

Чисельні  
характеристики  
побудови узагальнених  
FM-index для алфавітів  
різних розмірів

Підготував студент ІПЗ-4  
Фомін Володимир  
Науковий керівник: Зважій Д.В.

# Мета роботи

- Розробка ефективного алгоритму побудови узагальненого FM-index.
- Ретельне дослідження та аналіз чисельних характеристик побудови узагальненого FM-index залежно від розміру алфавіту.
- Визначення оптимального розміру алфавіту для ефективної роботи з узагальненим FM-index у веббраузері.

# Хід ВИКОНАННЯ

Оптимізація  
алгоритму

Аналіз  
чинників впливу на  
чисельні хар-ки

01 ..... 02 ..... 03 ..... 04 ..... 05

Вивчення  
класичного рішення

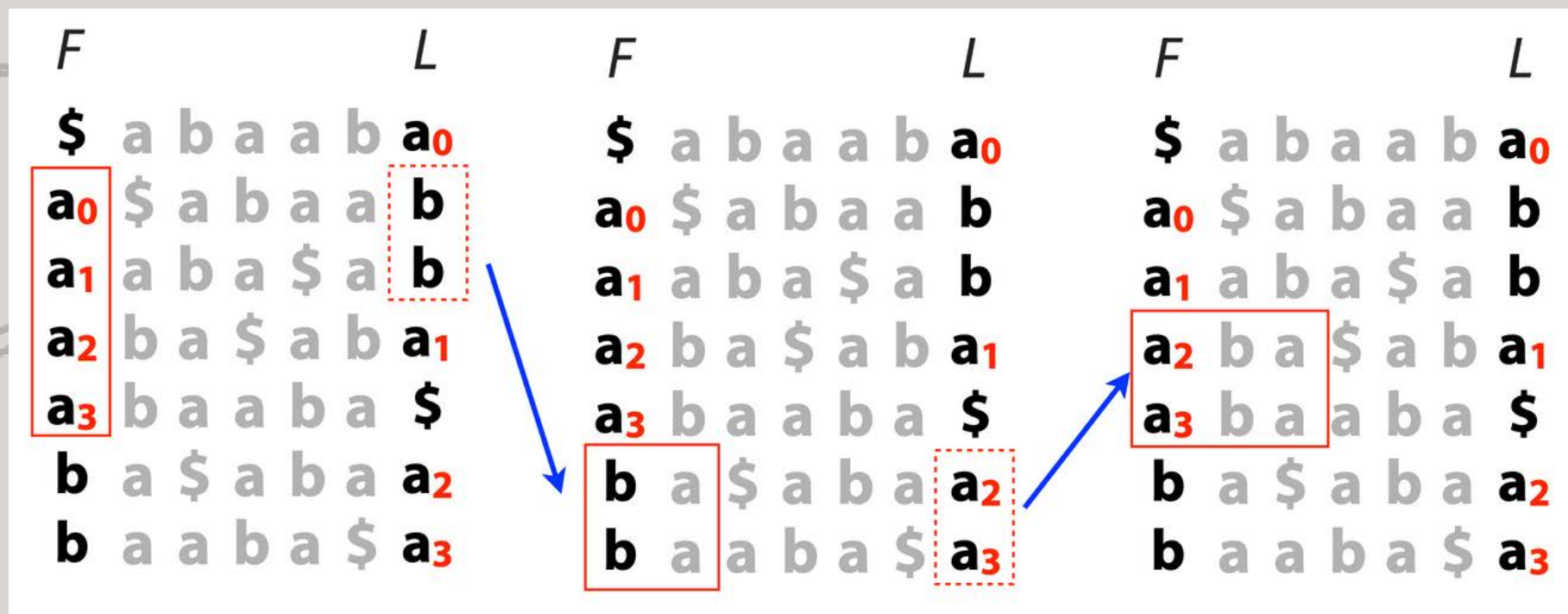
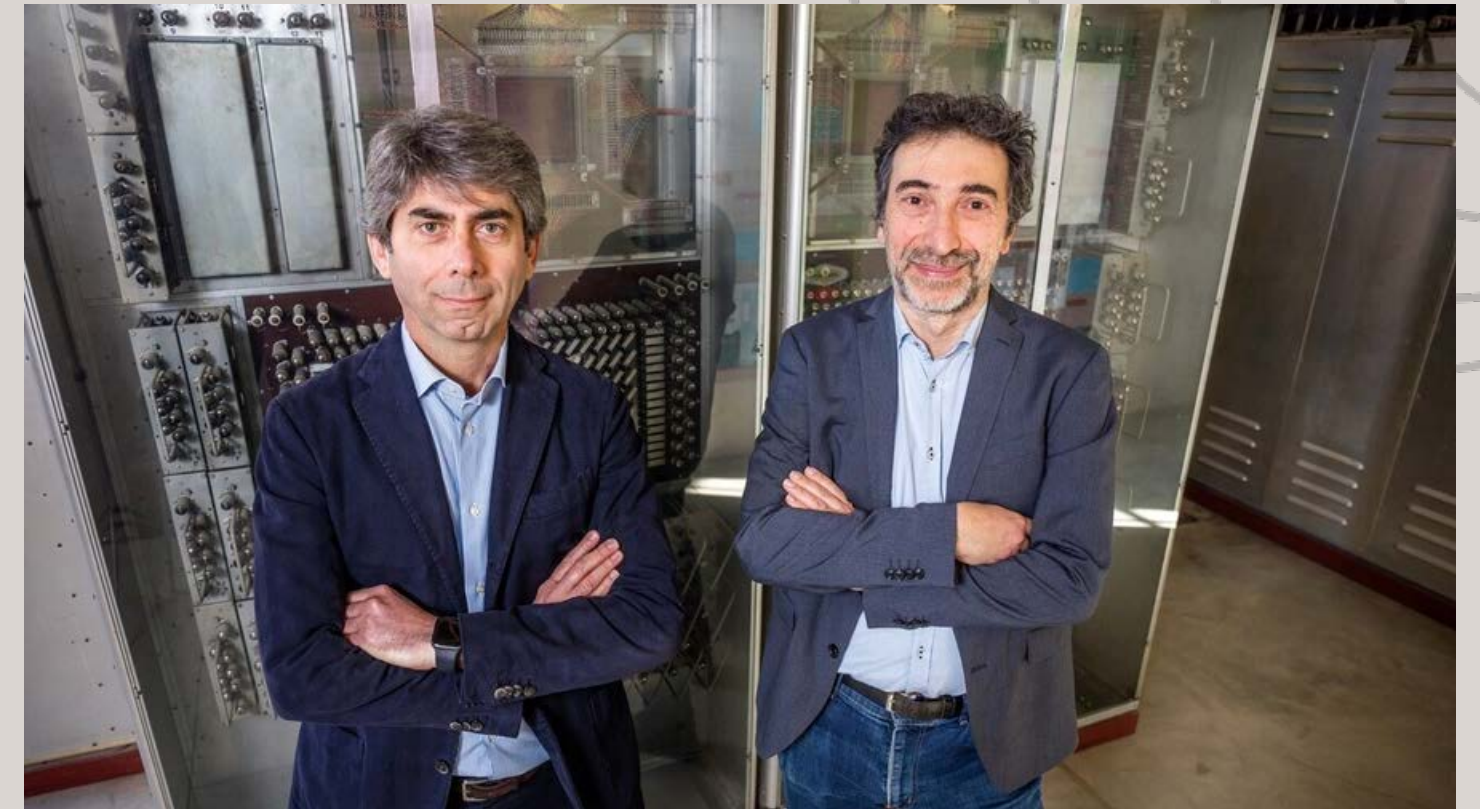
Модифікація  
під узагальнену  
структуру даних

Визначення  
оптимального  
розміру даних

# FM-index

Структура даних, яка дозволяє одночасно стискати вхідний текст та виконувати швидкий пошук підрядка. Вперше представлено Паоло Феррагіно та Джованні Манзіні у 2000 році.

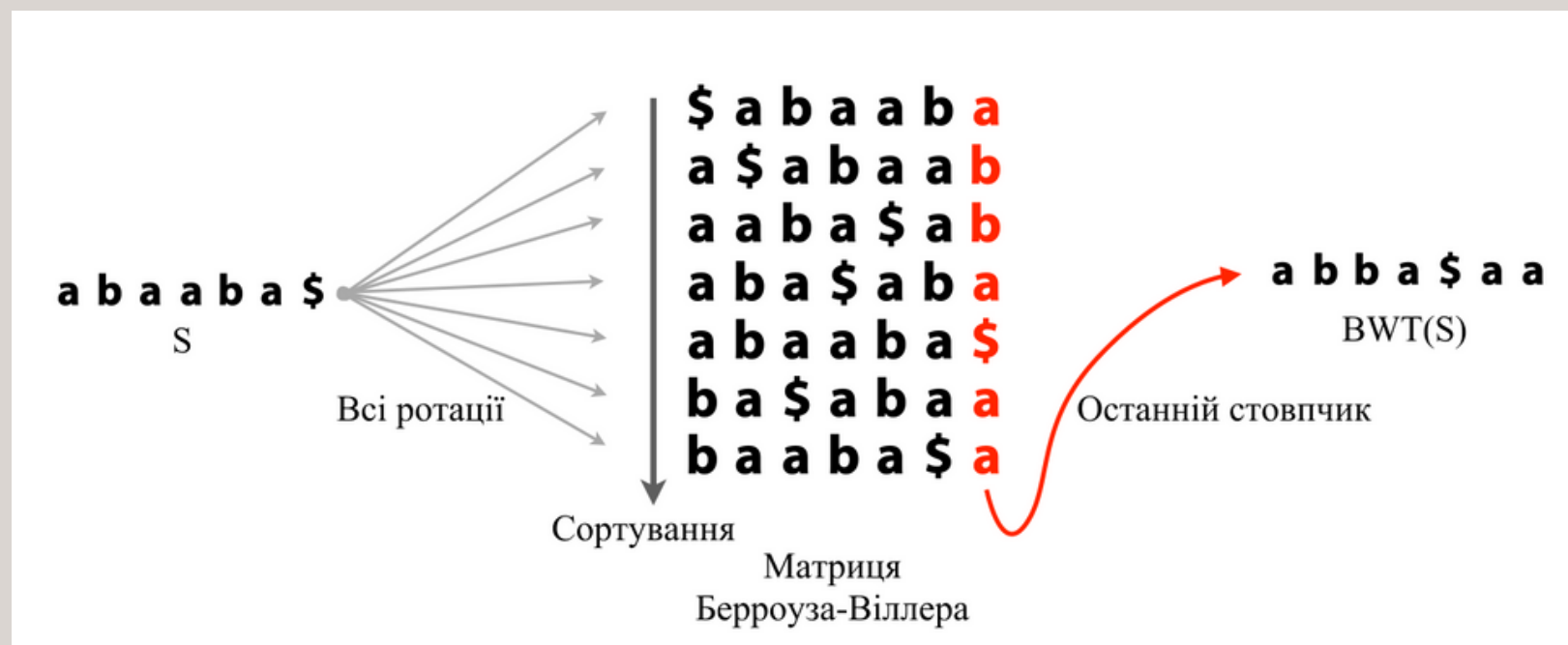
Широко використовується у біоінформатиці та інформаційному пошуку. Є дуже ефективним для обробки одного великого рядка з алфавітом невеликого розміру (ДНК).



Для побудови необхідно провести попередній аналіз тексту, що може вимагати багато ресурсів. Це є особливо критичним для тексту, який динамічно змінюється.

Менш ефективний для деяких типів запитів (пошук за виключенням).

# Перетворення Берроуза-Вілера



Перетворення Берроуза-Вілера (Burrows-Wheeler transform, BWT) — метод перестановки символів рядка у результаті якого отримується стрічка з якої можливо отримати початкову послідовність символів.

Алгоритм перетворення:

1. Створити  $n$  ротацій рядка ( $n$  – довжина рядка).
2. Створити матрицю BWM розташувавши кожен ротацію одна під одною.
3. Лексикографічно впорядкувати матрицю.
4. Скласти новий рядок за допомогою конкатенації останніх символів кожної ротації починаючи згори.

Перетворення відповідає новому рядку утвореному в останньому кроці.

# Суфіксні масиви

Суфіксний масив (Suffix array) — масив, де зберігається інформація про всі суфікси певного рядка відсортовані у лексикографічному порядку.

Суфіксний масив можна використати для здійснення перетворення Берроуза-Вілера завдяки їх схожим властивостям.

```

$ a b a a b a
a $ a b a a b
a a b a $ a b
a b a $ a b a
a b a a b a $
b a $ a b a a
b a a b a $ a

```

BWM(S)

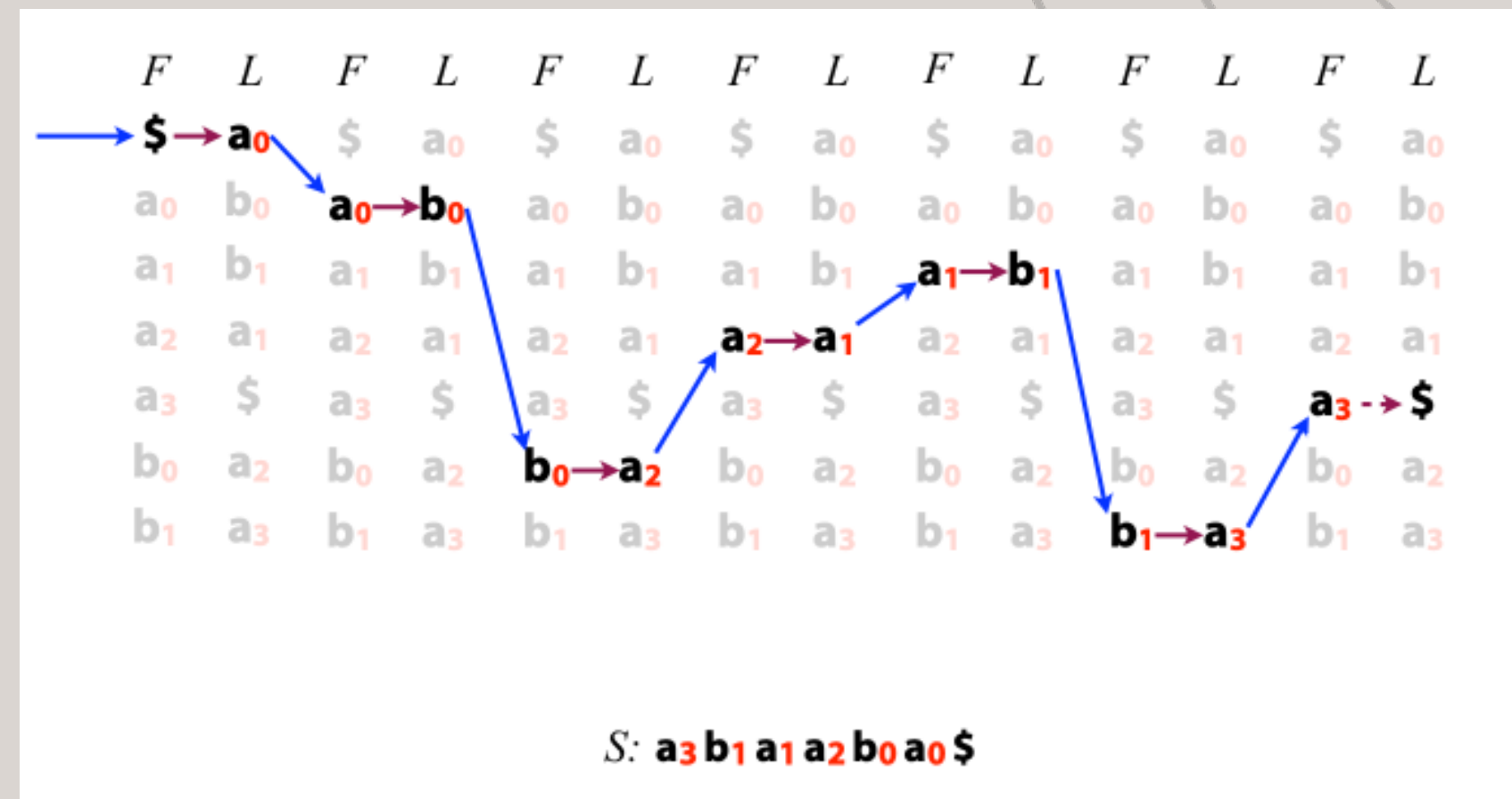
6	<b>\$</b>								
5	<b>a</b>	<b>\$</b>							
2	<b>a</b>	<b>a</b>	<b>b</b>	<b>a</b>	<b>\$</b>				
3	<b>a</b>	<b>b</b>	<b>a</b>	<b>\$</b>					
0	<b>a</b>	<b>b</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>a</b>	<b>\$</b>		
4	<b>b</b>	<b>a</b>	<b>\$</b>						
1	<b>b</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>a</b>	<b>\$</b>			

SA(S)

# LF-Відображення

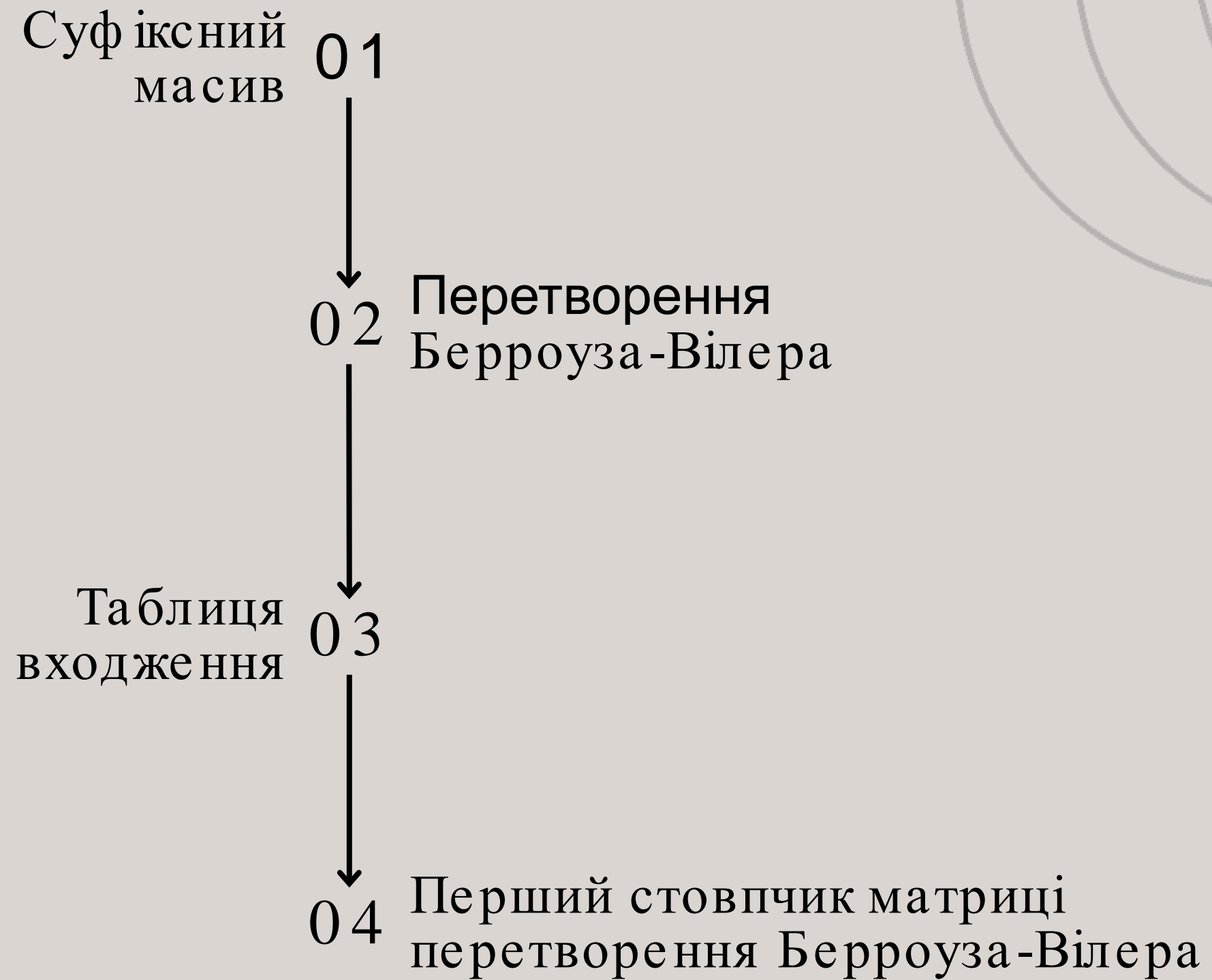
LF-Відображення (LF-Mapping) — операція, яка використовується для оберненого перетворення BWT та пошуку в FM-index.

<i>F</i>		<i>L</i>
\$	a b a a b	<b>a<sub>0</sub></b>
<b>a<sub>0</sub></b>	\$ a b a a	<b>b<sub>0</sub></b>
<b>a<sub>1</sub></b>	a b a \$ a	<b>b<sub>1</sub></b>
<b>a<sub>2</sub></b>	b a \$ a b	<b>a<sub>1</sub></b>
<b>a<sub>3</sub></b>	b a a b a	\$
<b>b<sub>0</sub></b>	a \$ a b a	<b>a<sub>2</sub></b>
<b>b<sub>1</sub></b>	a a b a \$	<b>a<sub>3</sub></b>



Ґрунтується на тому, що у матриці перетворення Берроуза-Вілера (BWM) зберігається відносний порядок входжень символів в першій та останній колонках.

# Алгоритм побудови



# Оптимізації

01

Суфіксний масив для перетворення Берроуза-Вілера

F	L	SA
\$ a b a a b a		6
a \$ a b a a b		
a a b a \$ a b		2
a b a \$ a b a		
a b a a b a \$		0
b a \$ a b a a		4
b a a b a \$ a		

F	L	SA
\$ a b a a b a		6
a \$ a b a a b		
a a b a \$ a b		2
a b a \$ a b a	a	
a b a a b a \$		0
b a \$ a b a a		4
b a a b a \$ a		

02

Часткове збереження суфіксного масиву та таблиці входження

L	a	b
a	1	0
b	1	1
b		
a	2	2
\$	2	2
a		
a	4	2

03

Компактне збереження першого стовпчика матриці BWM

F	L
\$ a b a a b a	a <sub>3</sub>
a \$ a b a a b	b <sub>1</sub>
a a b a \$ a b	b <sub>0</sub>
a b a \$ a b a	a <sub>1</sub>
a b a a b a \$	a <sub>3</sub>
b a \$ a b a a	b <sub>0</sub>
b a a b a \$ a	a <sub>2</sub>
b a a b a \$ a	b <sub>1</sub>

# Узагальнений FM-index

Узагальнений FM-index передбачає роботу з колекцією рядків, що вимагає певні модифікації:

- Для перетворення Берроуза-Вілера, необхідно конкатенувати всі рядки колекції.
- У таблиці входження необхідно зберігати інформацію про термінальний символ.

String1	MSBWM	MSBWT
ACCA\$	\$ACCA\$CAAA	A
CCA\$A	\$CAAA\$ACCA	A
CA\$AC	A\$ACCA\$CAA	A
A\$ACC	A\$CAAA\$ACC	C
\$ACCA	AA\$ACCA\$CA	A
	AAA\$ACCA\$C	C
	ACCA\$CAAA\$	\$
String2	CA\$CAAA\$AC	C
CAAA\$	CAAA\$ACCA\$	\$
AAA\$C	CCA\$CAAA\$A	A
AA\$CA		
A\$CAA		
\$CAAA		

# Реалізація та вимірювання

Choose file eng4.txt  
FmIndex built in 11146 milliseconds

The screenshot shows the Chrome DevTools Memory tab. The top bar indicates the file 'eng4.txt' and the construction time 'FmIndex built in 11146 milliseconds'. The main area shows the 'Constructor' call stack for the 'FmIndex' object. The object details on the right show the following properties:

Property	Value
cps	FmCheckpoints @43195
ssa	Object @43193
bwt	"blgcpqmbgenjlpidmngnImilneiqheodfjcnjllnojdampamccgegoahkdobdcmnfcgppqfbnkbnacaodedgfnifghaqcmloijcgfpqldjdaiffhjlobdpaehnnhqkppqiocbkqfaejgggimllhgbmhjlc"
first	Object @43107
__proto__	Object @43059
map	system / Map @43113

Choose file eng4.txt  
FmIndex built in 11146 milliseconds

The screenshot shows the Chrome DevTools Memory tab with a 'Statistics' view. A donut chart displays the memory usage breakdown for the 'FmIndex' object:

Category	Size
Code	159 kB
Strings	92 kB
JS arrays	2 kB
Typed arrays	24 kB
System objects	277,748 kB
<b>Total</b>	<b>278,995 kB</b>

```
class FmIndex {
  static downsampleSuffixArray(sa, n = 4) {
    const ssa = {};
    for (let i = 0; i < sa.length; i++) {
      if (sa[i] % n === 0) {
        ssa[i] = sa[i];
      }
    }
    return ssa;
  }

  constructor(t, cpIval = 4, ssaIval = 4) {
    const sa = suffixArray(t);
    this.bwt = bwtFromSa(t, sa);
    this.ssa = FmIndex.downsampleSuffixArray(sa, ssaIval);
    this.slen = this.bwt.length;
    this.cps = new FmCheckpoints(this.bwt, cpIval);
    const tots = {};
    for (const c of this.bwt) {
      tots[c] = (tots[c] || 0) + 1;
    }
    this.first = {};
    let totc = 0;
    for (const [c, count] of Object.entries(tots).sort()) {
      this.first[c] = totc;
      totc += count;
    }
  }

  count(c) {
    if (!(c in this.first)) {
      let cc;
      for (cc of Object.keys(this.first).sort()) {
        if (c < cc) return this.first[cc];
      }
      return this.first[cc];
    } else {
      return this.first[c] || 0;
    }
  }

  range(p) {
    let l = 0, r = this.slen - 1;
    for (let i = p.length - 1; i >= 0; i--) {
      const offset = this.count(p[i]);
      l = offset + this.cps.rank(this.bwt, p[i], l - 1);
      r = offset + this.cps.rank(this.bwt, p[i], r - 1);
      if (r < l) break;
    }
    return [l, r + 1];
  }

  resolve(row) {
    const stepLeft = (row) => {
      const c = this.bwt[row];
      return this.cps.rank(this.bwt, c, row - 1) + (this.count(c) || 0);
    };
    let nsteps = 0;
    while (!(row in this.ssa)) {
      row = stepLeft(row);
      nsteps++;
    }
    return this.ssa[row] + nsteps;
  }

  hasSubstring(p) {
    const [l, r] = this.range(p);
    return r > l;
  }

  occurrences(p) {
    const [l, r] = this.range(p);
    const occurrences = [];
    for (let x = l; x < r; x++) {
      occurrences.push(this.resolve(x));
    }
    return occurrences;
  }
}
```

# Результати

№	Назва	So	Σ	T	S	Опис
1	Eng	88 Мб	26	122,197 с	2719,84 Мб	8 млн рядків довжиною 5-15 символів англійського алфавіту
2	Eng2	88 Мб	8	117,27 с	1331,69 Мб	8 млн рядків довжиною 5-15 символів скороченого англійського алфавіту
3	Eng3	88 Мб	17	111,552 с	2001,22 Мб	8 млн рядків довжиною 5-15 символів скороченого англійського алфавіту
4	Eng4	11 Мб	17	11,146 с	264,55 Мб	1 млн рядків довжиною 5-15 символів скороченого англійського алфавіту
5	Ukr	839 Кб	34	269,619 с	17,82 Мб	40 тис. рядків довжиною 5-15 символів українського алфавіту
6	Kor	928 Кб	19	106,315 с	10,65 Мб	30 тис. рядків довжиною 5-15 символів корейського алфавіту

# Висновки

01

Дослідження  
Вивчення особливостей  
структури даних,  
класичного способу  
побудови та можливих  
варіантів оптимізації.

02

Реалізація  
Написання алгоритму  
побудови на JavaScript,  
модифікація під роботу з  
колекцією рядків, подальша  
оптимізація.

03

Вимірювання  
Тестування роботи на  
різних даних, визначення  
оптимальних розмірів  
даних, встановлення  
максимальної межі розміру.