



Composite repetition-aware data structures

Djamal Belazzougui¹, **Fabio Cunial**², Travis Gagie¹, Nicola Prezza³, Mathieu Raffinot⁴

(1) Department of Computer Science, University of Helsinki, Finland.

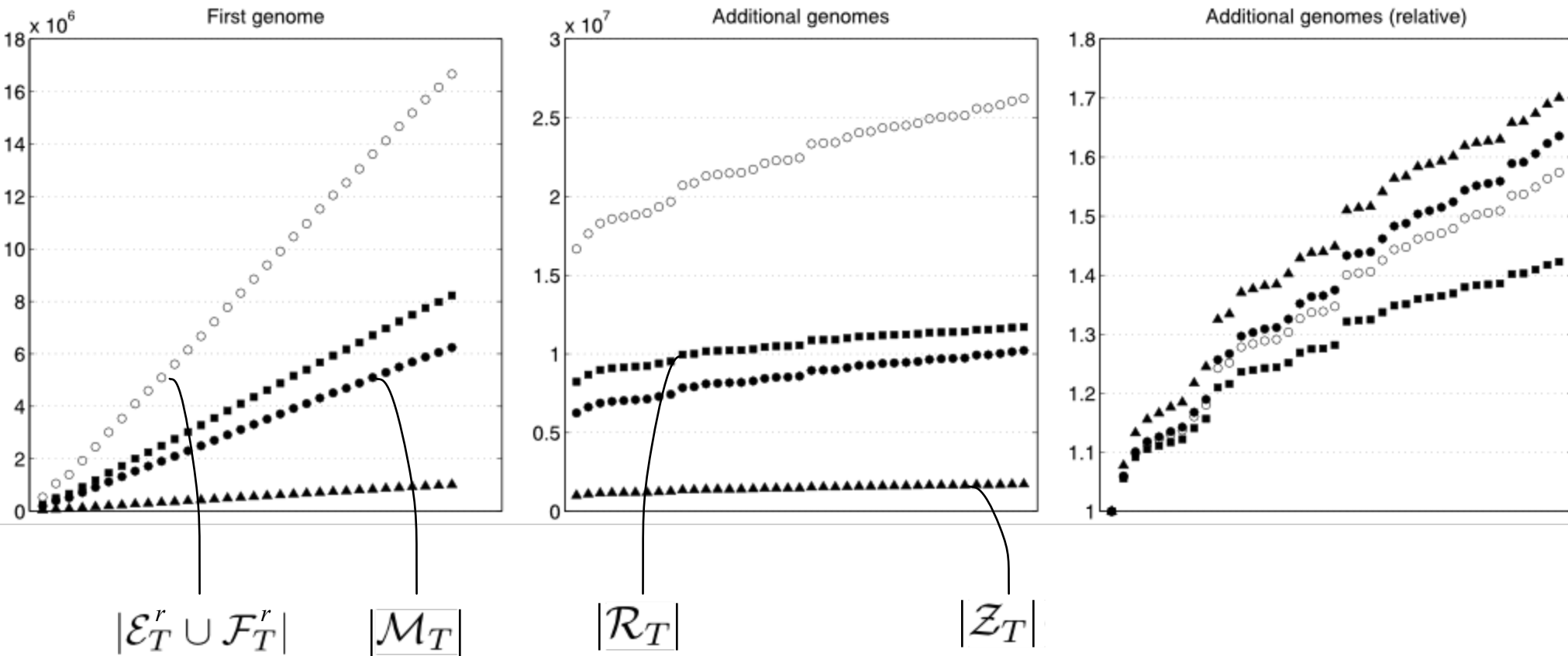
(2) Max Planck Institute for Molecular Cell Biology and Genetics, Dresden, Germany.

(3) Department of Mathematics and Computer Science, University of Udine, Italy.

(4) LIAFA, Paris Diderot University - Paris 7, France.

Highly-repetitive strings

39 *Saccharomyces cerevisiae* genomes



Distinct measures of repetition all grow sublinearly

Combining repetition-aware data structures



RLBWT_T

\mathcal{R}_T



LZ77 index

\mathcal{Z}_T

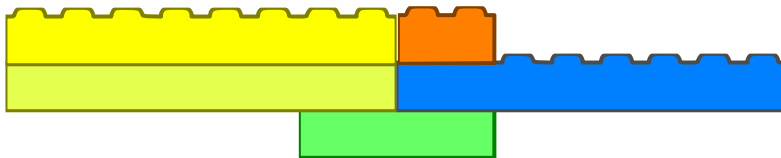


CDAWG_T

$\mathcal{M}_T, \mathcal{E}_T \cup \mathcal{F}_T$

Locating

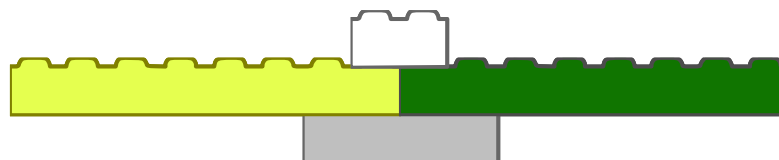
RLBWT _{\bar{T}}



Locating



**Suffix tree
representations**



Locating

Words:

RLBWT+CDAWG $O(|\mathcal{E}_T^r \cup \mathcal{F}_T^r|)$

RLBWT+LZ77 $O(|\mathcal{Z}_T| + |\mathcal{R}_T| + |\mathcal{R}_{\overline{T}}|)$

[1] $O(n/k + |\mathcal{R}_T|)$

Time:

RLBWT+CDAWG $O(m \log \log n + \text{occ})$

RLBWT+LZ77 $O(m(\log \log n + \log |\mathcal{Z}_T|) + p0\text{cc} \log^\epsilon |\mathcal{Z}_T| + s0\text{cc} \log \log n)$

[2] $O(m^2 h + (m + \text{occ}) \log |\mathcal{Z}_T|)$

[1] $O(m \log \log n + k \cdot \text{occ} \log \log n)$

[1] Veli Mäkinen, Gonzalo Navarro, Jouni Sirén, and Niko Välimäki. *Storage and retrieval of highly repetitive sequence collections*. Journal of Computational Biology, 17(3):281–308, 2010.

[2] Sebastian Kreft and Gonzalo Navarro. *On compressing and indexing repetitive sequences*. Theoretical Computer Science, 483:115–133, 2013.

Suffix tree representation

	stringDepth locateLeaf	isAncestor	parent nextSibling	child firstChild	suffixLink	weinerLink	edgeChar	nLeaves
1	$O(1)$	$O(1)$	$O(\log \log n)$	$O(1)$	$O(\log \log n)$	$O(\log \log n)$	$O(\log \log n)$	$O(1)$
2	$O(1)$	$O(1)$	$O(\log \log n)$	$O(1)$	$O(\log \log n)$	$O(\log \log n)$		$O(1)$
3	$O(1)$		$O(\log \log n)$	$O(1)$	$O(1)$			

Words: $O(|\mathcal{E}_T^r| + |\mathcal{F}_T^r| + |\mathcal{E}_T^\ell| + |\mathcal{F}_T^\ell|)$
 $O(|\mathcal{E}_T^r| + |\mathcal{F}_T^r|)$

Preliminaries

Maximal repeats

and the BWT

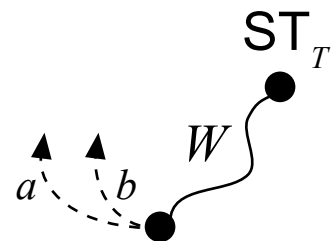
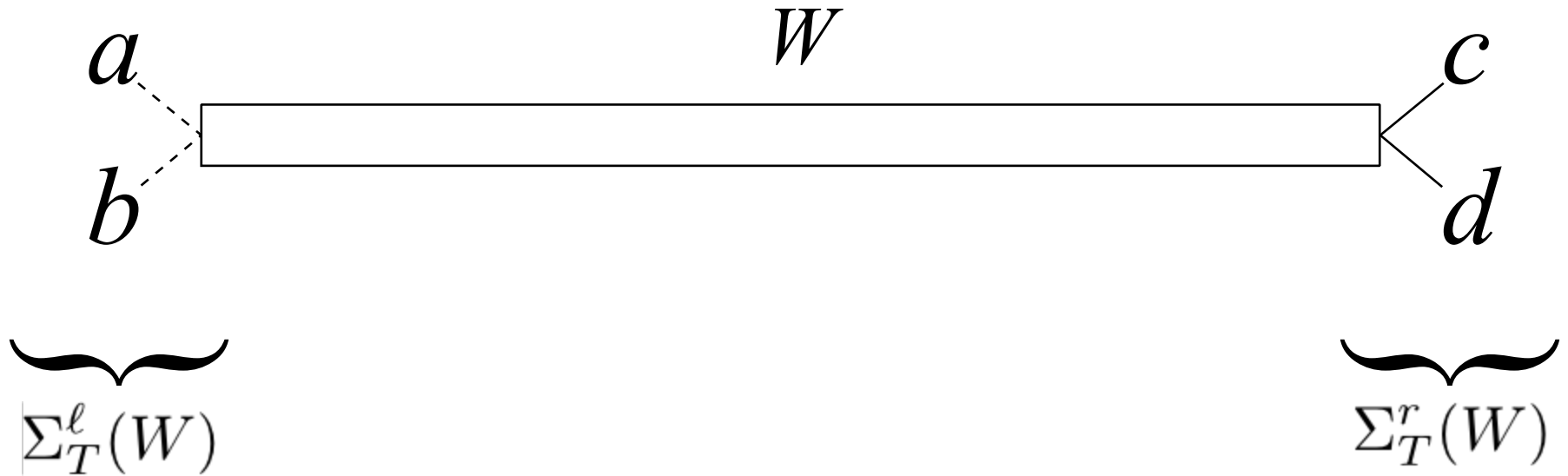
and the CDAWG

"Rightmost" maximal repeats

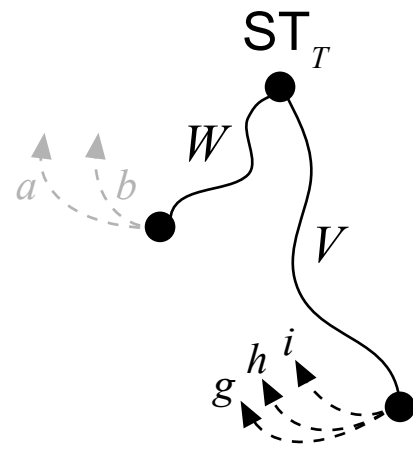
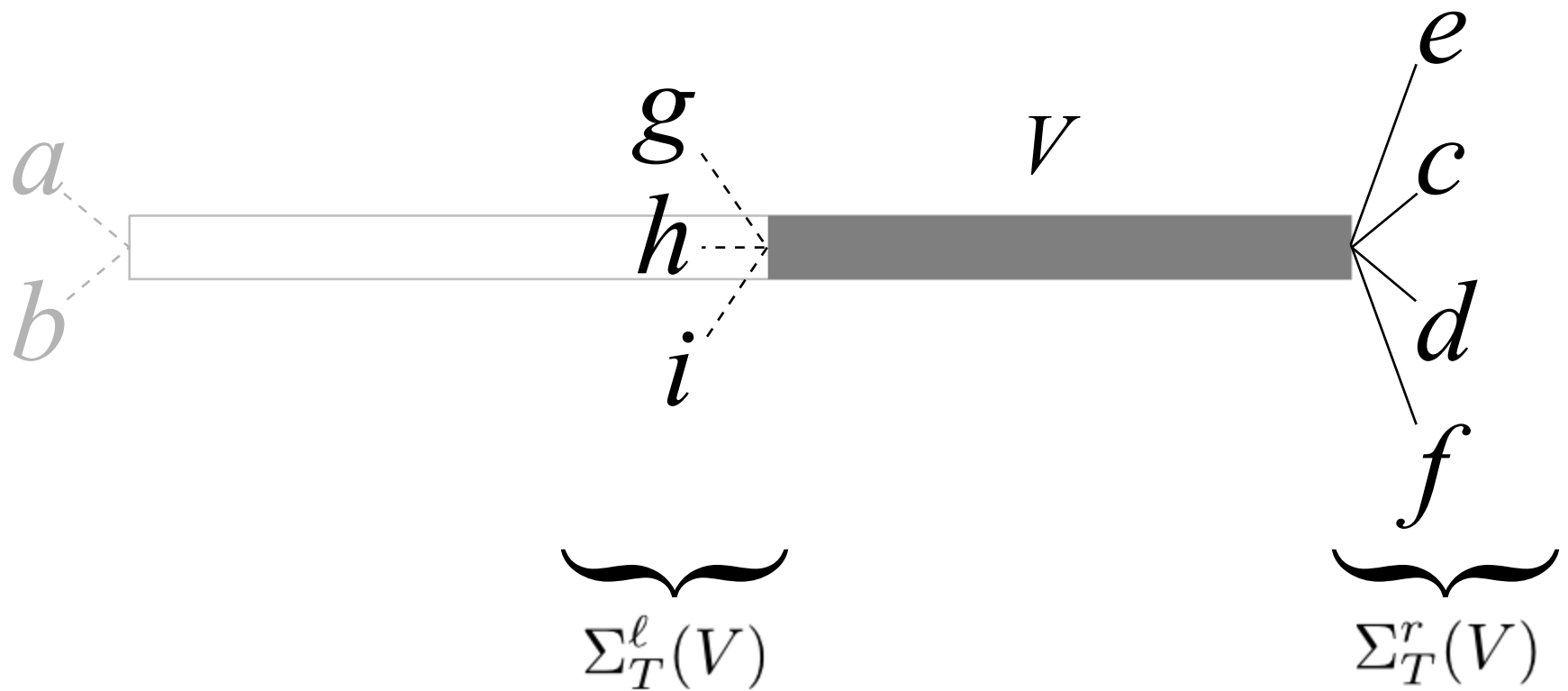
and the BWT

and LZ77 factors

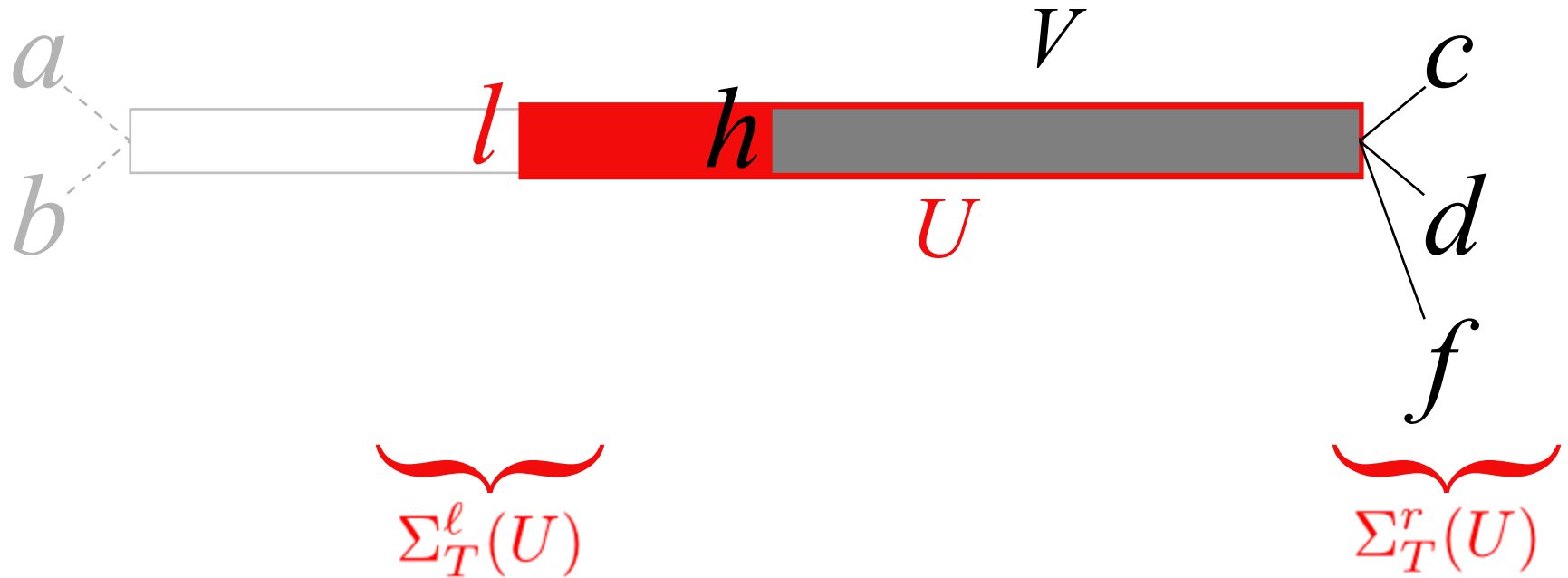
Maximal repeats



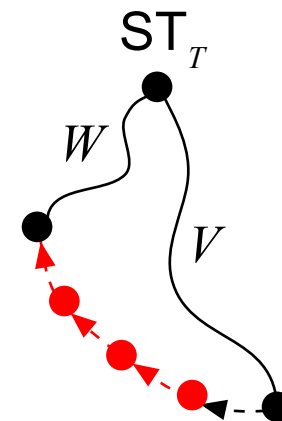
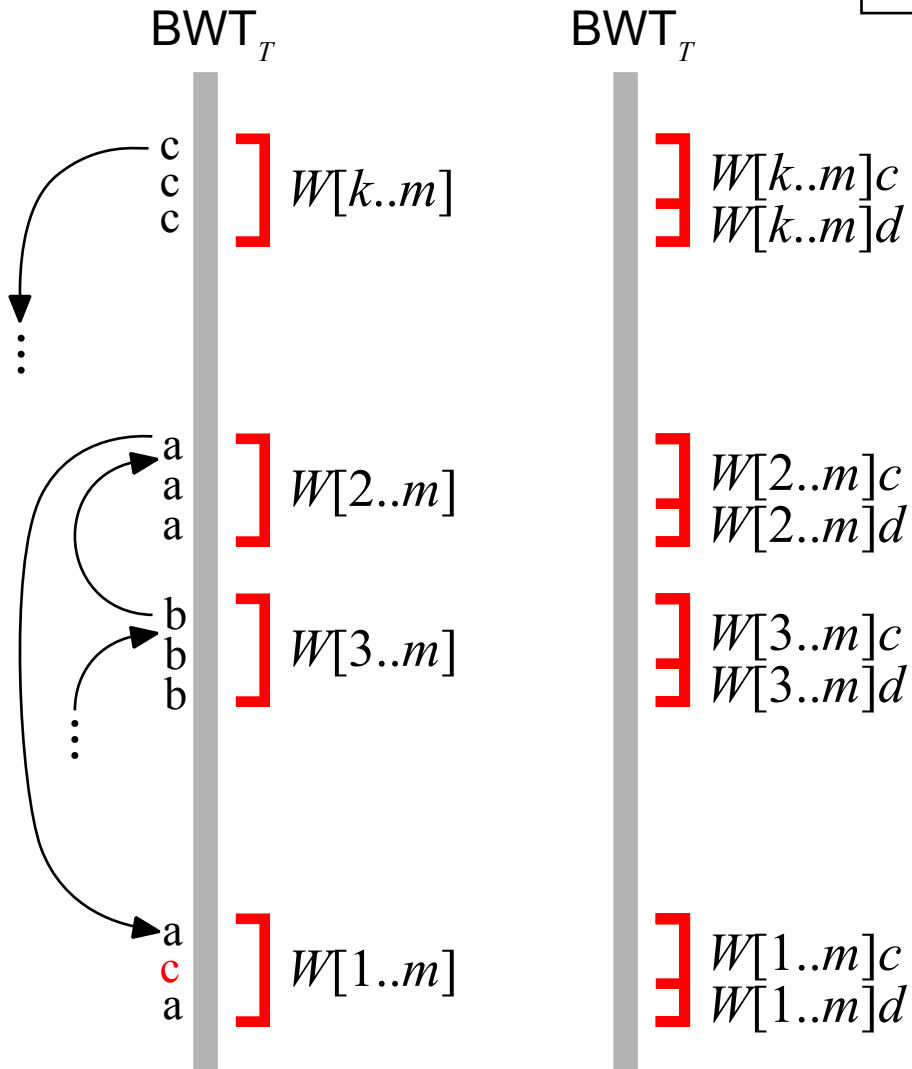
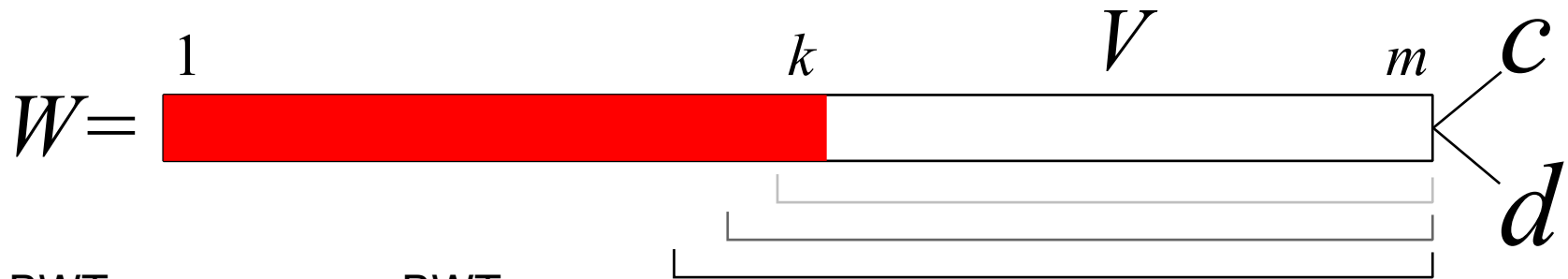
Maximal repeats



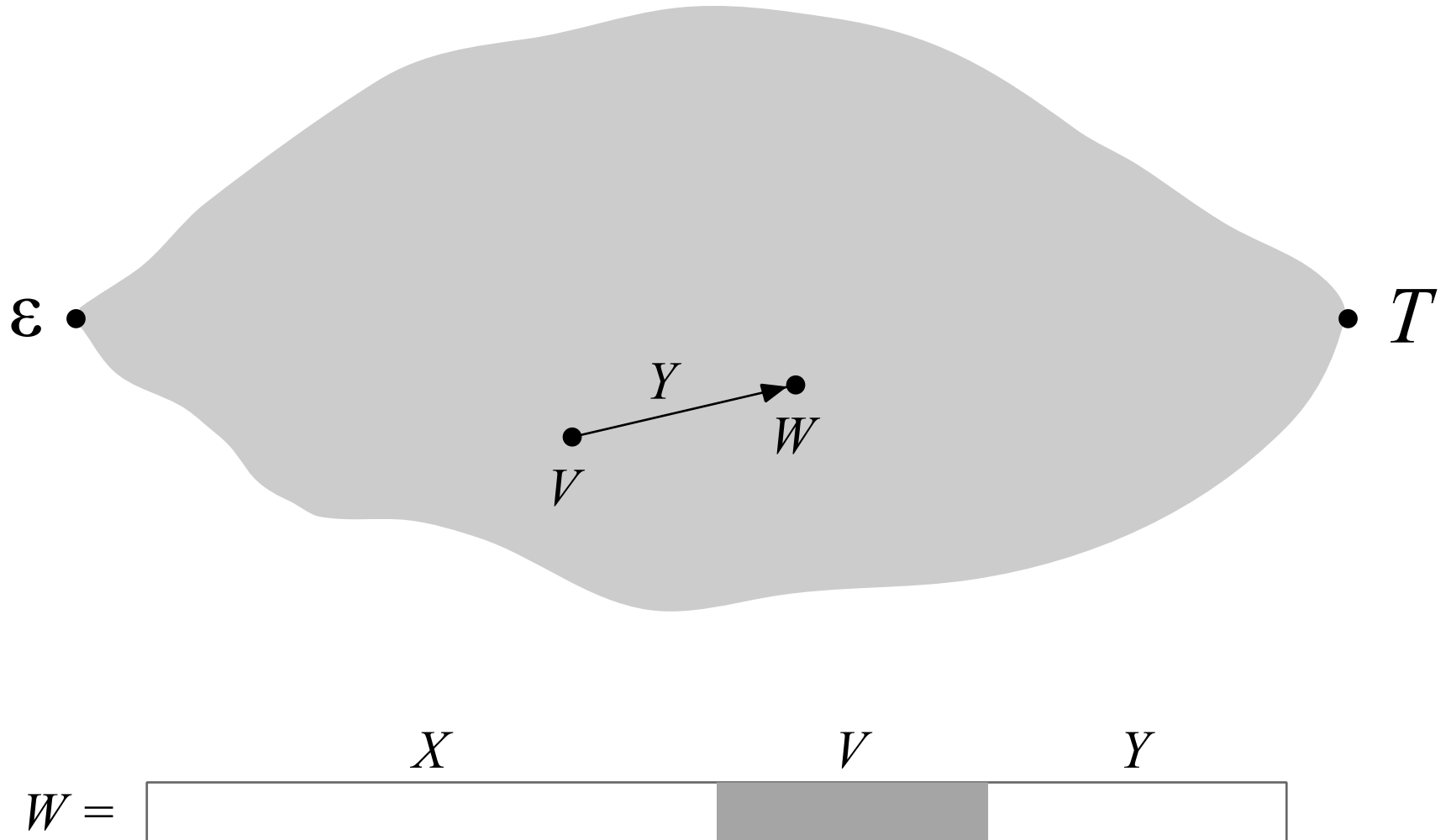
Maximal repeats



Maximal repeats and BWT



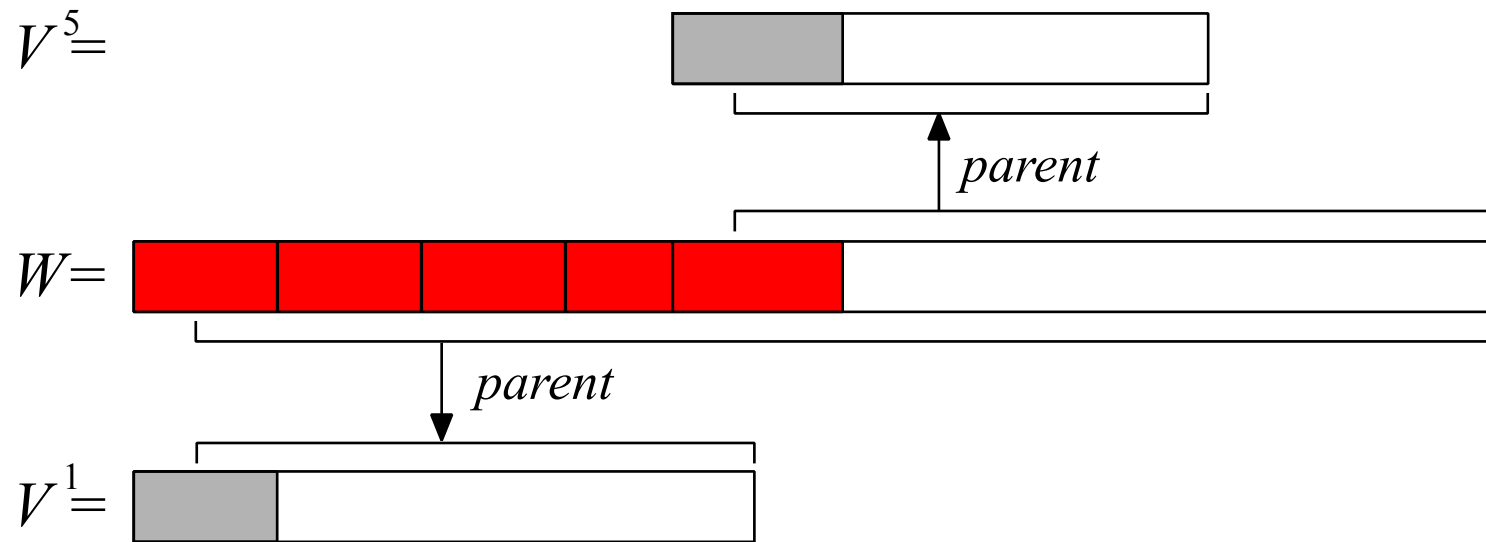
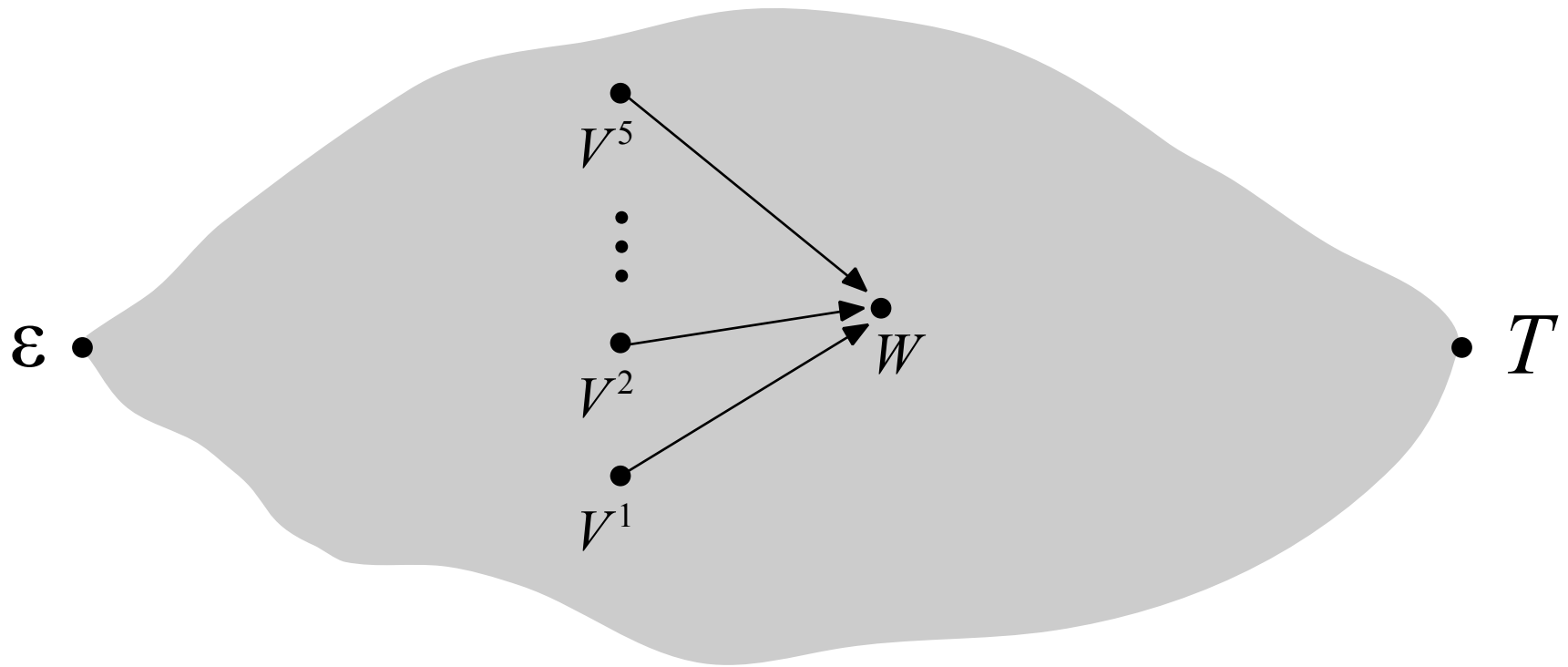
Compact Directed Acyclic Word Graph



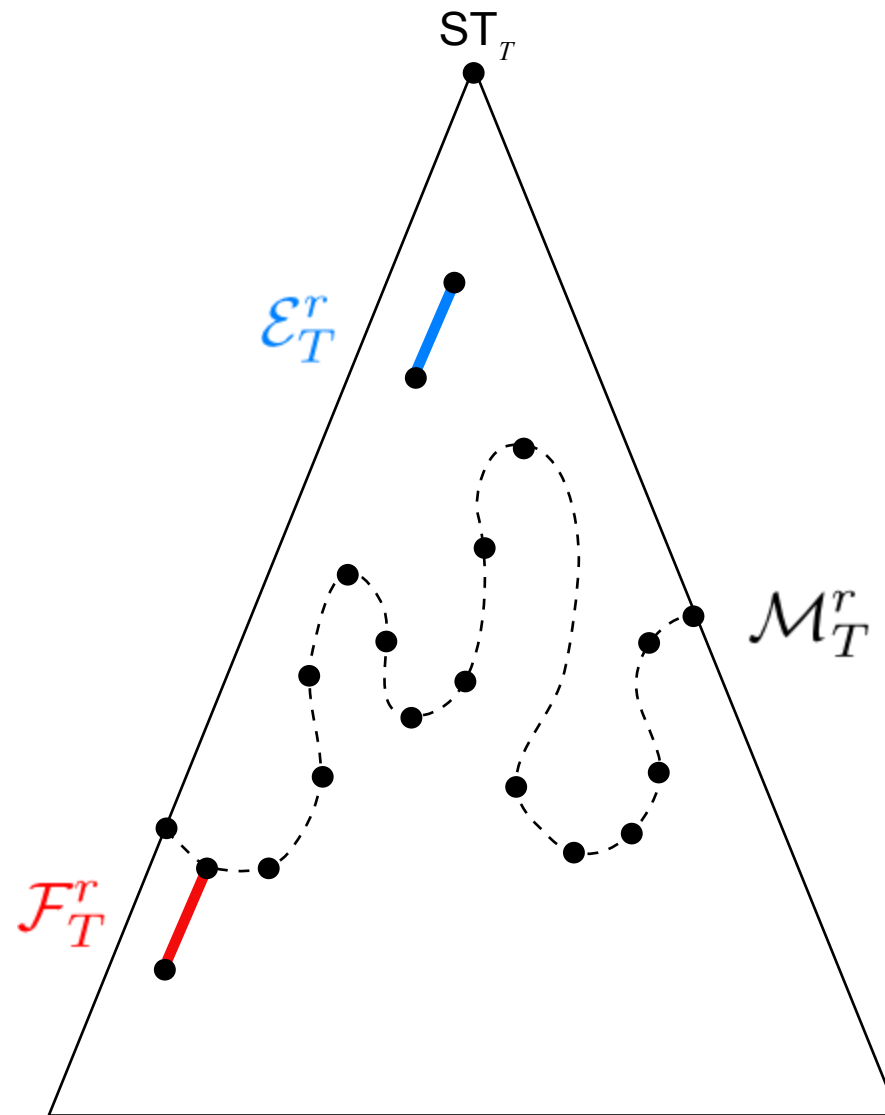
[1] Anselm Blumer, Janet Blumer, David Haussler, Ross McConnell, and Andrzej Ehrenfeucht. *Complete inverted files for efficient text retrieval and analysis*. Journal of the ACM, 34(3):578–595, 1987.

[2] Maxime Crochemore and Renaud V  rin. *Direct construction of compact directed acyclic word graphs*. In Alberto Apostolico and Jotun Hein, editors, CPM, volume 1264 of Lecture Notes in Computer Science, pages 116–129. Springer, 1997.

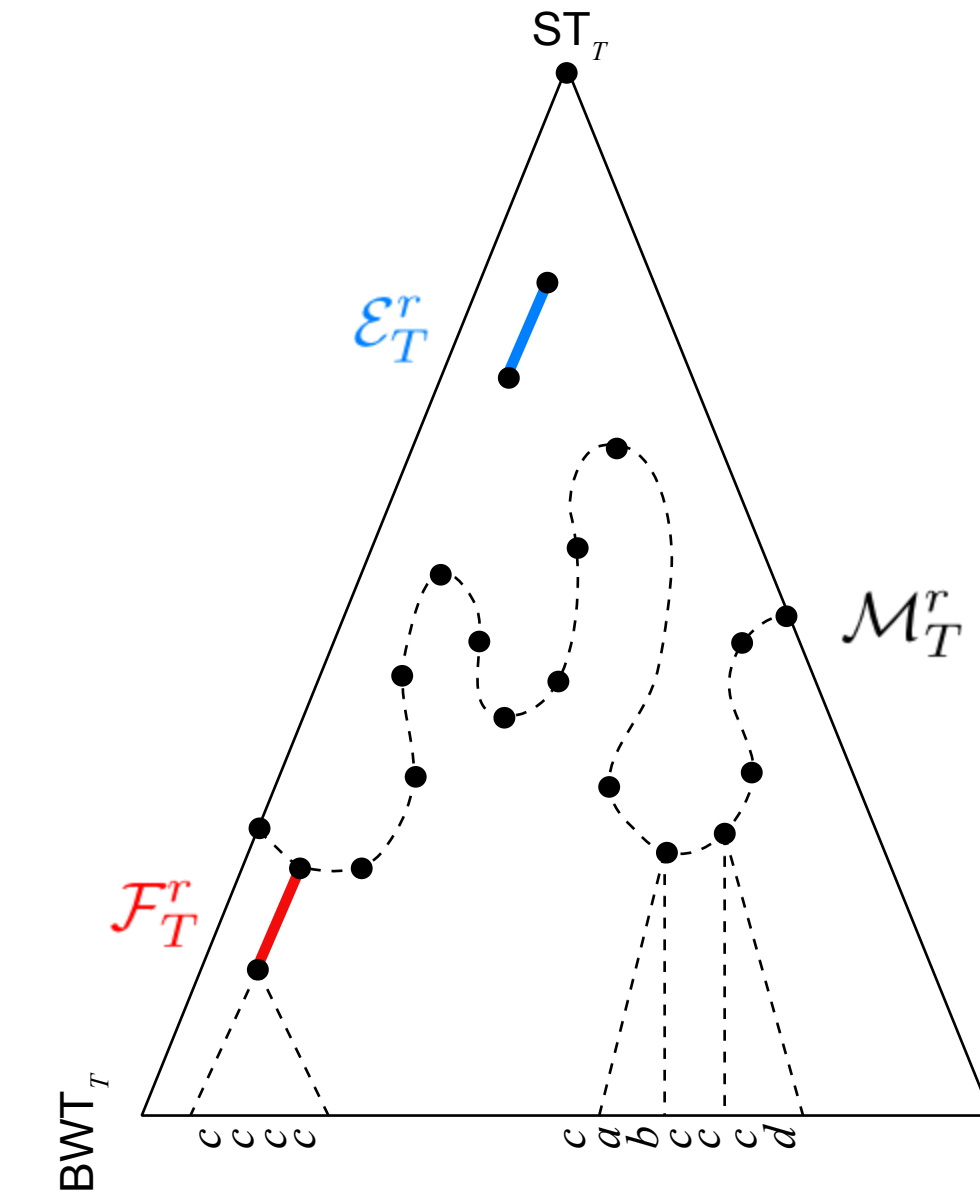
Maximal repeats and CDAWG



Rightmost maximal repeats



Rightmost maximal repeats and BWT runs



$$\sum_{W \in \mathcal{M}_T^r} |\Sigma_T^\ell(W)| - |\mathcal{M}_T^r| + 1 \leq |\mathcal{R}_T| \leq |\mathcal{F}_T^r|$$

Locating with RLBWT+LZ77



RLBWT_T

\mathcal{R}_T

Rank/select in $O(\log \log n)$
time, $O(|\mathcal{R}_T|)$ words
(predecessor data structure).



LZ77 index

\mathcal{Z}_T

Primary occurrences:
 $O((1 + \text{occ}) \log^\epsilon(|\mathcal{Z}_T|))$ time,
 $O(|\mathcal{Z}_T|)$ words
(4-sided range reporting).

Secondary occurrences:
 $O(\text{occ} \log \log n)$ time,
 $O(|\mathcal{Z}_T|)$ words
(2-sided range reporting).

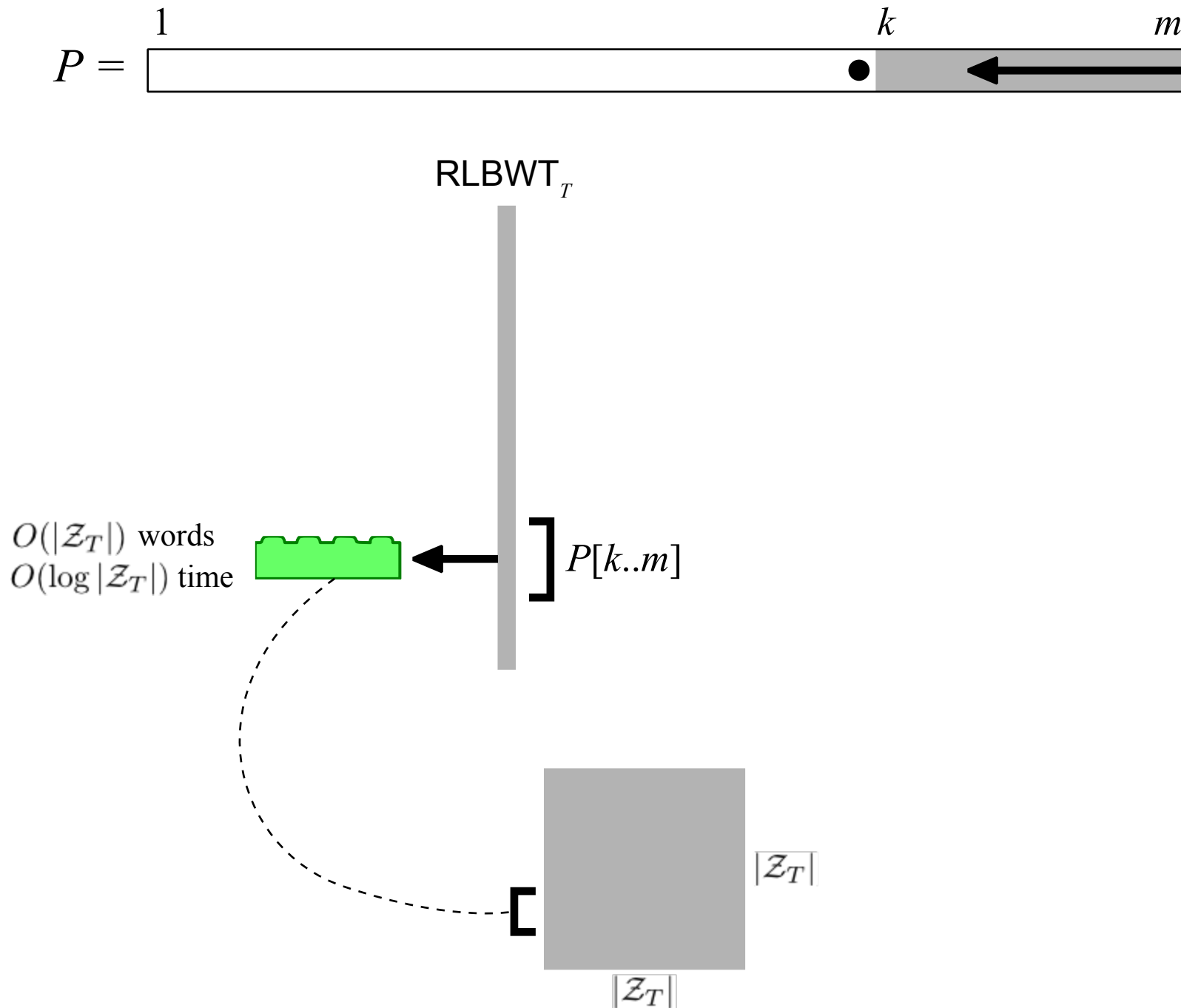


CDAWG_T

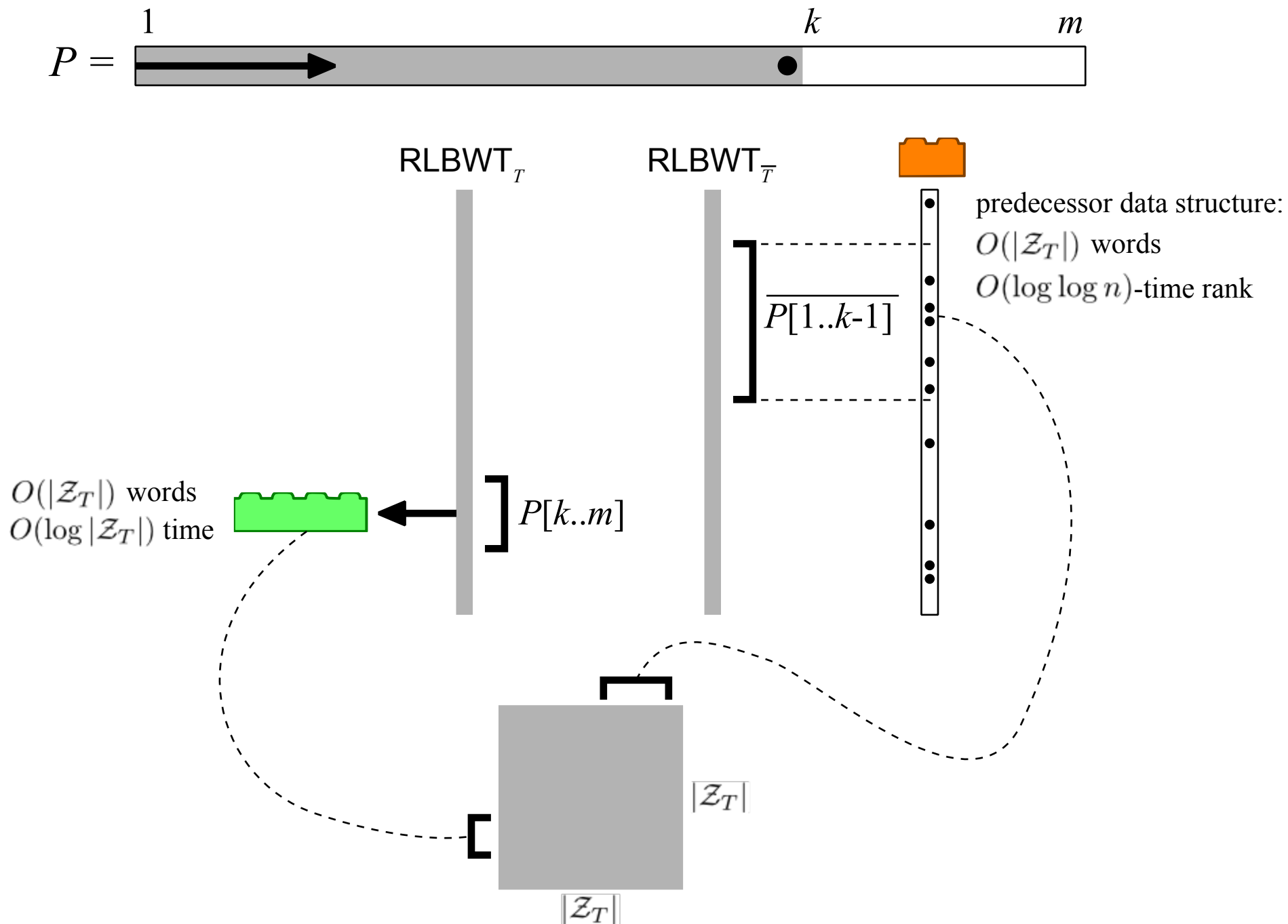
$\mathcal{M}_T, \mathcal{E}_T \cup \mathcal{F}_T$

- [1] Dan E Willard. *Log-logarithmic worst-case range queries are possible in space $\Theta(N)$* . Information Processing Letters, 17(2):81–84, 1983.
- [2] Timothy M. Chan, Kasper Green Larsen, and Mihai Pătraşcu. *Orthogonal range searching on the RAM, revisited*. In Proceedings of the Twenty-seventh Annual Symposium on Computational Geometry, pages 1–10. ACM, 2011.
- [3] Juha Kärkkäinen and Esko Ukkonen. *Lempel-Ziv parsing and sublinear-size index structures for string matching*. In Proc. 3rd South American Workshop on String Processing (WSP'96), pages 141–155, 1996.

Locating with RLBWT+LZ77



Locating with RLBWT+LZ77



Locating with RLBWT+LZ77

Words:

$$\text{RLBWT+LZ77} \quad O(|\mathcal{Z}_T| + |\mathcal{R}_T| + |\mathcal{R}_{\overline{T}}|)$$

$$[1] \quad O(n/k + |\mathcal{R}_T|)$$

Time:

$$\text{RLBWT+LZ77} \quad O(m(\log \log n + \log |\mathcal{Z}_T|) + p0cc \log^\epsilon |\mathcal{Z}_T| + s0cc \log \log n)$$

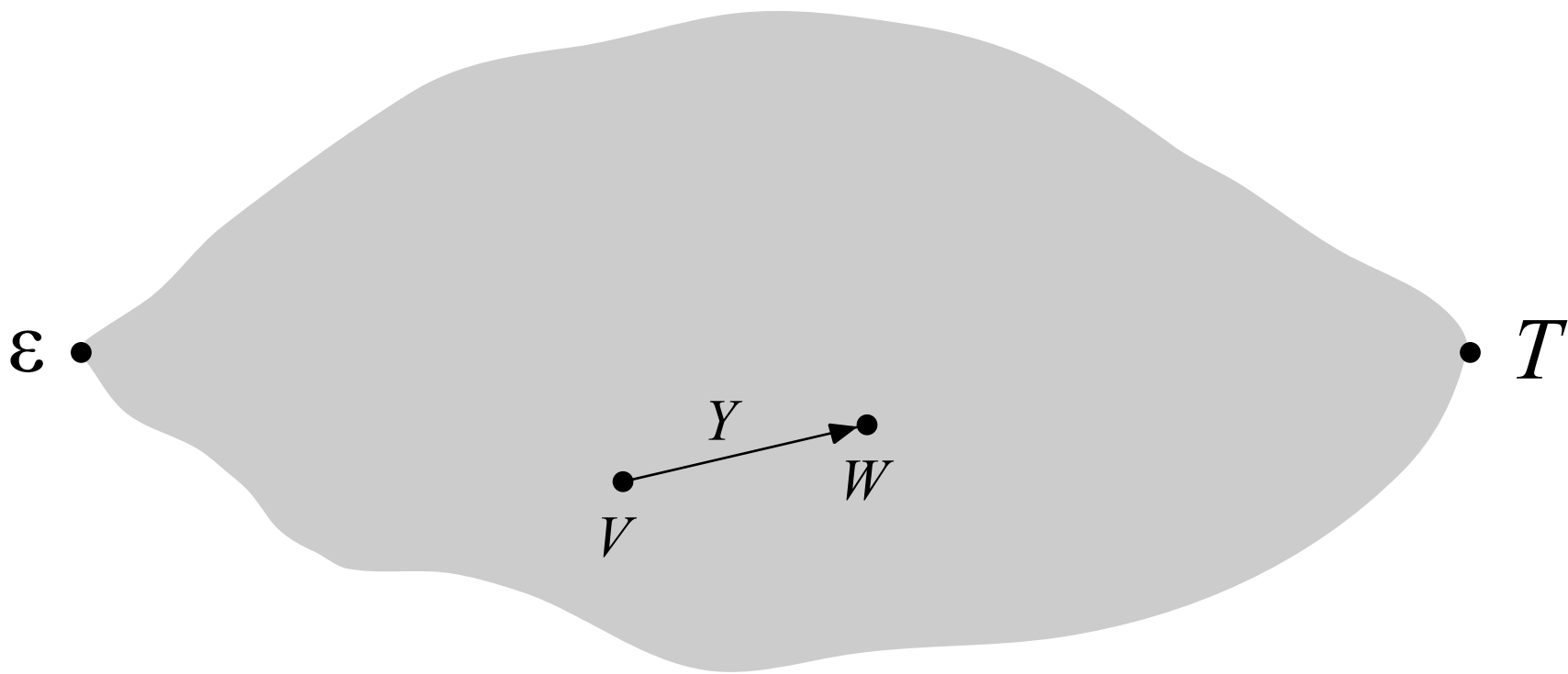
$$[2] \quad O(m^2 h + (m + occ) \log |\mathcal{Z}_T|)$$

$$[1] \quad O(m \log \log n + k \cdot occ \log \log n)$$

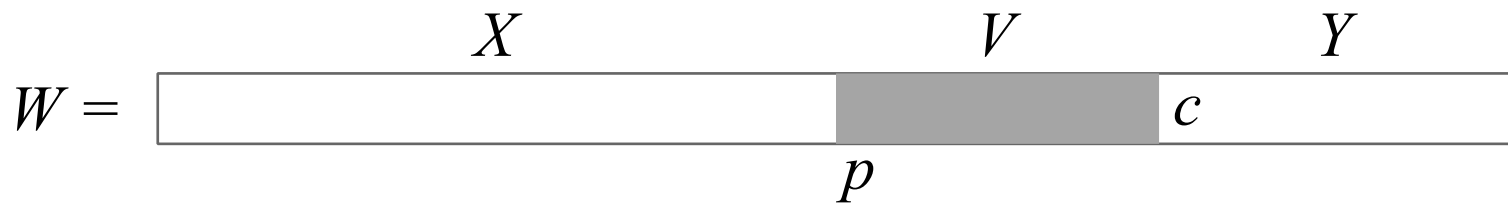
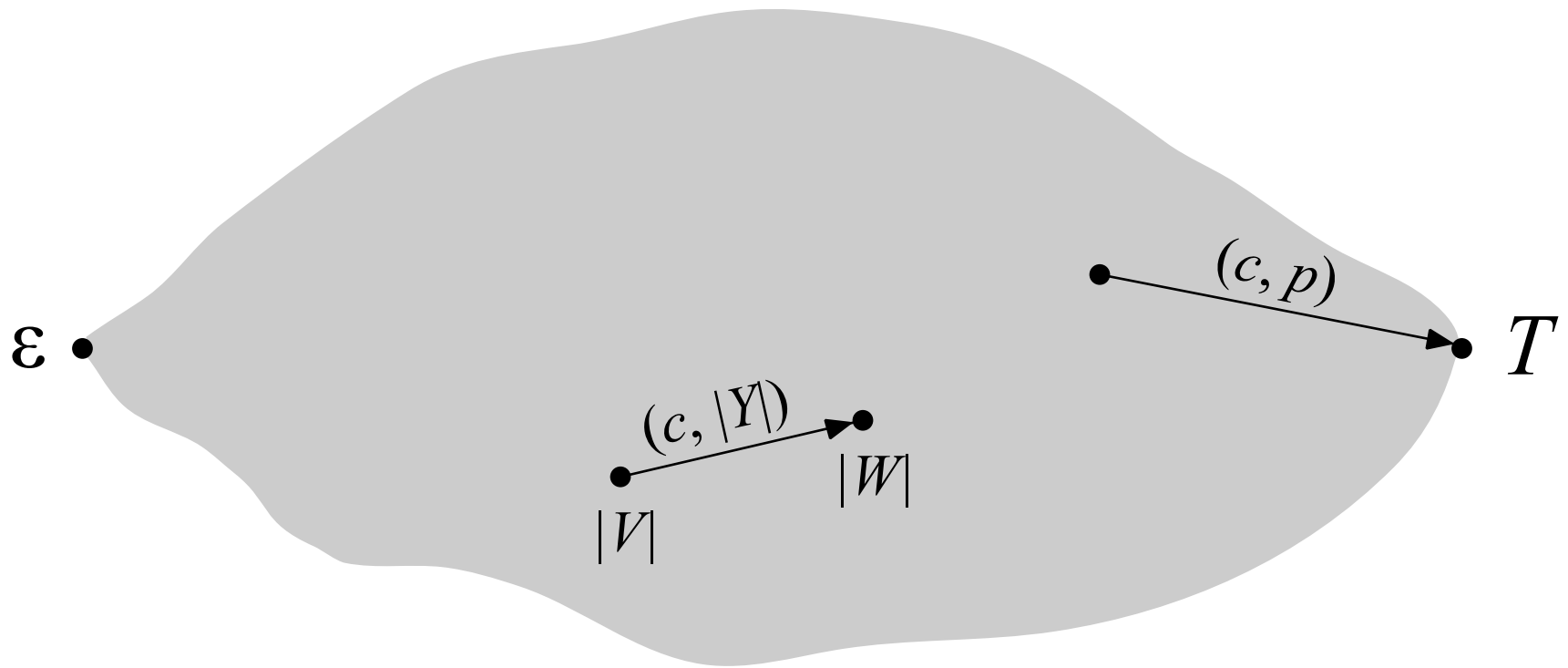
[1] Veli Mäkinen, Gonzalo Navarro, Jouni Sirén, and Niko Välimäki. *Storage and retrieval of highly repetitive sequence collections*. Journal of Computational Biology, 17(3):281–308, 2010.

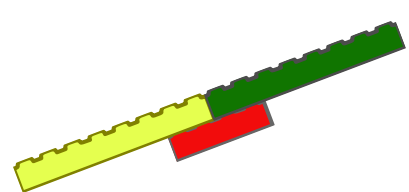
[2] Sebastian Kreft and Gonzalo Navarro. *On compressing and indexing repetitive sequences*. Theoretical Computer Science, 483:115–133, 2013.

Locating with RLBWT+CDAWG



Locating with RLBWT+CDAWG

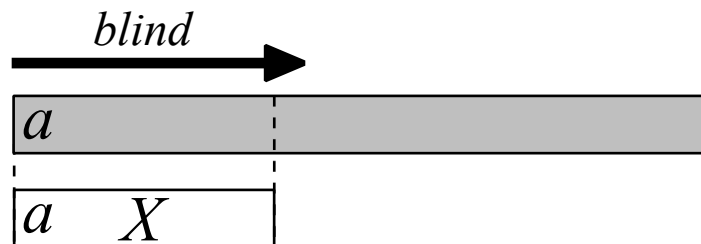




Locating with RLBWT+CDAWG

$$P =$$

$$W_1 =$$



RLBWT_T

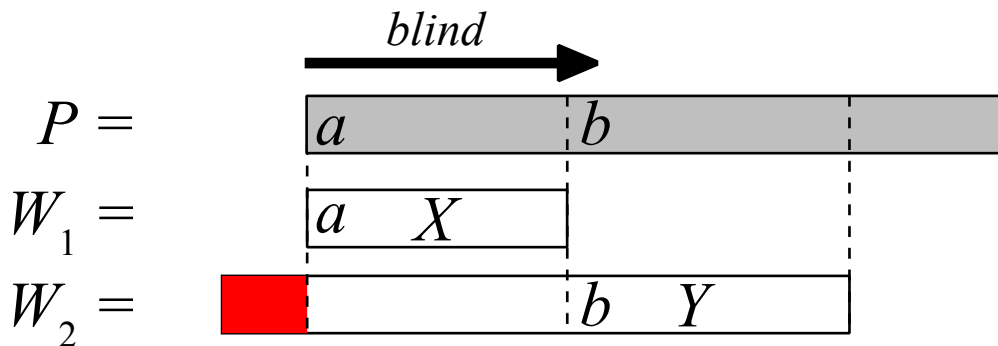
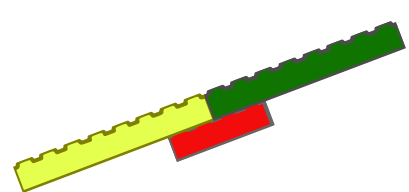
CDAWG_T

$(a, |X|)$

W_1

$]P$

Locating with RLBWT+CDAWG

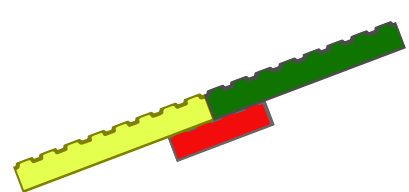


RLBWT_T

CDAWG_T



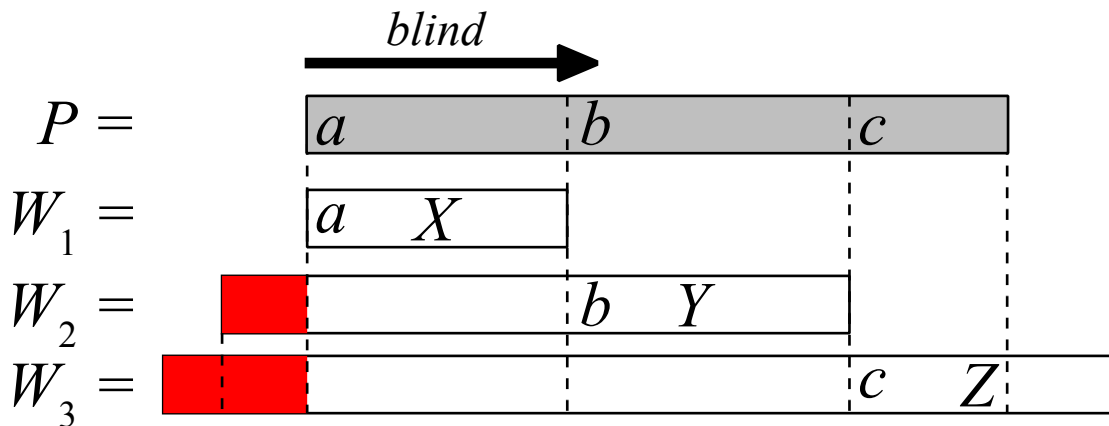
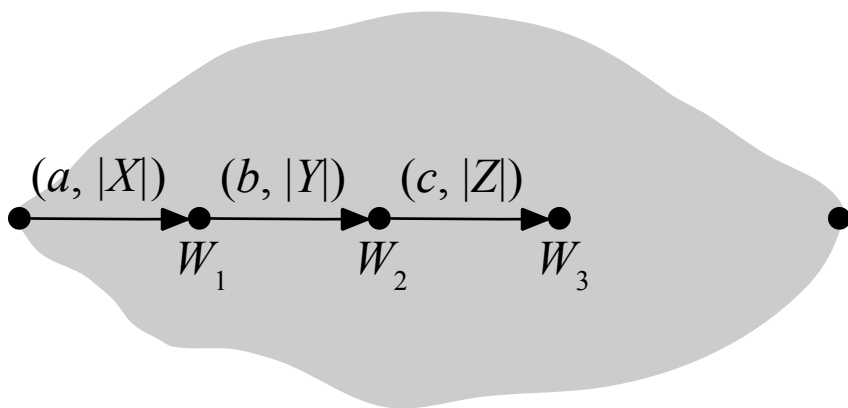
Locating with RLBWT+CDAWG



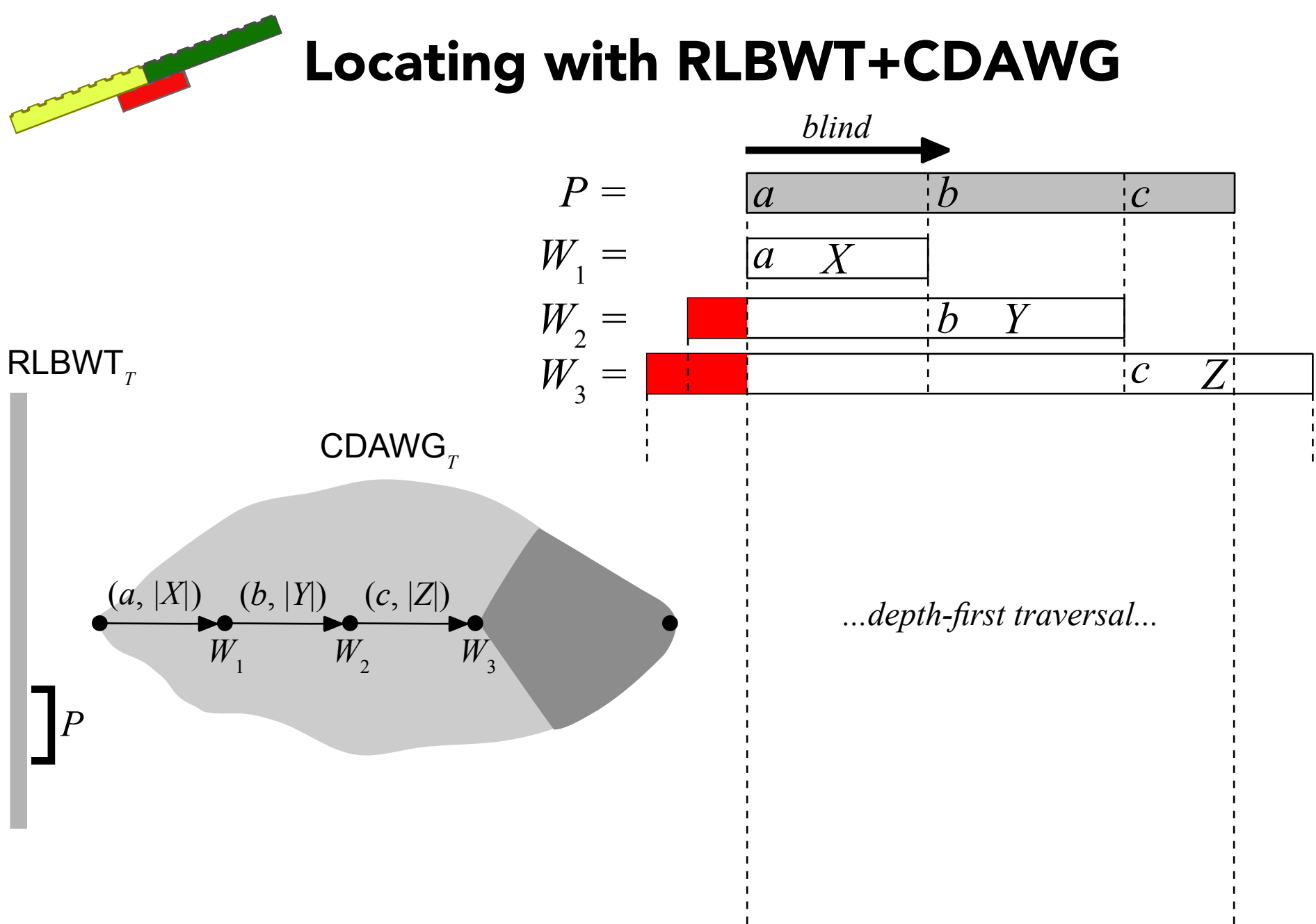
RLBWT_T



CDAWG_T



Locating with RLBWT+CDAWG



Locating with RLBWT+CDAWG

Words:

$$\text{RLBWT+CDAWG} \quad O(|\mathcal{E}_T^r \cup \mathcal{F}_T^r|)$$

$$\text{RLBWT+LZ77} \quad O(|\mathcal{Z}_T| + |\mathcal{R}_T| + |\mathcal{R}_{\overline{T}}|)$$

$$[1] \quad O(n/k + |\mathcal{R}_T|)$$

Time:

$$\text{RLBWT+CDAWG} \quad O(m \log \log n + \text{occ})$$

$$\text{RLBWT+LZ77} \quad O(m(\log \log n + \log |\mathcal{Z}_T|) + p_0 \text{occ} \log^\epsilon |\mathcal{Z}_T| + s_0 \text{occ} \log \log n)$$

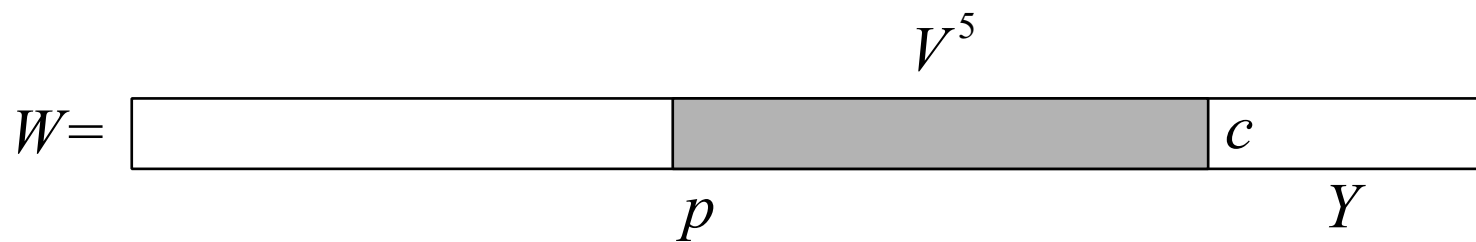
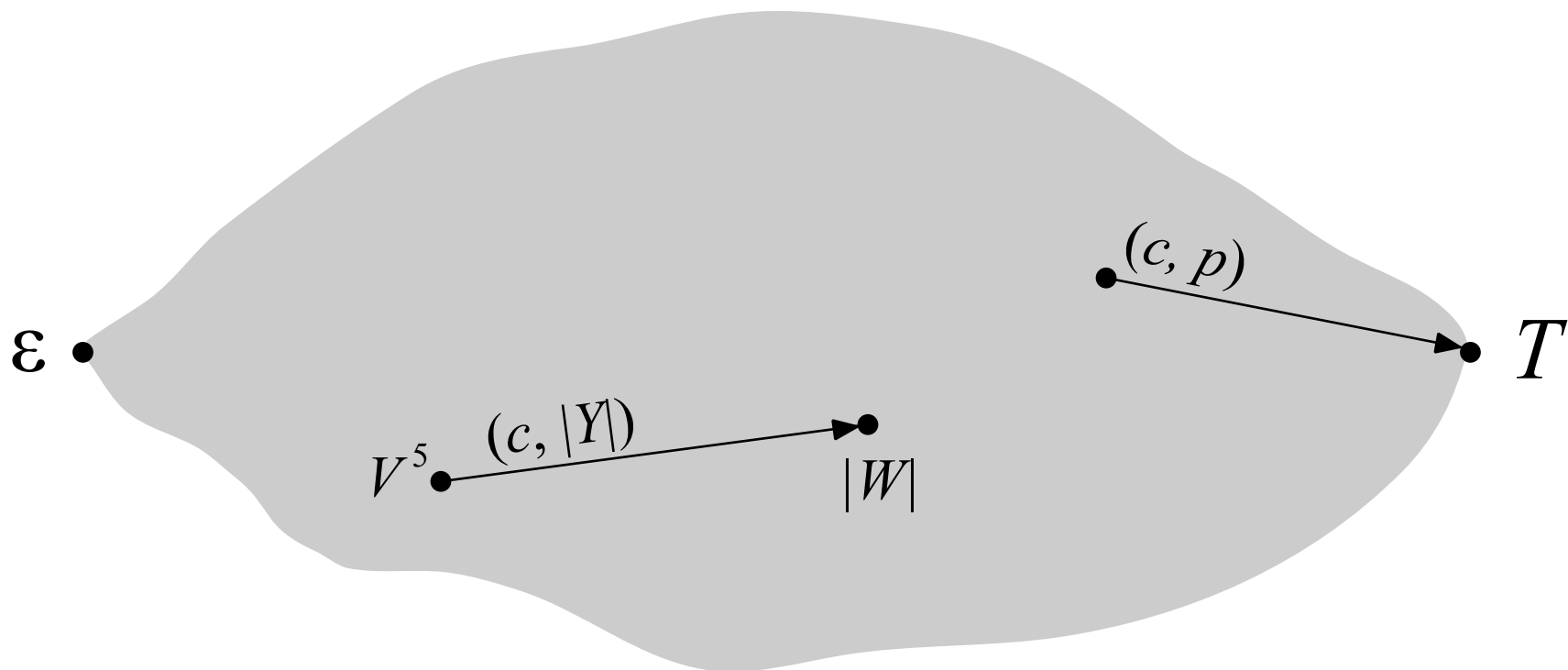
$$[2] \quad O(m^2 h + (m + \text{occ}) \log |\mathcal{Z}_T|)$$

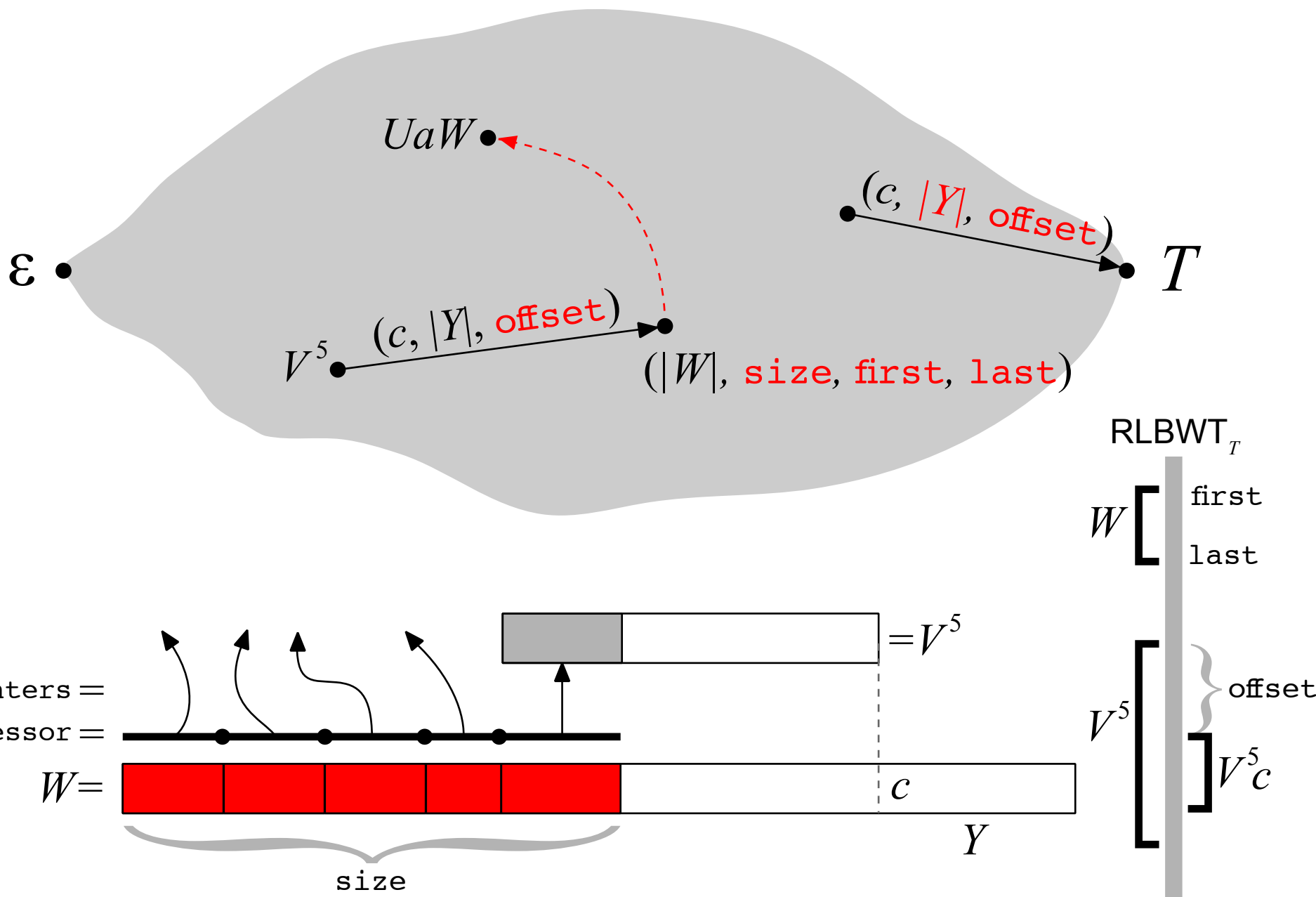
$$[1] \quad O(m \log \log n + k \cdot \text{occ} \log \log n)$$

[1] Veli Mäkinen, Gonzalo Navarro, Jouni Sirén, and Niko Välimäki. *Storage and retrieval of highly repetitive sequence collections*. Journal of Computational Biology, 17(3):281–308, 2010.

[2] Sebastian Kreft and Gonzalo Navarro. *On compressing and indexing repetitive sequences*. Theoretical Computer Science, 483:115–133, 2013.

CDAWG for locating





A node v of the suffix tree

$$\text{id}(v) = (\bullet, |\ell(v)|, \bullet, \bullet)$$

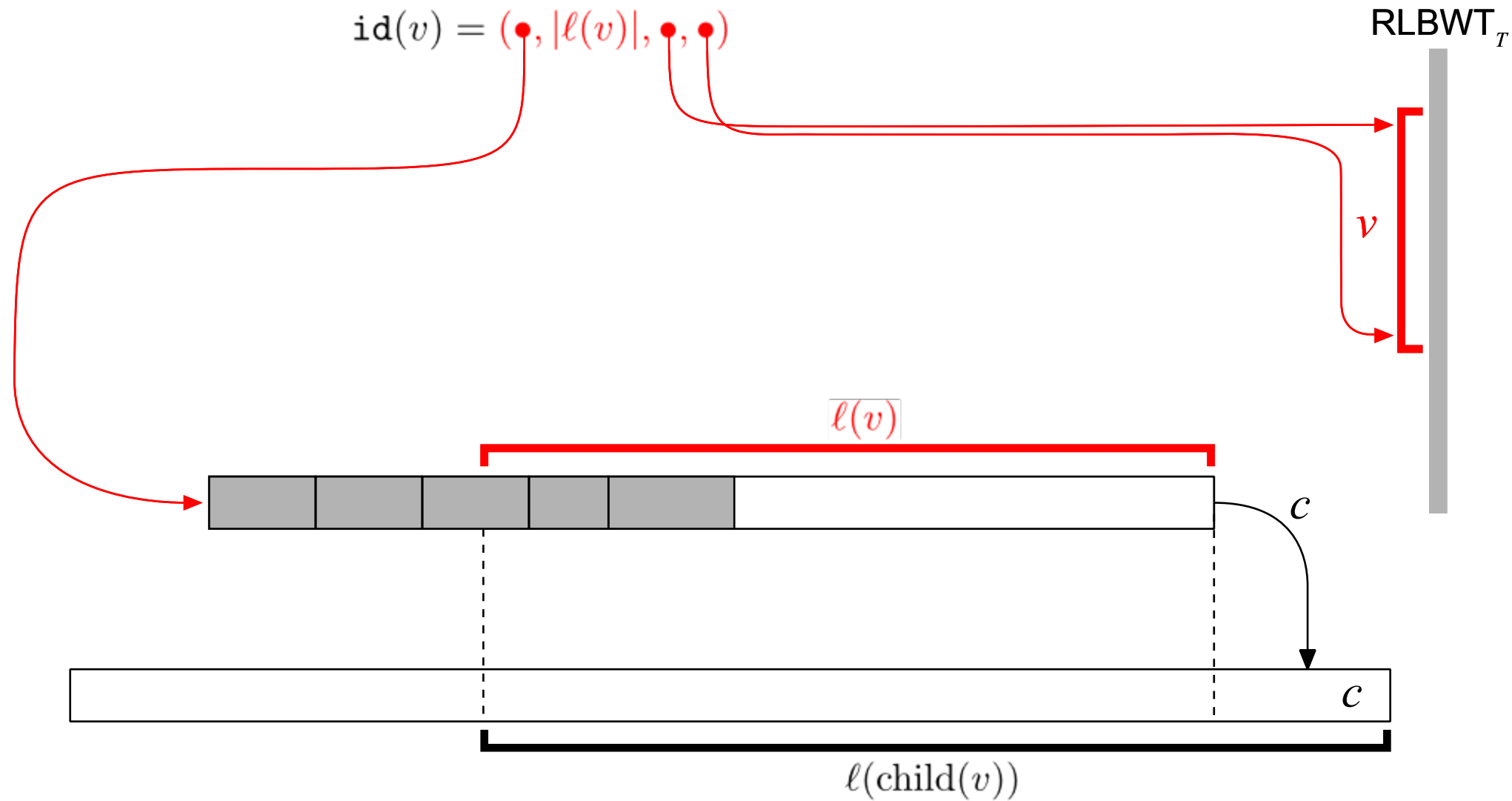
RLBWT_T

v

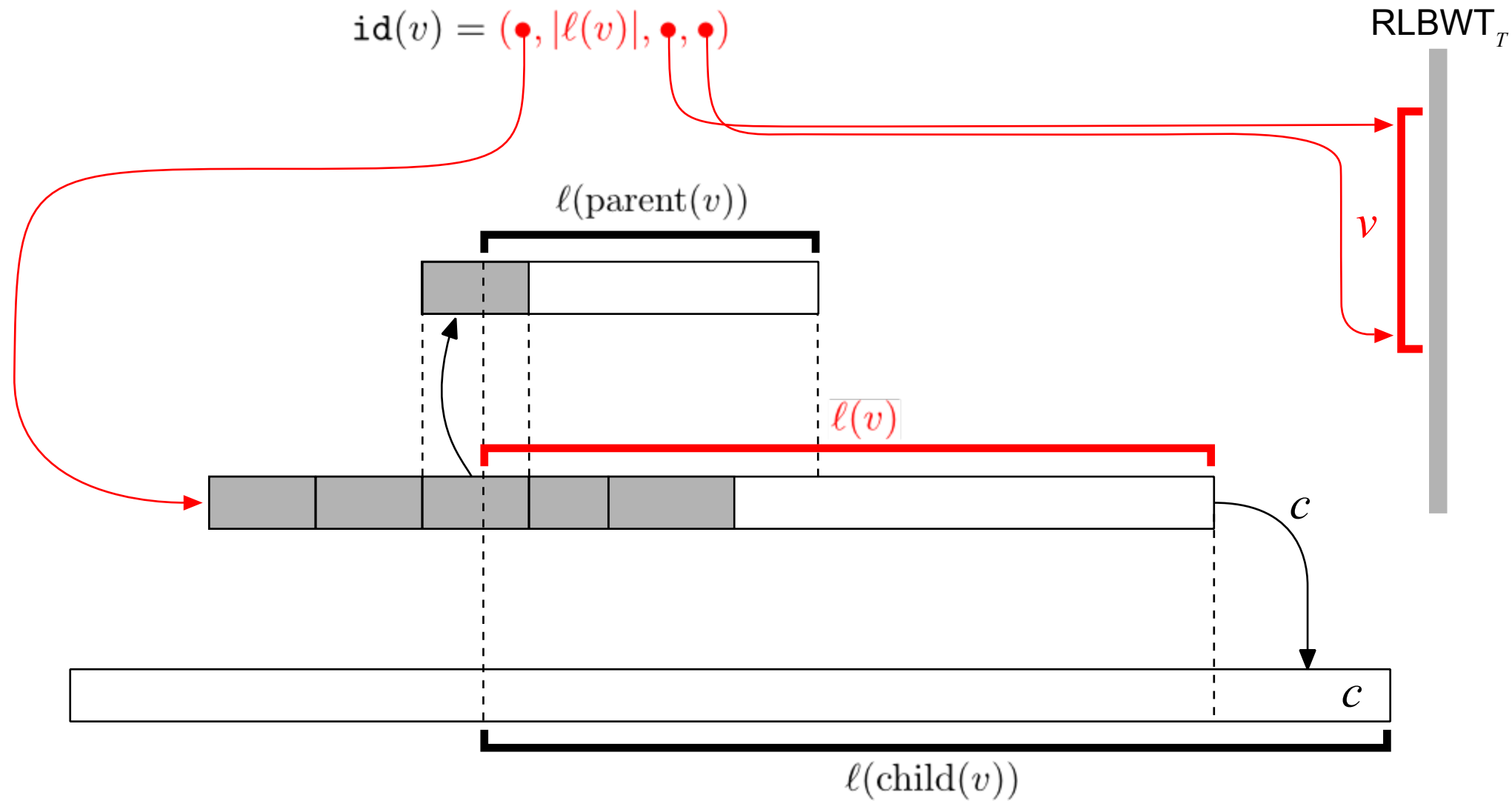
$\ell(v)$



A node v of the suffix tree



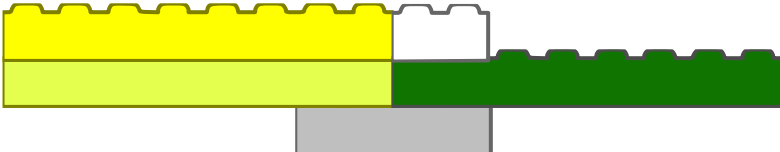
A node v of the suffix tree

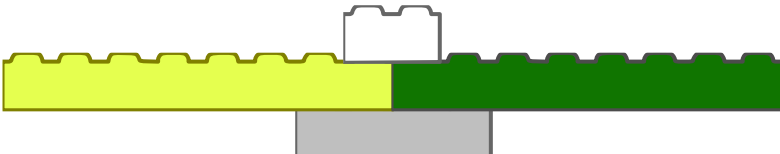


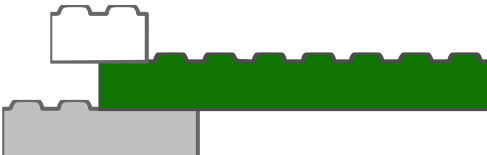
Suffix tree operations

Time:

	stringDepth locateLeaf	isAncestor	parent nextSibling	child firstChild	suffixLink	weinerLink	edgeChar	nLeaves
1	$O(1)$	$O(1)$	$O(\log \log n)$	$O(1)$	$O(\log \log n)$	$O(\log \log n)$	$O(\log \log n)$	$O(1)$
2	$O(1)$	$O(1)$	$O(\log \log n)$	$O(1)$	$O(\log \log n)$	$O(\log \log n)$		$O(1)$
3	$O(1)$		$O(\log \log n)$	$O(1)$	$O(1)$			

1)  $O(|\mathcal{E}_T^r| + |\mathcal{F}_T^r| + |\mathcal{E}_T^l| + |\mathcal{F}_T^l|)$

2)  $O(|\mathcal{E}_T^r| + |\mathcal{F}_T^r|)$

3)  $O(|\mathcal{E}_T^r| + |\mathcal{F}_T^r|)$

Suffix tree operations

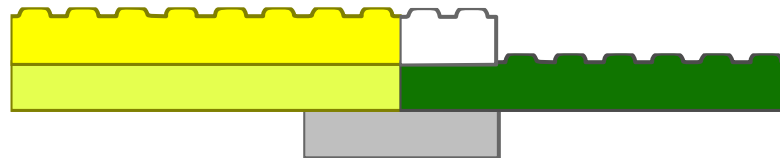
Time:

matching statistics

	stringDepth locateLeaf	isAncestor	parent nextSibling	child firstChild	suffixLink	weinerLink	edgeChar	nLeaves
1	$O(1)$	$O(1)$	$O(\log \log n)$	$O(1)$	$O(\log \log n)$	$O(\log \log n)$	$O(\log \log n)$	$O(1)$
2	$O(1)$	$O(1)$	$O(\log \log n)$	$O(1)$	$O(\log \log n)$	$O(\log \log n)$		$O(1)$
3	$O(1)$		$O(\log \log n)$	$O(1)$	$O(1)$			

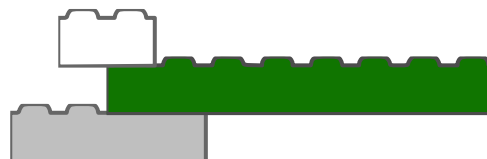
constant-space traversal

1)



$$O(|\mathcal{E}_T^r| + |\mathcal{F}_T^r| + |\mathcal{E}_T^\ell| + |\mathcal{F}_T^\ell|)$$

3)



$$O(|\mathcal{E}_T^r| + |\mathcal{F}_T^r|)$$



Composite repetition-aware data structures

Djamal Belazzougui¹, Fabio Cunial², Travis Gagie¹, Nicola Prezza³, Mathieu Raffinot⁴

(1) Department of Computer Science, University of Helsinki, Finland.

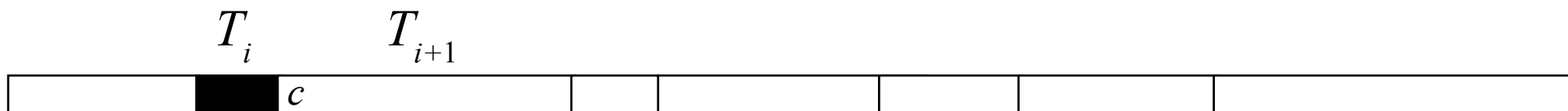
(2) Max Planck Institute for Molecular Cell Biology and Genetics, Dresden, Germany.

(3) Department of Mathematics and Computer Science, University of Udine, Italy.

(4) LIAFA, Paris Diderot University - Paris 7, France.

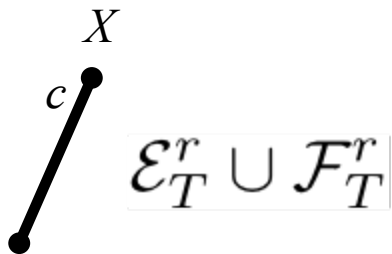
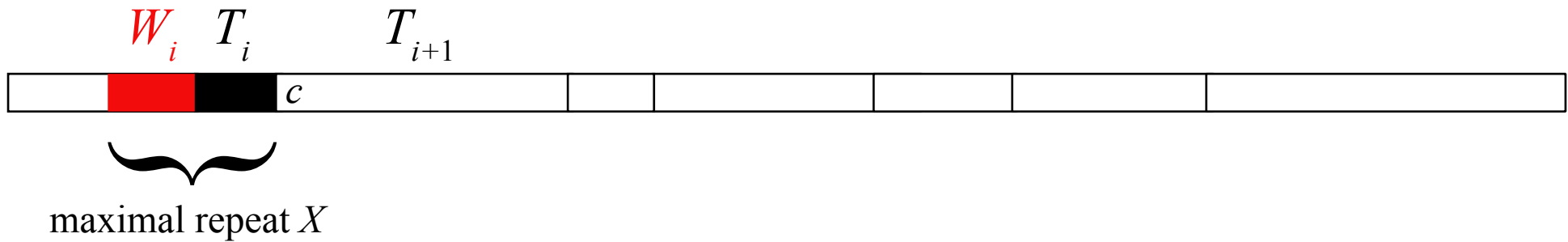
Rightmost maximal repeats and LZ factors

$$|\mathcal{Z}_T| \leq |\mathcal{E}_T^r \cup \mathcal{F}_T^r|$$



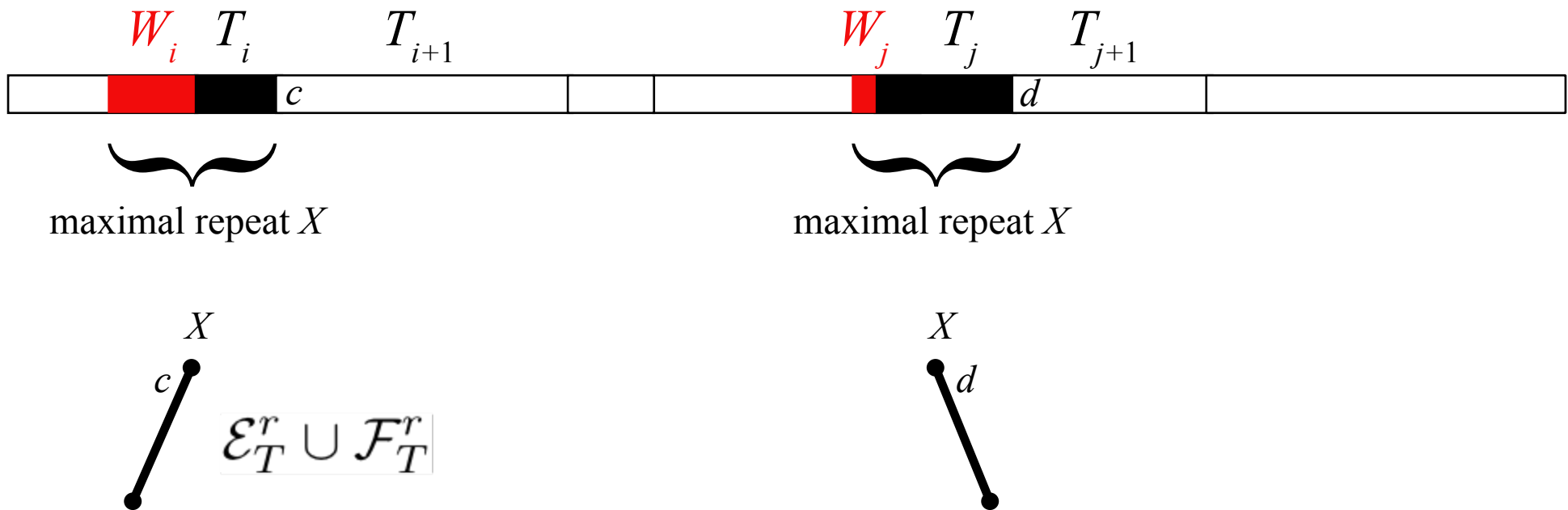
Rightmost maximal repeats and LZ factors

$$|\mathcal{Z}_T| \leq |\mathcal{E}_T^r \cup \mathcal{F}_T^r|$$



Rightmost maximal repeats and LZ factors

$$|\mathcal{Z}_T| \leq |\mathcal{E}_T^r \cup \mathcal{F}_T^r|$$



Measures of repetition

