

Методи оптимізації використання пам'яті для узагальнених суфіксних масивів

Підготувала студентка КН-4
Барабуха Марія
Науковий керівник: Глибовець А.М.

Мета роботи

1. Розробка алгоритму, що дозволить виконувати швидкий пошук по даним у вебсистемі.
2. Досягнення максимальної межі розміру даних, по яким виконуватиметься пошук.
3. Часова оптимізація та зменшення затрат пам'яті на побудову та зберігання структури.

Хід виконання

Вибір

оптимального
підходу

Застосування

методів
компресії

01 — 02 — 03 — 04 — 05

Дослідження

наявних рішень

Адаптація

під вебсистему
та узагальнені
структури

Встановлення

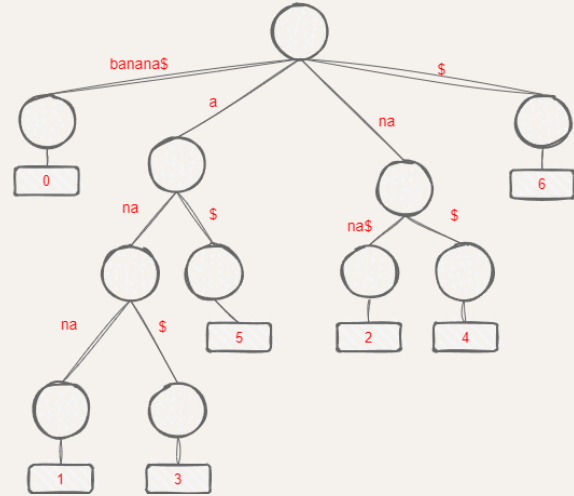
допустимих
меж розміру
даних

Класичне рішення

Стандартним рішенням для швидкого пошуку є використання *суфіксних дерев*.

Попри їх перевагу в швидкодії, вони мають недолік – велика *затрата* як *пам'яті*, так і часу на побудову, що ускладнює їхнє застосування у вебсистемі.

Стандартною альтернативою для суфіксних дерев, що використовує менше ресурсів, є *суфіксні масиви*.



Суфіксні масиви

Суфіксний масив (*suffix array*) – це масив, що складається зі стартових позицій всіх суфіксів певного рядка, впорядкованих в алфавітному порядку

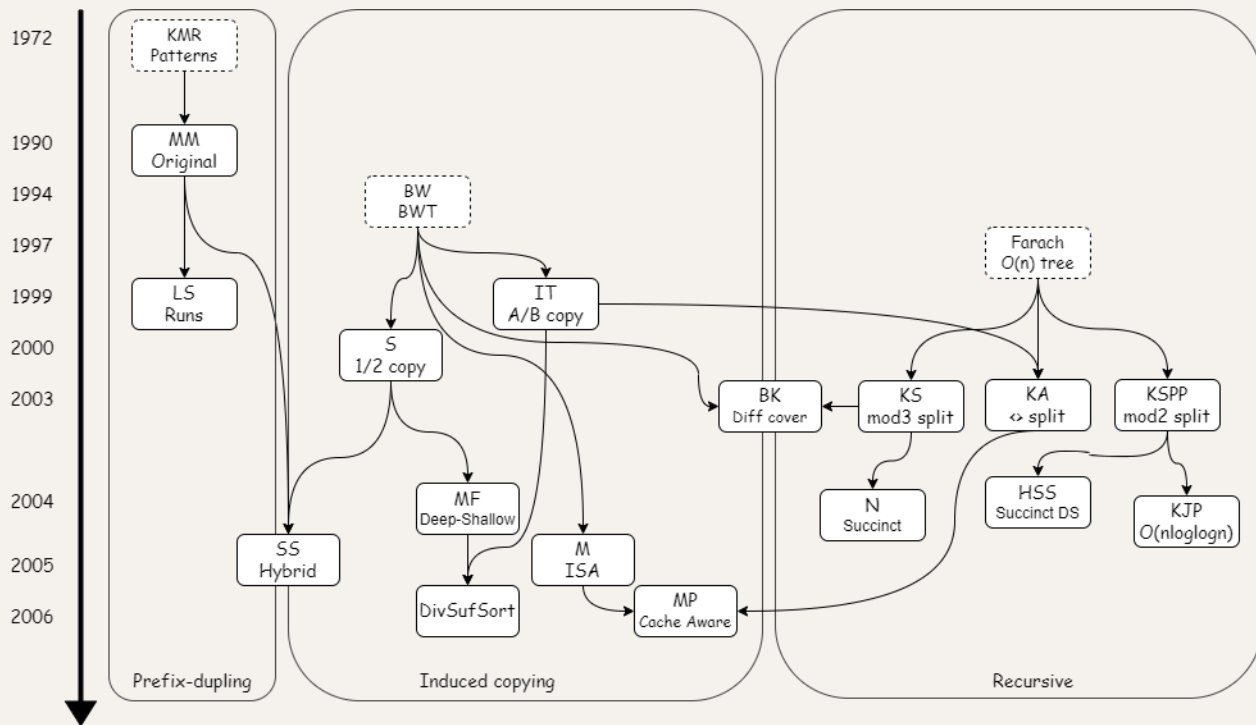
(0,0) banana\$	(1,0) ananas\$	(0,6) \$	(1,4) as\$
(0,1) anana\$	(1,1) nanas\$	(1,6) \$	(0,0) banana\$
(0,2) nana\$	(1,2) anas\$	(0,5) a\$	(0,4) na\$
(0,3) ana\$	(1,3) nas\$	(0,3) ana\$	(0,2) nana\$
(0,4) na\$	(1,4) as\$	(0,1) anana\$	(1,1) nanas\$
(0,5) a\$	(1,5) s\$	(1,0) ananas\$	(1,3) nas\$
(0,6) \$	(1,6) \$	(1,2) anas\$	(1,5) s\$

0 banana\$	6 \$
1 anana\$	5 a\$
2 nana\$	3 ana\$
3 ana\$	1 anana\$
4 na\$	0 banana\$
5 a\$	4 na\$
6 \$	2 nana\$

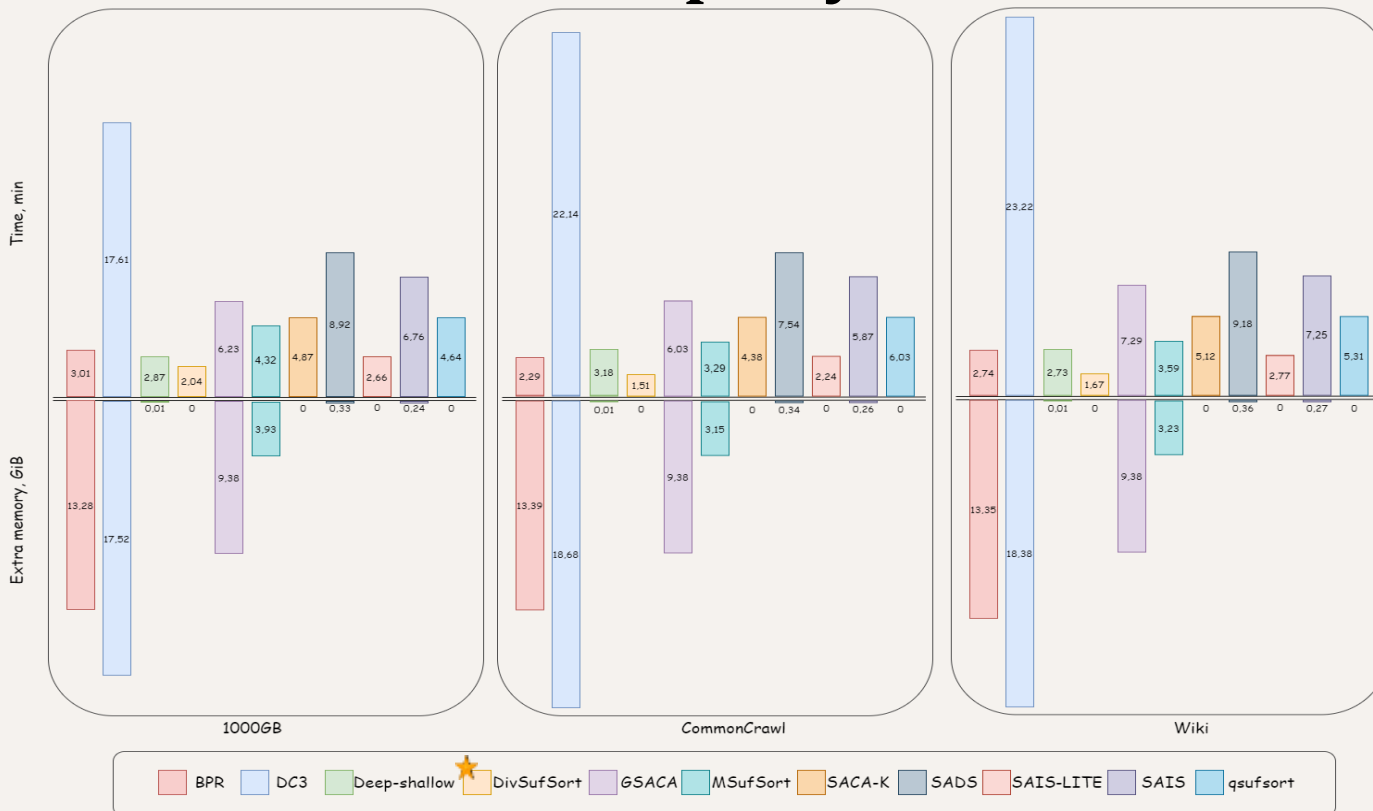
Для суфіксних масивів існує *узагальнений* варіант реалізації, що призначений для зберігання та обробки колекції рядків.

Основна проблема використання суфіксних масивів для пошуку полягає в побудові самої структури, а зокрема – в сортуванні суфіксів.

Найвні алгоритми



Показники продуктивності



DivSufSort

01

$O(n \log n)$

Часова складність для
побудови

02

$5n + O(1)$

Просторова
складність

03

$O(m \log n)$

Часова складність для
пошуку

04

С

Мова оригінального
алгоритму

05

Не підтримує
роботу з узагальненою
структурою

06

Найшвидший
наявний алгоритм

Модифікація DivSufSort

01 C → TypeScript

Перенесення коду алгоритму на TypeScript

02 Узагальнена структура

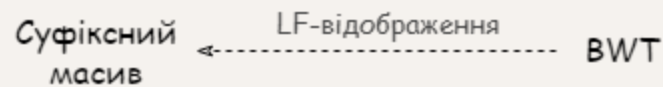
Алгоритм було модифіковано для роботи з колекцією рядків

03 Компресія

До отриманої структури були застосовані алгоритми компресії.

Компресія

Для оптимізації використання ресурсів було застосовано *LF-відображення*, що встановлює зв'язок між впорядкованими суфіксами та рядком, утвореним внаслідок *перетворення Берроуза-Вілера*.



banana\$	→	\$banana
anana\$b	→	a\$banan
nana\$ba	→	ana\$ban
ana\$ban	→	anana\$b
na\$bana	→	banana\$
a\$banan	→	na\$bana
\$banana	→	nana\$b

— Сортуємо —→

Перетворення Берроуза-Вілера: створюємо n ротацій рядка S , де n – довжина рядка S , впорядковуємо їх в лексикографічному порядку та складаємо останні літери кожної ротації в новий рядок.

Компресія

01 F-стовчик
Представлений у вигляді кількості входжень кожного символу алфавіту у тексти.

02 L-стовчик
Рядок, до якого було застосовано *BWT*.

03 Збереження частин масивів
Зберігаємо частини масивів рангів та суфіксів, де решту елементів можна відновити за $O(1)$ часу.

F	L	ranks			SA
		a	b	n	
1	a	1	0	0	6
3	n	-	-	-	-
1	n	-	-	-	3
2	b	1	1	2	-
	\$	-	-	-	0
	a	-	-	-	-
	a	3	1	2	-

Результати

	Назва	M_{orig} , Мб	$ \Sigma $	T, с	M, Мб	Опис
1	sudoku	164	10	75	986	Один мільйон партій в судоку
2	foodprices	87	42	52	407	Ціни на продукти харчування у різних країнах
3	sales	50	39	37	453	Інформація про продажі, близько 740 тис. записів

Висновки

Дослідження

01 було досліджено різні підходи до реалізації швидкого пошуку, а зокрема використання суфіксних дерев та їх аналогу – суфіксних масивів.

02 Модифікація

Перенесення на TypeScript, робота з колекцією рядків, компресія. Хоч алгоритм і поступається в швидкості, проте в результаті було отримано відносно компактну структуру.

03 Розмір даних

Визначення максимальної допустимої межі розміру даних