# Top-$k$ Frequent Patterns in Streams

# and Parameterized-Space LZ Compression

**Patrick Dinklage**     Johannes Fischer     Nicola Prezza

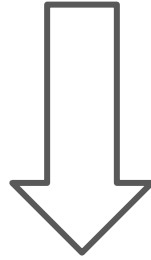technische universität
dortmund

Università
Ca'Foscari
Venezia

# Lempel-Ziv (LZ) Compression

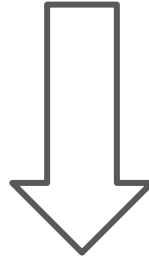panaman_ananas_banana_pancake

# Lempel-Ziv (LZ) Compression

panam<span style="color:blue">an</span>_ananas_banana_pancake

⬇

panam<span style="color:blue">(4,2)</span>_ananas_banana_pancake

# Lempel-Ziv (LZ) Compression

panam<u>an</u>_ananas_banana_pancake

⬇

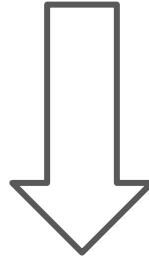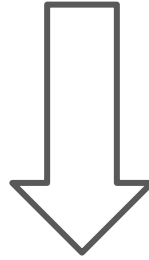p<u>an</u>am(4,2)_ananas_banana_pancake

# Lempel-Ziv (LZ) Compression

panaman_ananas_banana_pancake

panam(4,2)_(7,3)nas_banana_pancake

# Lempel-Ziv (LZ) Compression

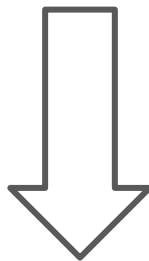panaman_ananas_banana_pancake

⬇

panam(4,2)_(7,3)(2,2)s_b(8,5)_(22,3)cake

# Lempel-Ziv (LZ) Compression

panaman_ananas_banana_pancake

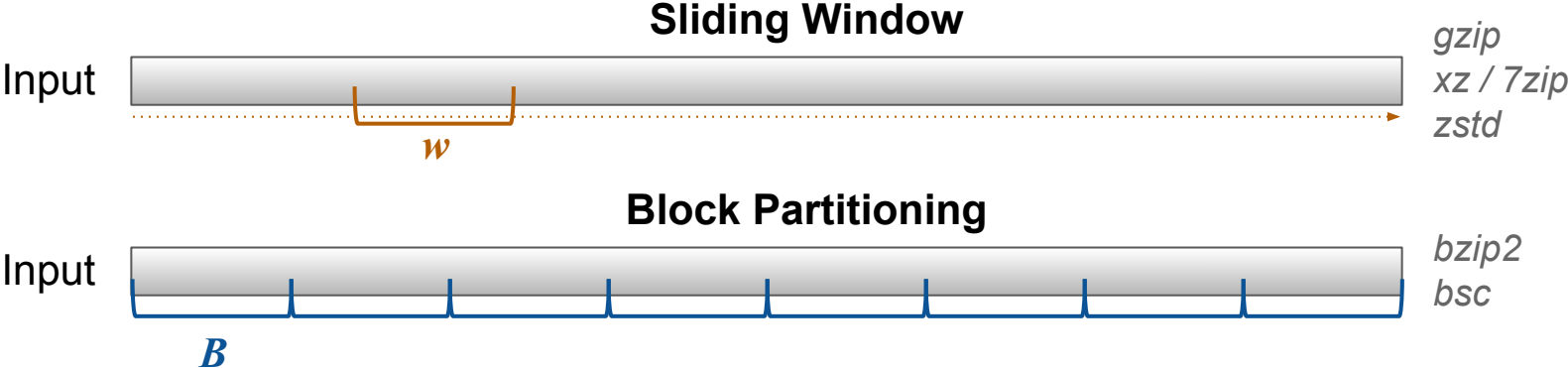panam(4,2)_(7,3)(2,2)s_b(8,5)_(22,3)cake
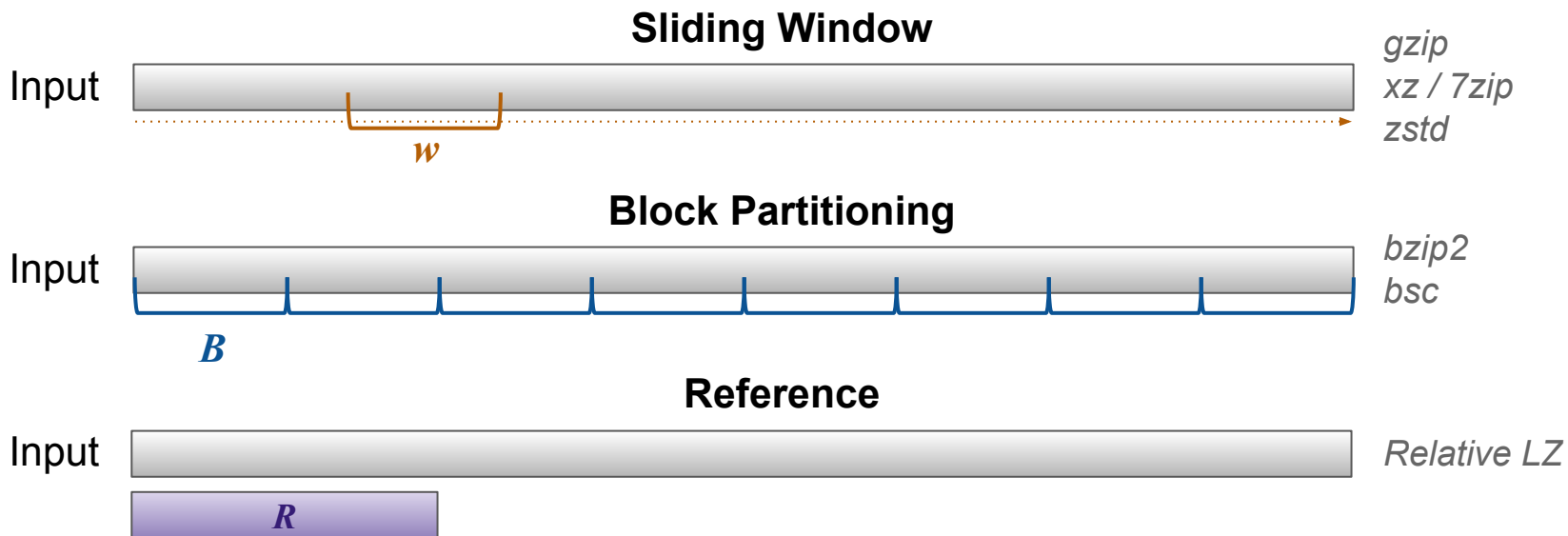
**LZ references**

# Compression in parameterized space
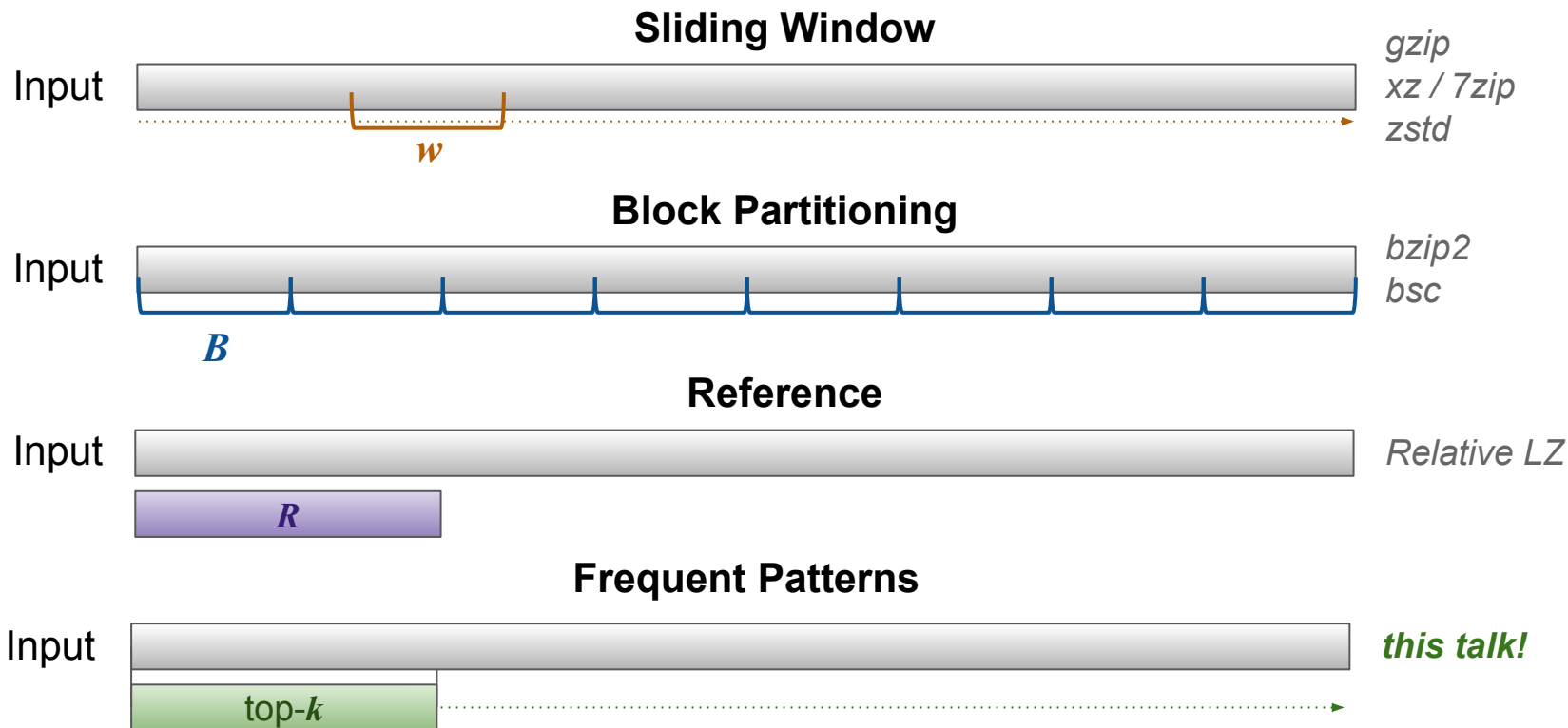
# Compression in parameterized space

**Sliding Window**

Input

$w$

*gzip*
*xz / 7zip*
*zstd*

# Compression in parameterized space

**Sliding Window**

Input

$w$

*gzip*
*xz / 7zip*
*zstd*

**Block Partitioning**

Input

$B$

*bzip2*
*bsc*

# Compression in parameterized space

**Sliding Window**

Input

*gzip*
*xz / 7zip*
*zstd*

$w$

**Block Partitioning**

Input

*bzip2*
*bsc*

$B$

**Reference**

Input

*Relative LZ*

$R$

# Compression in parameterized space

**Sliding Window**

Input

*gzip*
*xz / 7zip*
*zstd*

$w$

**Block Partitioning**

Input

*bzip2*
*bsc*

$B$

**Reference**

Input

*Relative LZ*

$R$

**Frequent Patterns**

Input

*this talk!*

top-$k$

# Frequent Patterns

→ We look for frequent **patterns** in a stream $S$ of characters

# Frequent Patterns

→ We look for frequent **patterns** (consecutive substrings) in a stream $S$ of characters

# Frequent Patterns

→ We look for **frequent patterns** (consecutive substrings) in a stream $S$ of characters

`panaman_ananas_banana_pancake`

# Frequent Patterns

→ We look for **frequent patterns** (consecutive substrings) in a stream $S$ of characters

p<u>ana</u>man_ananas_banana_pancake

# Frequent Patterns

→ We look for **frequent patterns** (consecutive substrings) in a stream $S$ of characters

p<u>an</u>am<span style="color:magenta">an_anan</span>as_b<span style="color:magenta">anan</span>a_p<span style="color:magenta">an</span>cake

# Compressing with Frequent Patterns

**Idea:**

- Estimate online which $k$ patterns are the most frequent

- Replace occurrences of frequent patterns by references

# Compressing with Frequent Patterns

**Idea:**

- Estimate online which $k$ patterns are the most frequent

- Replace occurrences of frequent patterns by references

$\rightarrow$ The number of candidate patterns is <span style="color:red">quadratic in $|S|$</span>

# Frequent Patterns    versus    LZ References

→ Frequent patterns are captured by LZ references

**Frequent Patterns**       versus       **LZ References**

→ Frequent patterns are captured by LZ references

**Frequent Patterns**          versus          **LZ References**

→ Patterns *not* captured by LZ references *cannot* be frequent (by any measure)

→ Frequent patterns are captured by LZ references

**Frequent Patterns**  versus  **LZ References**

→ Patterns *not* captured by LZ references *cannot* be frequent (by any measure)

→ The number of LZ references is bounded by $|S|$

# top-k LZ78

*k=6*

( 0 )

panamanananasbananapancake

# top-k LZ78

*k=6*

( 0 )

panamanananasbananapancake

# top-k LZ78

*k=6*

0

p

1 |

panamanananasbananapancake

(0,p)

# top-k LZ78

*k=6*



**0**

p

**1** |

pa namanananasbananapancake

(0,p)

# top-k LZ78

*k=6*

```
        0
     a /   \ p
      /     \
    (2)      (1)
     |        |
```

p<u>a</u>namanananasbananapancake

(0,p)<u>(0,a)</u>

# top-k LZ78

*k=6*

0

a    n    p

2  |    3  |    1  |

pa<u>n</u>amanananasbananapancake

(0,p)(0,a)<u>(0,n)</u>

# top-k LZ78

*k=6*



pan<u>a</u>manananasbananapancake

(0,p)(0,a)(0,n)

# top-k LZ78

pan<u>a</u>manananasbananapancake

(0,p)(0,a)(0,n)

# top-k LZ78

*k=6*

0

a     n     p

2  ||    3  |    1  |

pana<u>m</u>anananasbananapancake

(0,p)(0,a)(0,n)

# top-k LZ78



panam<u>an</u>anananasbananapancake

(0,p)(0,a)(0,n)<u>(2,m)</u>

# top-k LZ78

*k=6*

panam**a**nananasbananapancake

(0,p)(0,a)(0,n)(2,m)

# top-k LZ78

*k=6*

```
       0
    a  | n   p
   2 |||  3 |   1 |
 m
4 |
```

panama<u>n</u>ananasbananapancake
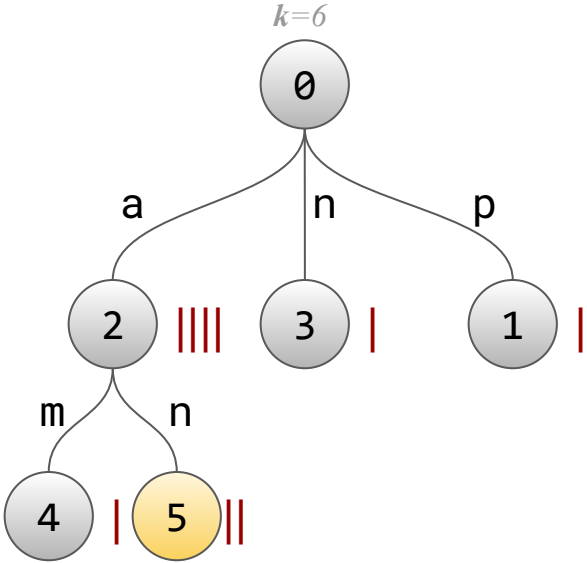
(0,p)(0,a)(0,n)(2,m)

# top-k LZ78

*k=6*



panama<u>n</u>ananasbananapancake

(0,p)(0,a)(0,n)(2,m)(2,n)

# top-k LZ78



panamana<u>a</u>nanasbananapancake

(0,p)(0,a)(0,n)(2,m)(2,n)

# top-k LZ78



*k=6*

panamana<u>n</u>anasbananapancake

(0,p)(0,a)(0,n)(2,m)(2,n)

38

# top-k LZ78

*k=6*



panamanan<u>a</u>nasbananapancake

(0,p)(0,a)(0,n)(2,m)(2,n)
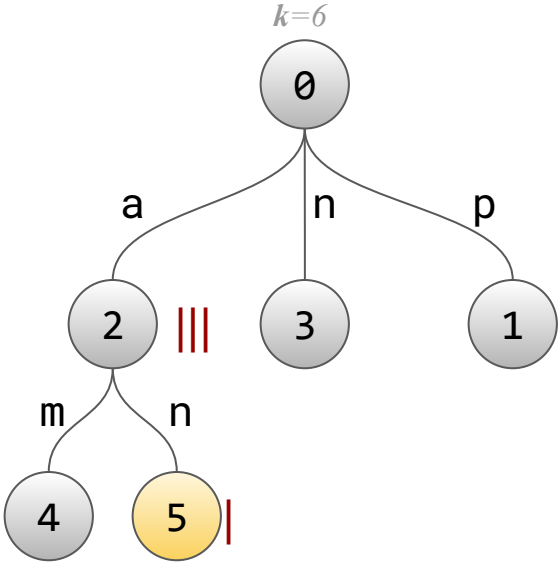
# top-k LZ78

*k=6*



panamanan<u>a</u>nasbananapancake

trie is full, nothing to "recycle"

(0,p)(0,a)(0,n)(2,m)(2,n)<u>(5,a)</u>

# top-k LZ78

*k=6*

0

a     n     p

2  |||    3         1

m     n

4     5  |

trie is full, nothing to "recycle" → decrement **all** counters

panamanana nasbananapancake

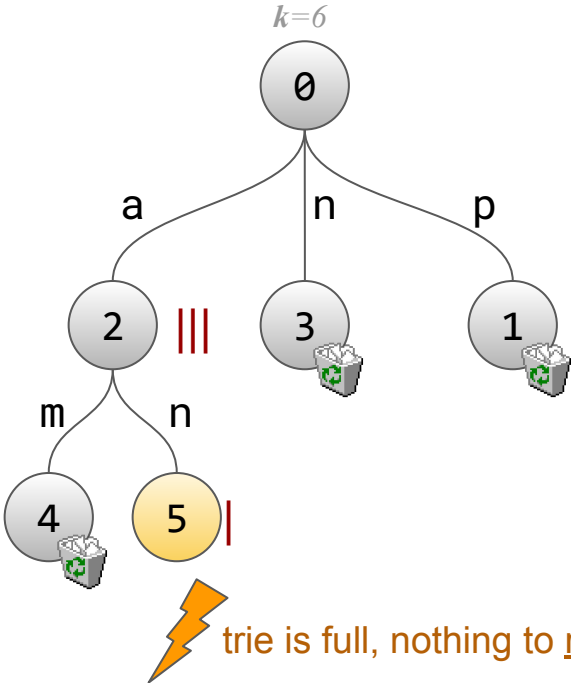(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)

# top-k LZ78

*k=6*

panamanana<u>a</u>nasbananapancake

trie is full, nothing to <u>recycle</u> → decrement **all** counters
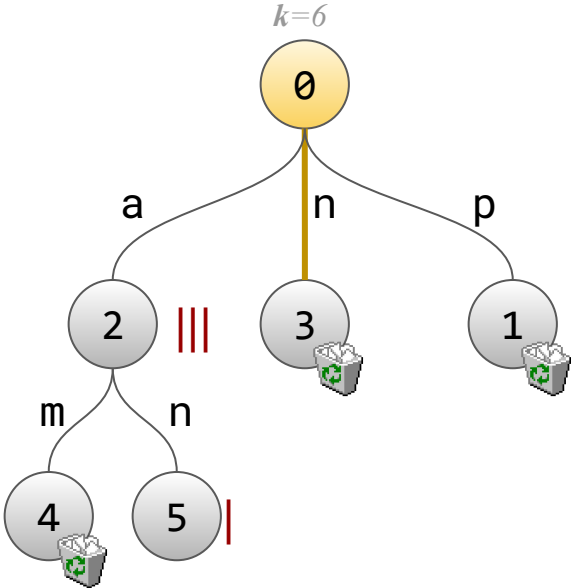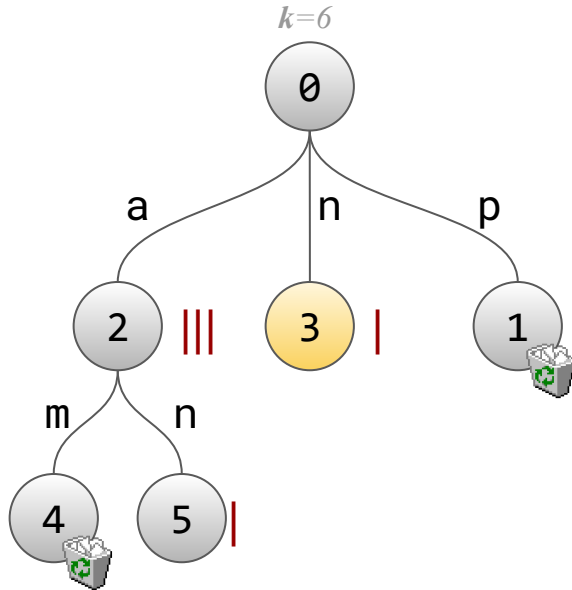
(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)

# top-k LZ78

*k=6*

panamanana<u>n</u>asbananapancake

(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)
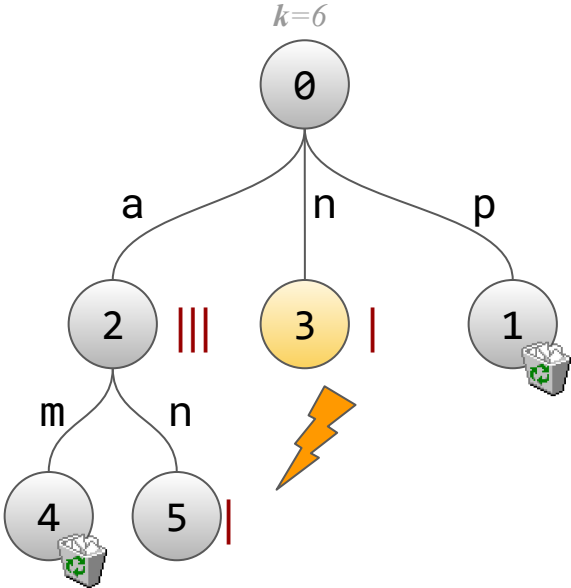
# top-k LZ78

*k=6*



panamanana<u>a</u>sbananapancake

(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)

# top-k LZ78

*k=6*

panamananan**a**sbananapancake

(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)(3,a)

# top-k LZ78



*k=6*

panamananan**a**sbananapancake

(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)(3,a)
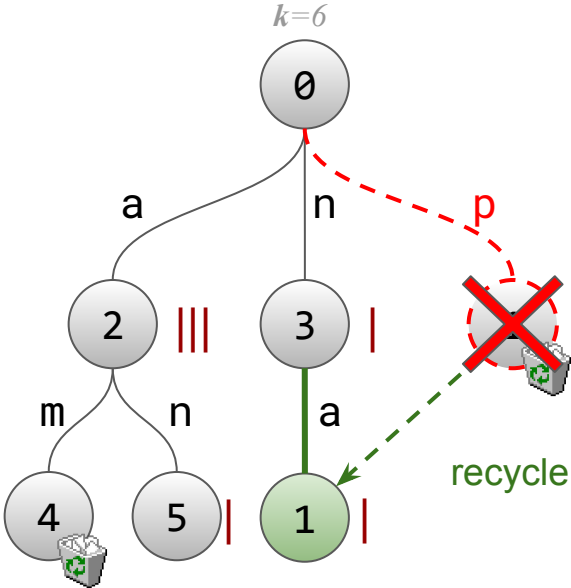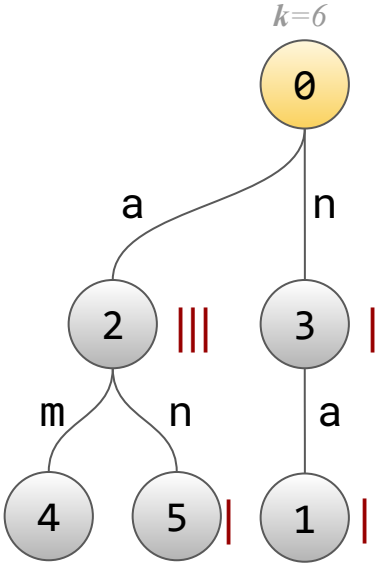
# top-k LZ78

*k=6*

0

a          n

2  ‖‖          3  |
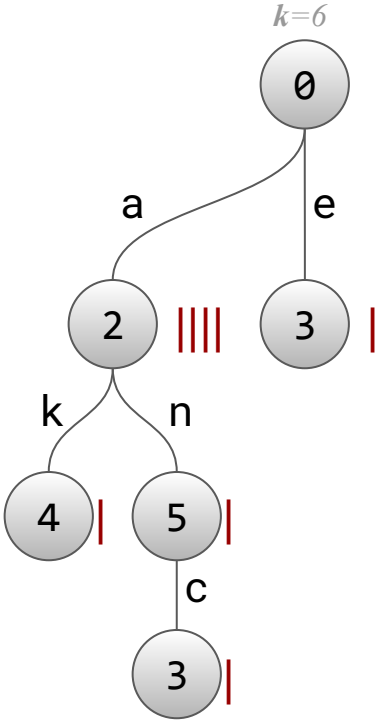
m      n          a

4      5 |      1 |

panamanananasbananapancake

## … and so forth …

(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)(3,a)

# top-k LZ78



*k=6*

panamanananasbananapancake

(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)(3,a)(0,s)(0,b)(5,a)(1,p)(5,c)(2,k)(0,e)

# top-k LZ78

(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)(3,a)(0,s)(0,b)(5,a)(1,p)(5,c)(2,k)(0,e)

→ **time-dependent** top-k trie references

# top-k LZ78

`(0,p)(0,a)(0,n)(2,m)(2,n)(5,a)(3,a)(0,s)(0,b)(5,a)(1,p)(5,c)(2,k)(0,e)`

→ **time-dependent** top-k trie references

→ hybrid of **LZ78** compression and **Misra-Gries** sketch

→ simulate decrements & manage garbage queue using [Metwally et al., 2005]

→ constant amortized time per input character

# LZ78 vs. LZ77

Input:    $a^n$ = aaaaaaaaaaaaaaa …

LZ78:

LZ77:

# LZ78 vs. LZ77

Input:   $a^n$ = aaaaaaaaaaaaaaa …

LZ78:   `(0,a)(1,a)(2,a)(3,a)`  …        $\rightarrow \Theta(\sqrt{n})$ phrases

LZ77:

# LZ78 vs. LZ77

Input:     $a^n$ = aaaaaaaaaaaaaa …

LZ78:     `(0,a)(1,a)(2,a)(3,a)` …          $\rightarrow \Theta(\sqrt{n})$ phrases

LZ77:     `a(1,`*n-1*`)`                              $\rightarrow$ 2 phrases

# LZ78 vs. LZ77

Input:      $a^n$ = aaaaaaaaaaaaaaaa …

LZ78:      `(0,a)(1,a)(2,a)(3,a)` …          $\rightarrow \Theta(\sqrt{n})$ phrases

    ⊕ simple algorithm (compressed space by default)

    ⊕ straightforward efficient encoding

LZ77:      a`(1,`*n-1*`)`                    $\rightarrow$ 2 phrases

    ⊖ much harder to compute

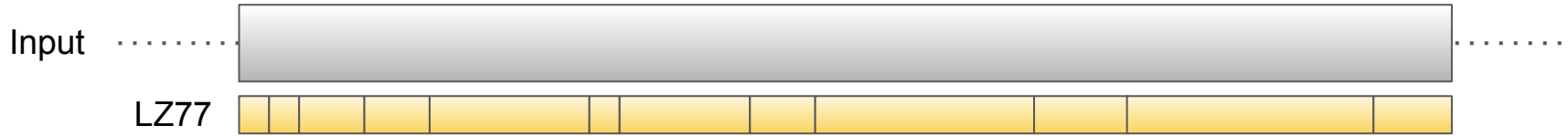    ⊖ harder to encode efficiently

$\rightarrow$ best of both worlds?

# top-k LZ77

**Algorithm Sketch:**

1. Maintain a top-k trie of phrases (like in top-k LZ78)

2. Partition the input into blocks of size $B = O(k)$

3. Compute the block's LZ77 parsing in time and space $O(B) = O(k)$

4. Greedily pick the next LZ77 phrase or matching pattern from the trie
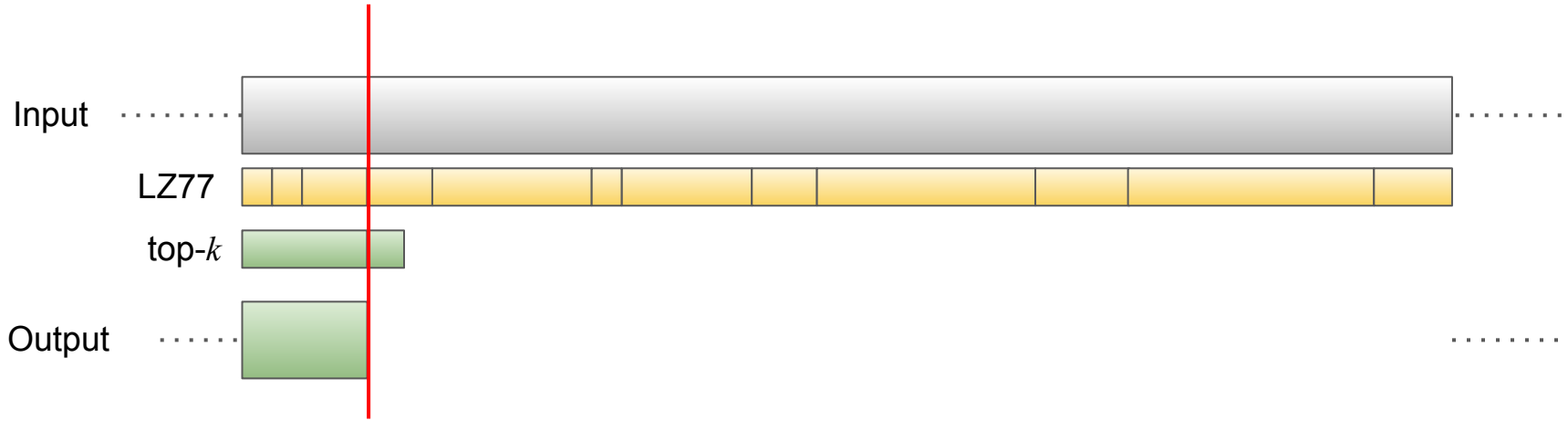
# top-k LZ77

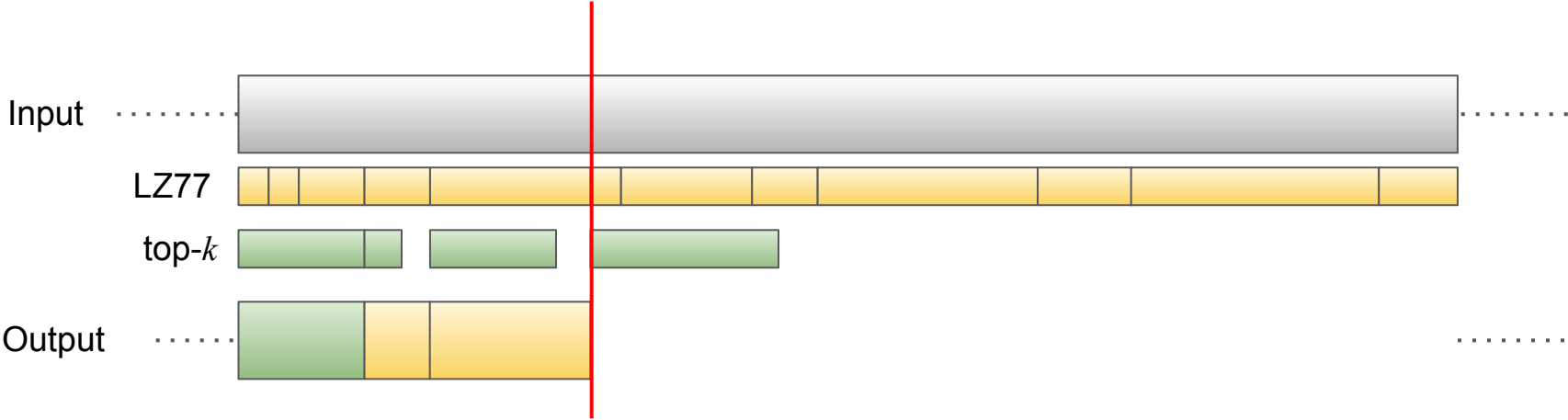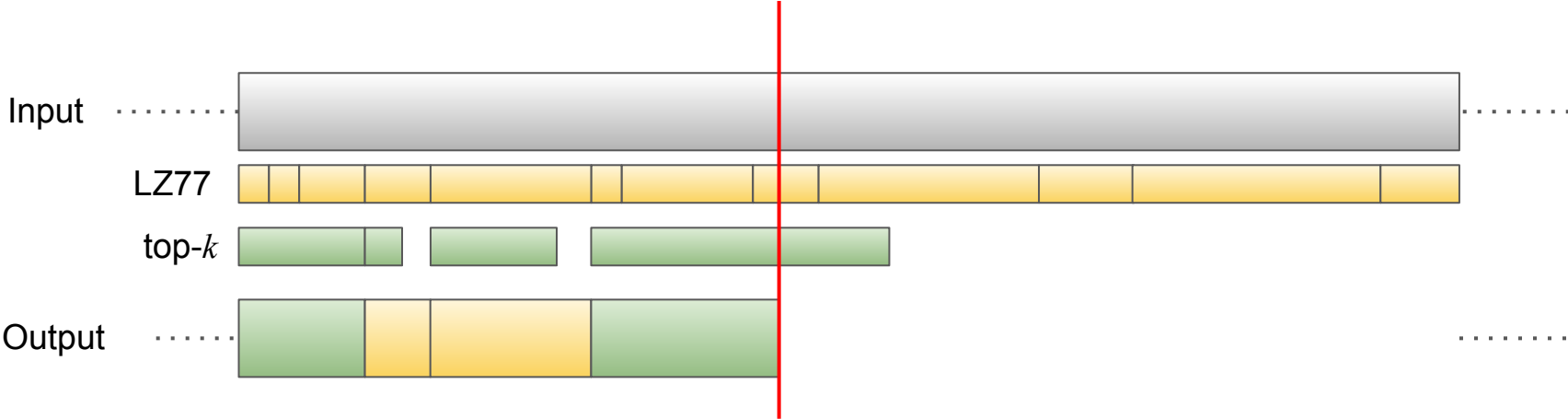Input · · · · · · · · · ·

LZ77

# top-k LZ77

Input

LZ77

top-$k$

# top-k LZ77

Input

LZ77

top-$k$

Output

# top-k LZ77

Input

LZ77

top-$k$

Output

# top-k LZ77

Input

LZ77

top-$k$

Output

# top-k LZ77

Input

LZ77

top-$k$

Output

# top-k LZ77

# top-k LZ77



Input

LZ77

top-$k$

Output

CommonCrawl (100 GiB web text crawl)



comp. ratio [%]

time [h]

max. RAM: 64 GB

COMMONCRAWL (100 GiB web text crawl)

topk-lz77
topk-lz78

max. RAM: 64 GB

comp. ratio [%]

time [h]
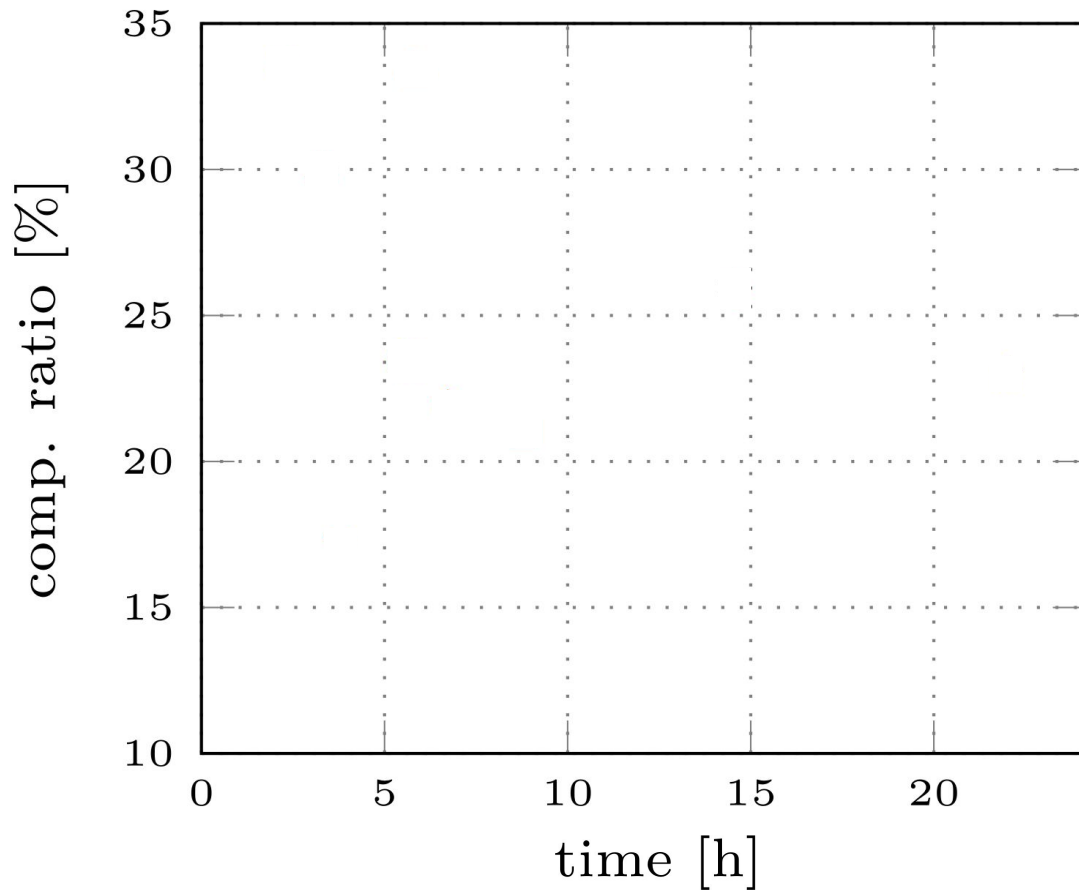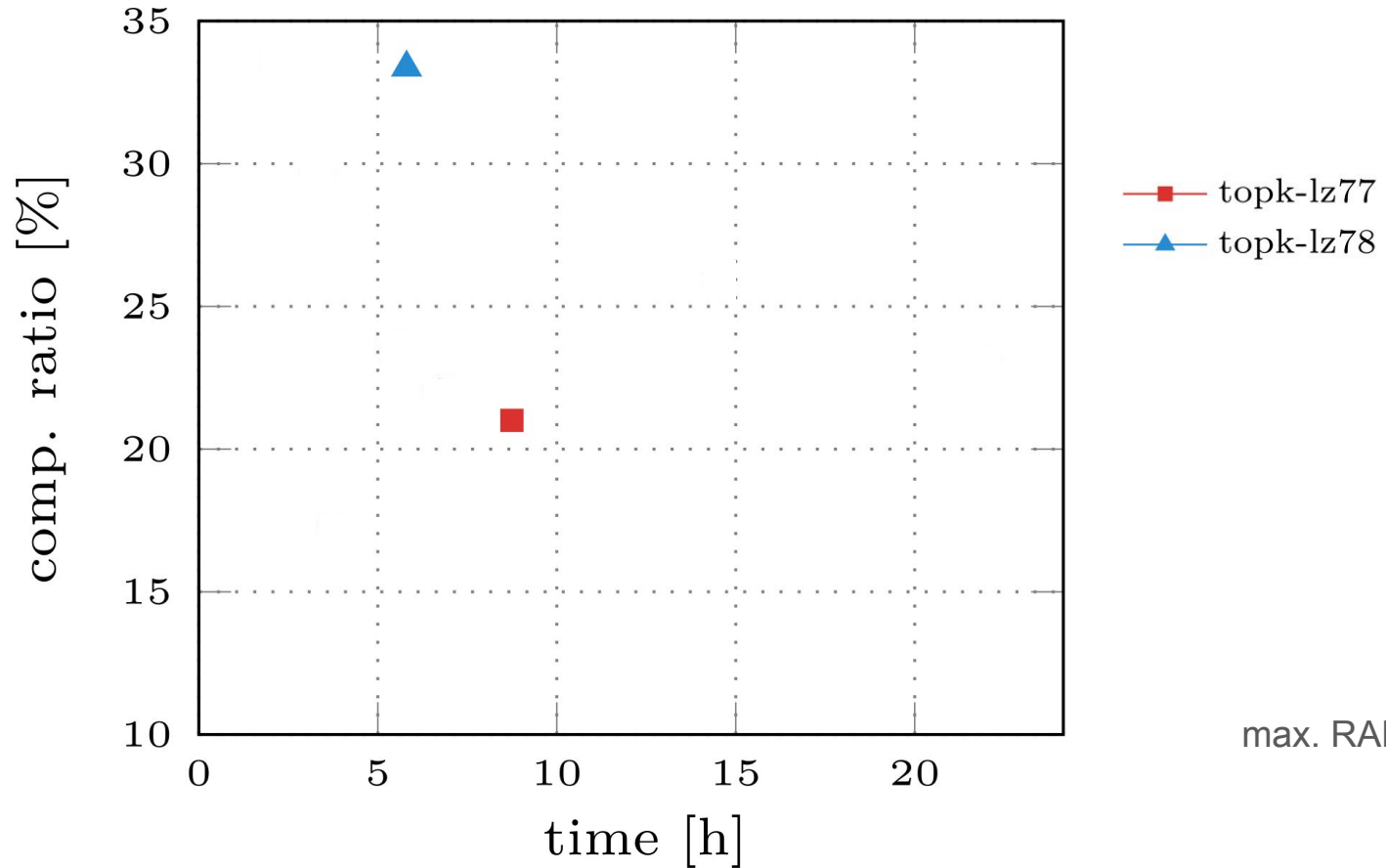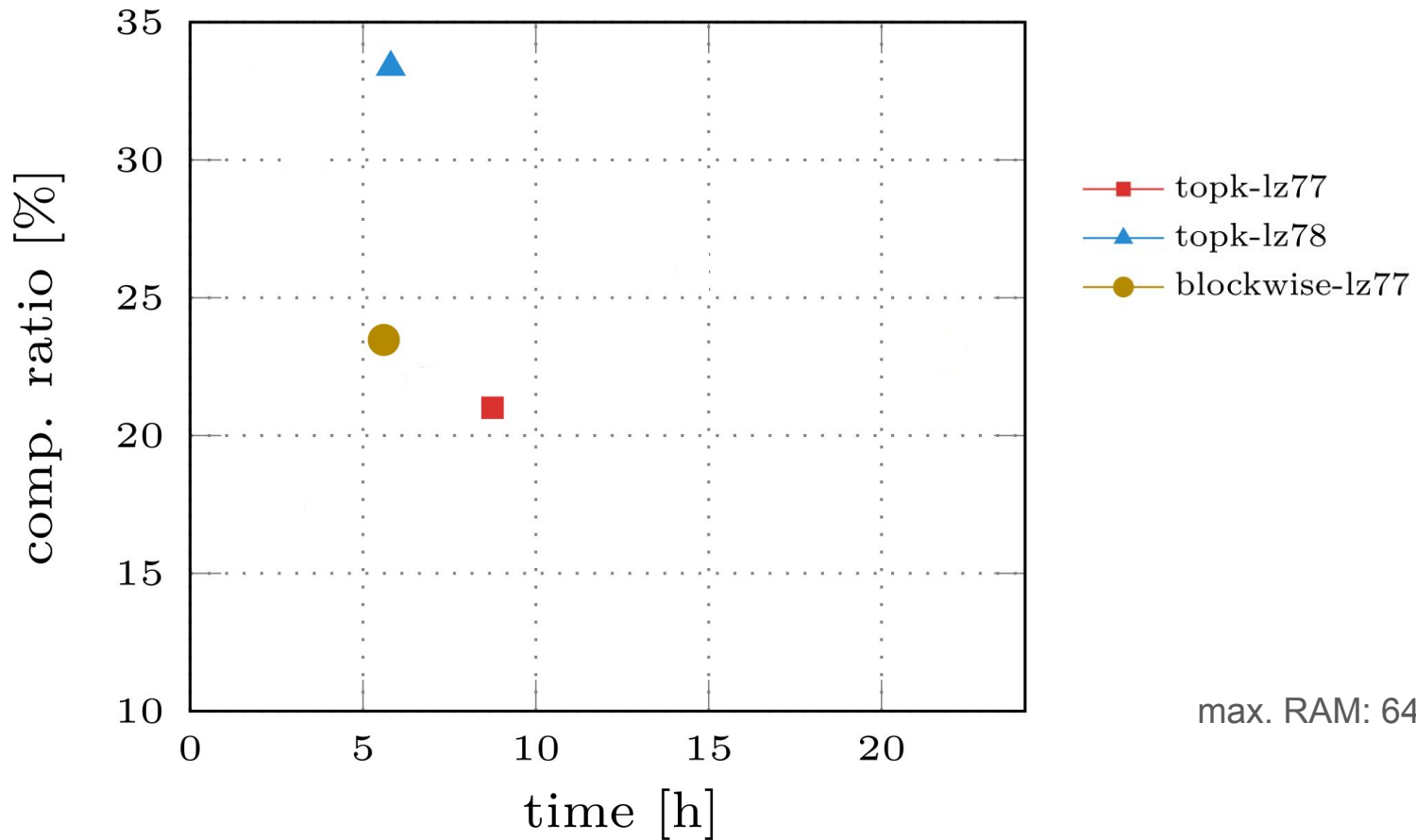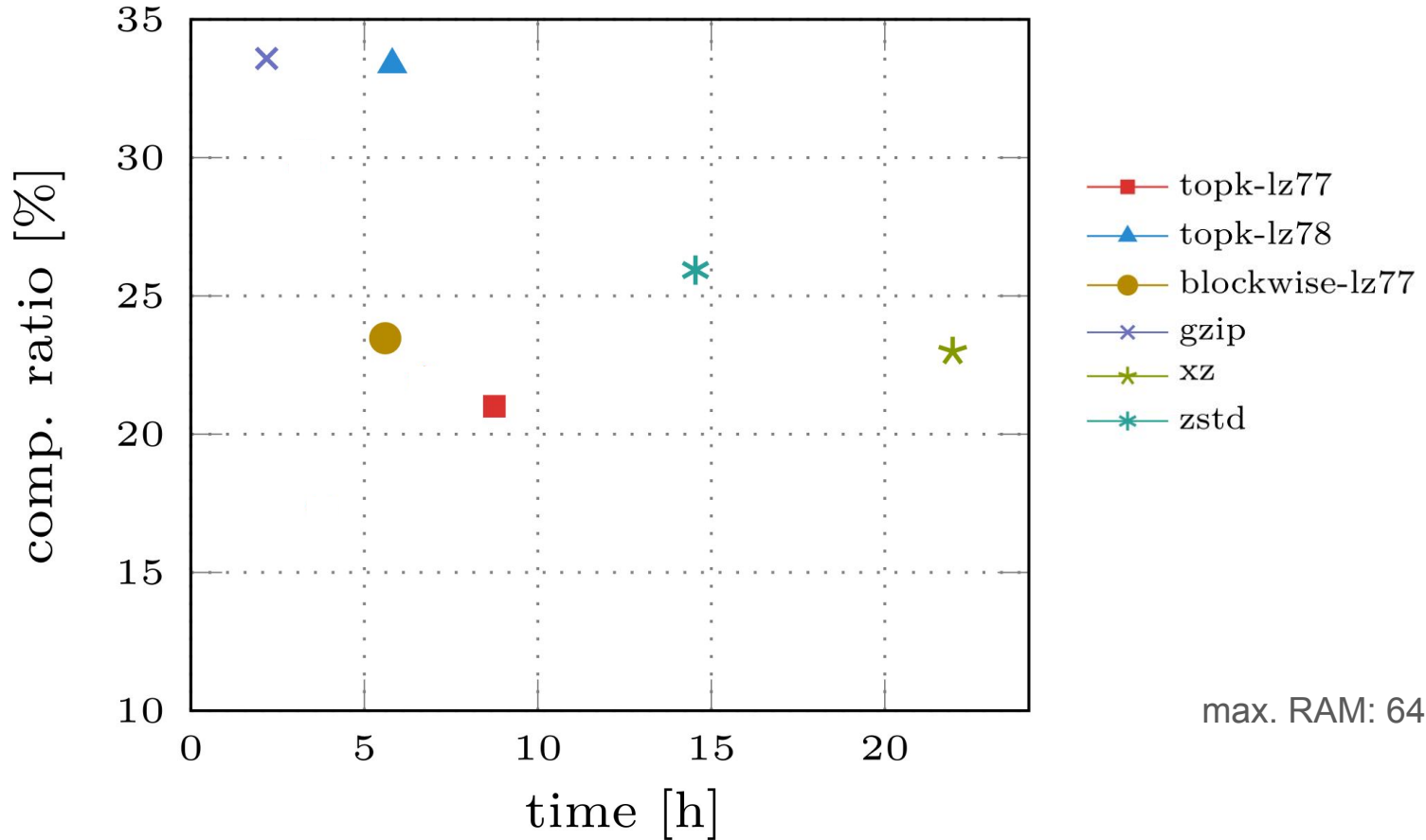
# CommonCrawl (100 GiB web text crawl)



max. RAM: 64 GB

# CommonCrawl (100 GiB web text crawl)
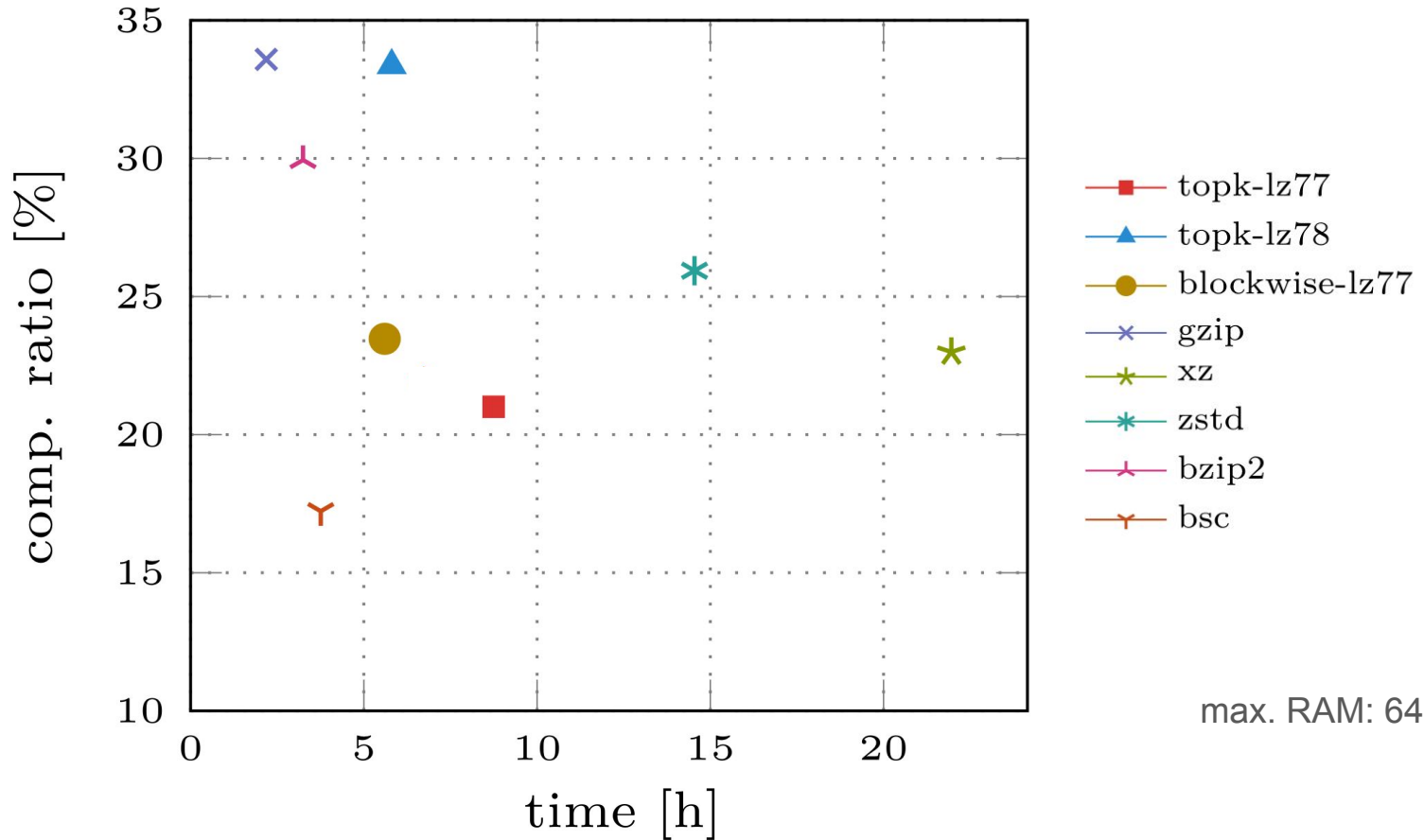


max. RAM: 64 GB

CommonCrawl (100 GiB web text crawl)

comp. ratio [%] vs time [h]

Legend:
- topk-lz77
- topk-lz78
- blockwise-lz77
- gzip
- xz
- zstd
- bzip2
- bsc

max. RAM: 64 GB

68

# Future Work

- Parallel computation

- Statistical encoding of trie references

- Dynamic string attractors? (as opp. to tries)

- **Precompression** of long repetitions

- Random Access (e.g., like [Arz & Fischer, 2018])

- Information Retrieval

- Your idea here!