

# Efficient Parallel Output-Sensitive Edit Distance

Xiangyun Ding, Xiaojun Dong, Yan Gu, Youzhe Liu, Yihan Sun  
University of California, Riverside

Full version: <https://arxiv.org/pdf/2306.17461.pdf>

Code: <https://github.com/ucrparlay/Edit-Distance>



# Problem Definitions: Edit Distance

- **Input:**  $A[1..n]$  and  $B[1..m]$  over a fixed alphabet  $\Sigma$  (varieties of characters)
- **Output: Levenshtein distance**
  - unit-cost, single-character
  - insertions, deletions, and substitutions (replacement).

a	m	s	t	e	r	d	a	m
---	---	---	---	---	---	---	---	---

Input:

a	s	t	e	r	i	s	m
---	---	---	---	---	---	---	---

# Problem Definitions: Edit Distance

- **Input:**  $A[1..n]$  and  $B[1..m]$  over a fixed alphabet  $\Sigma$  (varieties of characters)
- **Output: Levenshtein distance**
  - unit-cost, single-character
  - insertions, deletions, and substitutions (replacement).

a m s t e r d a m

Input:

a s t e r i s m

Output: The edit distance  $k = 3$



# The classic algorithm

## Dynamic programming

$$T[i, j] = \min \begin{cases} T[i-1, j] + 1 \\ T[i, j-1] + 1 \\ T[i-1, j-1] + \delta(x[i-1], y[j-1]) \end{cases}$$

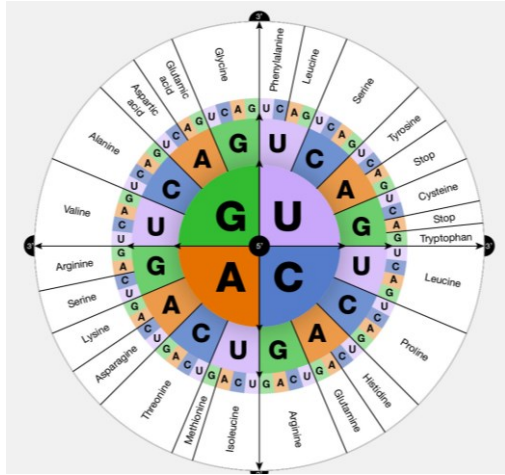
$= 0$  If  $x[i-1] = y[j-1]$   
 $= 1$  If  $x[i-1] \neq y[j-1]$

size of the two sequences  $\leftarrow$

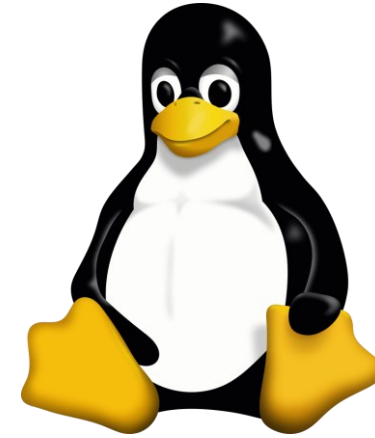
**Work (time complexity):**  $O(mn)$

		A	s	t	e	r	i	s	m
A	0	1	2	3	4	5	6	7	8
m	1	0	1	2	3	4	5	6	7
s	2	1	1	2	3	4	5	6	6
t	3	2	1	2	3	4	5	5	6
e	4	3	2	1	2	3	4	5	6
r	5	4	3	2	1	2	3	4	5
d	6	5	4	3	2	1	2	3	4
a	7	6	5	4	3	2	2	3	4
m	8	7	6	5	4	3	3	3	4
	9	8	7	6	5	4	4	4	3

# The sizes of $m$ and $n$ are quite large!



Bacteria\*: about **4.6 million** base pairs



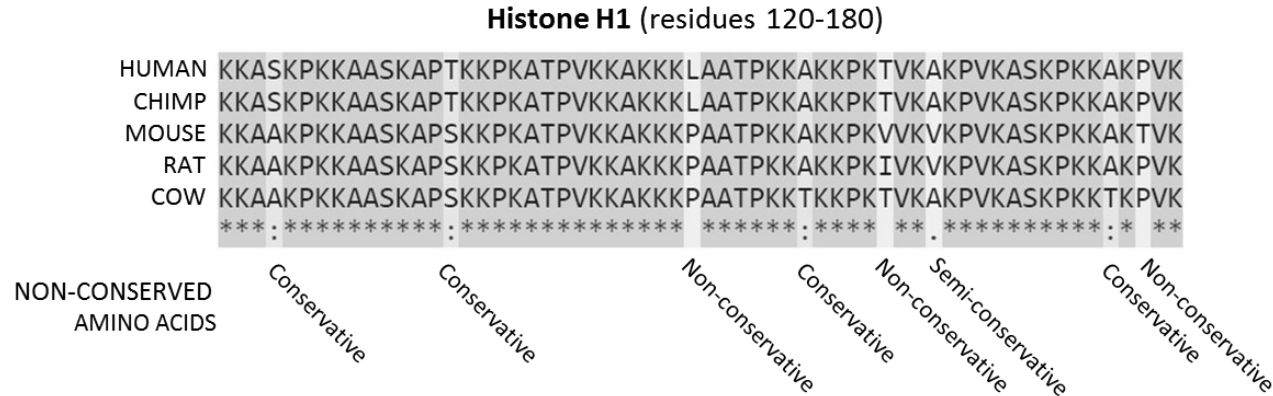
Linux kernel\*: around **27 million** lines of code

**Input: millions-scale or even larger**

**Better solutions for large-scale applications in practice?**

\* Figure from Wikipedia

# The two sequences are very similar!



**Alignment:** edits are **small**  
(we care alignments on similar sequence pairs)

Treat **different sizes of changes** in different ways

The two sequences are generally very **similar!**  $\rightarrow$  #edits:  $k = o(n)$

**Output-sensitive** algorithms!

$\Rightarrow$  When  $k = o(n)$ , the cost is  $o(nm)$

# Another resource: parallelism

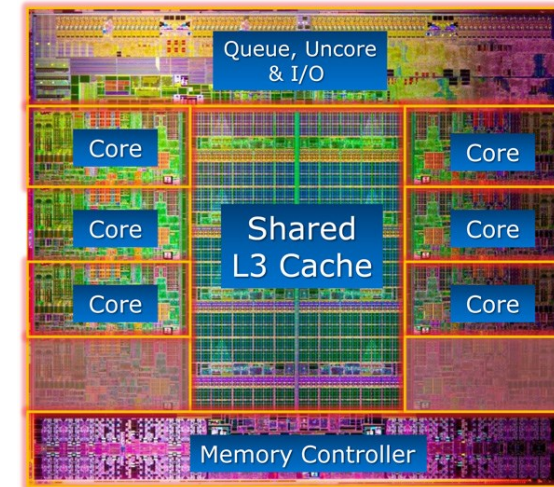
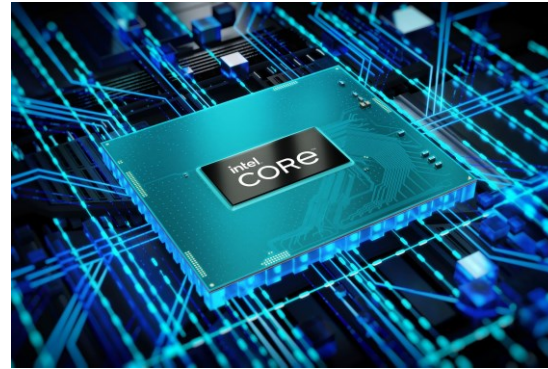
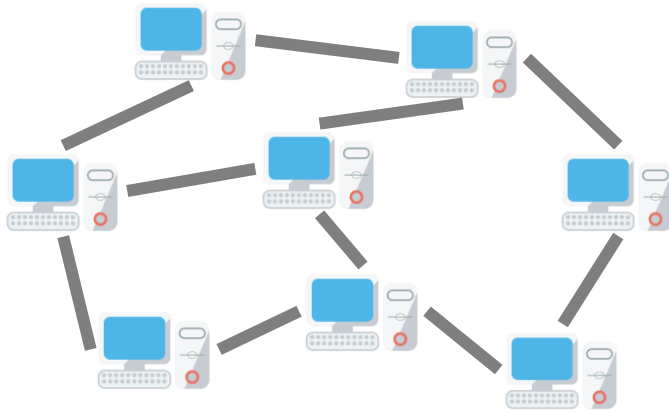
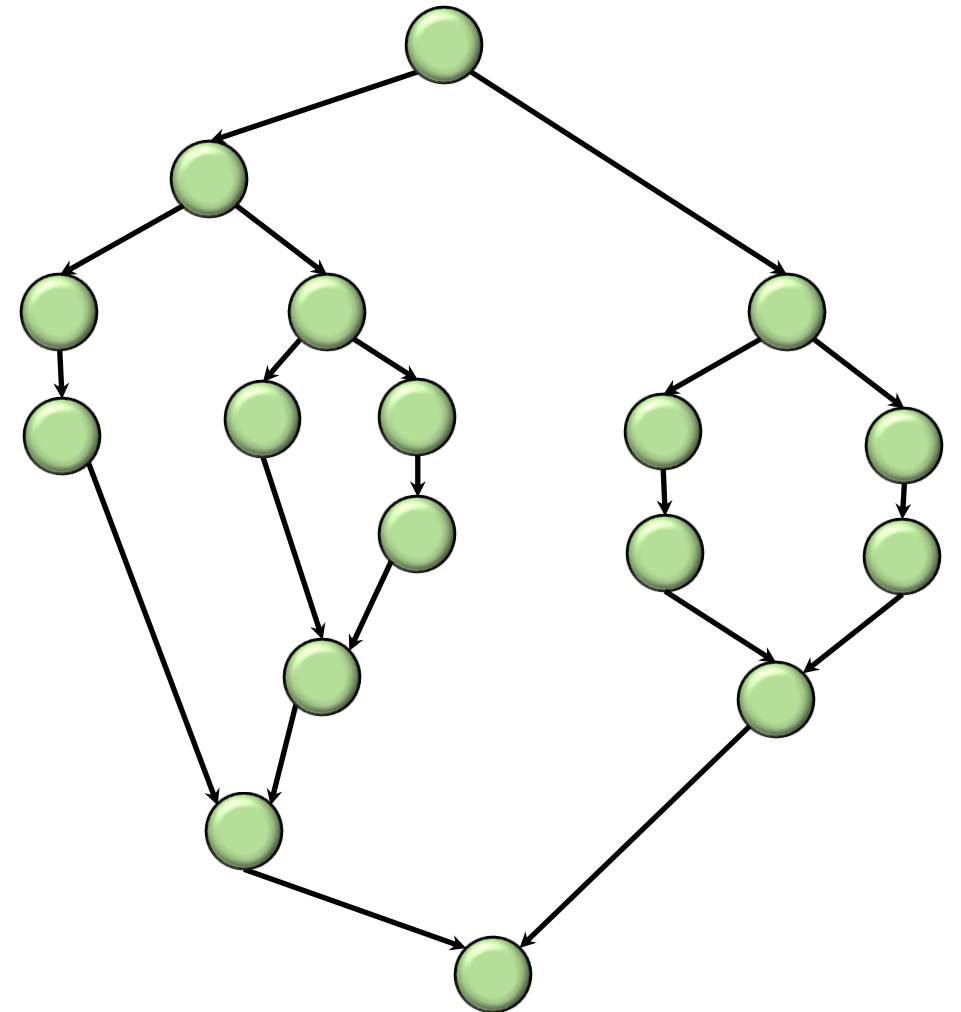


Figure credit: Wikipedia

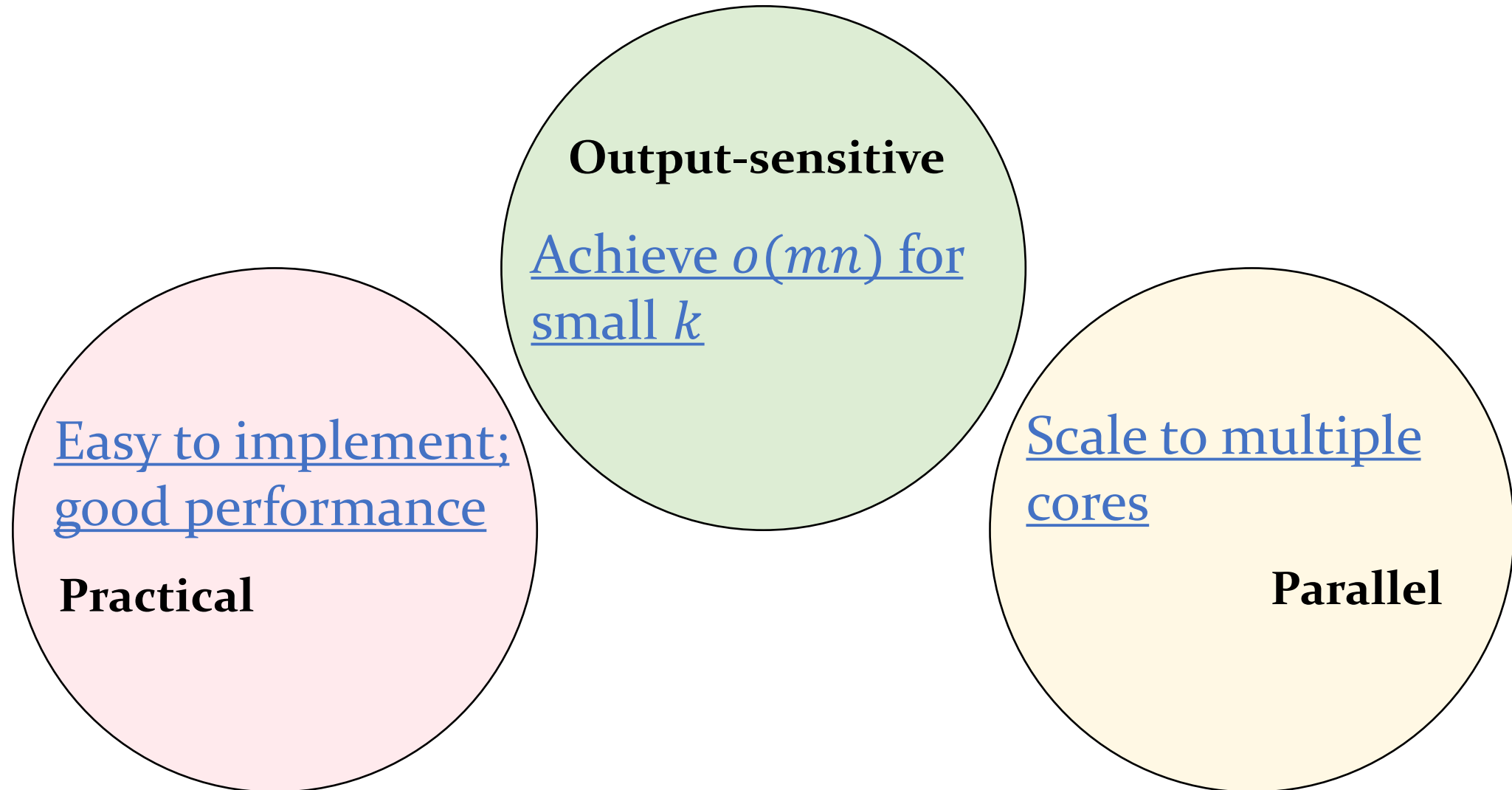
**Parallelism** is ubiquitous and can be used to accelerate algorithms

# Computational Model

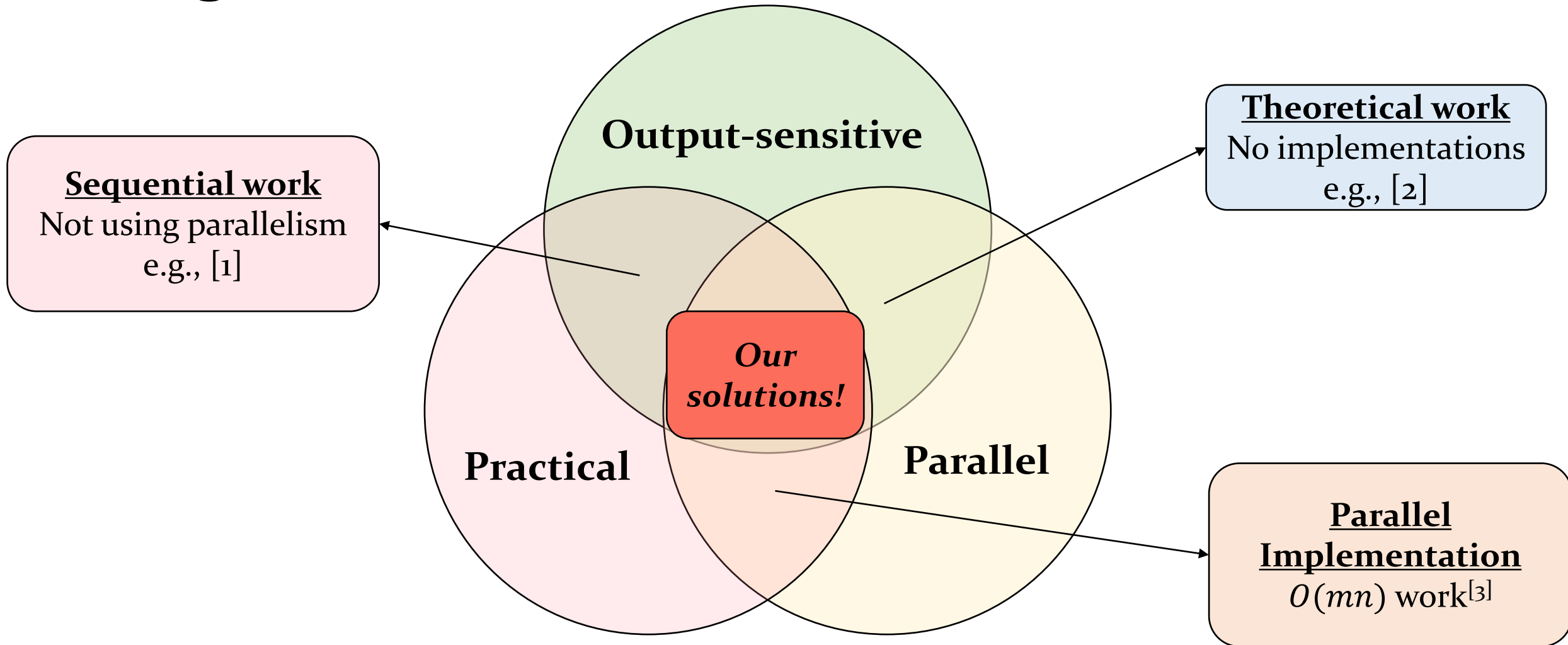
- Shared-memory multi-core setting
- **Work**: the total number of operations
  - = running time on one core
- **Span (depth)**: the length of the longest dependency chain
  - = running time on an infinite number of cores



# Our goals



# Our goals

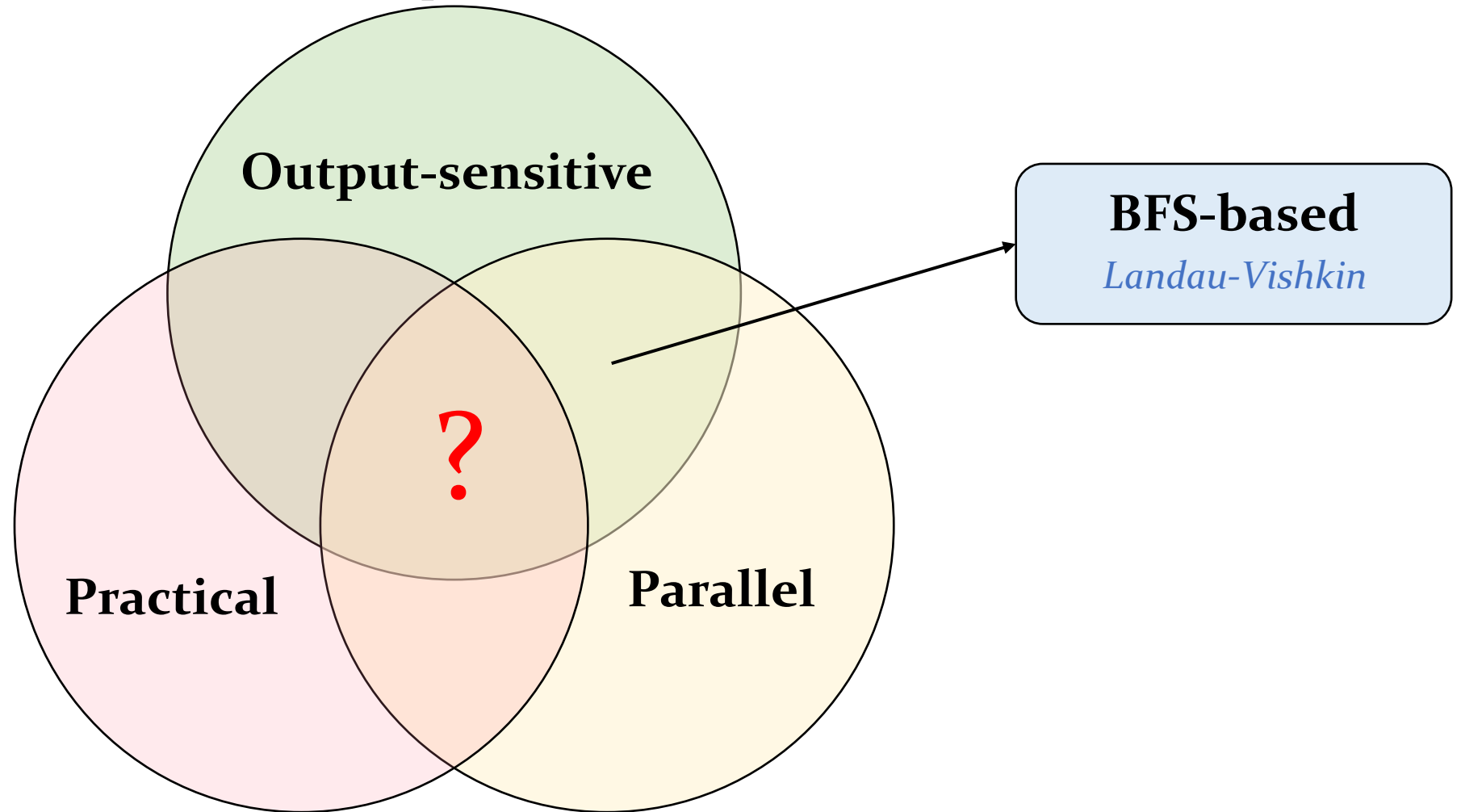


[1] Marçais, DeBlasio, Pandey, and Kingsford. Locality-sensitive hashing for the edit distance. Bioinformatics, 2019.

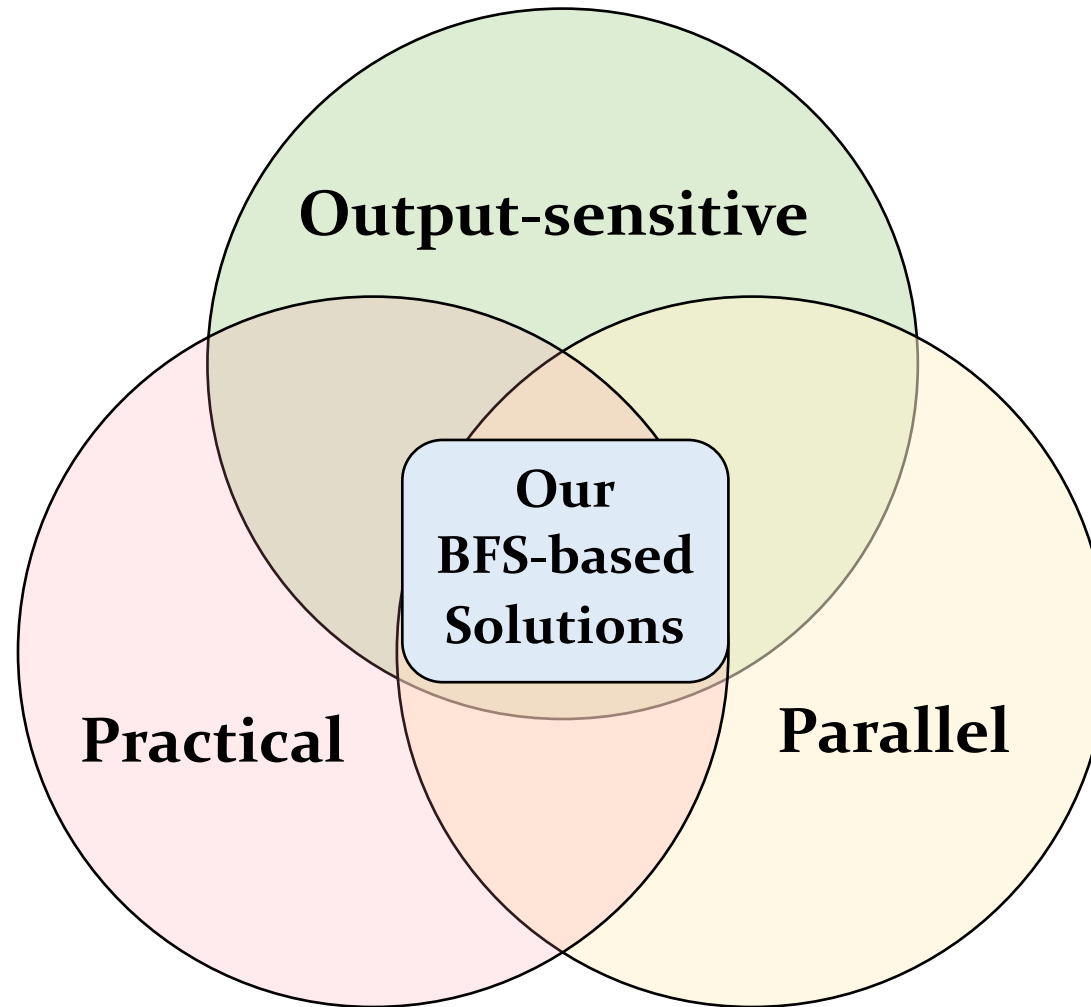
[2] Landau and Vishkin. Fast parallel and serial approximate string matching. J. Algorithms, 1989.

[3] Blelloch, Anderson, and Dhulipala. "ParlayLib-a toolkit for parallel algorithms on shared-memory multicore machines

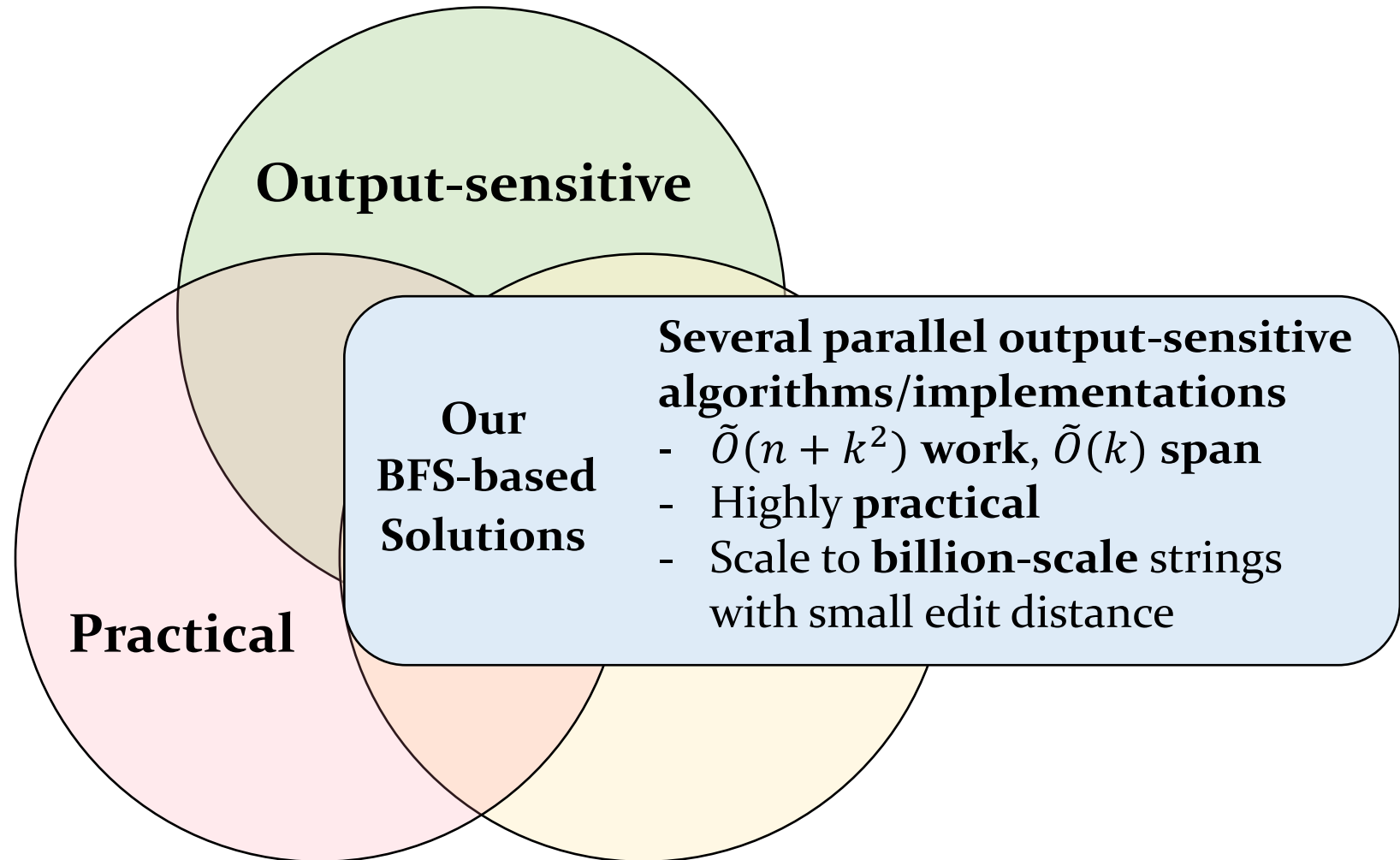
# Our contributions: in practice



# Our contributions: in practice



# Our contributions: in practice

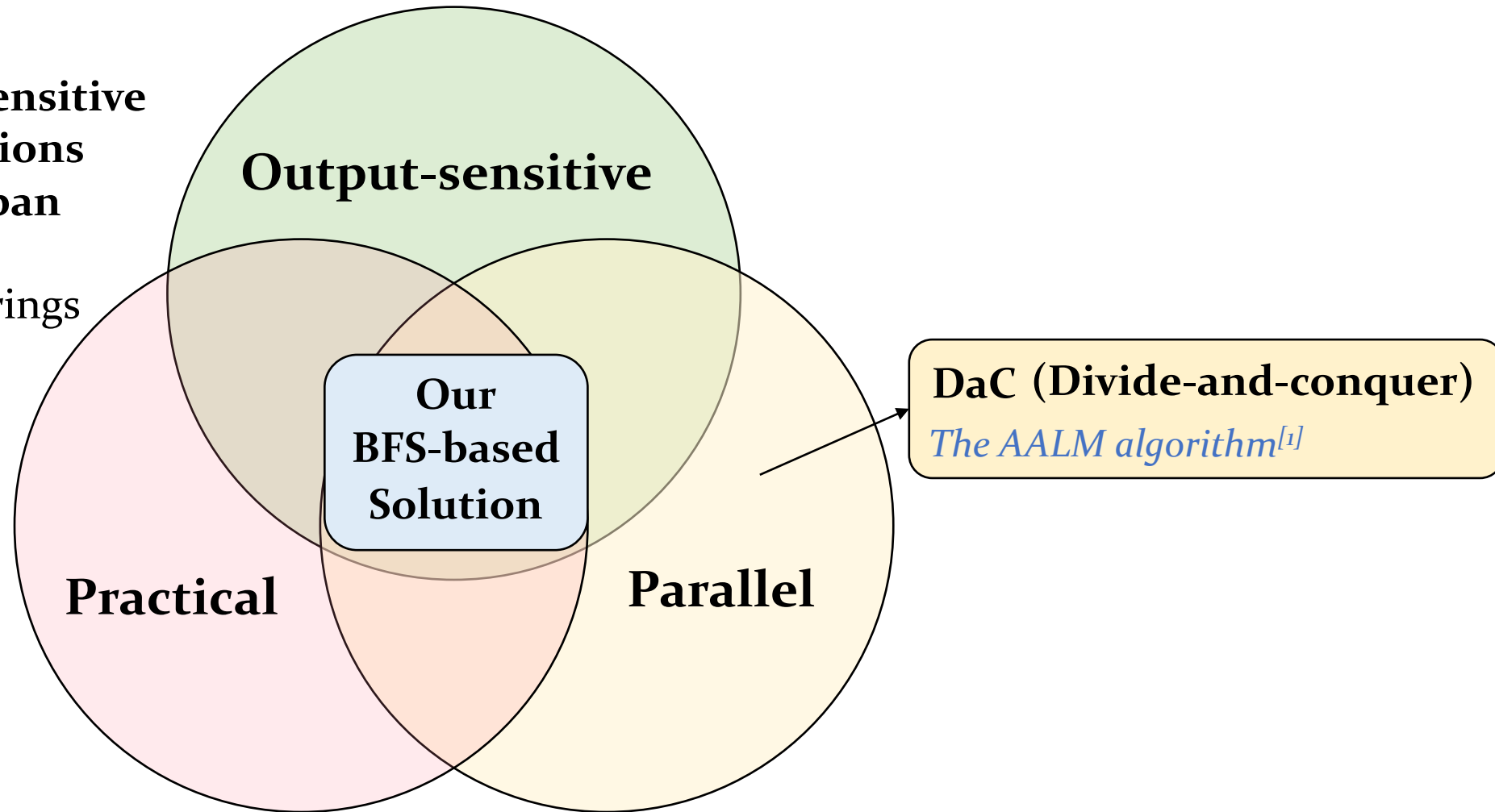


# Our contributions: in theory

## In practice

Several parallel output-sensitive algorithms/implementations

- $\tilde{O}(n + k^2)$  work,  $\tilde{O}(k)$  span
- Highly practical
- Scale to **billion-scale** strings with small distance

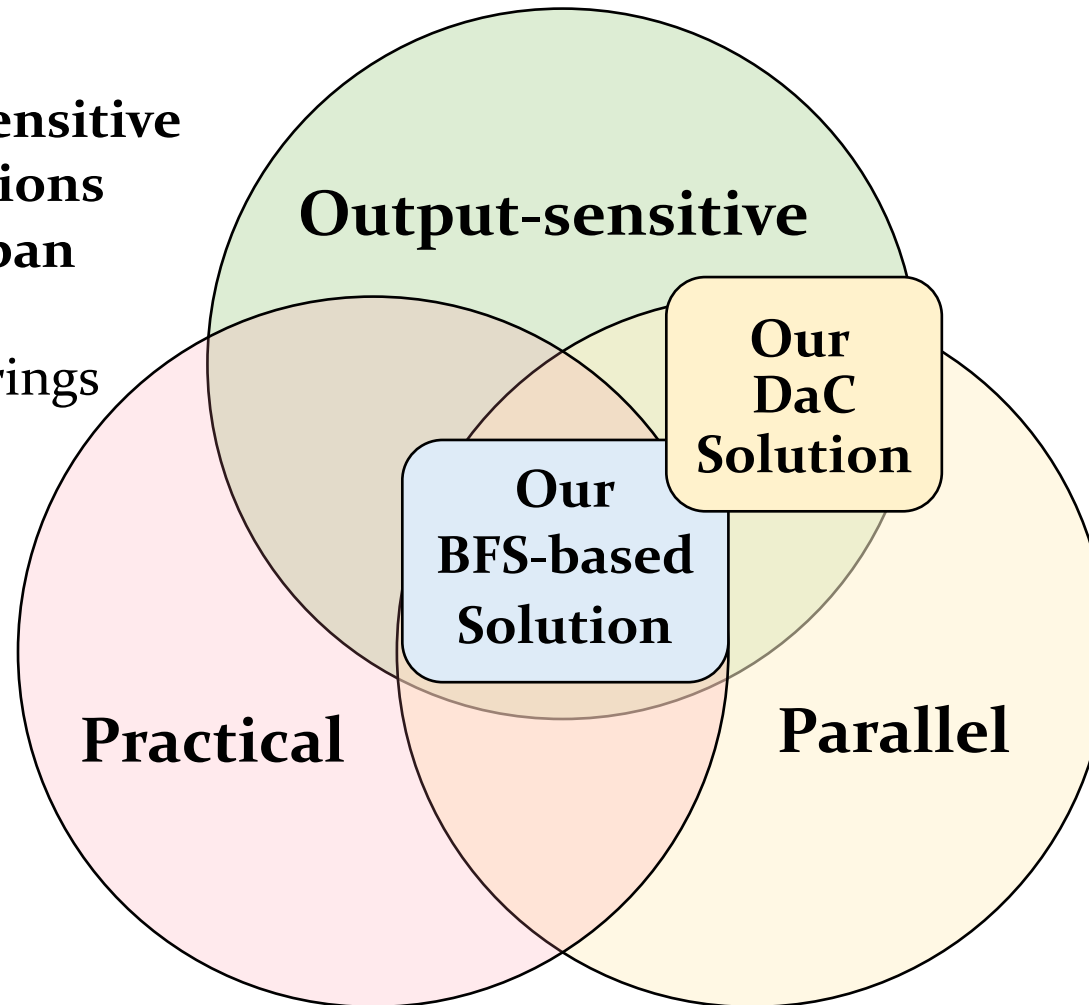


# Our contributions: in theory

## In practice

Several parallel output-sensitive algorithms/implementations

- $\tilde{O}(n + k^2)$  work,  $\tilde{O}(k)$  span
- Highly practical
- Scale to **billion-scale** strings with small distance

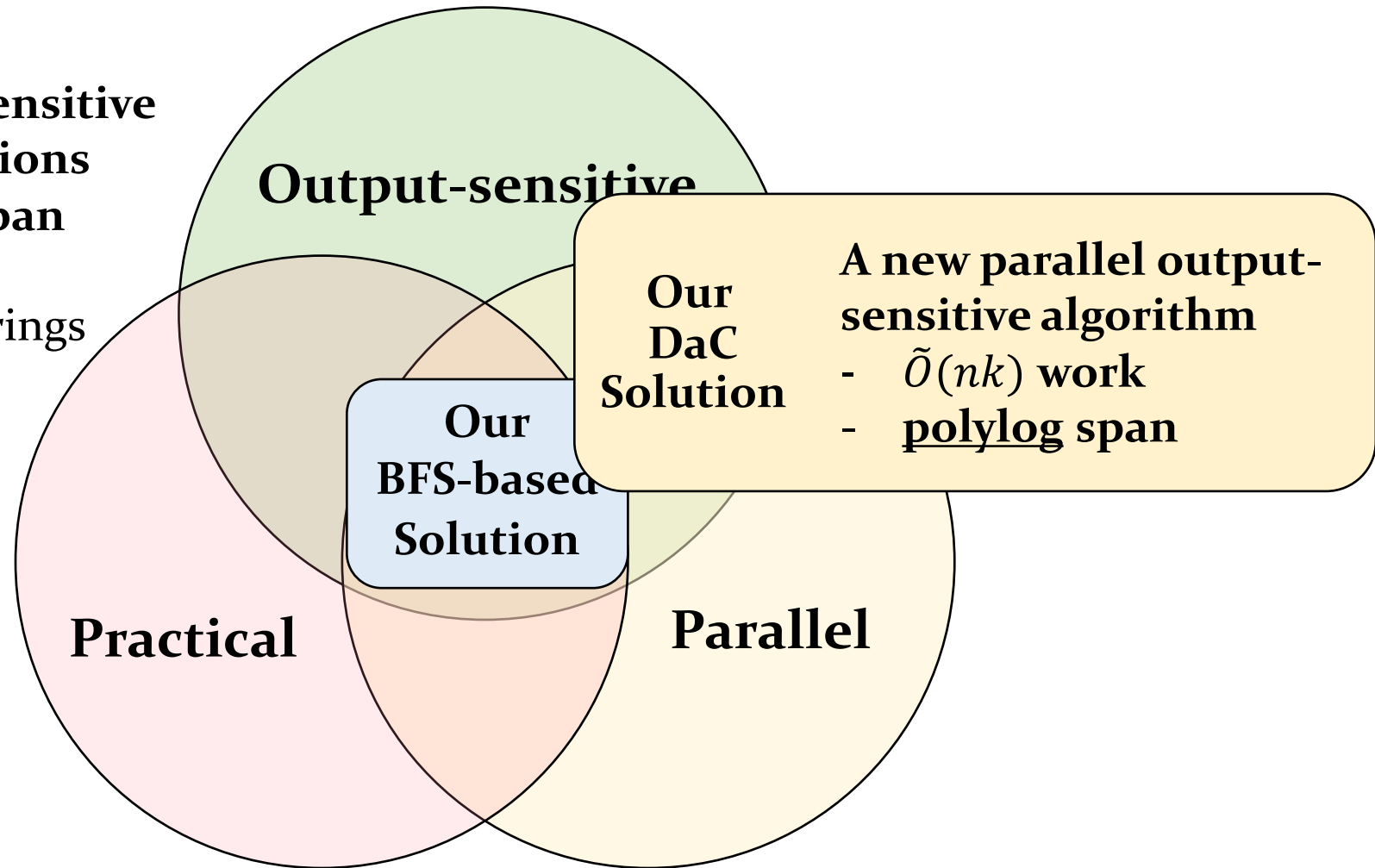


# Our contributions: in theory

## In practice

Several parallel output-sensitive algorithms/implementations

- $\tilde{O}(n + k^2)$  work,  $\tilde{O}(k)$  span
- Highly practical
- Scale to **billion-scale** strings with small distance

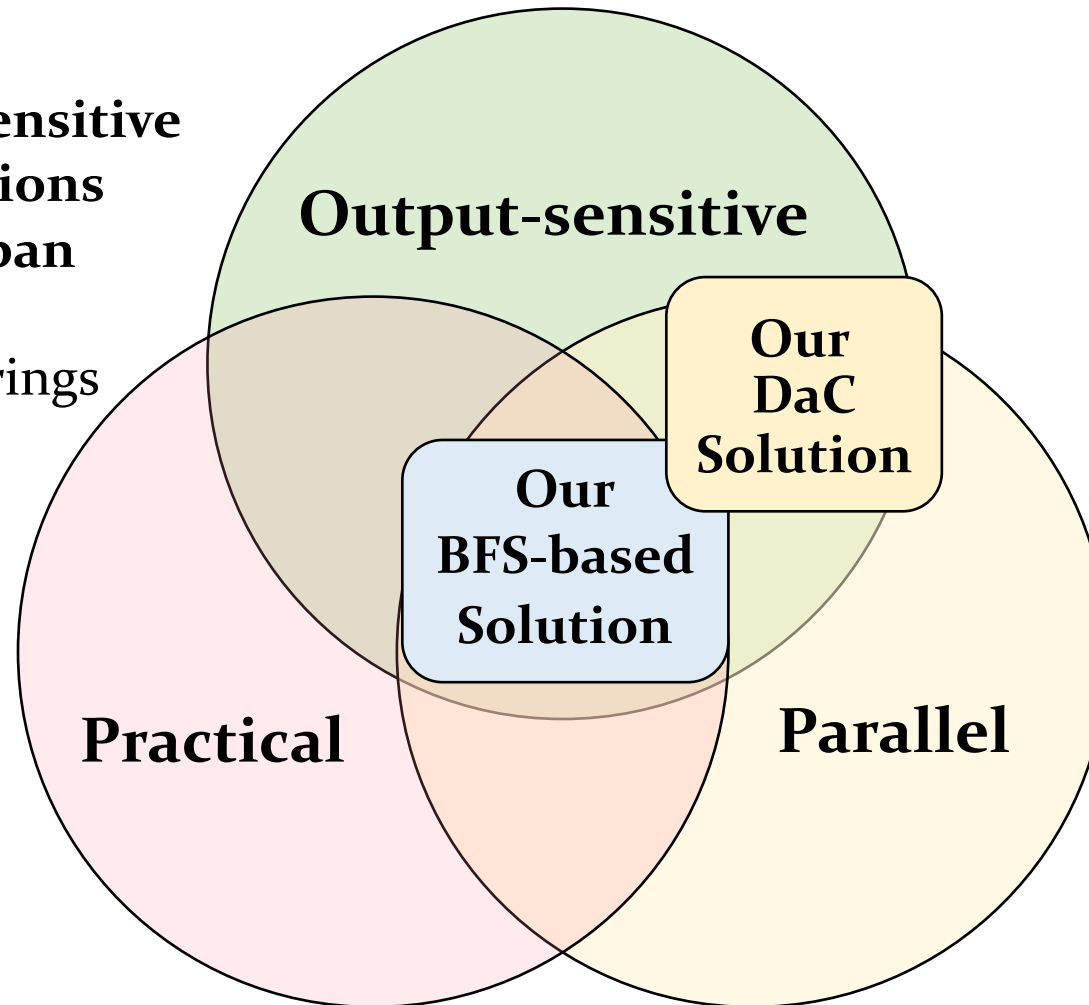


# Our contributions

## In practice

Several parallel output-sensitive algorithms/implementations

- $\tilde{O}(n + k^2)$  work,  $\tilde{O}(k)$  span
- Highly practical
- Scale to **billion-scale** strings with small distance



## In theory

A new parallel output-sensitive algorithm

- $\tilde{O}(nk)$  work
- polylog span

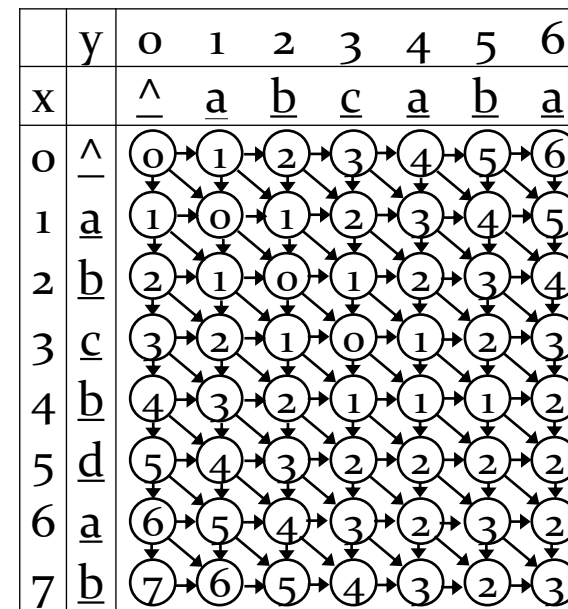
In practice

# Parallel BFS-based algorithms

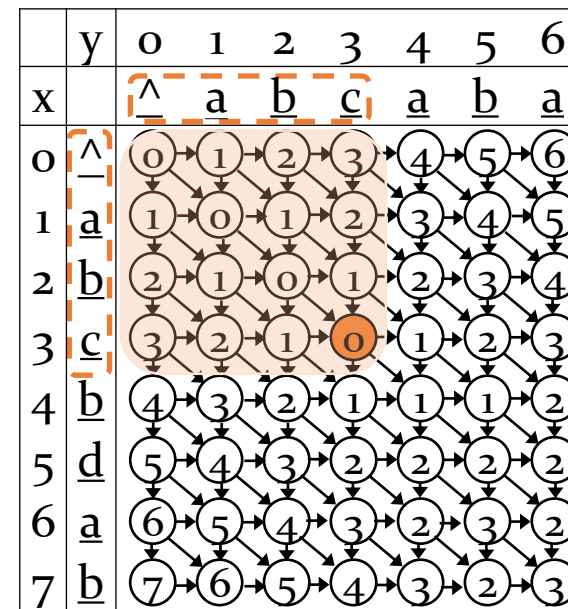
# BFS on the DAG

	y	0	1	2	3	4	5	6
x		<u>^</u>	<u>a</u>	<u>b</u>	<u>c</u>	<u>a</u>	<u>b</u>	<u>a</u>
0	<u>^</u>	0	1	2	3	4	5	6
1	<u>a</u>	1	0	1	2	3	4	5
2	<u>b</u>	2	1	0	1	2	3	4
3	<u>c</u>	3	2	1	0	1	2	3
4	<u>b</u>	4	3	2	1	1	1	2
5	<u>d</u>	5	4	3	2	2	2	2
6	<u>a</u>	6	5	4	3	2	3	2
7	<u>b</u>	7	6	5	4	3	2	3

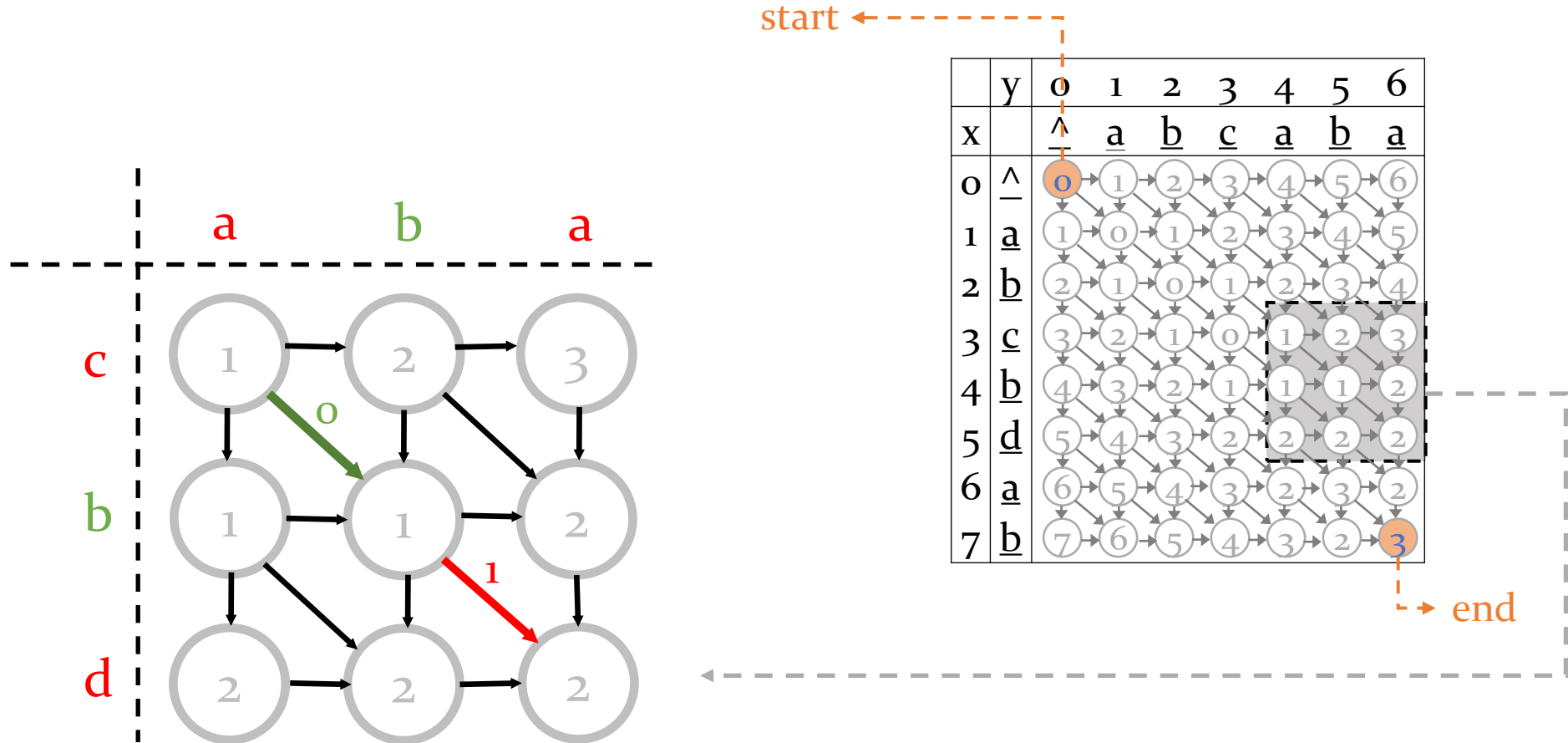
# BFS on the DAG



# BFS on the DAG



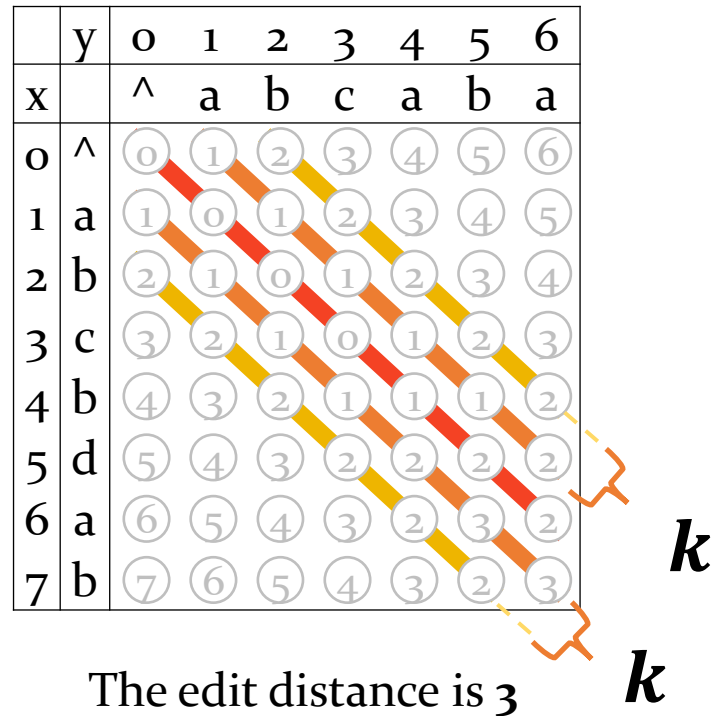
# BFS on the DAG



# BFS on the DAG: two key observations

	y	0	1	2	3	4	5	6
x		<u>^</u>	<u>a</u>	<u>b</u>	<u>c</u>	<u>a</u>	<u>b</u>	<u>a</u>
0	<u>^</u>	0	1	2	3	4	5	6
1	<u>a</u>	1	0	1	2	3	4	5
2	<u>b</u>	2	1	0	1	2	3	4
3	<u>c</u>	3	2	1	0	1	2	3
4	<u>b</u>	4	3	2	1	1	1	2
5	<u>d</u>	5	4	3	2	2	2	2
6	<u>a</u>	6	5	4	3	2	3	2
7	<u>b</u>	7	6	5	4	3	2	3

# BFS on the DAG: two key observations



# BFS on the DAG: two key observations

Vertices (states) with  $|x - y| > k$  will **NOT** be touched

**Work:**  $O(kn)$



**Output-sensitive !**

	y	o	1	2	3	4	5	6
x		^	a	b	c	a	b	a
0	^	0	1	2	3	4	5	6
1	a	1	0	1	2	3	4	5
2	b	2	1	0	1	2	3	4
3	c	3	2	1	0	1	2	3
4	b	4	3	2	1	1	1	2
5	d	5	4	3	2	2	2	2
6	a	6	5	4	3	2	3	2
7	b	7	6	5	4	3	2	3

The edit distance is 3

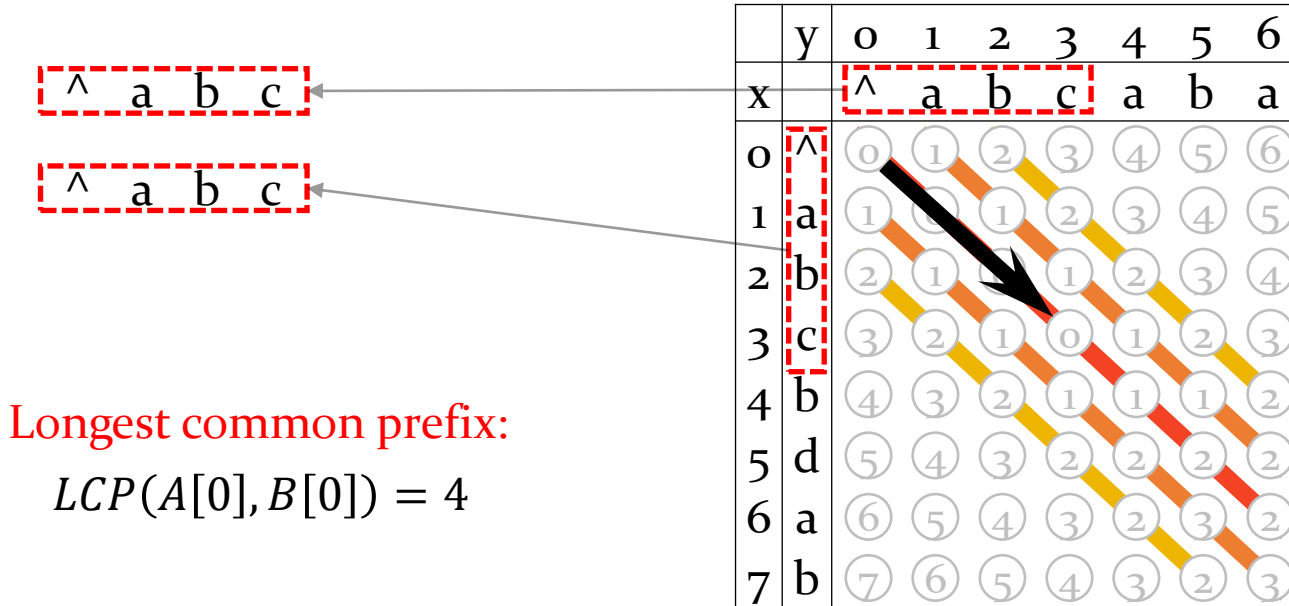
# BFS on the DAG: two key observations

Shown by **Ukkonen** in his paper:  
*Algorithms for approximate string matching. Information and Control, 1985.*

	y	o	1	2	3	4	5	6
x		^	a	b	c	a	b	a
0	^	0	1	2	3	4	5	6
1	a	1	0	1	2	3	4	5
2	b	2	1	0	1	2	3	4
3	c	3	2	1	0	1	2	3
4	b	4	3	2	1	1	1	2
5	d	5	4	3	2	2	2	2
6	a	6	5	4	3	2	3	2
7	b	7	6	5	4	3	2	3

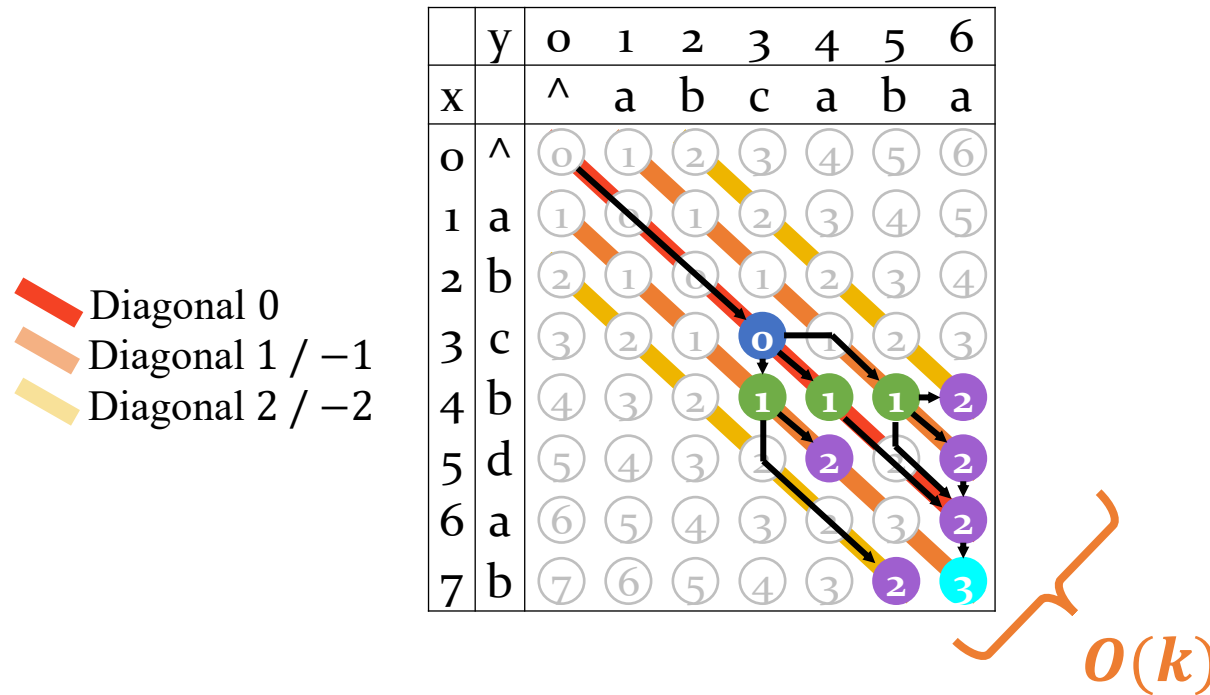
The edit distance is 3

# BFS on the DAG: two key observations



The edit distance is 3

# BFS on the DAG: two key observations



1. **Restricted computing area** →  $k$  rounds

2. **LCP computing instead** →  $O(k)$  vertices per round

$O(k^2)$  states

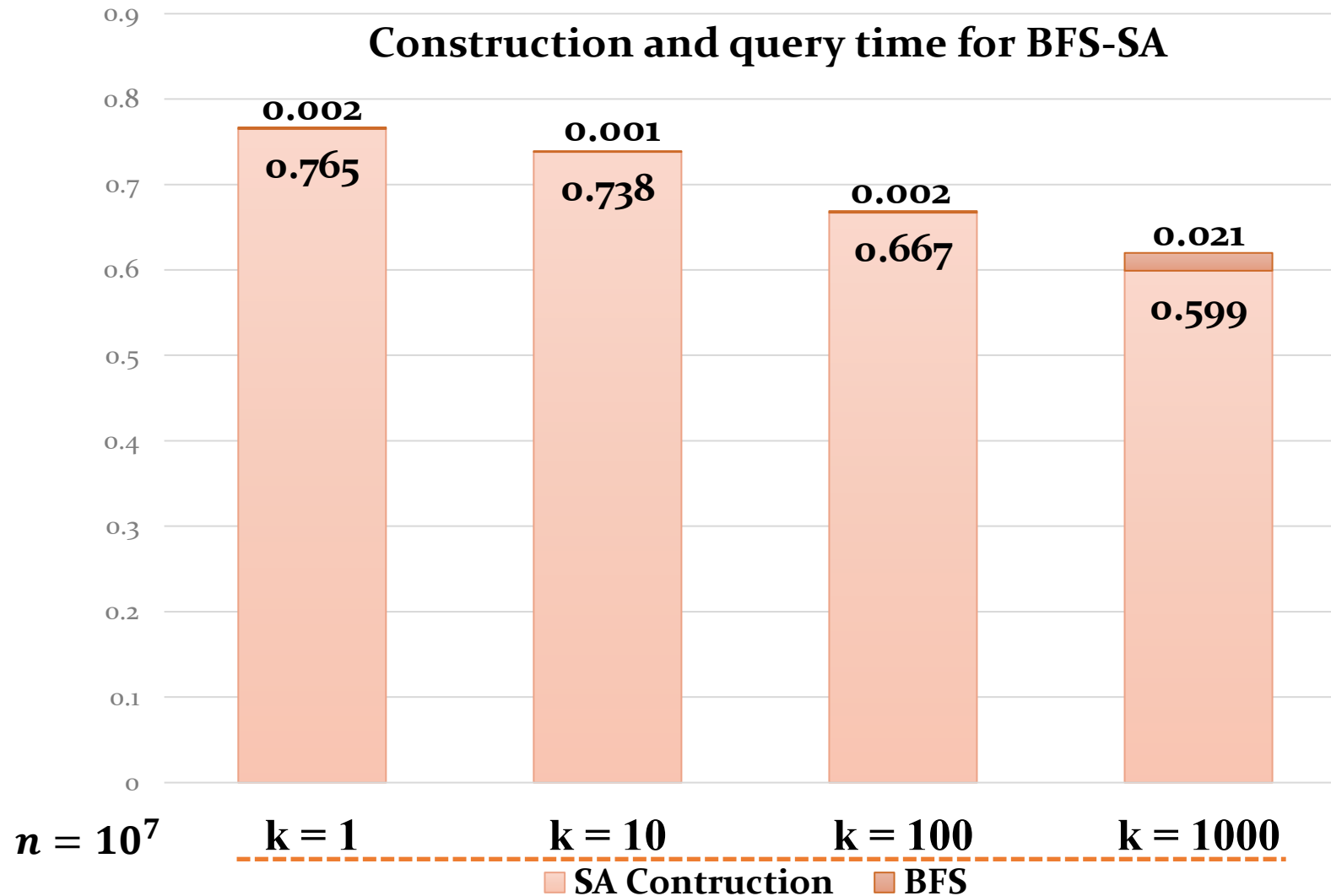
# How to compute the LCP

- **Suffix tree** by Landau and Vishkin<sup>[1]</sup>
  - Hard to implement, especially in parallel settings
  - the large constant of its work leads to a poor performance in practice
- **Our improvement & implementation: BFS-SA**
  - Alternative data structure: **Suffix array**<sup>[2]</sup>
  - Easy to implement in parallel
  - Achieve better performance

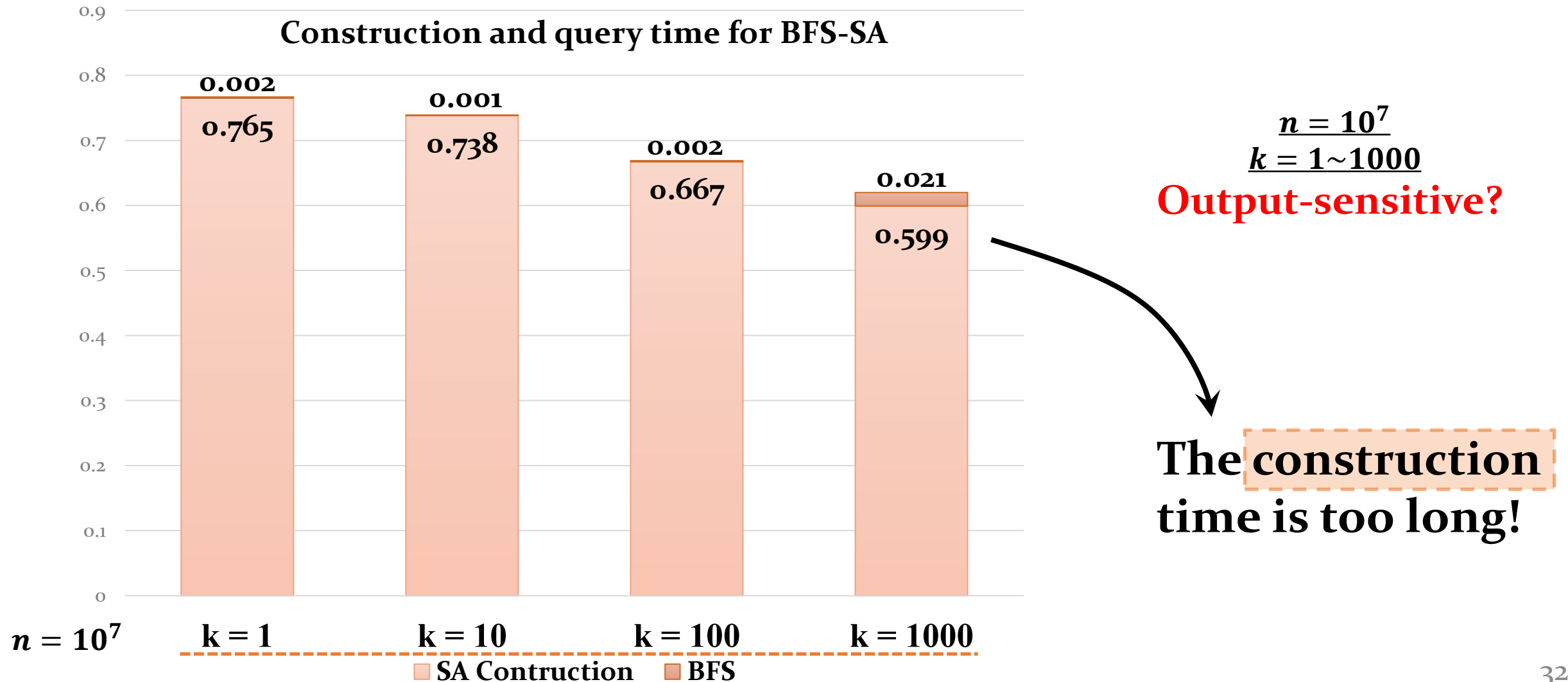
[1] Gad M Landau and Uzi Vishkin. Fast parallel and serial approximate string matching. J. Algorithms, 1989.

[2] Juha Kärkkäinen and Peter Sanders. Simple linear work suffix array construction, ICALP, 2003.

# BFS-SA performance at a glance



# BFS-SA performance at a glance



# Solutions with more practical design?

- Look at the cost of BFS-SA

- **Work:**  $O(n + k^2)$

- **Span:**  $\tilde{O}(k)$

SA Construction → BFS:  $O(1) \times O(k^2)$

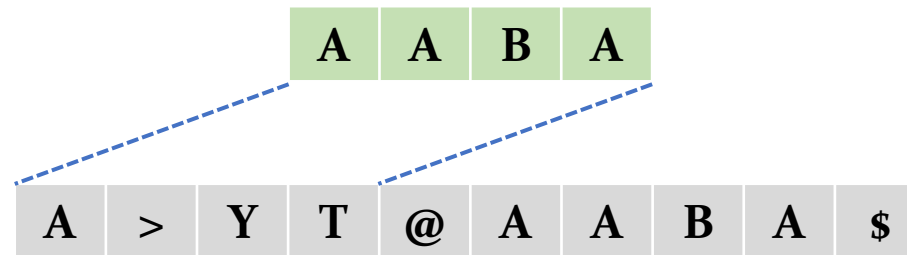
- Expecting trade-offs for better solutions
  - **Accelerate** the **construction** speed as much as possible
  - Acceptable **sacrifice** on **BFS** query work

Our contribution in practice:  
faster algorithm using hashing

# Parallel BFS-based algorithms

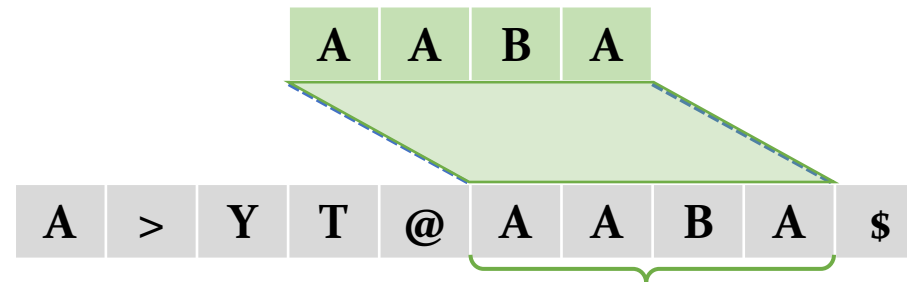
# Prefix table for edit distance

- **String hashing** for string matching
  - Compare the (sub)string based on hashing values



# Prefix table for edit distance

- **String hashing** for string matching
  - Compare the (sub)string based on hashing values



# Prefix table for edit distance

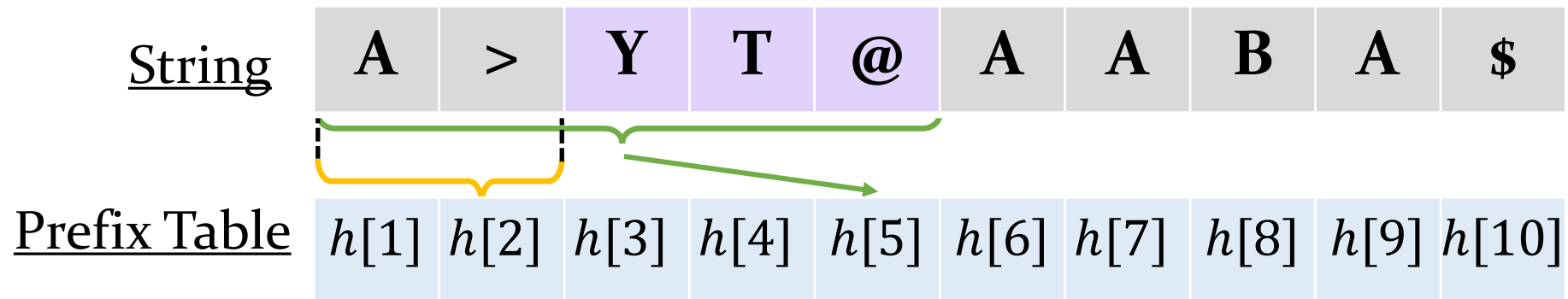
- The “string hashing arithmetic”
  - Preprocessing a “prefix table”

<u>String</u>	A	>	Y	T	@	A	A	B	A	\$
<u>Prefix Table</u>	$h[1]$	$h[2]$	$h[3]$	$h[4]$	$h[5]$	$h[6]$	$h[7]$	$h[8]$	$h[9]$	$h[10]$

$$h[4] = h(A[1..4]) = h[3] \oplus h(A[4])$$

# Prefix table for edit distance

- The “string hashing arithmetic”
  - Preprocessing a “prefix table”
  - Extracting the hashing value of a substring in  $O(1)$  work



$$h(A[2..5]) = h[5] \ominus h[2]$$

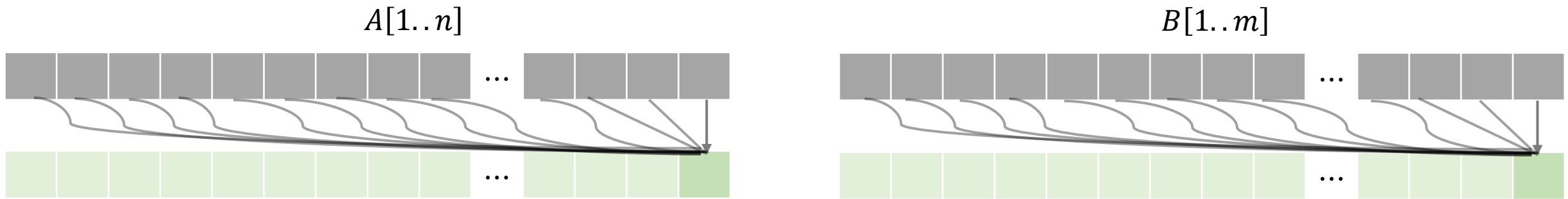
# Prefix table construction

- Prefix tables construction for fast LCP query



# Prefix table construction

- Prefix tables construction for fast LCP query



**Prefix sum in parallel**

# Query based on the prefix tables

- Prefix tables construction for fast LCP query



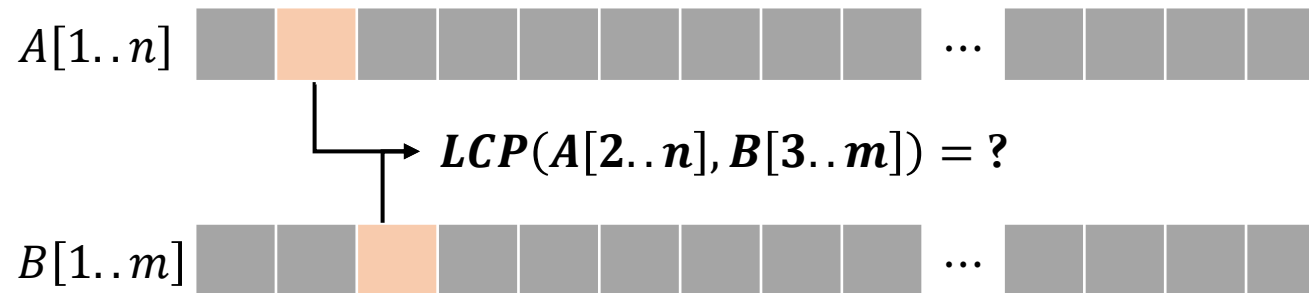
**For construction:**

**Work :**  $O(n)$

**Span:**  $O(\log n)$

# Query based on the prefix tables

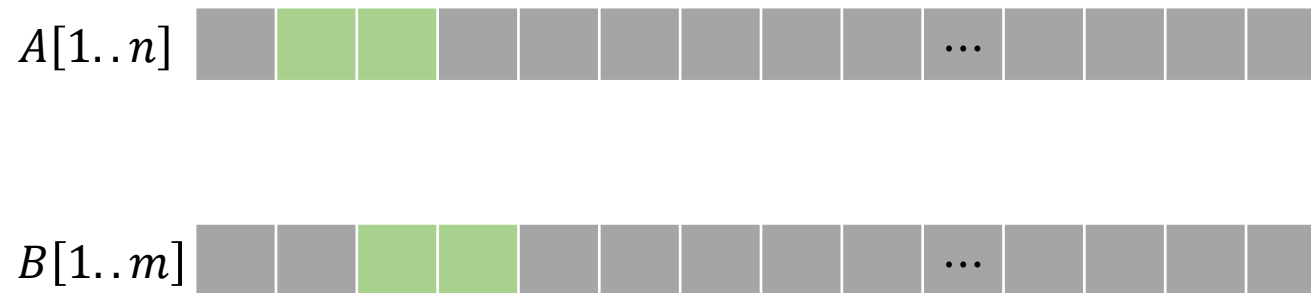
- LCP query based on the prefix tables



Dual binary search: to find the longest common prefix

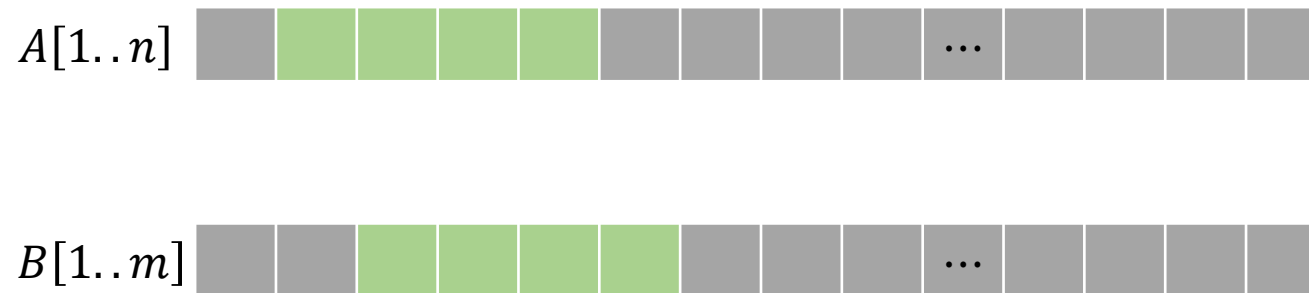
# Query based on the prefix tables

- LCP query based on the prefix tables



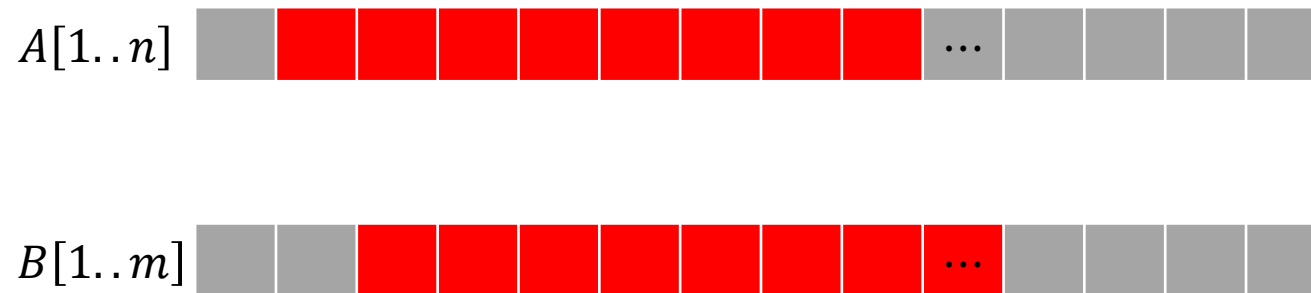
# Query based on the prefix tables

- LCP query based on the prefix tables



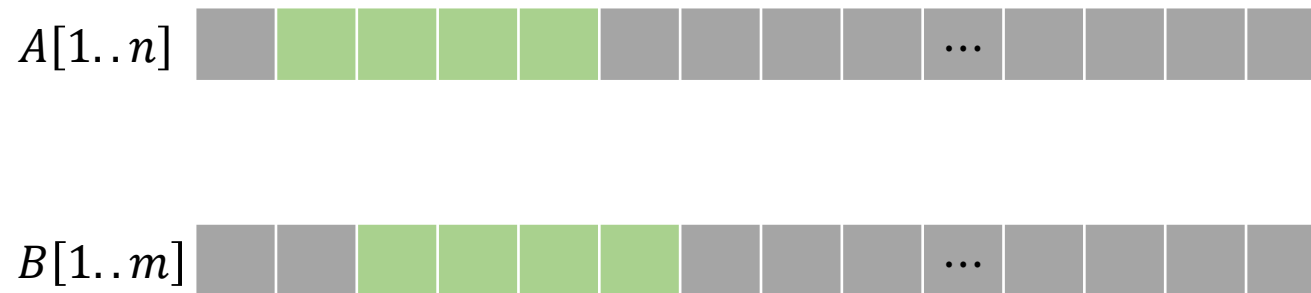
# Query based on the prefix tables

- LCP query based on the prefix tables



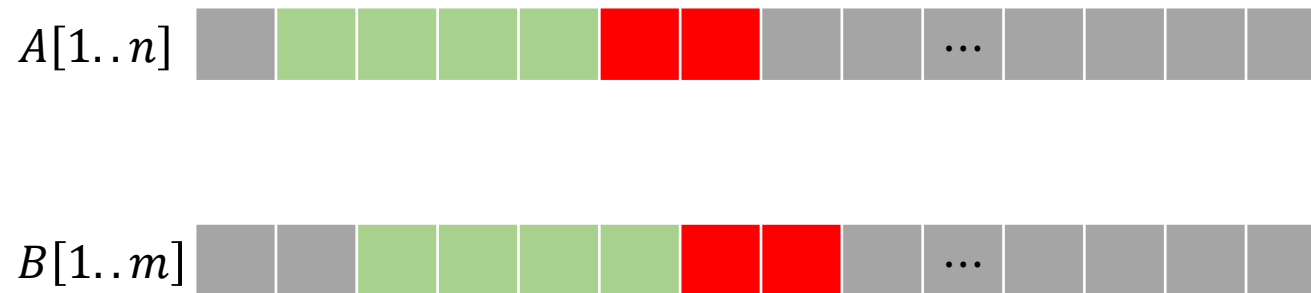
# Query based on the prefix tables

- LCP query based on the prefix tables



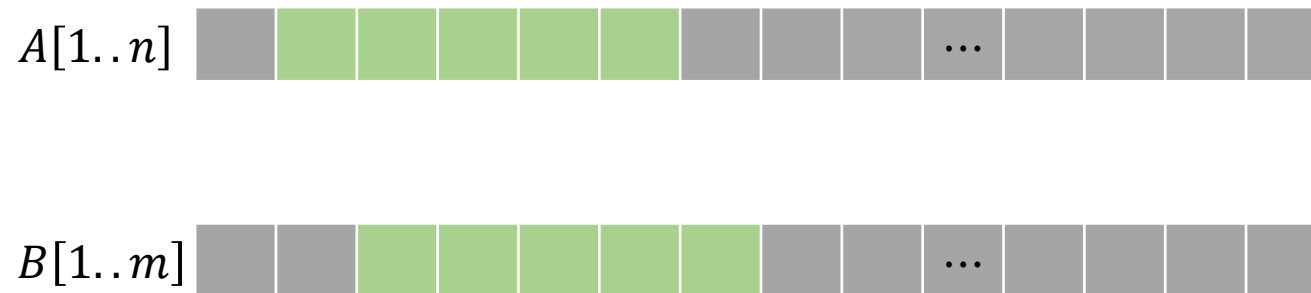
# Query based on the prefix tables

- LCP query based on the prefix tables



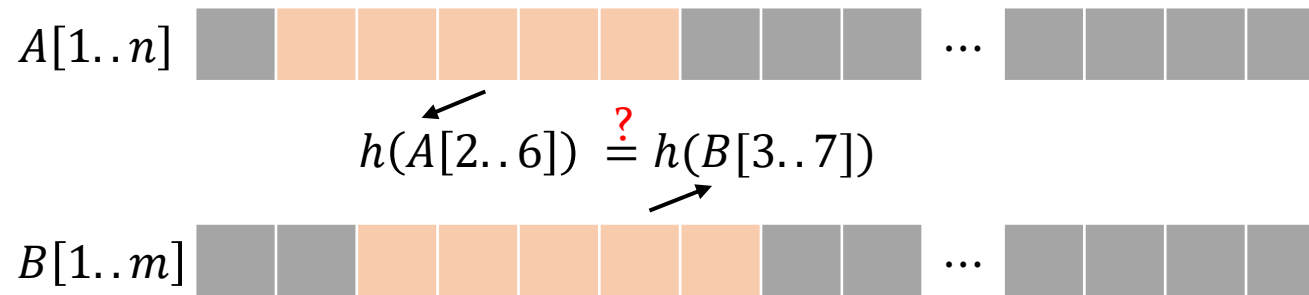
# Query based on the prefix tables

- LCP query based on the prefix tables



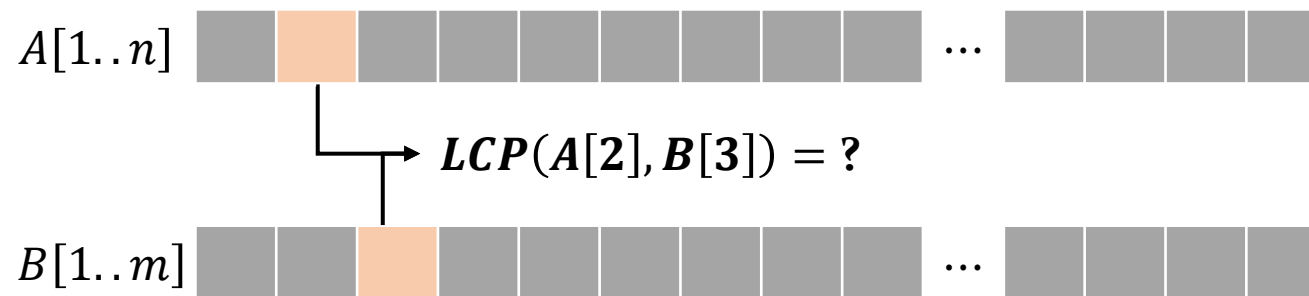
# Query based on the prefix tables

- LCP query based on the prefix tables



# Query based on the prefix tables

- LCP query based on the prefix tables



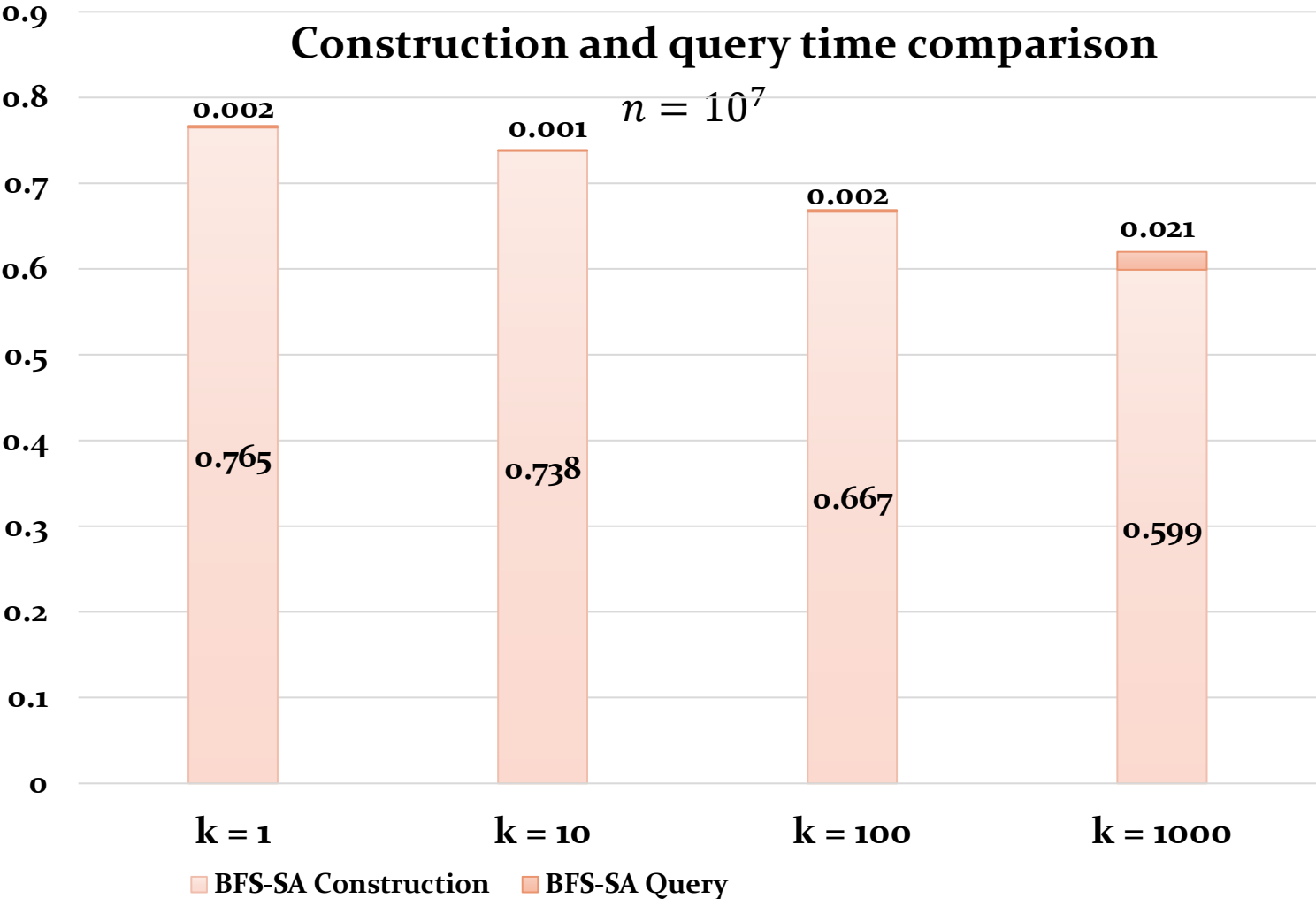
**BFS-Hash**

**Work for construction:**  $O(n)$

**Work for each query:**  $O(\log n)$

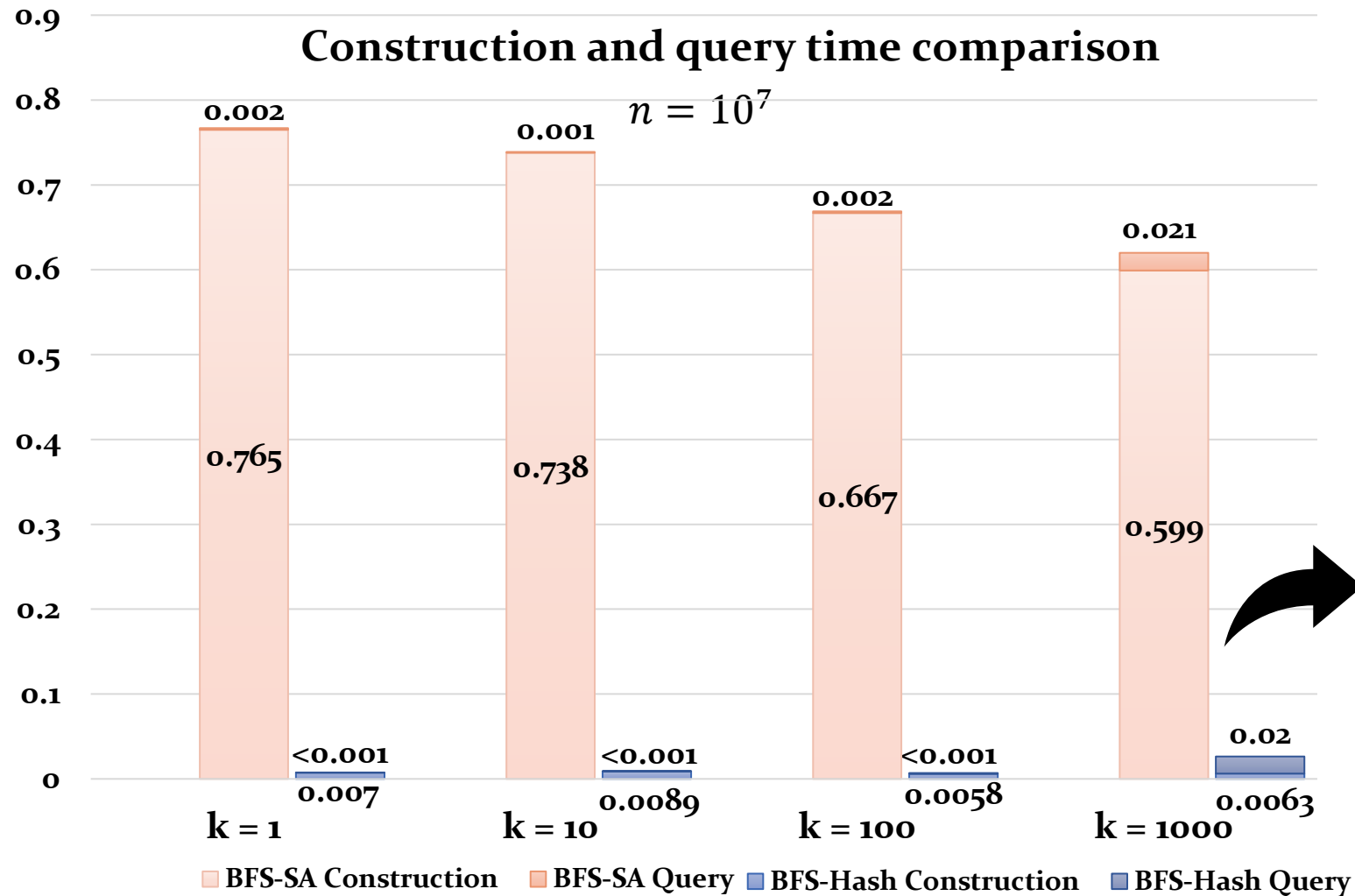
**Total work:**  $O(n + k^2 \log n)$

# Did BFS-Hash actually improve the performance?

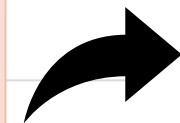


Lower is better

# BFS-Hash is significantly better than BFS-SA

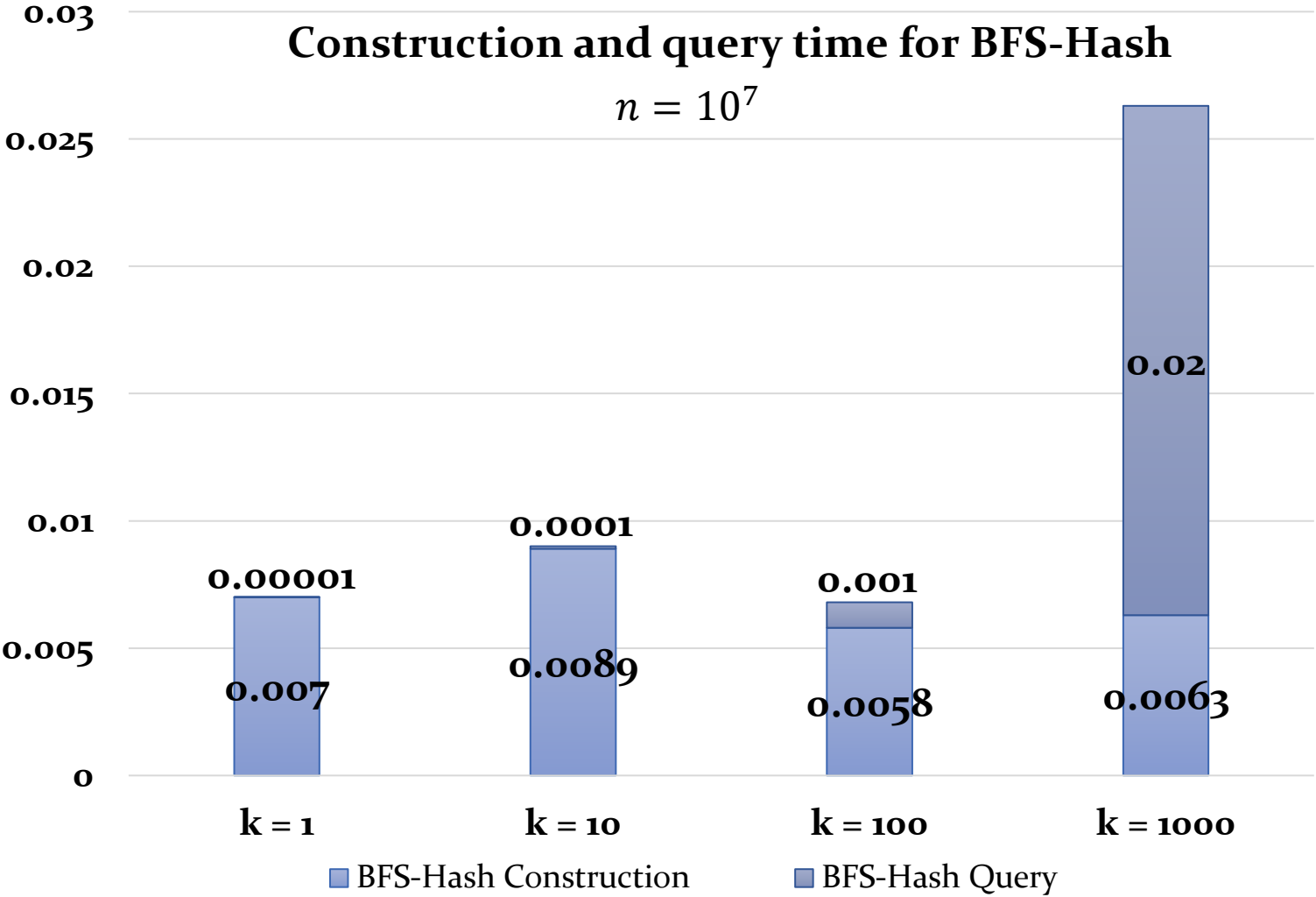


Lower is better



BFS-Hash is 24 × better than BFS-SA

# BFS-Hash is significantly better than BFS-SA



# What about the space of prefix tables?



We use the **same** size for hash coding preprocessing  $\longrightarrow O(n)$

Do we have ideas to **reduce the space usage**?

# What about the space of prefix tables?

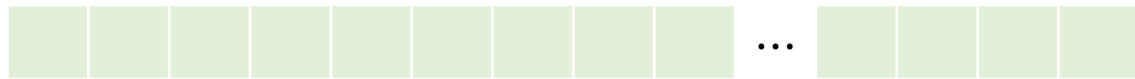
$A[1..n]$



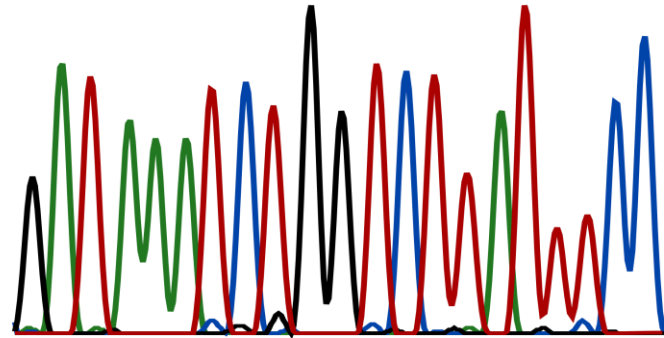
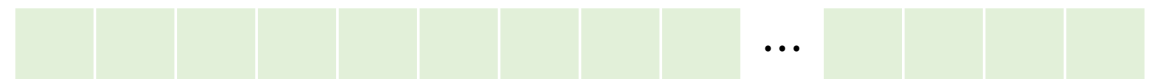
$B[1..m]$



Prefix table for  $A[1..n]$



Prefix table for  $B[1..m]$



120 130  
GAT AAAT CT GGTCTT ATTTCC

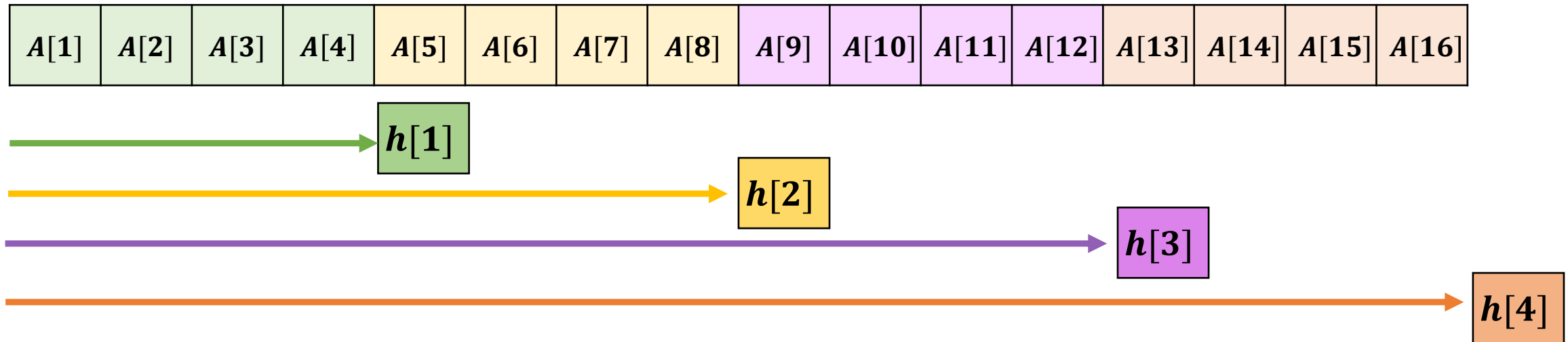
# BFS-B-Hash: block-version prefix table

$A[1]$	$A[2]$	$A[3]$	$A[4]$	$A[5]$	$A[6]$	$A[7]$	$A[8]$	$A[9]$	$A[10]$	$A[11]$	$A[12]$	$A[13]$	$A[14]$	$A[15]$	$A[16]$
--------	--------	--------	--------	--------	--------	--------	--------	--------	---------	---------	---------	---------	---------	---------	---------

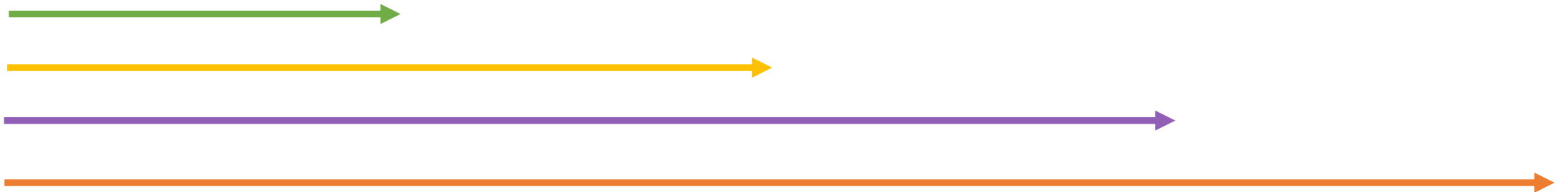
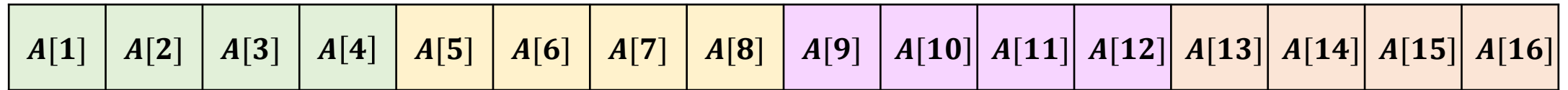
# BFS-B-Hash: block-version prefix table

$A[1]$	$A[2]$	$A[3]$	$A[4]$	$A[5]$	$A[6]$	$A[7]$	$A[8]$	$A[9]$	$A[10]$	$A[11]$	$A[12]$	$A[13]$	$A[14]$	$A[15]$	$A[16]$
--------	--------	--------	--------	--------	--------	--------	--------	--------	---------	---------	---------	---------	---------	---------	---------

# BFS-B-Hash: block-version prefix table



# BFS-B-Hash: block-version prefix table



Block size: 4

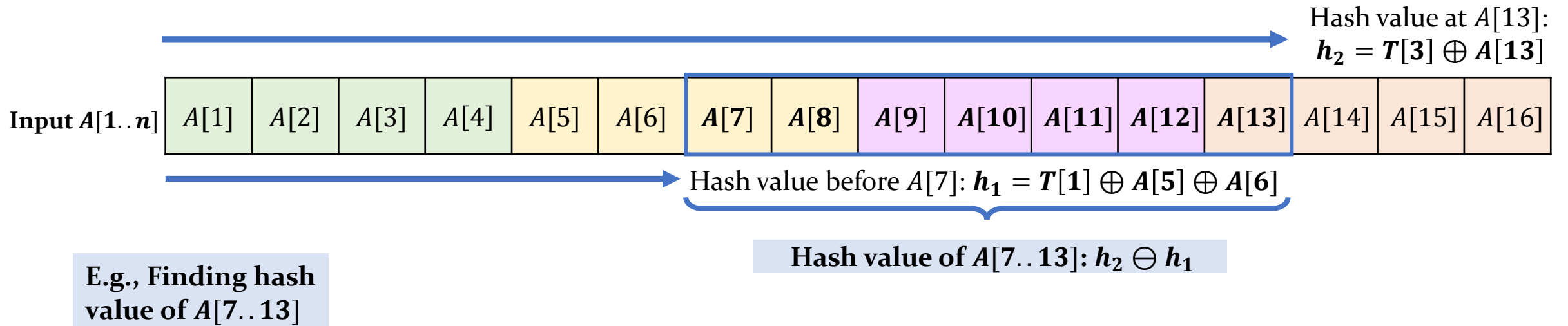
Sequence size: 16

B-Hash Prefix table size:  $16 / 4 = 4$

The B-Hash Prefix Table 

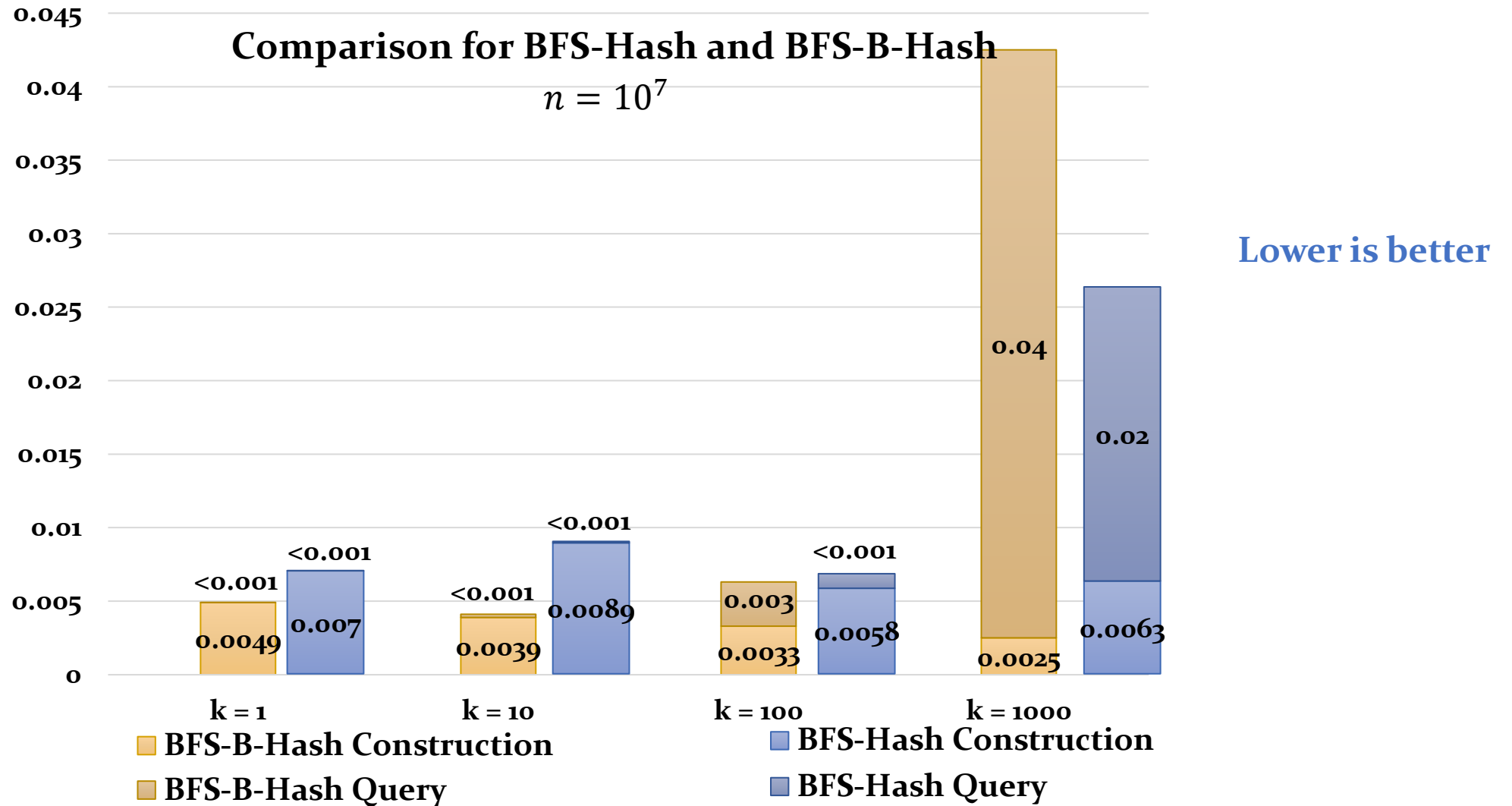
$h[1]$	$h[2]$	$h[3]$	$h[4]$
--------	--------	--------	--------

# BFS-B-Hash: Query



Work for one query:  $O(b \log n)$

# BFS-B-Hash: block-version prefix table

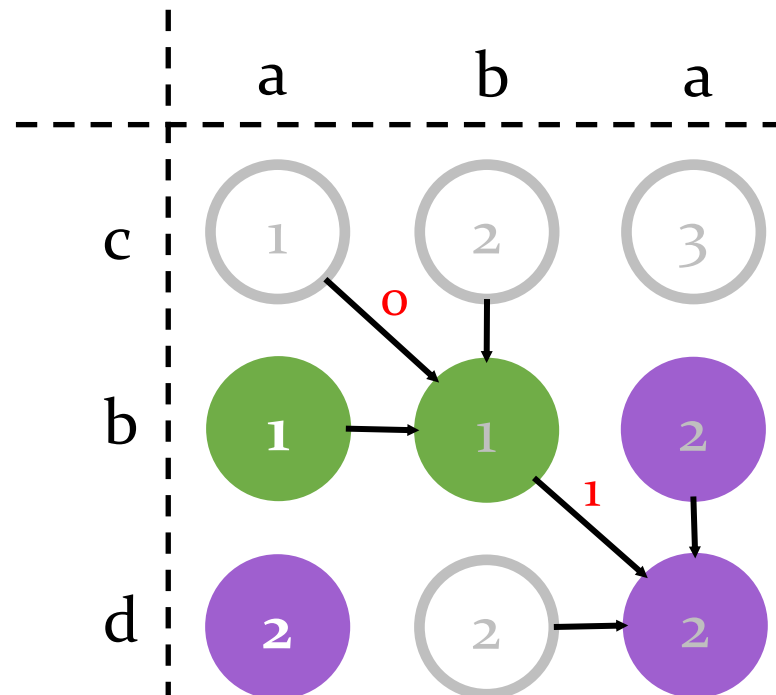


*Our contribution in theory*

**Parallel Divide-and-conquer-based algorithms**

# Revisit the DP matrix

- DP matrix -> **Directed acyclic graph (DAG)**
- The solution of edit distance is the shortest distance from  $(0, 0)$  to  $(n, m)$



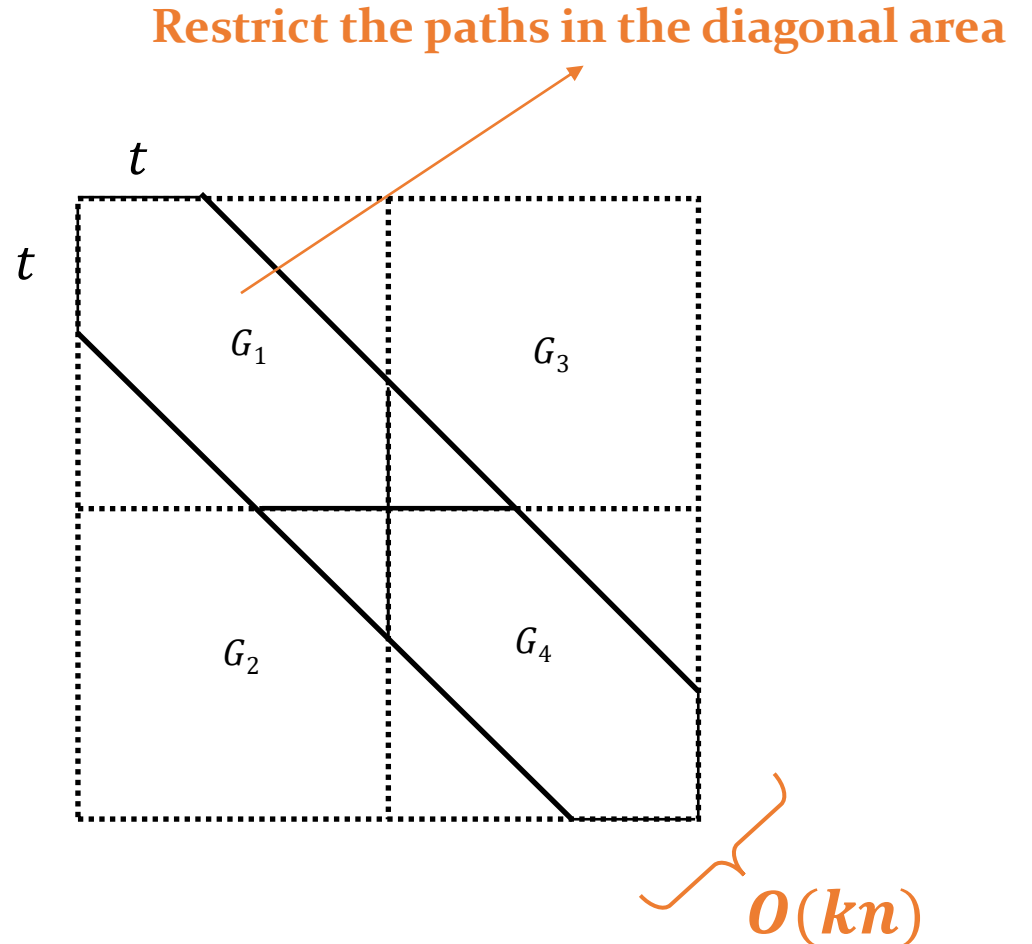
	y	0	1	2	3	4	5	6
x	^	a	b	c	a	b	a	
0	^	0	1	2	3	4	5	6
1	a	1	0	1	2	3	4	5
2	b	2	1	0	1	2	3	4
3	c	3	2	1	0	1	2	3
4	b	4	3	2	1	0	1	2
5	d	5	4	3	2	1	0	1
6	a	6	5	4	3	2	1	0
7	b	7	6	5	4	3	2	1

- Frontier 0
- Frontier 1
- Frontier 2
- Frontier 3

# Our exploration of the DaC-based approach

- The AALM algorithm<sup>[1]</sup>
  - The **monotonicity** of the paths
  - Theoretically, the AALM can achieve
    - Work:  $O(n^2 \log n)$
    - Span:  $O(\log^3 n)$
- Our algorithm **DAC-SD**
  - **Work:**  $O(nk \log k)$
  - **Span:**  $O(\log n \log^3 k)$

best so far under output-sensitive settings



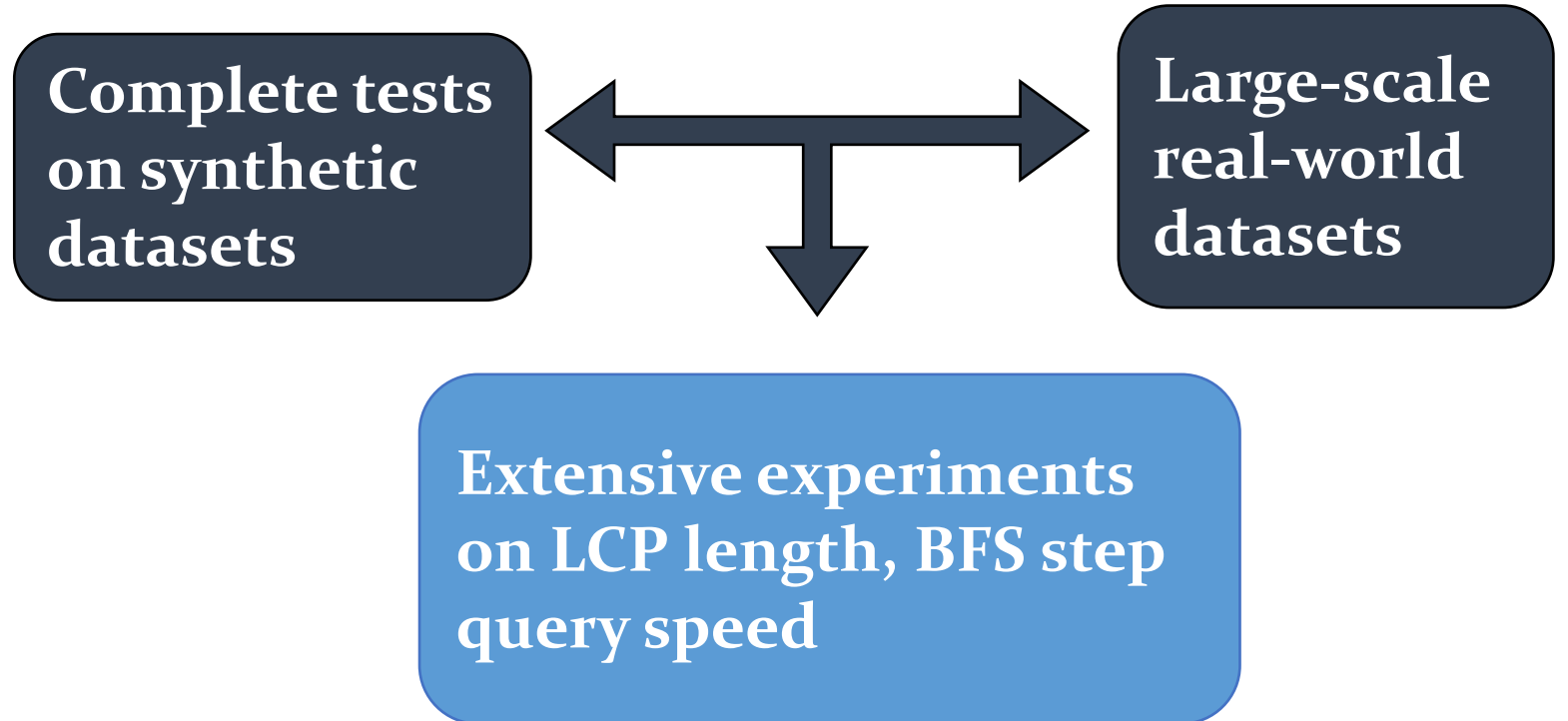
# Our exploration of the DaC-based approach

- Our algorithm **DAC-SD**
  - **Work:**  $O(nk \log k)$
  - **Span:**  $O(\log n \log^3 k)$
- Carefully implemented **DAC-SD** with careful engineering details
  - Relatively worse performance than BFS-based algorithms in practice

# Overview of the experiments

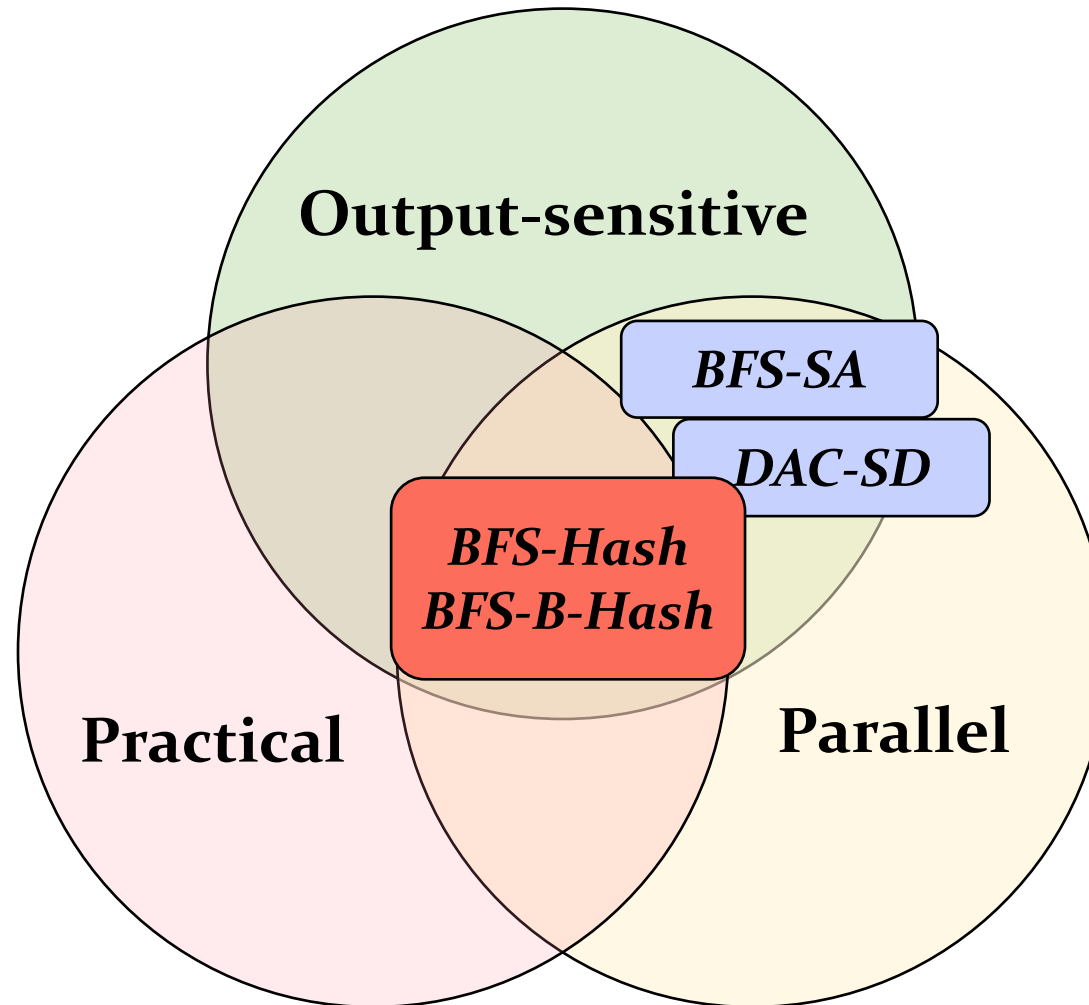
## Experiments setup

- ParlayLib<sup>[1]</sup> for **fork-join parallelism** and **primitives**
- **96-core (192 hyperthreads)** machine with four Intel Xeon Gold 6252 CPUs



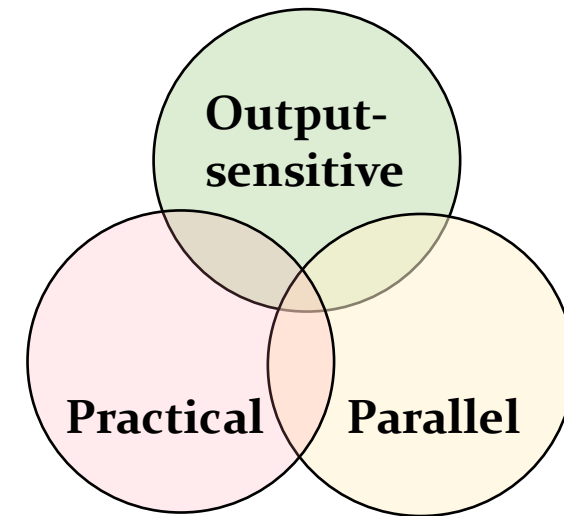
# Conclusions

- Summary



# Conclusions

	<u>Work</u>	<u>Span</u>	<u>Space</u>
<div data-bbox="76 496 738 568" style="border: 1px solid blue; padding: 2px;">Our efficient implementations</div> <div data-bbox="76 682 738 811" style="border: 1px solid purple; padding: 2px;">Our theoretical improvement: Improved theoretical bound</div>	<div data-bbox="828 406 1141 502" style="border: 1px solid blue; border-radius: 10px; padding: 5px; text-align: center;"><i><b>BFS-SA</b></i></div> <div data-bbox="828 564 1141 659" style="border: 1px solid blue; border-radius: 10px; padding: 5px; text-align: center;"><i><b>DAC-SD</b></i></div>	$O(n + k^2)$ $\tilde{O}(k)$	$O(n)$ $O(nk)$
<div data-bbox="76 892 738 1013" style="border: 1px solid red; padding: 2px;">Our improved algorithms: Best performance in practice</div>	<div data-bbox="802 849 1151 1053" style="border: 1px solid red; border-radius: 10px; padding: 5px; text-align: center;"><i><b>BFS-Hash</b></i> <i><b>BFS-B-Hash</b></i></div>	$O(n + k^2 \log n)$ $O(n + k^2 b \log n)$	$\tilde{O}(k)$ $\tilde{O}\left(\frac{n}{b} + k\right)$ $O(n)$ $O(n/b + k)$



- Code on GitHub:
  - <https://github.com/ucrparlay/Edit-Distance>
- Contact
  - Youzhe Liu ([yliu908@ucr.edu](mailto:yliu908@ucr.edu))