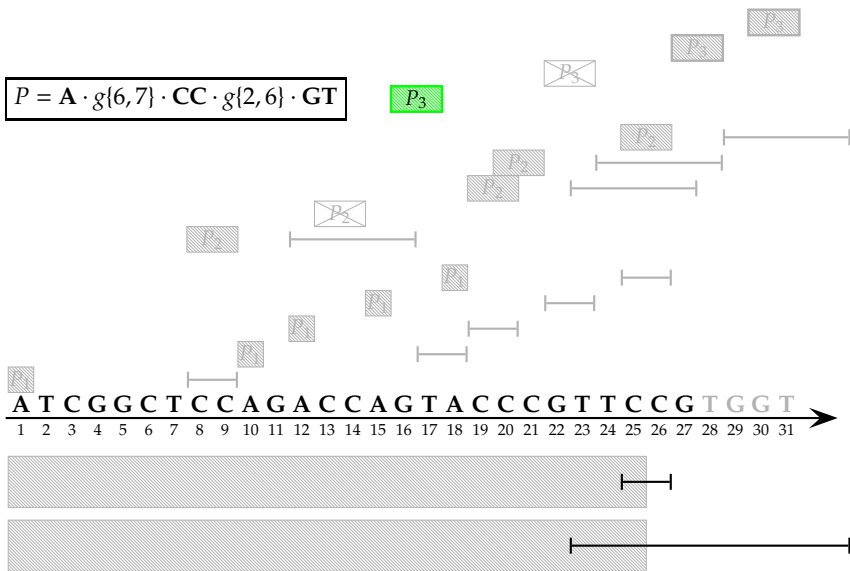
P_3



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

P_3



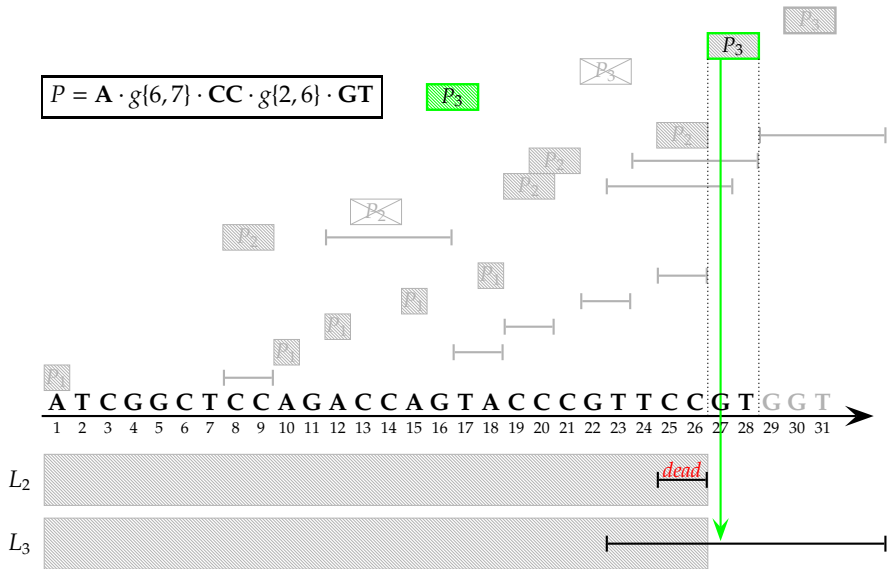
L_2

L_3

Illustrating the Algorithm

$$P = \mathbf{A} \cdot g\{6,7\} \cdot \mathbf{CC} \cdot g\{2,6\} \cdot \mathbf{GT}$$

P_3



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

P_3

P_3

P_3

P_3

P_2

P_2

P_2

P_2

P_1

P_1

P_1

P_1

P_1

A T C G G C T C C A G A C C A G T A C C C G T T C C G T G G T

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31

L_2

dead

L_3



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$

P_3

P_3

P_3

P_3

P_2

P_2

P_2

P_2

P_2

P_1

P_1

P_1

P_1

P_1

A T C G G C T C C A G A C C A G T A C C C G T T C C G T G G T

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31

L_2

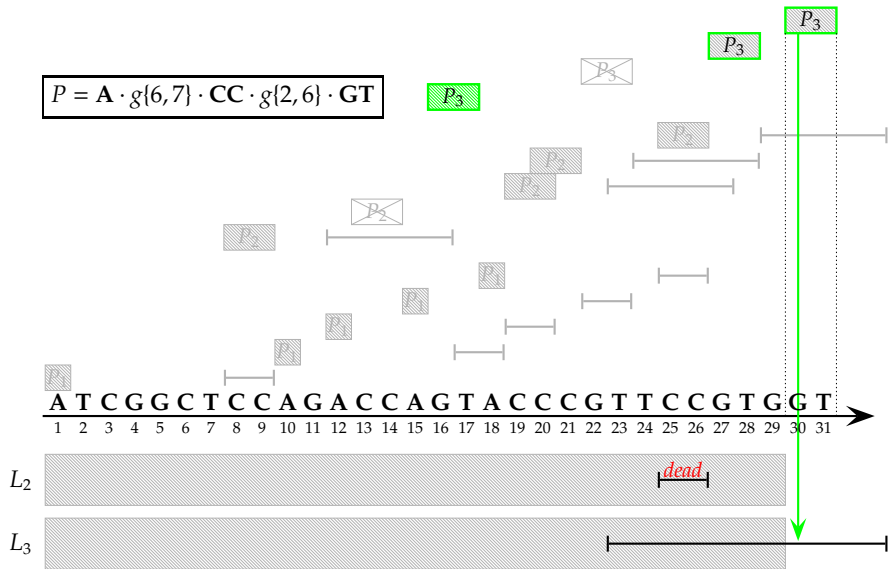
dead

L_3



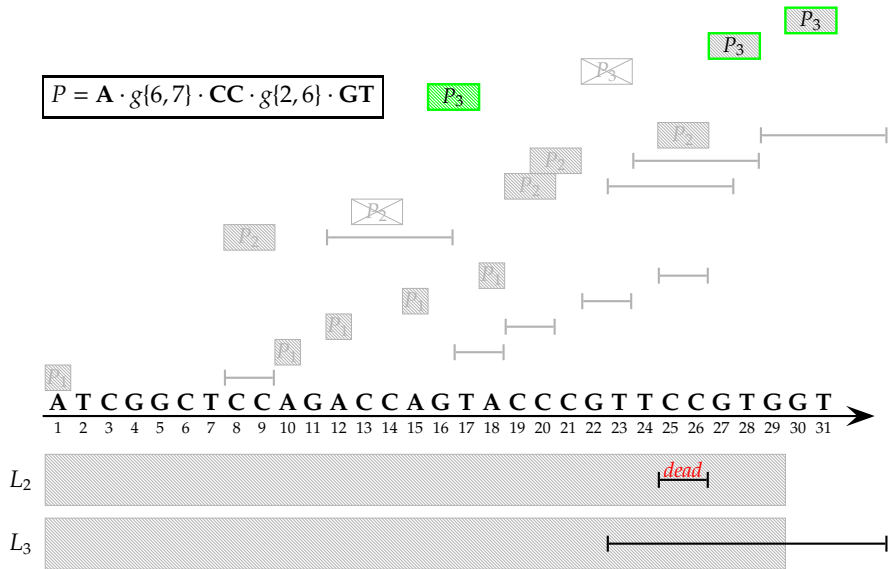
Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



Illustrating the Algorithm

$$P = A \cdot g\{6,7\} \cdot CC \cdot g\{2,6\} \cdot GT$$



Time and Space

Claim: The algorithm runs in $O((n + m) \log k + \alpha)$ time and uses $O(m + A)$ space.

Time and Space

Claim: The algorithm runs in $O((n + m) \log k + \alpha)$ time and uses $O(m + A)$ space.

Time

- ▶ Processing T using AC automaton takes $O((n + m) \log k + \alpha)$ time.
- ▶ At most α ranges are added and removed, so $O(\alpha)$ extra time is spent maintaining the lists.

Time and Space

Claim: The algorithm runs in $O((n + m) \log k + \alpha)$ time and uses $O(m + A)$ space.

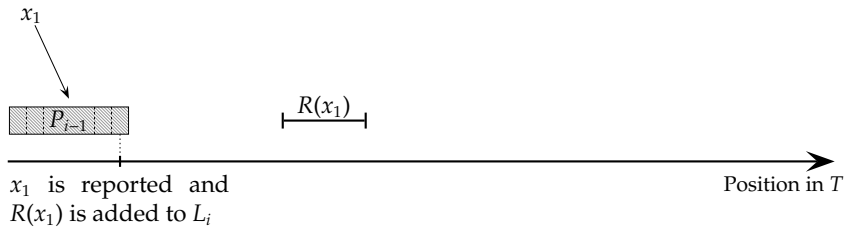
Time

- ▶ Processing T using AC automaton takes $O((n + m) \log k + \alpha)$ time.
- ▶ At most α ranges are added and removed, so $O(\alpha)$ extra time is spent maintaining the lists.

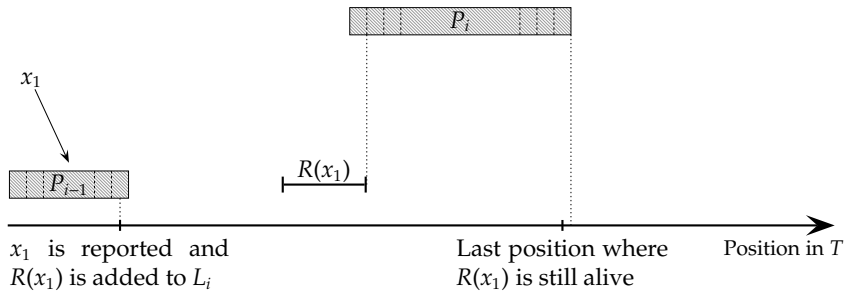
Space

- ▶ AC automaton takes $O(m)$ space.
- ▶ How much space is used by L_2, \dots, L_k ?

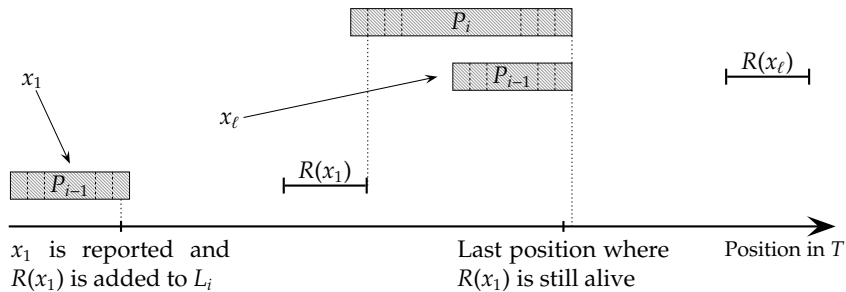
Maximum Size of L_i



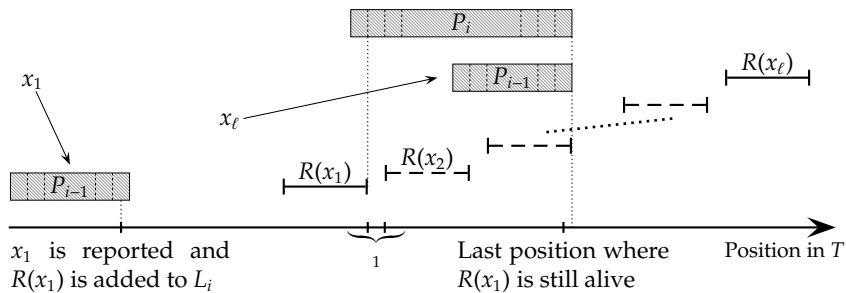
Maximum Size of L_i



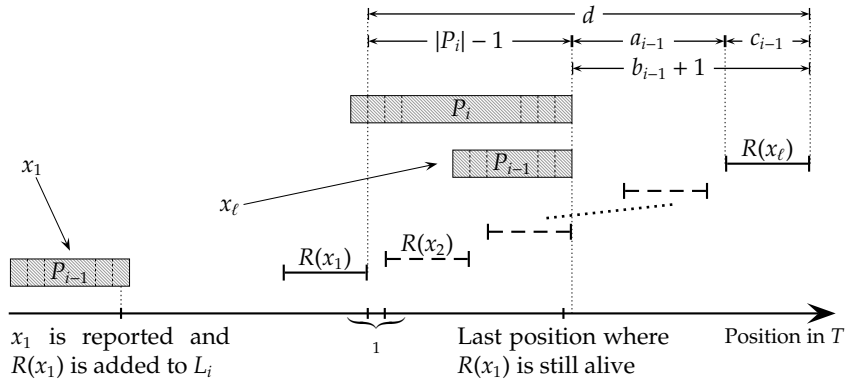
Maximum Size of L_i



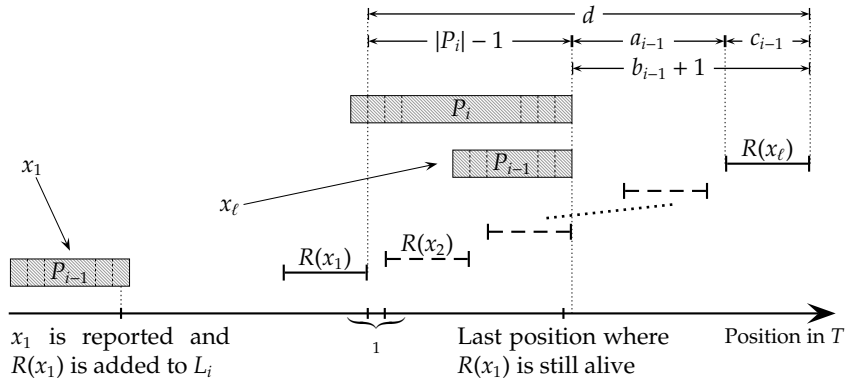
Maximum Size of L_i



Maximum Size of L_i

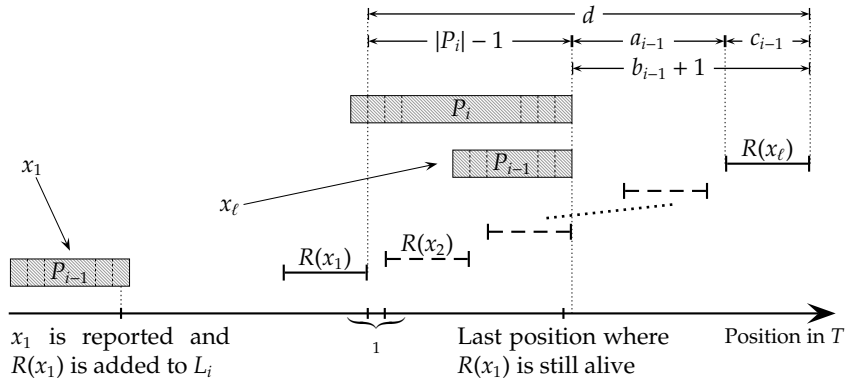


Maximum Size of L_i



$$|L_i| \leq \left\lfloor \frac{d}{c_{i-1} + 1} \right\rfloor + 1 = \left\lfloor \frac{2c_{i-1} + |P_i| + a_{i-1}}{c_{i-1} + 1} \right\rfloor = O(|P_i| + a_{i-1}).$$

Maximum Size of L_i



$$|L_i| \leq \left\lfloor \frac{d}{c_{i-1} + 1} \right\rfloor + 1 = \left\lfloor \frac{2c_{i-1} + |P_i| + a_{i-1}}{c_{i-1} + 1} \right\rfloor = O(|P_i| + a_{i-1}).$$

Total space: $\sum_{i=2}^k |L_i| = O\left(\sum_{i=2}^k |P_i| + \sum_{i=1}^{k-1} a_i\right) = O(m + A)$