



# Optimally Computing Compressed Indexing Arrays Based on the Compact Directed Acyclic Word Graph

Hiroki Arimura<sup>1</sup>

Shunsuke Inenaga<sup>2</sup>

Yasuaki Kobayashi<sup>1</sup>

Yuto Nakashima<sup>2</sup>

Mizuki Sue<sup>1</sup>

1) Graduate School of IST,  
Hokkaido University, Japan

2) Department of Informatics,  
Kyushu University, Japan

The longer version: <https://arxiv.org/abs/2308.02269>

This slide pdf: <https://ikndeva.github.io> (or arXiv entry's "Code, Data, Media/Paper with Code" section)

This work is partly supported by MEXT Grant-in-Aid for Basic Research A, 2000-2004, Japan

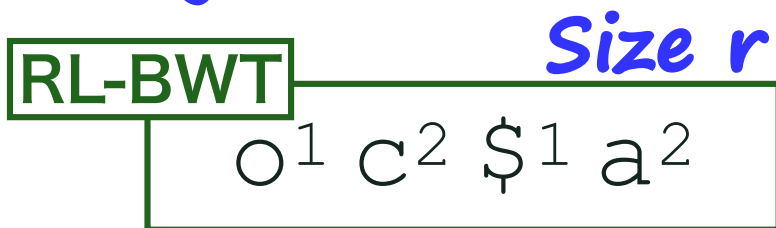


- Increasing amount and types of **repetitive texts**
  - Markup texts (Wikipedia), Genome sequences
- Development of **compressed index structures** for these **repetitive texts** attracts much attention.
- Such indices can compress highly-repetitive texts beyond the entropy bounds up to “compression parameters” – the sizes of indices
- We focus on the relationship between compression structures

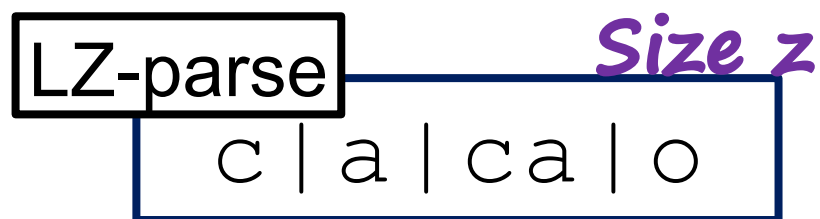
# Three compressed index structures

index	0	1	2	3	4	5
Text $T$	c	a	c	a	o	\$

- **RL-BWT** is obtained from SA by taking the **preceding letter** and **run-length encoded**

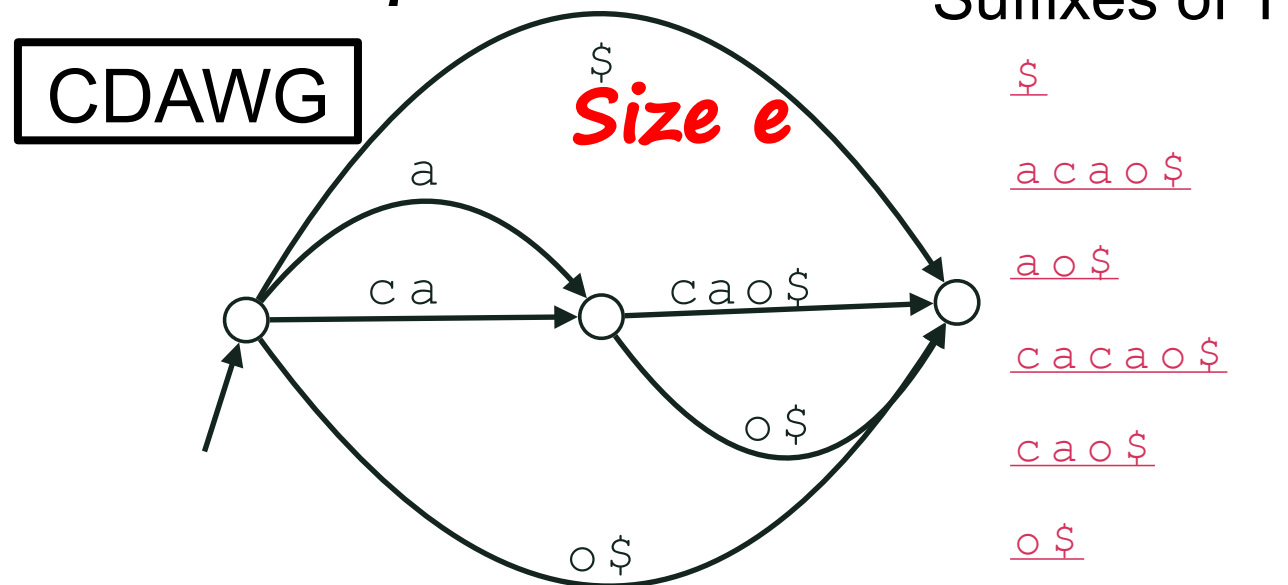


- The **LZ-parse** is obtained by partitioning  $T$  into the **longest previous factors** (PLFs)



- The **CDAWG** (Compact Directed Word Graphs) for a text  $T$  is an **automata-based index in a DAG form**

- It is obtained from the **Suffix Tree** of  $T$  by merging isomorphic subtrees



# Three compressed index structures

index	0	1	2	3	4	5
Text $T$	c	a	c	a	o	\$

- The **CDAWG** (Compact Directed Word Graphs) for a text  $T$  is an automata-based index in a DAG form

- **RL-BWT** is obtained from SA by taking the run-length

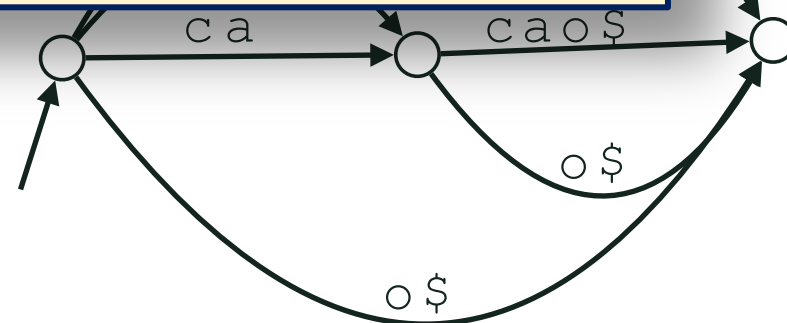
RL-BWT

We are interested in the size-relation and computational complexities of **conversion** between them.

the Suffix

Suffixes of  $T$

\$  
a c a o \$  
a o \$  
c a c a o \$  
c a o \$  
o \$



- The LZ-parse partitioning previous factors (PLFs)

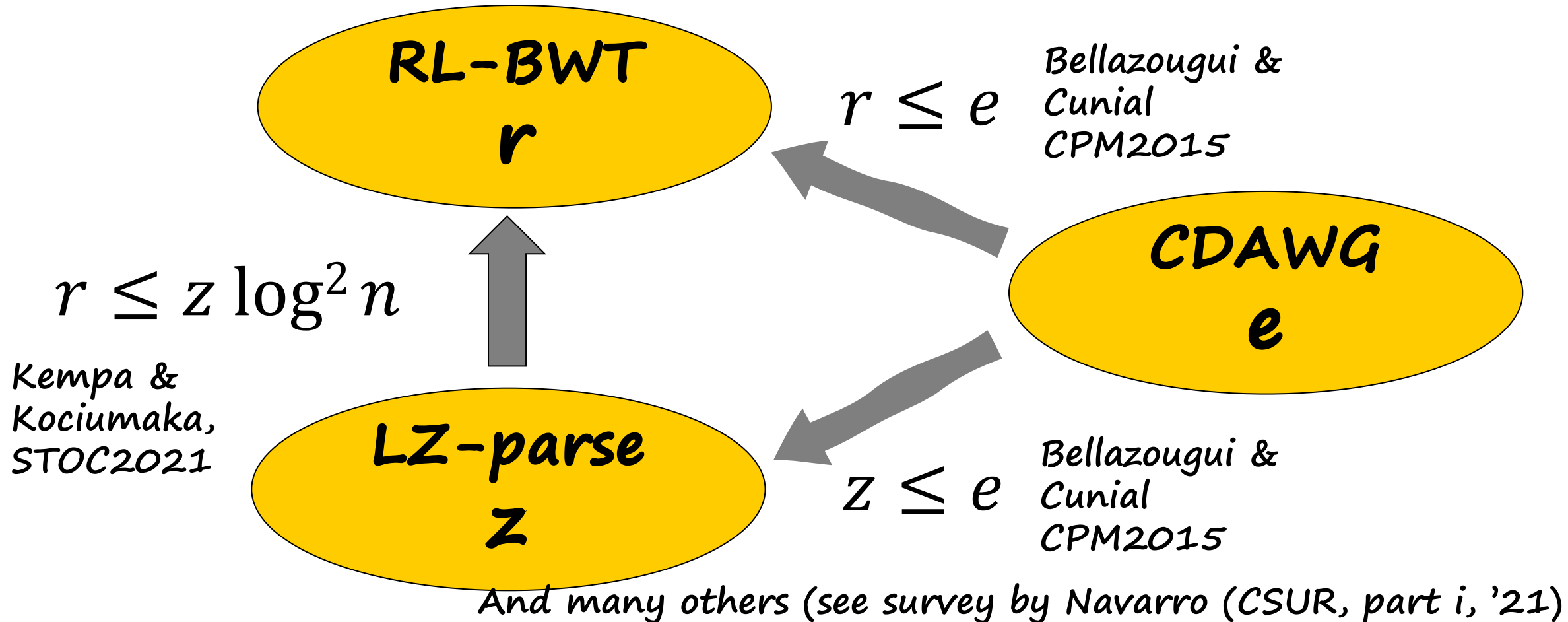
LZ-parse

Size  $z$

c | a | c a | o

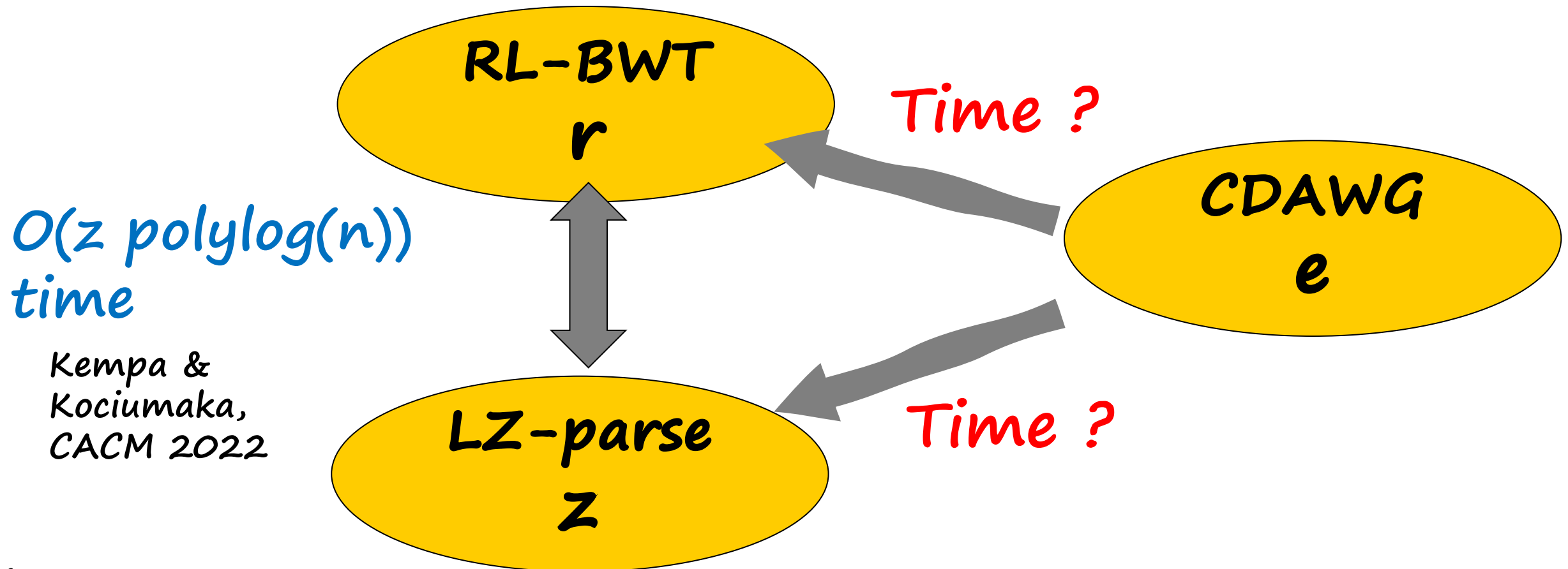


- Consider the relationship between their sizes
  - has been studied so far.



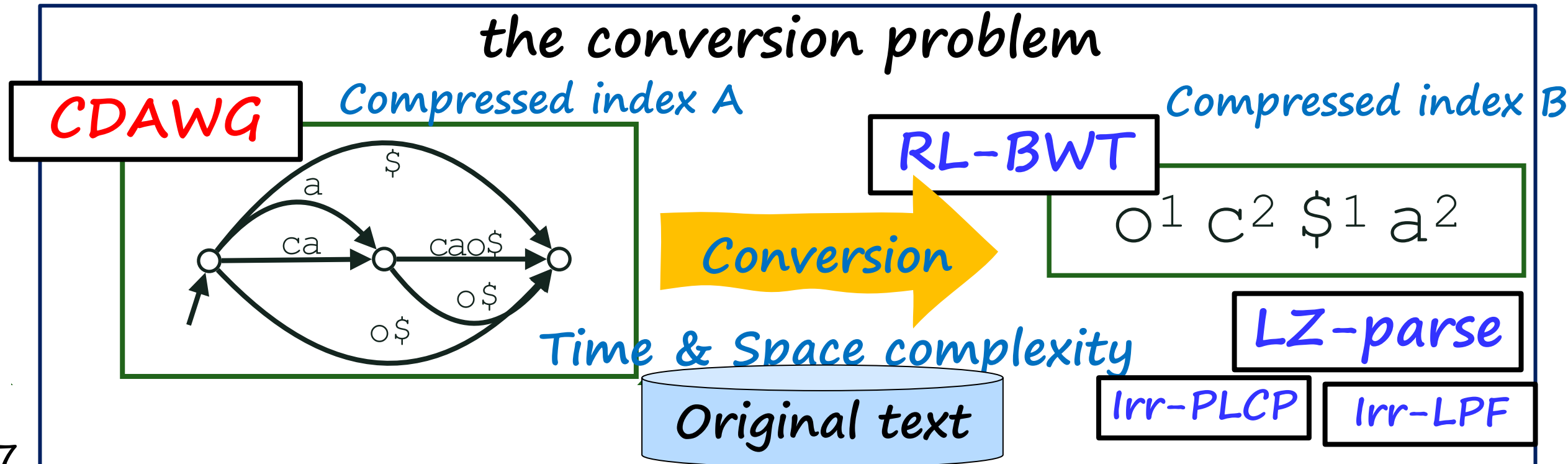


- The time and space complexities of conversions
  - have not been studied very much



# Research Goal:

We devise efficient algorithms that solves the **conversion problem** from the **CDAWG** for a text  $T$  into various compressed indexes for  $T$  in linear time and space in the combined input/output sizes







## Sublinear time and space conversion between two indices

### ■ Kempa [SODA'19]

- Converting an RL-BWT-based index into the irreducible PLCP, CSA, and LZ-parse for a text  $T$  of length  $n$  in  $O(n / \log_\sigma n + r \text{ polylog } n)$  time and  $O(r)$  space.

### ■ Kempa & Kociumaka [STOC'21, CACM'22]

- Converting the LZ77-parse of a text  $T$  into the RL-BWT for  $T$  in  $O(z \text{ polylog } n)$  time and space.
- This work solved a long-standing open problem

### ■ Bannai et al. [CPM'13]

- Converting an SLP of size  $g$  into LZ78-parse of size  $z_{78}$  in  $O(g + z_{78} \log z_{78})$  time and space.
- Combined with Belazzougui & Cunial [CPM'15], we obtain the conversion from the CDAWG for  $T$  into LZ78-parse in  $O(e + z_{78} \log z_{78})$  time and space.





Thm (4.1, 5.1, 5.2): For any integer alphabet  $\Sigma$ , we can convert **the CDAWG  $G$  of size  $e$  for a text  $T$**  into the following compressed indexing structures for  $T$  in  $O(e)$  deterministic time and words of space:

- **The RL-BWT** (run-length BWT) of size  $r$
- The irreducible PLCP (permuted LCP) array of size  $r$
- The quasi-irreducible LPF (longest previous factor) array of size  $e$  (def. Sec. 2 of this paper)
- The Lex-parse of size  $2r = O(r)$
- **The LZ-parse** of size  $z$

$G$  is given in either

- the CDAWG of size  $e$  with the read only text of length  $n$ ,
- the self-index version of CDAWG of size  $O(e)$  without a text

# Algorithms



Coming back to the relationship  
between the sizes ...

**Observation:** The proof by  
Bellazougui & Cunial (2015)  
is done by relating “ $r$ ” and  
“ $z$ ” to  $O(e)$  secondary  
incoming/ outgoing edges of  
 $CDAWG(T)$

LZ-parse  
 $z$

$$r \leq e$$

Bellazougui &  
Cunial  
CPM2015

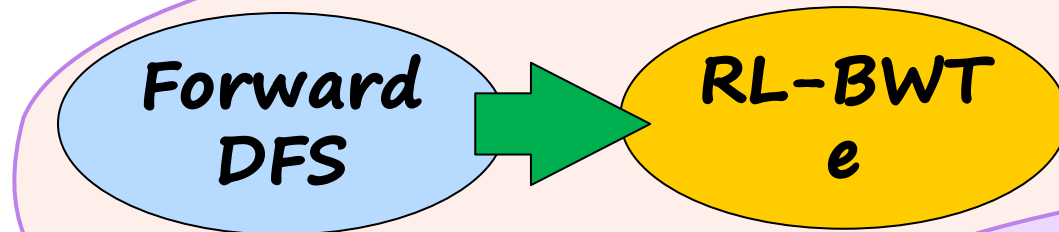
**CDAWG**  
 $e$

$$z \leq e$$

Bellazougui &  
Cunial  
CPM2015

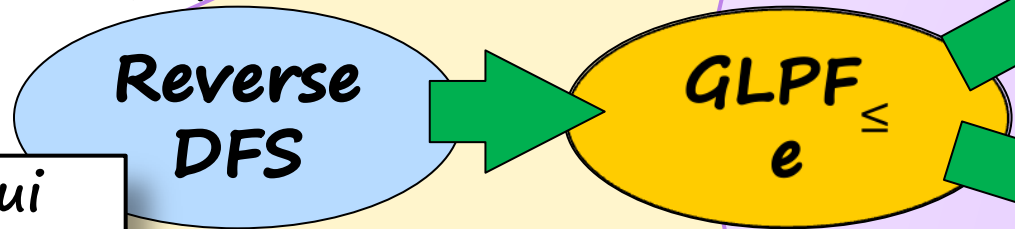
- We use **two orders** of paths
- Order for defining **2ndary edges**
- One for **traversal** of CDAWG

- Order for traversal



Ordered DFS from the source in the lexicographic order

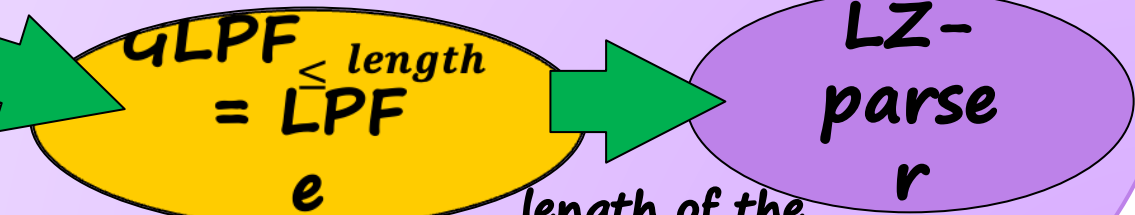
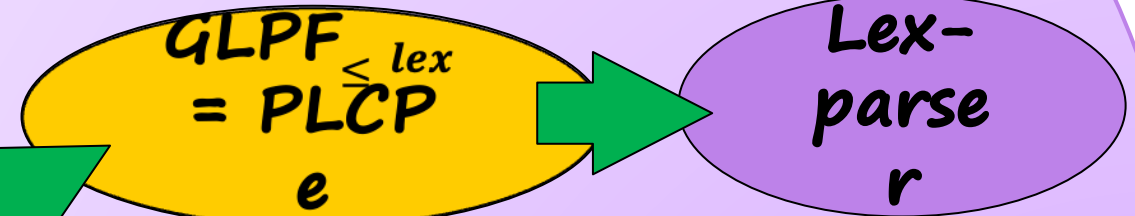
2ndary edge  $\approx$  same-letter run



Ordered DFS from the sink in the text order

Generalized Longest Previous Factor Array [This work]

- Order for 2ndary edges



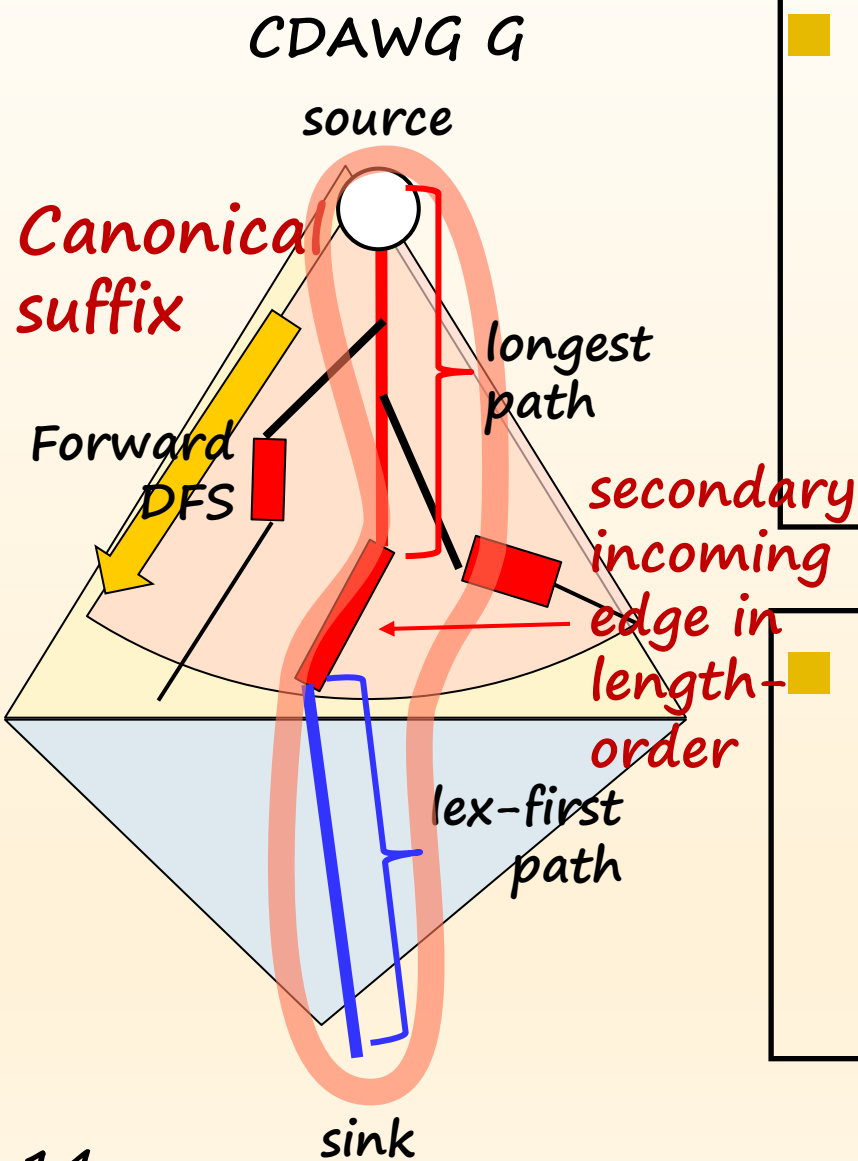
length of the longest upper path  $\approx$  irreducible GLPF-value

Bellazougui & Cunial CPM2015

Navarro, Ochoa, & Prezza (Trans. Inf. Theory, '20).

- We generalize PLCP & LPF into GLPF by the framework of (NOP'20)

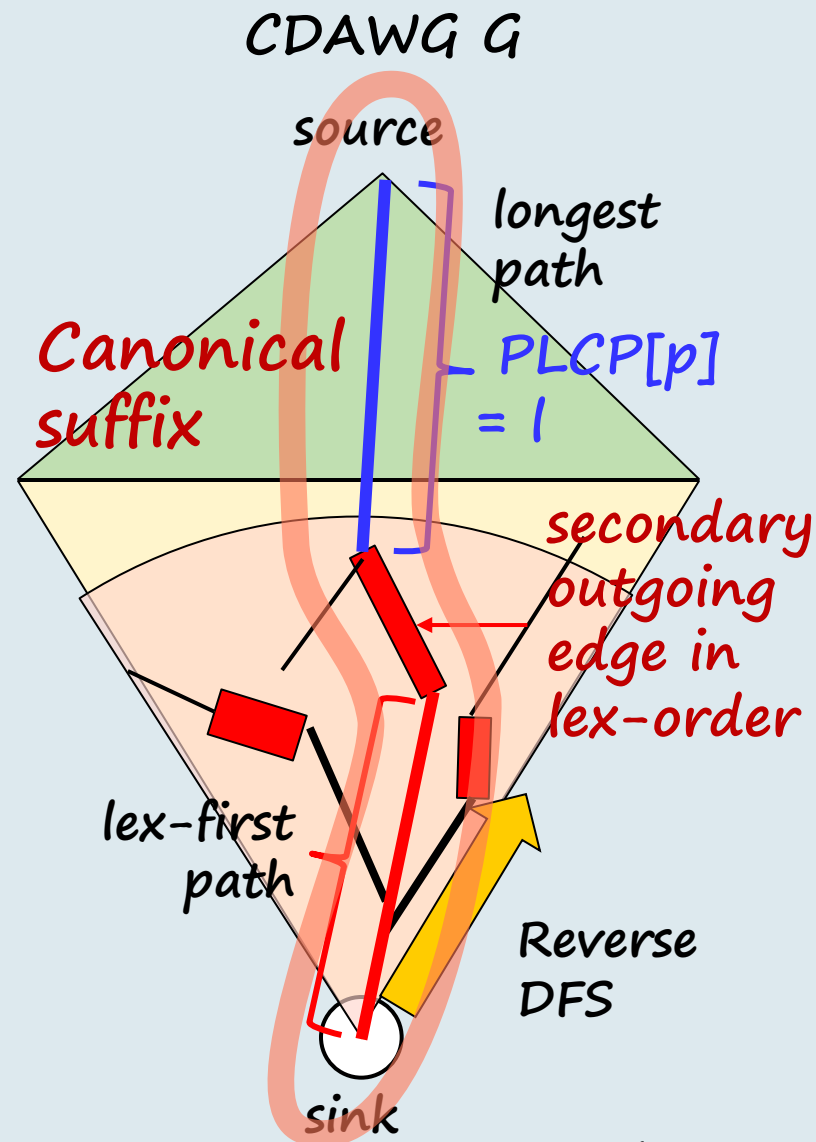
# Sec4. Computing RL-BWT in $O(e)$ time&space



■ Observation A1:  $O(e)$  secondary incoming edges of  $CDAWG(T)$  under the length-order correspond to subintervals of the same-letter runs of the BWT.  
(this is because such a search path defines a non-left-maximal factor in  $T$ )

■ Observation A2:  $O(e)$  incoming edges of  $CDAWG(T)$  can be enumerated in the lexicographic order of its “canonical suffix” by the forward DFS from the source.

# Sec5. Computing PLCP in $O(e)$ time&space



■ Observation A1:  $O(e)$  secondary outgoing edge of  $CDAWG(T)$  under the length-order determines the irreducible value  $PLCP[p] = l$  by the length  $l$  of the longest path from the source to the corresponding branching node

■ Observation A2:  $O(e)$  secondary outgoing edges can be enumerated in the text order of its "canonical suffix" by the reverse DFS from the sink.

We can extend the above result from PLCP to PLPF by employing the definition of 2ndary outgoing edges in length-order





# Conclusions

- Conversion problem from the CDAWG into other compressed indices for highly-repetitive texts:
  - Input: either the CDAWG of a text  $T$  or its self-index
  - Output: RL-BWT, irreducible PLCP and LPF, Lex- and LZ-parse
- We obtained **Optimal  $O(e)$  time and space conversion algorithms** for the above indices:
  - **Effective version** of the result by Belazzougui & Cunial (CPM'15) that  $r \leq e$  and  $z \leq e$  to actual conversion.
- Techniques:
  - Characterization of the “irreducible values” by **secondary edges**.
  - **Forward/reverse DFS** under the lexicographic/text order
- Future Work:
  - Conversion from RL-BWT & LZ-parse into CDAWG in  $O(e)$  time & space; Extension of the techniques to other indexing structures



*Thank you!*