

# Computing Minimal Absent Words and Extended Bispecial Factors with CDAWG Space

---

Shunsuke Inenaga<sup>1</sup>, Takuya Mieno<sup>2</sup>, Hiroki Arimura<sup>3</sup>,  
Mitsuru Funakoshi<sup>1</sup>, Yuta Fujishige<sup>4</sup>

<sup>1</sup> Kyushu University

<sup>2</sup> University of Electro-Communications

<sup>3</sup> Hokkaido University

<sup>4</sup> Fujitsu

# Minimal Absent Words (MAWs) [1/2]

- A string  $w$  over an alphabet  $\Sigma$  is called a **Minimal Absent Word (MAW)** for a string  $S$ , if:
  1.  $w$  is a character from  $\Sigma$  not occurring in  $S$ , or
  2.  $w = aub$  ( $a, b \in \Sigma, u \in \Sigma^*$ ) does not occur in  $S$ , but both  $au$  and  $ub$  occur in  $S$ .

## Example

$w = b \ a \ b$

$S = a \ b \ a \ a \ b$

# Minimal Absent Words (MAWs) [1/2]

- A string  $w$  over an alphabet  $\Sigma$  is called a **Minimal Absent Word (MAW)** for a string  $S$ , if:
  1.  $w$  is a character from  $\Sigma$  not occurring in  $S$ , or
  2.  $w = aub$  ( $a, b \in \Sigma, u \in \Sigma^*$ ) does not occur in  $S$ , but both  $au$  and  $ub$  occur in  $S$ .

## Example

$w =$  b a b

$S =$  a b a a b



$\text{bab}$  is a MAW for  $\text{abaab}$

# Minimal Absent Words (MAWs) [2/2]

- MAW( $S$ ) denotes the set of MAWs for a string  $S$ .

## Example

$S = \text{a b a a b}$

$\Sigma = \{a, b, c\}$

$\text{MAW}(S) = \{\text{aaa}, \text{aaba}, \text{bab}, \text{bb}, c\}$

- The number  $|\text{MAW}(S)|$  of MAWs for a string  $S$  of length  $n$  over an alphabet of size  $\sigma$  is  $O(\sigma n)$ , and there is a matching lower bound [Crochemore et al. 1998].

# Motivations for Computing MAWs

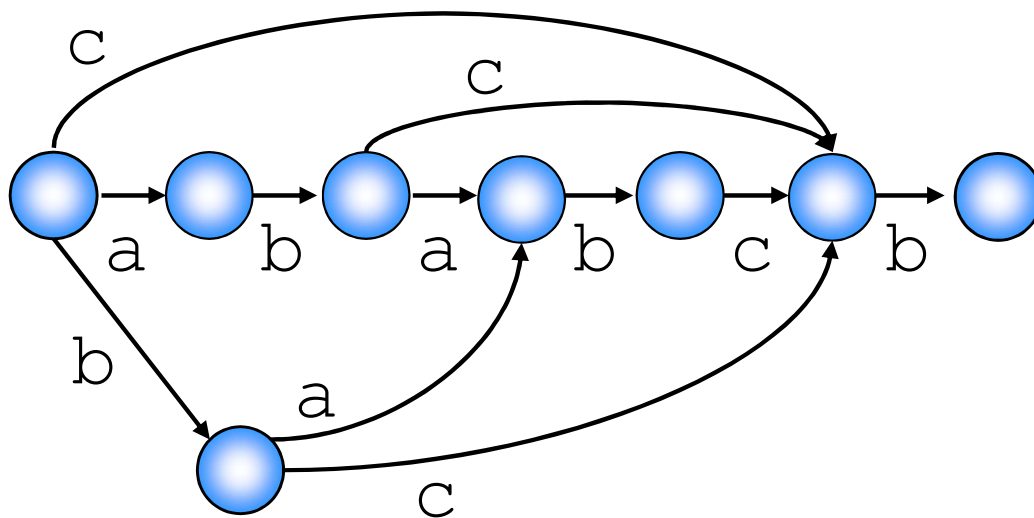
- Computing (minimal) absent words for a given string has various applications including:
  - ◆ **Data compression** (anti-dictionary)  
[Crochemore et al. 2000] [Crochemore & Navarro 2002]  
[Ayad et al. 2021]
  - ◆ **Sequence comparison**  
[Chairungsee & Crochemore 2012]  
[Charalampopoulos et al. 2018]
  - ◆ **Bioinformatics**  
[Almirantis et al. 2017] [Koulouras & Frith 2019]  
[Pratas & Silva 2020]

# Computing MAWs with DAWG [1/3]

- Previous algorithms [Crochemore et al. 1998, Fujishige et al. 2023] for computing MAWs for a string  $S$  of length  $n$  use **DAWG (Directed Acyclic Word Graph)** for  $S$ , which is an  $O(n)$ -size automaton representing all substrings of  $S$ .

[Blumer et al. 1985]

E.g.  $S = a b a b c b$



Substrings of  $S$  are represented by the same node of  $\text{DAWG}(S)$  iff they have the same ending position(s) in  $S$ .

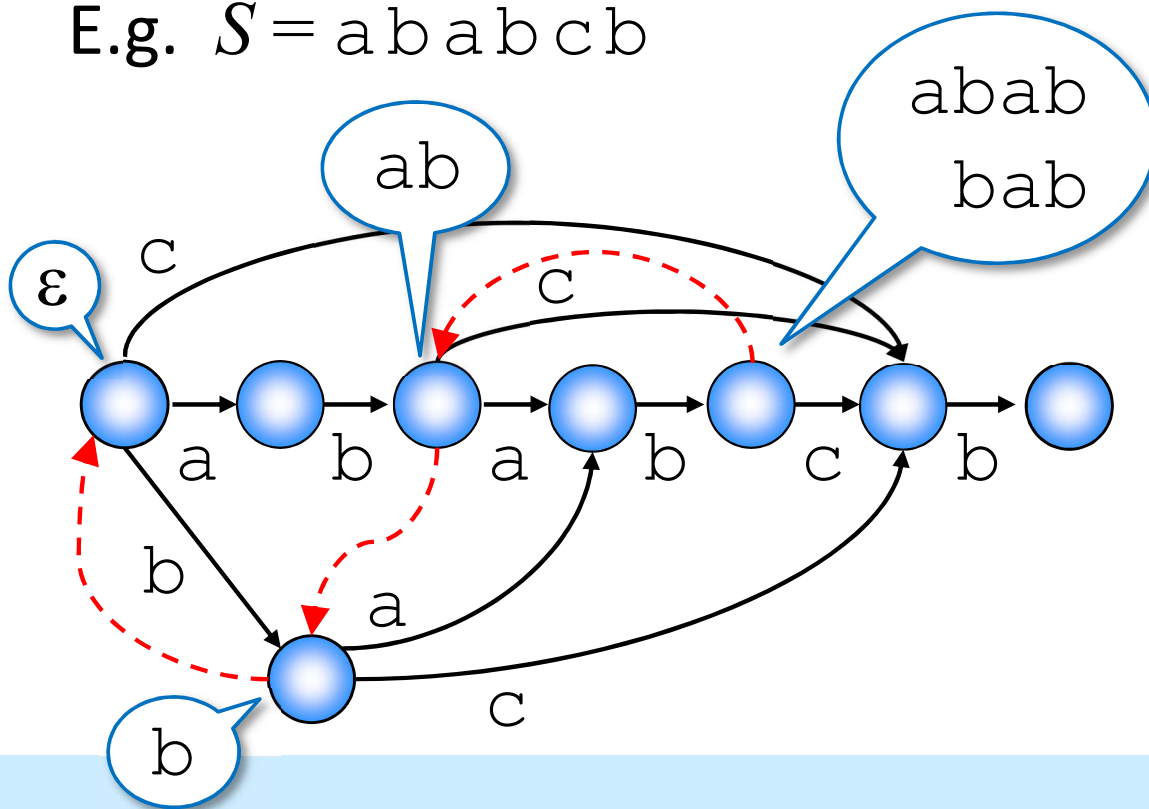
-----> suffix link

# Computing MAWs with DAWG [1/3]

- Previous algorithms [Crochemore et al. 1998, Fujishige et al. 2023] for computing MAWs for a string  $S$  of length  $n$  use **DAWG (Directed Acyclic Word Graph)** for  $S$ , which is an  $O(n)$ -size automaton representing all substrings of  $S$ .

[Blumer et al. 1985]

E.g.  $S = a b a b c b$



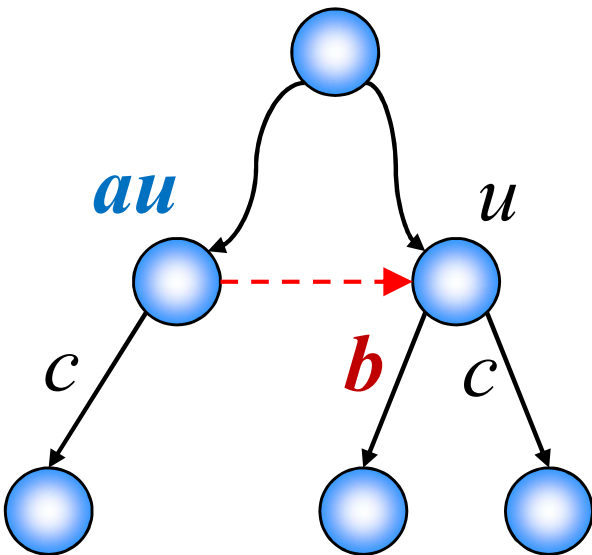
Substrings of  $S$  are represented by the same node of  $\text{DAWG}(S)$  iff they have the same ending position(s) in  $S$ .

---> suffix link

# Computing MAWs with DAWG [2/3]

- If the edges of DAWG are sorted, then one can compute  $\text{MAW}(S)$  in  $O(n + |\text{MAW}(S)|)$  time [Fujishige et al. 2023].

DAWG for string  $S$

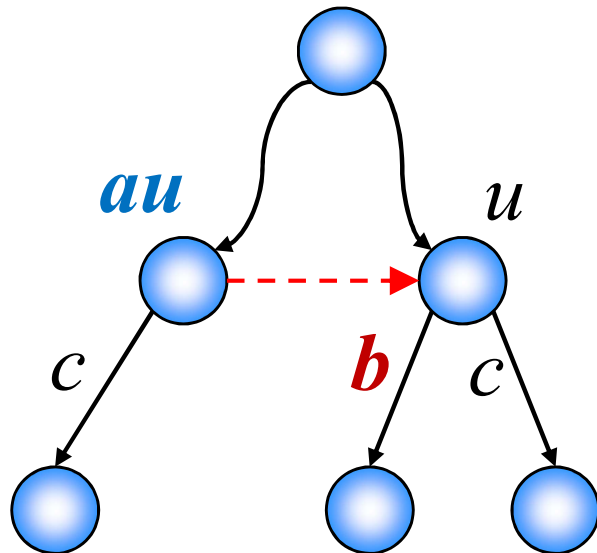


- Consider each pair of nodes  $au$  and  $u$  which are connected by a suffix link, where  $a$  is a character and  $u$  is a string.
- Compare the labels of the out-edges of nodes  $au$  and  $u$  in sorted order.
  - ◆ For  $b$ :  $au$  has no out-edge with  $b$ , but  $u$  has an out-edge with  $b$ .  
→  $aub$  is a MAW for the input string  $S$ .
  - ◆ For  $c$ : both  $au$  and  $u$  have out-edges with  $c$ .  
→  $auc$  is not a MAW for the input string  $S$ , but this cost of character comparisons can be charged to this out-edge of  $au$  labeled  $c$ .



# Computing MAWs with DAWG [3/3]

DAWG for string  $S$

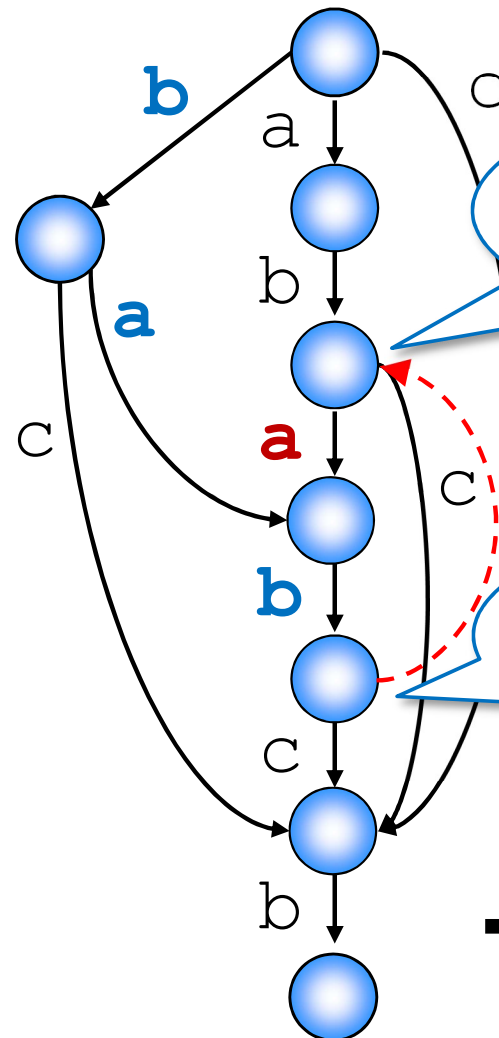


→  $aub$  is a MAW



suffix link

E.g.  $S = ababcb$



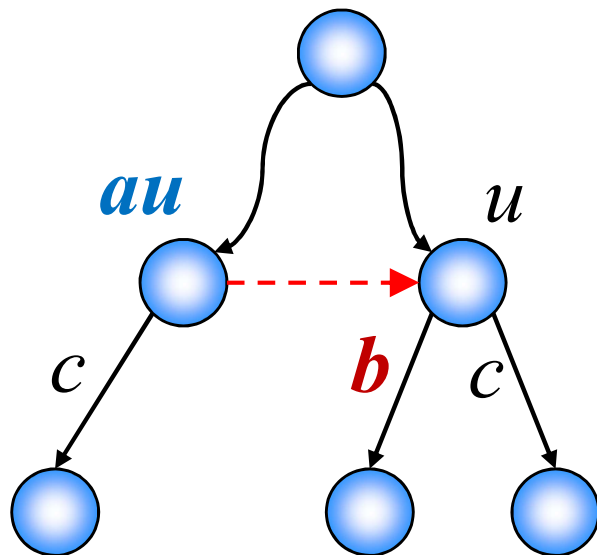
Longest string in this node is  $ab$

Shortest string in this node is  $bab$

→  $baba$  is a MAW

# Redundancy in MAW computation with DAWG

DAWG for string  $S$

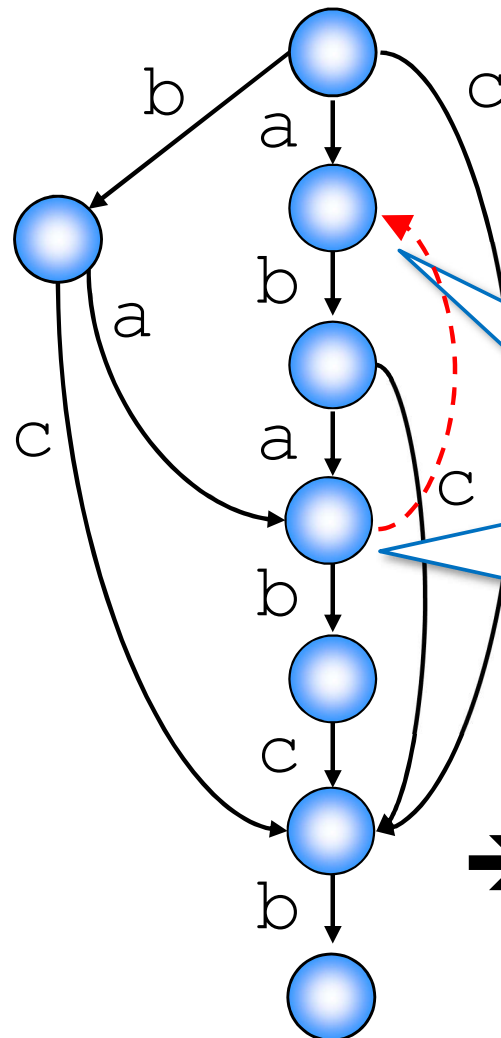


→  $aub$  is a MAW



suffix link

E.g.  $S = ababcb$

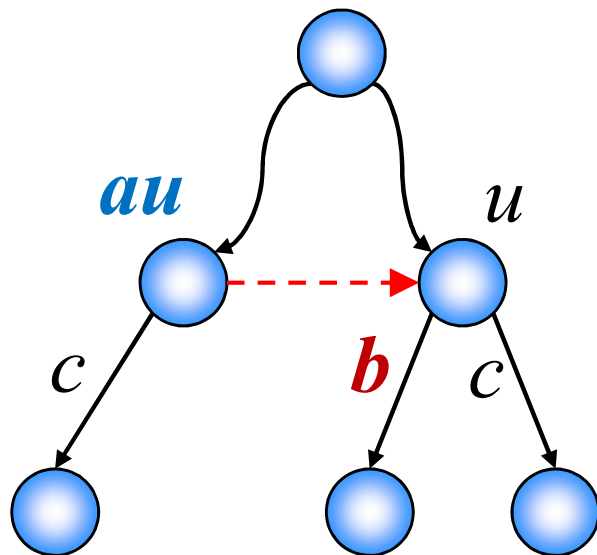


Both these nodes have only one child

→ No MAWs to output

# Redundancy in MAW computation with DAWG

DAWG for string  $S$

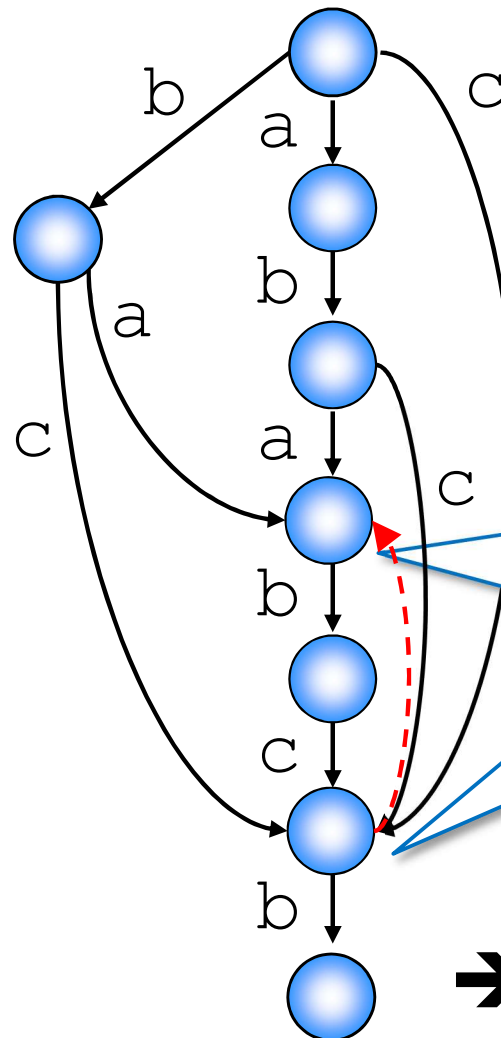


→  $aub$  is a MAW



suffix link

E.g.  $S = ababcb$



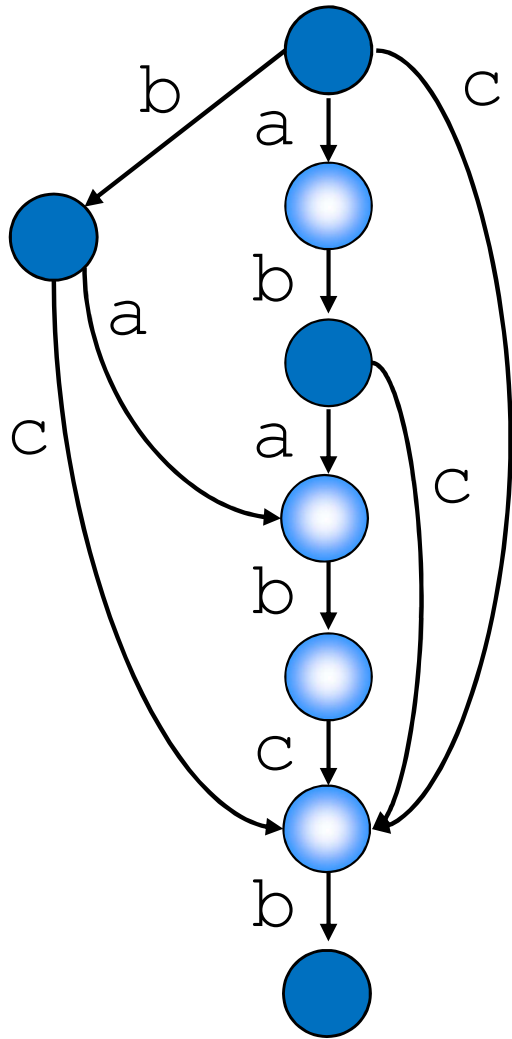
Both these nodes have only one child

→ No MAWs to output

# From DAWG to CDAWG (Compact DAWG)

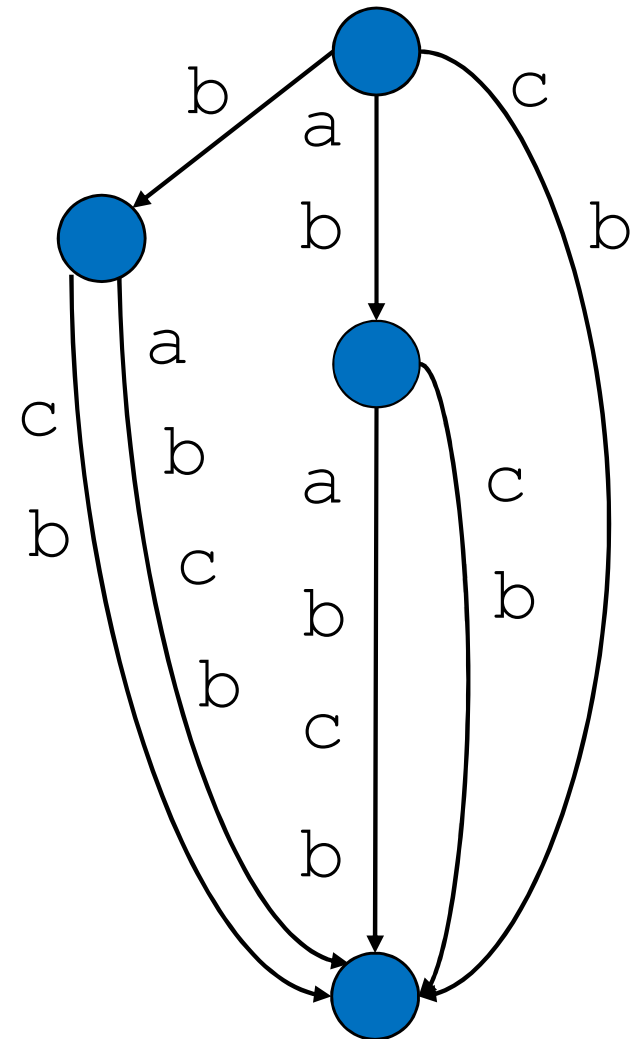
E.g.  $S = a b a b c b$

DAWG [Blumer et al. 1985]

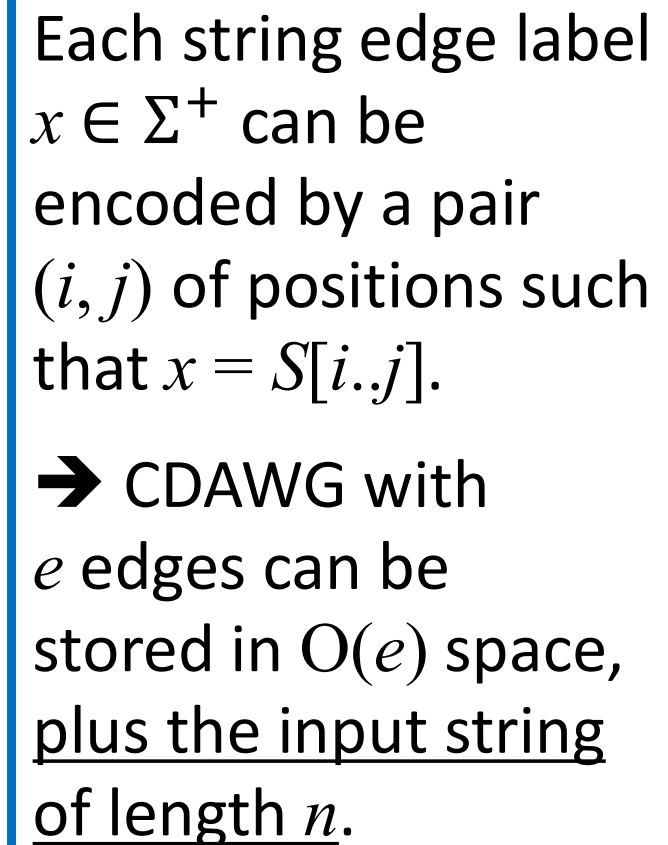


- (1) Keep only branching nodes + source & sink
- (2) Path contraction

CDAWG [Blumer et al. 1987]

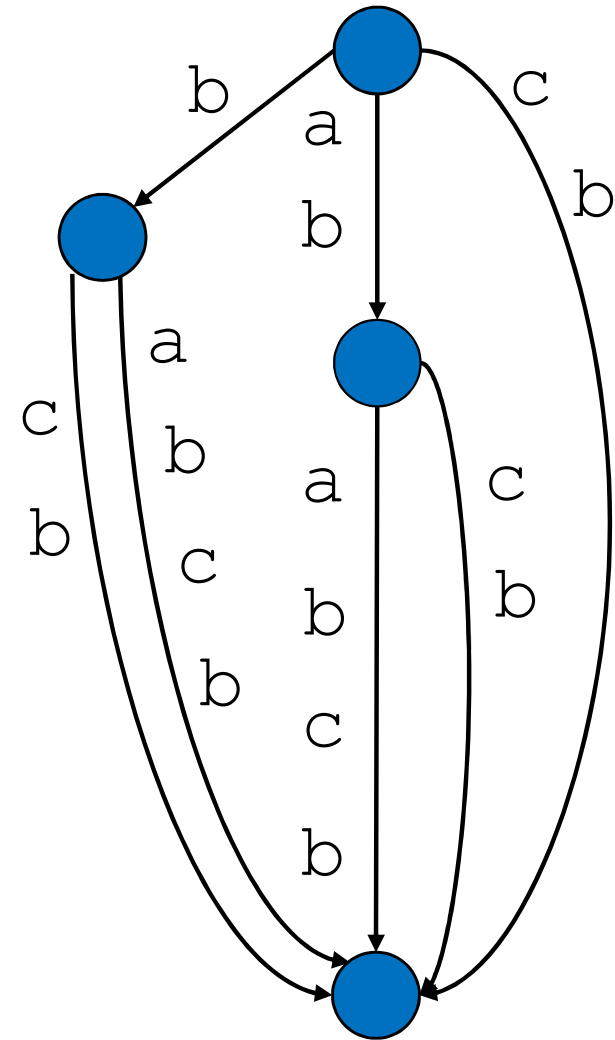


CDAWG for  $S = \begin{smallmatrix} & 1 & 2 & 3 & 4 & 5 & 6 \\ a & b & a & b & c & b \end{smallmatrix}$



# Encoding CDAWG in $O(e)$ space [2/2]

$S =$ ~~ababcb~~



## Theorem 1

[Belazzougui & Cunial 2015, Inenaga 2024]

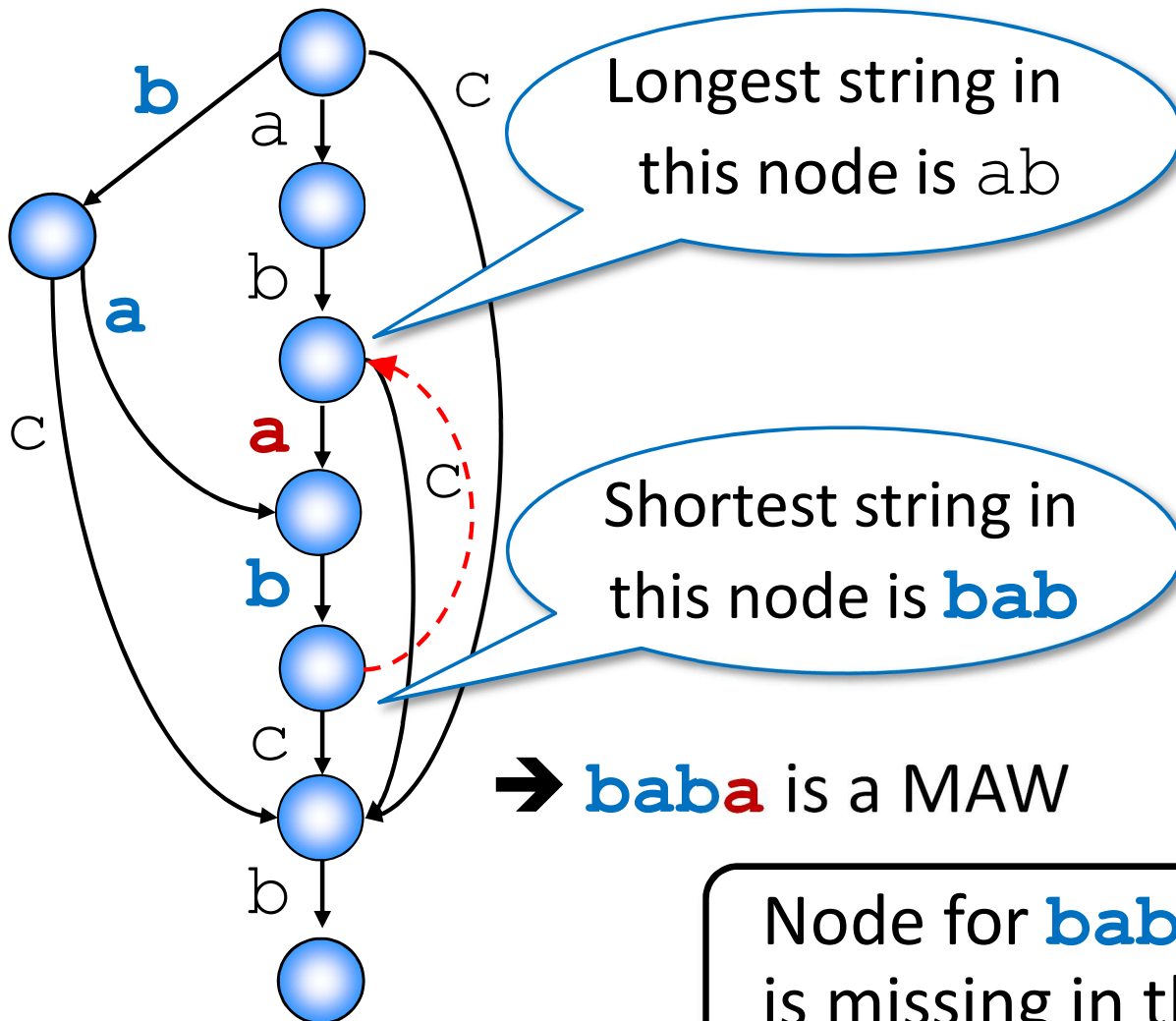
It is possible to store the edge labels of the CDAWG in  $O(e)$  total space, so that each edge string label  $x$  can be retrieved in  $O(|x|)$  time, where  $e$  is the number of edges in the CDAWG.

Note:  $e$  can be as small as  $O(\log n)$  for some highly repetitive strings.

# Computing MAWs with CDAWG [1/3]

E.g.  $S = a b a b c b$

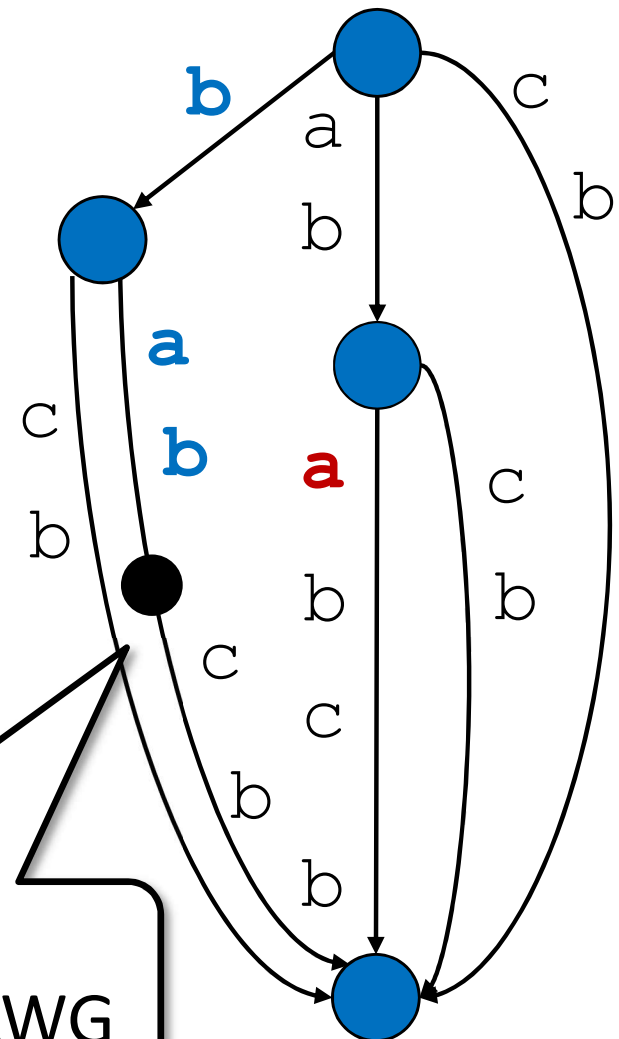
DAWG



→ **bab**a is a MAW

Node for **bab**  
is missing in the CDAWG

CDAWG



# Computing MAWs with CDAWG [2/3]

Theorem 2 [Belazzougui & Cunial 2015]

The number of “missing” nodes in the CDAWG of string  $S$  for computing  $\text{MAW}(S)$  is at most the number  $\bar{e}$  of edges of the CDAWG for the *reversed string*  $\bar{S}$ .

Theorem 3 [Belazzougui & Cunial 2015]

There exists a data structure of  $O(e + \bar{e})$  space that can output a representation of  $\text{MAW}(S)$  for the input string  $S$  in  $O(e + \bar{e} + |\text{MAW}(S)|)$  time.

How much is the gap between  $e$  and  $\bar{e}$  ?



# Computing MAWs with CDAWG [3/3]

How much is the gap between  $e$  and  $\bar{e}$  ?

Theorem 4 [Karkkainen 2017, Inenaga 2024]

There exists a family of strings of length  $n$  such that  $\bar{e}/e = \Omega(\sqrt{n})$ .

Theorem 5 [Inenaga, Kosolobov, Zuba 2024 (unpublished)]

For any string of length  $n$ ,  $\bar{e}/e = O(\sqrt{n})$ .

# Our Contribution

## Theorem 6

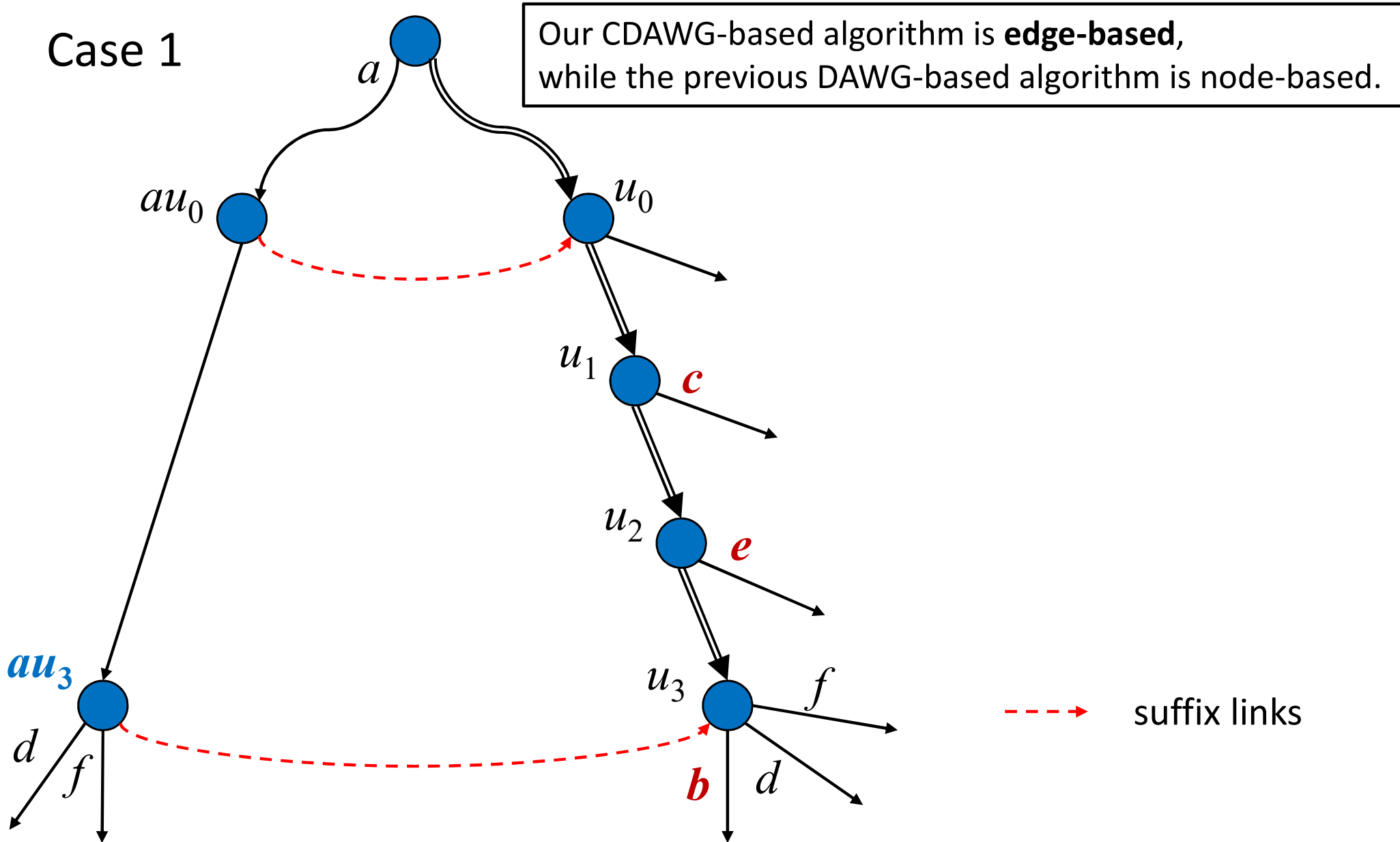
There exists a data structure of  $O(e_{\min})$  space that can output a representation of  $\text{MAW}(S)$  for string  $S$  in  $O(|\text{MAW}(S)|)$  time, where  $e_{\min} = \min\{e, \bar{e}\}$ .

algorithm	time	space
Crochemore et al. 1998	$O(\sigma n)$	$O(n)$
Fujishige et al. 2023	$O(n +  \text{MAW}(S) )$	$O(n)$
Belazzougui & Cunial 2015	$O(e + \bar{e} +  \text{MAW}(S) )$	$O(e + \bar{e})$
<b>Ours</b>	<b><math>O( \text{MAW}(S) )</math></b>	<b><math>O(e_{\min})</math></b>

Note:  $e$  and  $\bar{e}$  are at most  $2n-1$ , and can be as small as  $O(\log n)$  for some highly repetitive strings (e.g. Fibonacci strings).

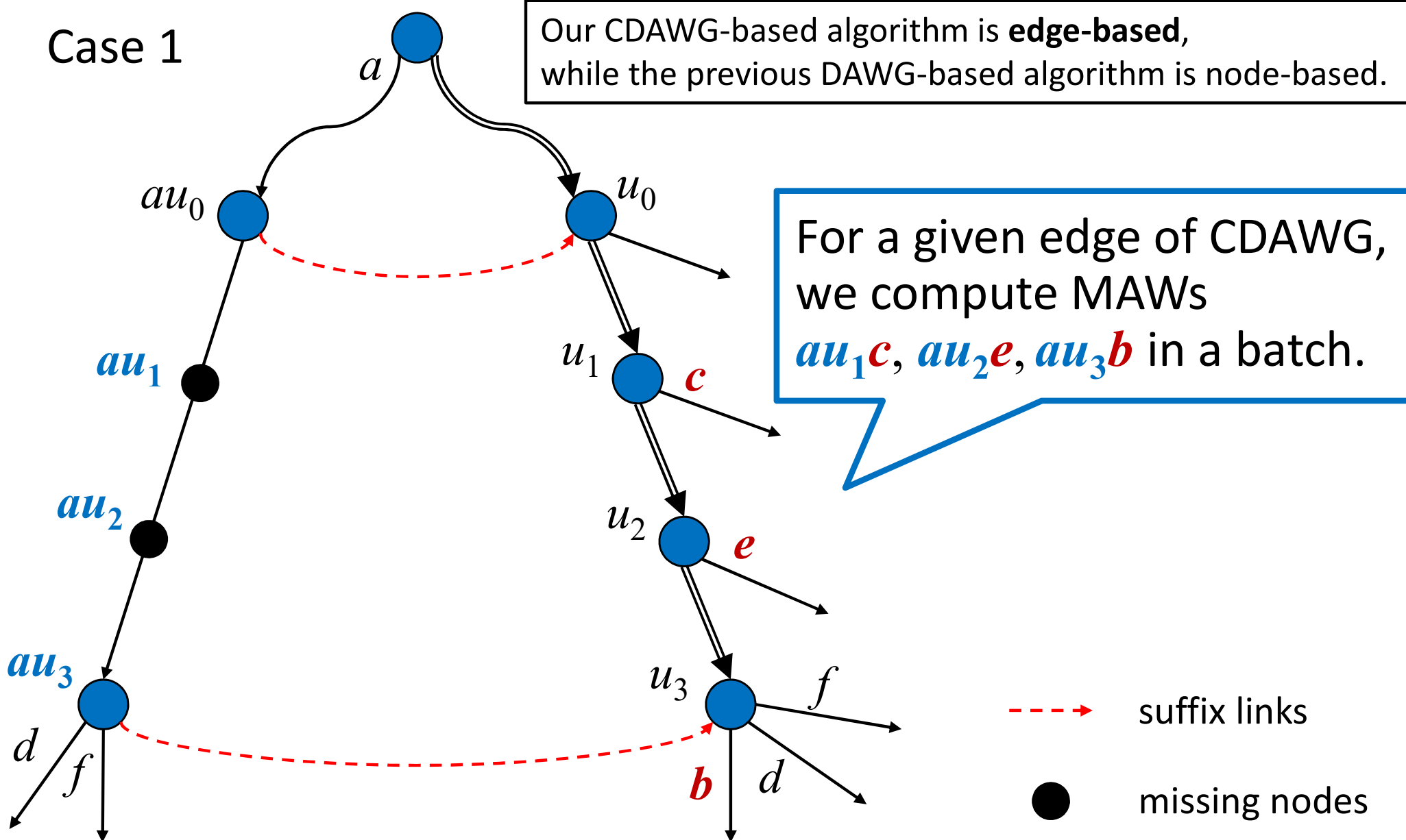
# Computing MAWs in $O(e_{\min})$ space [1/4]

Case 1



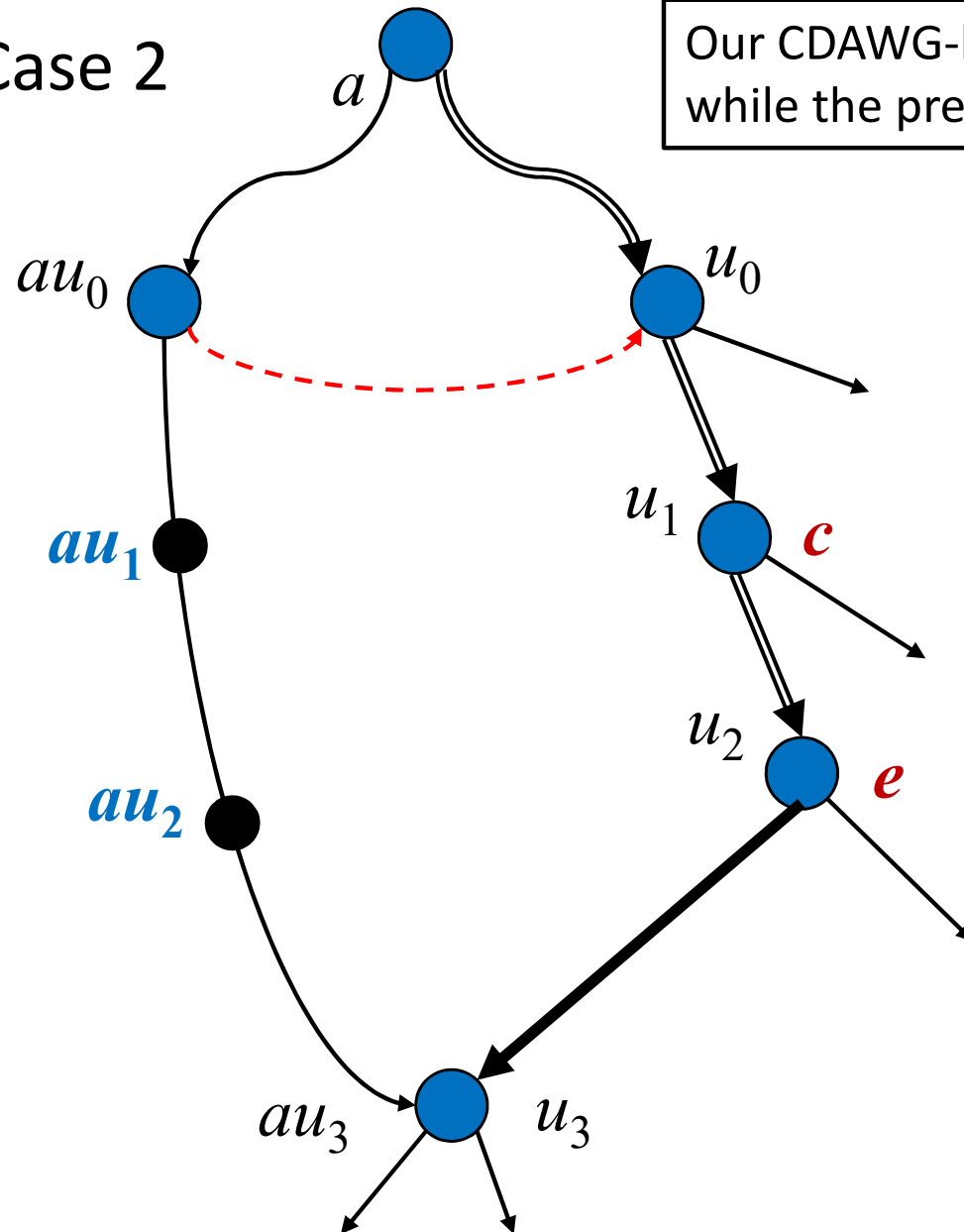
# Computing MAWs in $O(e_{\min})$ space [1/4]

Case 1



# Computing MAWs in $O(e_{\min})$ space [2/4]

Case 2



Our CDAWG-based algorithm is **edge-based**, while the previous DAWG-based algorithm is node-based.

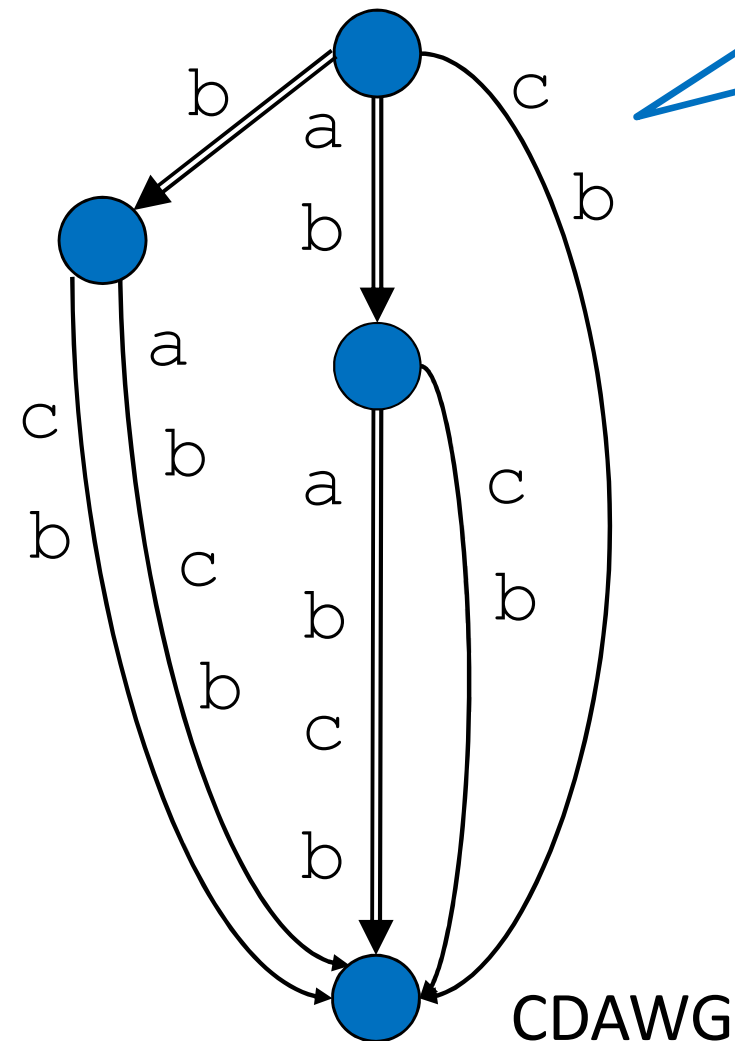
For a given edge of CDAWG, we compute MAWs  $au_1c$  and  $au_2e$  in a batch.

---> suffix links  
● missing nodes

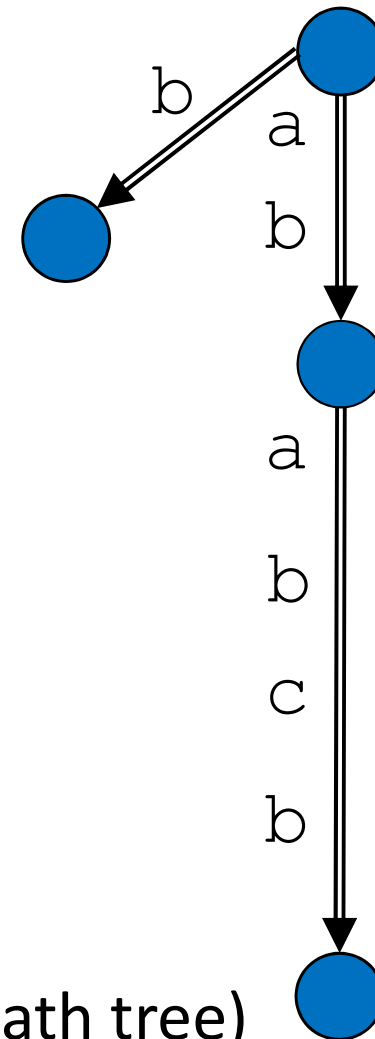
# Computing MAWs in $O(e_{\min})$ space [3/4]

$S = a b a b c b$

The **primary edges** are those that are in the **longest paths** from the source to the nodes.



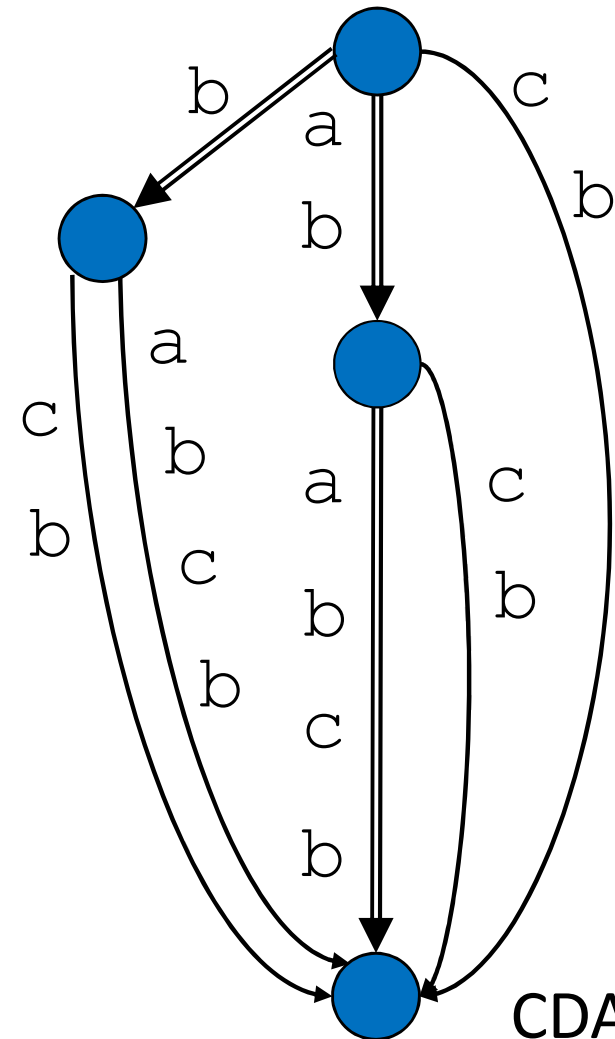
LPT  
(longest path tree)



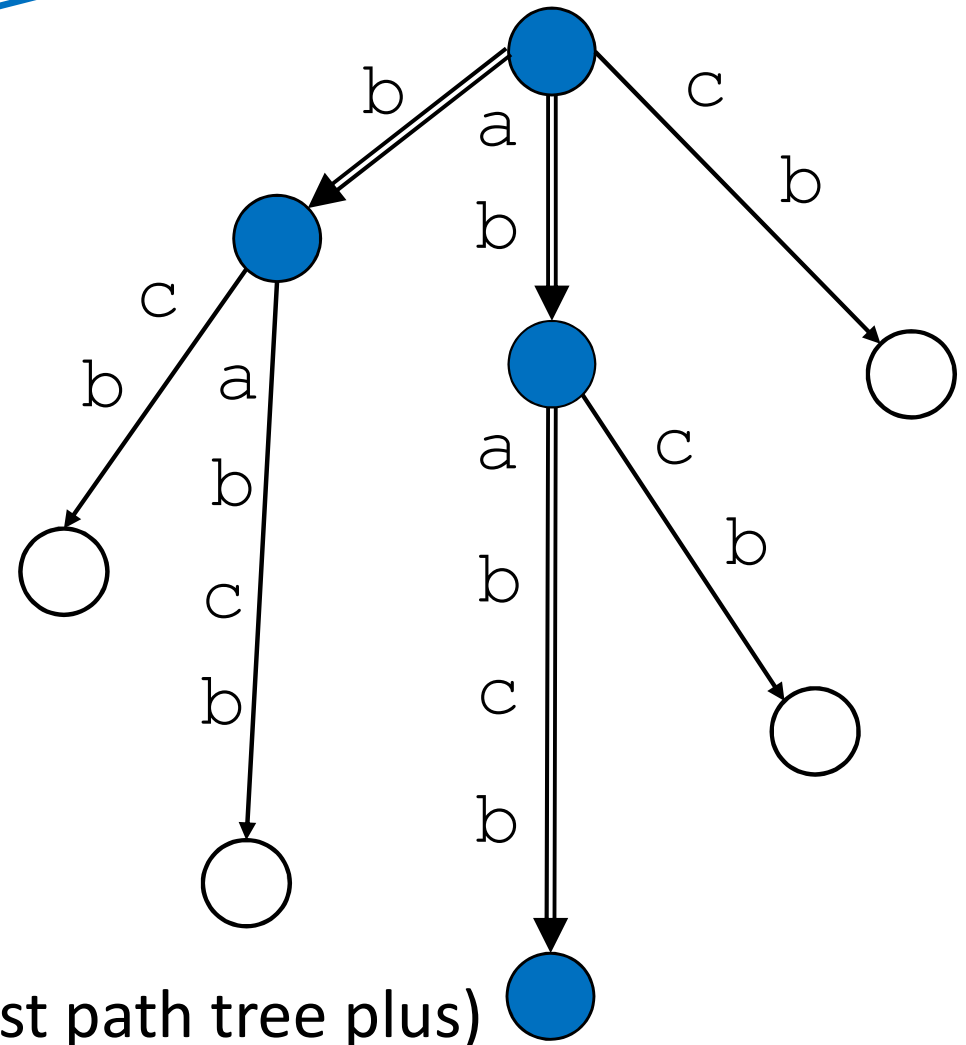
# Computing MAWs in $O(e_{\min})$ space [3/4]

$S = a b a b c b$

The **primary edges** are those that are in the **longest paths** from the source to the nodes.



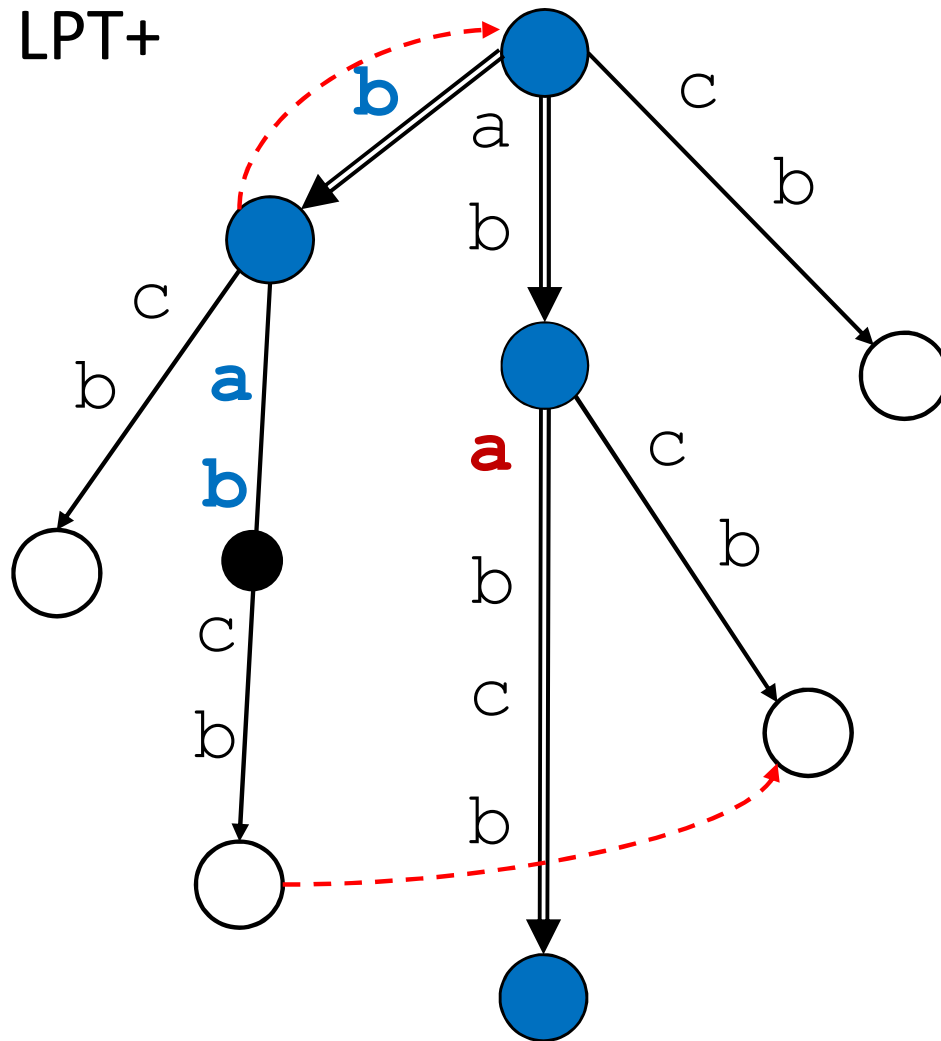
equivalent



# Computing MAWs in $O(e_{\min})$ space [4/4]

$S = a b a b c b$

LPT+



$\Rightarrow$  primary edges

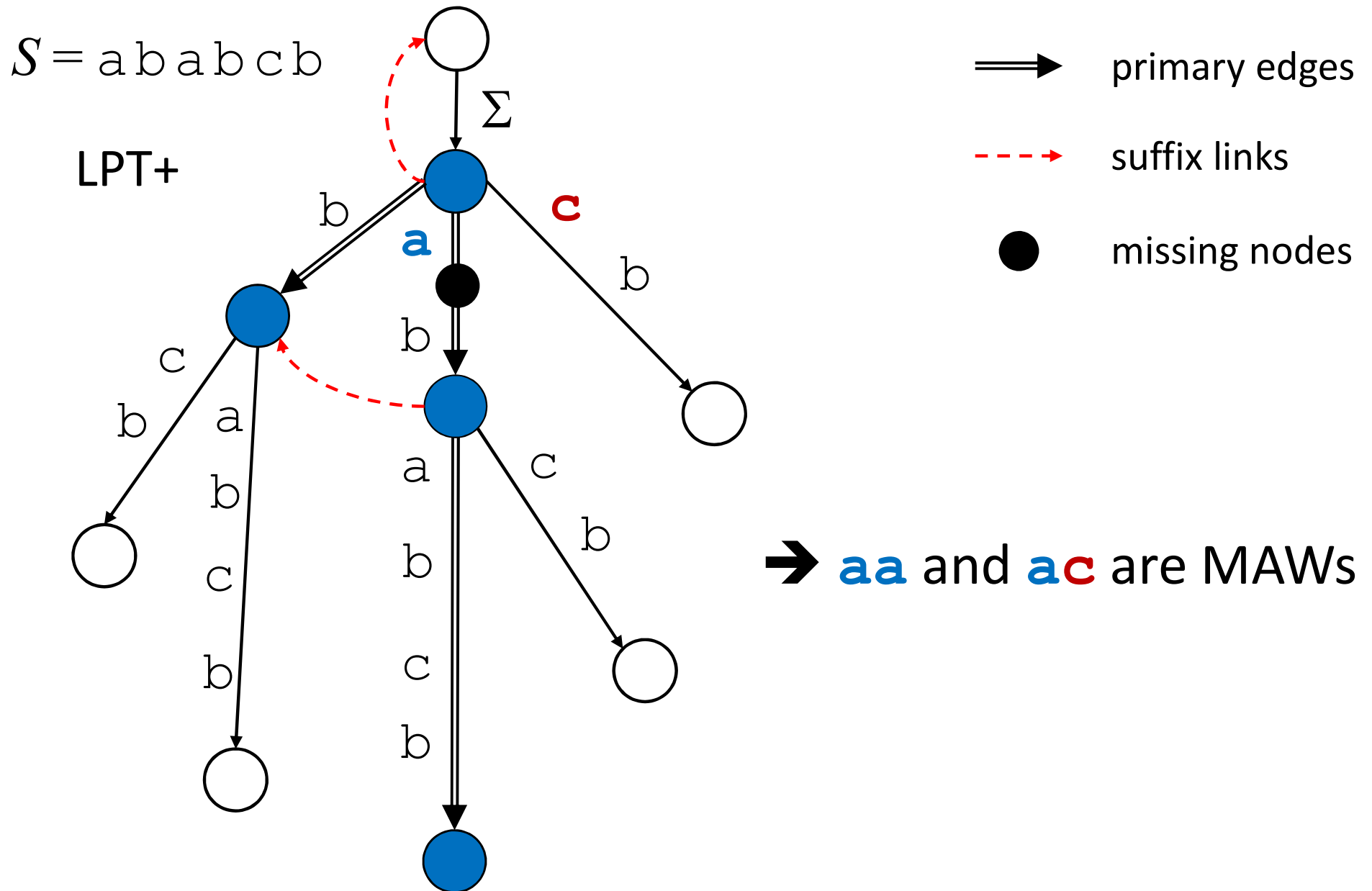
$\cdots \rightarrow$  suffix links

● missing nodes

$\Rightarrow$  **b****a****b****a** is a MAW



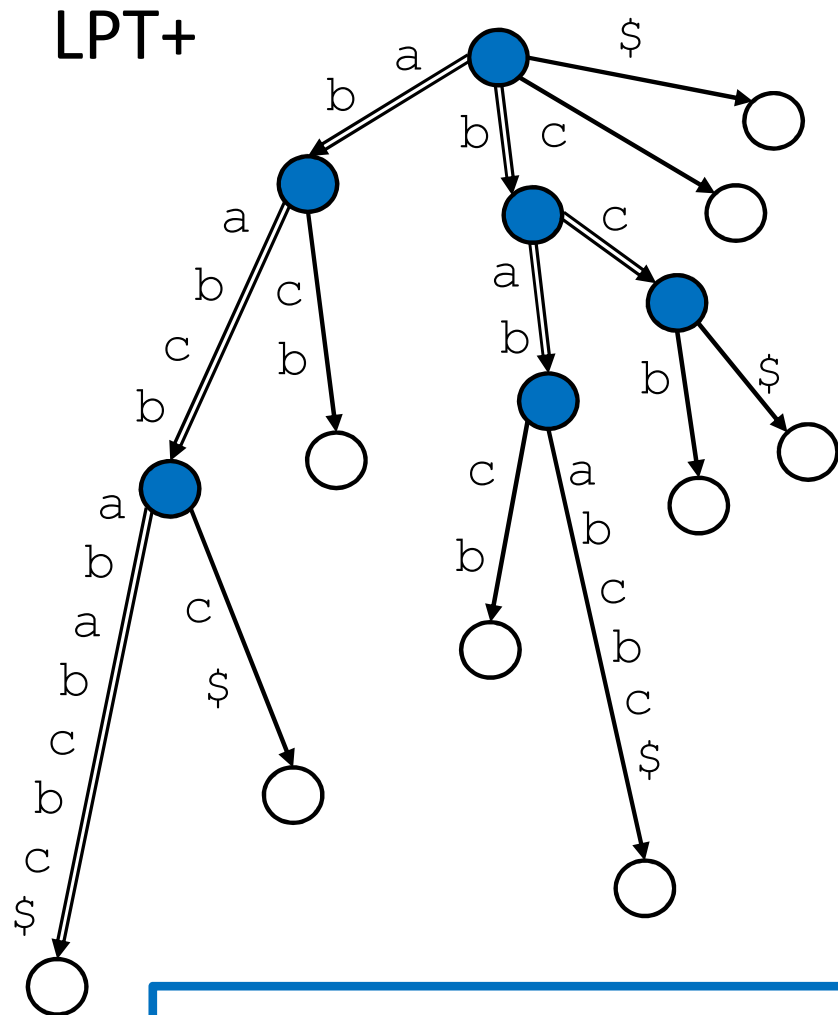
# Computing MAWs in $O(e_{\min})$ space [4/4]



# LPT+ $\neq$ Suffix Tree

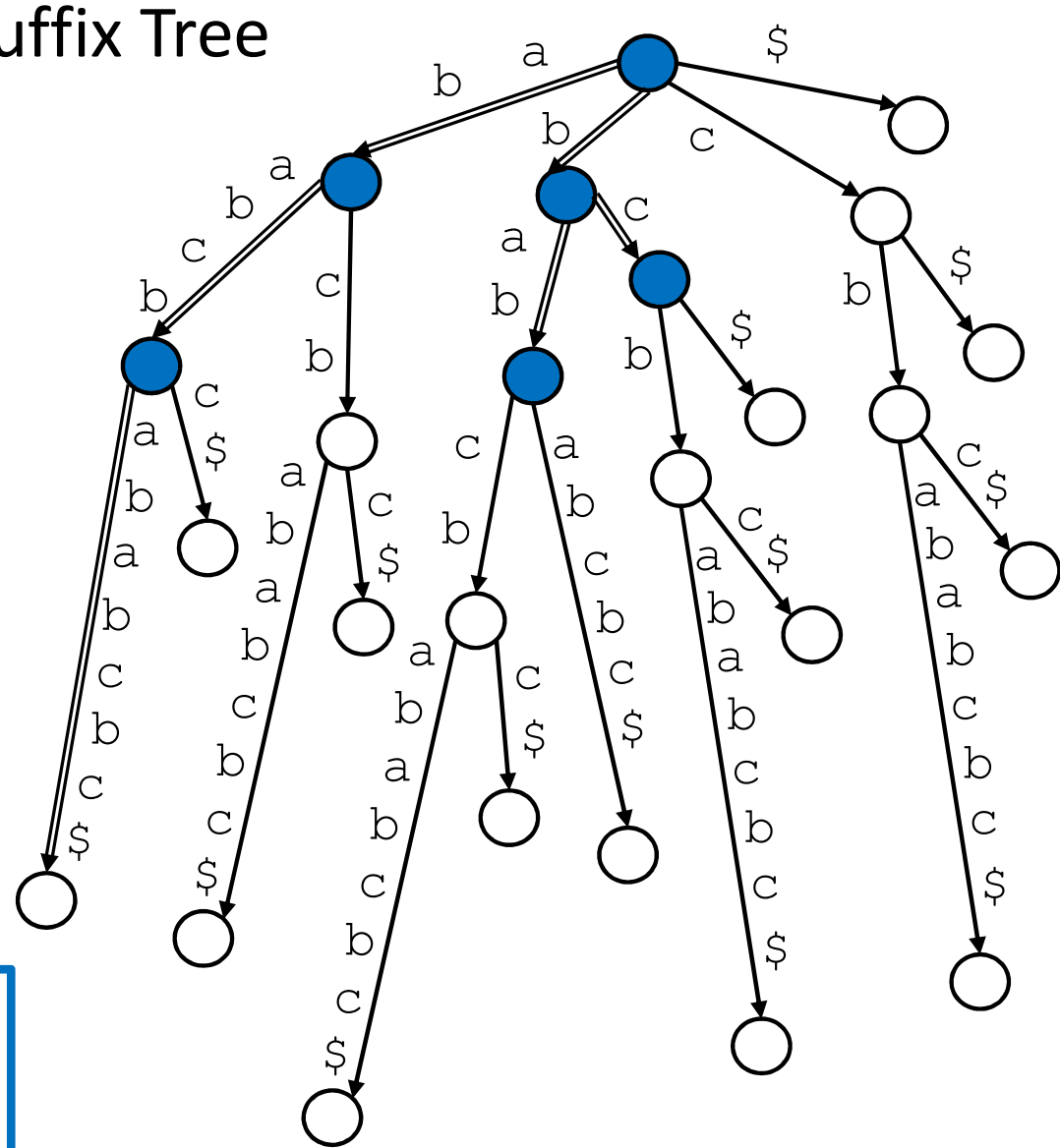
$$S = \text{a b a b c b a b a b c b c \$}$$

LPT+



LPT+ is of size  $\Theta(e)$  while  
Suffix Tree is of size  $\Theta(n)$ .

# Suffix Tree



# Conclusions

- We have shown a **CDAWG**-based data structure of  $O(e)$  space that can report all **MAWs** in optimal time.
- Our data structure can be extended to report **minimal rare words (MRW)** in optimal time:  
A string  $aub$  is a  $k$ -MRW of a string  $S$  if  $aub$  occurs  $k$  times in  $S$ , while  $au$  and  $ub$  occur more than  $k$  times in  $S$ .
  - ◆ **0-MRWs are MAWs.**
  - ◆ **1-MRWs are MUSs** (minimal unique substrings) [Ilie & Smyth 2011]

# Future Work

- Is there a data structure of  $O(r)$  space that reports all MAWs/MRWs in optimal time, where  $r$  is the number of equal-letter runs in the **BWT** (Burrows-Wheeler transform)?
- It is known that  $r \leq e$  [Belazzougui & Cunial 2017].
- There is a data structure of  $O(r \text{ polylog}(n))$  space that can do this in  $O((e + \bar{e})\text{polylog}(n) + \text{output})$  time with  $O(\sigma \log n)$  working space (under review).
- We want to do this in optimal  $O(\text{output})$  time with  $O(r + \bar{r})$  space.