

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ
ДЕРЖАВНЕ НЕКОМЕРЦІЙНЕ ПІДПРИЄМСТВО
ДЕРЖАВНИЙ УНІВЕРСИТЕТ «КИЇВСЬКИЙ АВІАЦІЙНИЙ ІНСТИТУТ»
ФАКУЛЬТЕТ КОМП'ЮТЕРНИХ НАУК ТА ТЕХНОЛОГІЙ
КАФЕДРА КОМП'ЮТЕРНИХ ІНФОРМАЦІЙНИХ ТЕХНОЛОГІЙ

ДОПУСТИТИ ДО ЗАХИСТУ
Завідувач випускової кафедри
_____ Аліна САВЧЕНКО
«___» _____ 2024 р.

КВАЛІФІКАЦІЙНА РОБОТА

(ПОЯСНЮВАЛЬНА ЗАПИСКА)

ВИПУСКНИКА ОСВІТНЬОГО СТУПЕНЯ МАГІСТРА
ЗА ОСВІТНЬО-ПРОФЕСІЙНОЮ ПРОГРАМОЮ
«ІНФОРМАЦІЙНІ ТЕХНОЛОГІЇ ПРОЕКТУВАННЯ»

Тема: «Метод оптимізації алгоритмів сортування для великих наборів даних»

Виконавець:	Богдан ГРИЦАК
Керівник:	к.т.н., доцент Вікторія СИДОРЕНКО
Нормоконтролер:	к.т.н., доцент Олена ТОЛСТІКОВА

ДЕРЖАВНЕ НЕКОМЕРЦІЙНЕ ПІДПРИЄМСТВО
ДЕРЖАВНИЙ УНІВЕРСИТЕТ «КИЇВСЬКИЙ АВІАЦІЙНИЙ ІНСТИТУТ»
Факультет комп'ютерних наук та технологій
Кафедра комп'ютерних інформаційних технологій
Спеціальність 122 «Комп'ютерні науки»
Освітньо-професійна програма «Інформаційні технології проектування»

ЗАТВЕРДЖУЮ

Завідувач кафедри КІТ

Аліна САВЧЕНКО

(підпис)

«___» _____ 2024 р.

ЗАВДАННЯ

на виконання кваліфікаційної роботи

Грицак Богдан Ярославович

(ПІБ випускника)

1. Тема кваліфікаційної роботи: «Метод оптимізації алгоритмів сортування для великих наборів даних» затверджена наказом ректора № 1782/ст від 06.09.2024р.
2. Термін виконання роботи: з 26 серпня 2024 року по 03 грудня 2024 року.
3. Вихідні дані до роботи: Алгоритм методу оптимізації сортування для великих наборів даних
4. Зміст пояснювальної записки: 1. Огляд предметної області. 2. Оптимізація алгоритмів сортування для великих наборів даних. 3. Реалізація оптимізації алгоритму.
5. Перелік обов'язкового ілюстративного матеріалу: 1. Приклади методів. 2. Скріншоти реалізації. 3. Скріншоти функцій бібліотек.

6. Календарний план-графік

№ з/п	Завдання	Термін виконання	Підпис керівника
1	Огляд та аналіз предметної області. Написання 1 розділу, представлення керівнику.	26.08.2024 - 10.09.2024	
2	Вибір та опис використаних технологій. Написання 2 розділу, представлення керівнику.	11.09.2024 - 29.09.2024	
3	Написання 3 розділу, представлення керівнику.	03.10.2024 – 17.11.2024	
4	Загальне редагування та друк пояснювальної записки.	18.11.2024 – 21.11.2024	
5	Проходження нормоконтролю, перепліт пояснювальної записки.	22.11.2024- 27.11.2024	
6	Розробка тексту доповіді. Оформлення графічного матеріалу для презентації	28.11.2024- 03.12.2024	

7. Дата видачі завдання 26.08.2024р.

Керівник кваліфікаційної роботи _____ Вікторія СИДОРЕНКО
(підпис керівника)

Завдання прийняв до виконання _____ Богдан ГРИЦАК
(підпис випускника)

РЕФЕРАТ

Пояснювальна записка до кваліфікаційної роботи на тему: «Метод оптимізації алгоритмів сортування для великих наборів даних» містить: 94 сторінки, 35 рисунків, 22 інформаційних джерел.

Об'єкт дослідження – процес сортування великих наборів даних.

Предмет дослідження – це методи, та засоби оптимізації алгоритмів сортування, які спрямовані на підвищення їх продуктивності та зменшення використання ресурсів під час роботи з великими обсягами даних

Мета роботи – розробка та аналіз ефективного методу оптимізації алгоритмів сортування, який забезпечить високу продуктивність та оптимальне використання ресурсів при сортуванні великих наборів даних

Методи дослідження – логічний, алгоритмічний аналіз, порівняльний, аналіз інформаційних джерел, моделювання та симуляція.

PYTHON, АЛГОРИТМ, СОРТУВАННЯ, MAPREDUCE,
ОПТИМІЗАЦІЯ, HADOOP, SORT, MERGE, ДАНІ, КЕШ, CHECKPOINT

ЗМІСТ

ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ, СКОРОЧЕНЬ, ТЕРМІНІВ	6
ВСТУП.....	7
РОЗДІЛ 1 ОГЛЯД ПРЕДМЕТНОЇ ОБЛАСТІ.....	9
1.1. Огляд алгоритмів сортування	9
1.2. Класифікація алгоритмів за складністю	15
1.3. Аналіз теоретичної складності	23
1.4. Аналіз продуктивності алгоритмів сортування	30
1.5. Переваги та недоліки основних алгоритмів сортування	35
1.6. Модифікації класичних алгоритмів сортування	37
1.7. Висновки до розділу 1	40
РОЗДІЛ 2 Оптимізація алгоритмів сортування для великих наборів даних...	42
2.1. Аналіз проблем сортування великих наборів даних	42
2.2. Підходи до оптимізації.....	46
2.3. Ефективне використання кеш-пам'яті	50
2.4. Зовнішні алгоритми сортування	53
2.5. Оптимізація розподілу ресурсів.....	57
2.6. Оцінка продуктивності та експерименти	60
2.7. Висновки до розділу 2	63
РОЗДІЛ 3 РЕАЛІЗАЦІЯ ОПТИМІЗАЦІЇ АЛГОРИТМУ	65
3.1. Огляд алгоритму для оптимізації	65
3.2. Основні концепції MapReduce	66
3.3. Основні недоліки Map Reduce та можливі вирішення	68
3.4. Висновки до розділу 3	91
ВИСНОВКИ	93
СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ	95

ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ, СКОРОЧЕНЬ, ТЕРМІНІВ

БД	–	База даних
DAG	–	Directed Acyclic Graph
GPU	–	Graphics processing unit
I/O	–	Input / Output
SSD	–	Solid-state drive

ВСТУП

Сучасна епоха характеризується стрімким зростанням обсягів даних, які потребують обробки та аналізу. Зокрема, у сфері обчислювальної техніки, великих даних та машинного навчання, сортування є фундаментальною задачею. Традиційні алгоритми сортування, такі як сортування бульбашкою чи швидке сортування, не завжди здатні ефективно справлятися з масивами даних, які перевищують доступні обчислювальні ресурси. Це створює потребу у дослідженні нових підходів до оптимізації алгоритмів, здатних обробляти великі набори даних швидко, точно та з мінімальними витратами ресурсів.

Актуальність даної роботи зумовлена збільшенням ролі великих даних у сучасному суспільстві та бізнесі. Застосування великих даних у науці, економіці та технологіях вимагає ефективних способів обробки інформації. Проте традиційні методи часто не враховують обмеження пам'яті, значні витрати на введення/виведення даних та складність паралельних обчислень. Вивчення методів оптимізації алгоритмів сортування, таких як паралельні обчислення або використання інструментів MapReduce, дозволяє вирішити ці проблеми, надаючи нові підходи до обробки великих обсягів інформації.

Метою роботи є розробка та аналіз методу оптимізації алгоритмів сортування для великих наборів даних, який забезпечить високу продуктивність, оптимальне використання обчислювальних ресурсів та мінімізацію часу обробки. Основну увагу зосереджено на адаптації алгоритмів до середовищ, що використовують розподілені обчислення, а також на впровадженні практичних рішень для підвищення їхньої ефективності.

Об'єктом дослідження виступає процес сортування великих наборів даних, а предметом — методи і засоби оптимізації алгоритмів сортування, що спрямовані на підвищення їхньої продуктивності.

Для досягнення мети використано такі методи дослідження: логічний і алгоритмічний аналіз, порівняння різних алгоритмів сортування,

моделювання процесів оптимізації та симуляція роботи алгоритмів у середовищах великих даних.

У роботі вперше запропоновано метод оптимізації сортування великих даних на основі адаптації традиційних алгоритмів до розподілених обчислень із використанням MapReduce.

Запропоновані рішення допоможуть зменшити вплив I/O-операцій, оптимізувати процес обробки в паралельному середовищі та покращити балансування навантаження між вузлами кластеру, що має практичне значення для обробки великих даних у різних галузях — від аналізу фінансових потоків до розробки систем штучного інтелекту. Реалізація методу оптимізації алгоритмів сортування сприятиме зниженню витрат на обчислювальні ресурси та підвищенню швидкості обробки даних.

РОЗДІЛ 1

ОГЛЯД ПРЕДМЕТНОЇ ОБЛАСТІ

1.1. Огляд алгоритмів сортування

Основні поняття сортування.

Сортування – це один з найважливіших процесів в інформатиці, який полягає у впорядкуванні набору даних за певним критерієм, зазвичай за зростанням або спаданням. Цей процес фундаментальний у багатьох задачах і алгоритмах, оскільки впорядковані дані значно спрощують подальші операції з ними. Відсортовані дані полегшують пошук інформації, дозволяють оптимізувати процеси злиття різних наборів, прискорюють доступ до конкретних елементів та є основою для багатьох інших алгоритмів, зокрема для пошуку, обробки і аналізу великих обсягів інформації [1].

Сортування впливає не тільки на ефективність роботи з даними, але й на загальну продуктивність систем, оскільки багато задач і програм спираються на упорядковані структури. Наприклад, бази даних використовують сортування для індексації, що прискорює виконання запитів; пошукові системи застосовують сортування для ранжування результатів; системи управління файлами сортують дані для полегшення навігації користувачем. Це пояснює, чому розуміння сортування і вибір оптимального алгоритму є критично важливим для будь-якого розробника.

Упорядкування даних може базуватися на різних критеріях: числа можуть сортуватися за значенням, рядки – за алфавітом, а об'єкти – за будь-якими властивостями, залежно від потреб конкретної задачі. Основний критерій сортування визначається правилами порівняння між елементами. Наприклад, при сортуванні чисел ми порівнюємо їх між собою на основі

Кафедра КІТ				ДНП ДУ КАІ 24 16 03 000 ПЗ						
	ПІБ			РОЗДІЛ 1. ОГЛЯД ПРЕДМЕТНОЇ ОБЛАСТІ			Літ.	Аркуш	Аркушів	
Розроб.	Грицак Б. Я.								9	33
Керівник	Сидоренко В. М.						М-122-23-1-ТП			
Н. Контр.	Толстікова О.В.									

величини: більше або менше. При сортуванні рядків критерій може бути заснований на лексикографічному порядку, де символи порівнюються за їх позицією в алфавіті.

Сортування можна застосовувати до найрізноманітніших типів даних: чисел, символів, рядків, складних структур, як-от об'єктів, списків, масивів, файлів тощо. Це робить сортування надзвичайно універсальною задачею, яка має величезне значення в різних галузях, починаючи від наукових розрахунків і закінчуючи електронною комерцією. Наприклад, онлайн-магазини часто застосовують сортування для впорядкування товарів за ціною, популярністю, датою додавання. Ця функція не тільки полегшує користувачам пошук потрібних товарів, але й підвищує ефективність роботи системи.

Важливим аспектом сортування є його ефективність. Вона вимірюється не тільки швидкістю виконання, але й використанням ресурсів системи, таких як оперативна пам'ять і процесорний час. В сучасних умовах, коли обсяги даних можуть бути колосальними, ефективність алгоритмів сортування стає критично важливою. Неправильно обраний або неоптимізований алгоритм може призвести до значного зниження продуктивності системи, особливо на великих масивах даних, що містять мільйони або навіть мільярди елементів. У зв'язку з цим, сортування є не лише базовою задачею програмування, але й дослідницькою темою, де пошук нових методів оптимізації та поліпшення продуктивності залишається актуальним.

Основна складність сортування полягає в тому, що для його виконання зазвичай потрібно порівнювати кожен елемент набору даних з іншими, що може вимагати значної кількості операцій. У найпростішому випадку, коли елементів мало, цей процес може бути виконаний досить швидко навіть за допомогою найпростіших алгоритмів. Проте, коли кількість елементів збільшується, кількість необхідних порівнянь і обмінів зростає в геометричній прогресії, що робить сортування дуже затратним за часом і ресурсами.

Для вимірювання ефективності сортування використовують два ключові параметри: часову складність і просторову складність. Часова складність

визначає, скільки часу потрібно алгоритму для завершення роботи в залежності від кількості елементів, що сортуються. Просторова складність визначає, скільки додаткової пам'яті потрібно алгоритму для виконання своїх операцій. Наприклад, деякі алгоритми виконують сортування "на місці", тобто вони не потребують додаткової пам'яті для зберігання проміжних результатів, в той час як інші алгоритми можуть потребувати значних ресурсів пам'яті для підтримки роботи.

Інший важливий аспект сортування – це стабільність. Стабільний алгоритм сортування гарантує, що порядок однакових елементів у вихідному наборі буде збережений після сортування. Це є важливим у тих випадках, коли сортування виконується за кількома критеріями. Наприклад, при сортуванні списку студентів спочатку за прізвищем, а потім за іменем, стабільний алгоритм дозволяє зберегти початковий порядок тих студентів, у яких прізвища однакові. Це може бути важливо для задач, де поряд з основним критерієм є додаткові властивості, які мають значення для користувача.

Сортування також може бути адаптивним, тобто таким, що враховує вже існуючий порядок у наборі даних і оптимізує свою роботу в залежності від цього. Наприклад, якщо дані вже частково відсортовані, адаптивний алгоритм може виконати сортування значно швидше, ніж якби він сортував зовсім неупорядкований набір. Це дуже корисно у випадках, коли дані оновлюються або додаються поступово, і немає потреби сортувати їх повністю кожен раз.

Ефективність і точність сортування можуть значно вплинути на роботу системи загалом. Від того, як швидко й ефективно буде виконано сортування, може залежати швидкість виконання багатьох інших алгоритмів, оскільки сортування зазвичай є попередньою операцією перед іншими важливими задачами, як-от пошук, фільтрація або злиття даних.

Види алгоритмів сортування

Алгоритми сортування є однією з найважливіших тем в інформатиці, оскільки вони вирішують задачу впорядкування даних, тобто процесу, який надзвичайно важливий для оптимізації багатьох інших задач, таких як пошук,

фільтрація, злиття даних та їх аналіз. Існує безліч алгоритмів сортування, кожен з яких має свої унікальні характеристики та застосовується в залежності від конкретних вимог, розміру набору даних, обмежень по часу і ресурсах.

Розглянемо детальніше основні види алгоритмів сортування, зосереджуючи увагу на їх характеристиках, сильних та слабких сторонах, а також на те, в яких ситуаціях кожен з них може бути найбільш ефективним.

Один з найпростіших і водночас найвідоміших алгоритмів сортування – це сортування бульбашкою. Його принцип заснований на поступовому "виштовхуванні" найбільших (або найменших) елементів до кінця списку через послідовні порівняння та обміни сусідніх елементів. Якщо поточний елемент більше за наступний, вони міняються місцями. Такий процес повторюється, поки всі елементи не будуть на своїх місцях. Сортування бульбашкою легко зрозуміти і реалізувати, проте його ефективність на великих наборах даних залишає бажати кращого: у найгіршому випадку алгоритм має часову складність $O(n^2)$, що робить його надзвичайно повільним при великій кількості елементів. Однак для невеликих або частково відсортованих масивів, а також в навчальних цілях, цей алгоритм залишається популярним через свою простоту.

Дуже схожим за ідеєю є сортування вибором, яке також відносно просте, але менш ефективне на великих наборах даних. Принцип цього алгоритму полягає в тому, що на кожному кроці вибирається мінімальний (або максимальний) елемент з неупорядкованої частини масиву і ставиться на своє місце. Це дозволяє поступово "збирати" відсортований масив зліва направо, обираючи найменший елемент серед залишкових. Як і сортування бульбашкою, цей алгоритм має складність $O(n^2)$, що робить його непридатним для обробки великих наборів даних, проте його використовують в тих випадках, коли пам'ять є критичним ресурсом, оскільки сортування вибором є алгоритмом сортування на місці і не потребує додаткової пам'яті для проміжних результатів.

Наступним простим, але трохи ефективнішим варіантом є сортування вставками. Цей алгоритм працює за аналогією з тим, як ми сортуємо карти в руках: на кожному кроці алгоритм бере один елемент з невпорядкованої частини і вставляє його на відповідне місце у впорядкованій частині масиву. Хоча його найгірша складність також $O(n^2)$, він набагато ефективніший на невеликих наборах даних або на вже частково відсортованих масивах, оскільки може виконати менше операцій порівняння і обміну. Сортування вставками є адаптивним алгоритмом, що означає, що його швидкість значно підвищується, якщо масив вже частково відсортований. Це робить його корисним для випадків, коли необхідно обробляти невеликі або майже відсортовані масиви даних.

Одним з найвідоміших і найефективніших алгоритмів є швидке сортування (quicksort). Цей алгоритм використовує підхід розділяй і володарюй, при якому масив поділяється на дві частини на основі вибору опорного елемента. Елементи, менші за опорний, переміщуються в ліву частину, а ті, що більші – в праву. Потім цей процес повторюється рекурсивно для кожної частини, доки масив не буде повністю відсортований. Швидке сортування є одним з найбільш ефективних алгоритмів для великих наборів даних, оскільки в середньому його складність становить $O(n \log n)$. Однак, у найгіршому випадку (коли опорний елемент вибирається невдало), складність може вирости до $O(n^2)$. Незважаючи на це, швидке сортування залишається найпопулярнішим алгоритмом завдяки своїй високій ефективності на практиці. Алгоритм легко модифікувати для вирішення конкретних завдань, що робить його надзвичайно гнучким. До того ж він є сортуванням на місці, що дозволяє економити пам'ять.

Як і швидке сортування, сортування злиттям використовує принцип розділяй і володарюй. Спочатку масив ділиться на дві частини, кожна з яких рекурсивно сортується, після чого два відсортовані підмасиви зливаються в один. Сортування злиттям є стабільним алгоритмом і має гарантовану складність $O(n \log n)$, навіть у найгіршому випадку. Однак його недоліком є

те, що він потребує додаткової пам'яті для зберігання проміжних результатів, що робить його менш підходящим для ситуацій, коли оперативна пам'ять обмежена. Це робить сортування злиттям чудовим вибором для обробки великих масивів даних, особливо тоді, коли стабільність є важливим фактором.

Пірамідальне сортування (heapsort) – це ще один алгоритм сортування, який має складність $O(n \log n)$ і працює на основі структури даних, відомої як піраміда або купа. Спочатку масив перетворюється в максимальну купу, де кожен батьківський елемент більше за своїх нащадків. Після цього найбільший елемент (корінь купи) переміщується в кінець масиву, а процес перебудови купи повторюється для залишкових елементів. Пірамідальне сортування є сортуванням на місці, що означає відсутність потреби в додатковій пам'яті, і воно завжди має гарантовану часову складність $O(n \log n)$, навіть у найгіршому випадку. Це робить його гарним вибором для випадків, коли потрібен надійний і ефективний алгоритм, здатний працювати в обмежених умовах пам'яті.

Окрім порівняльних алгоритмів сортування, існують так звані непорівняльні алгоритми, які не покладаються на пряме порівняння елементів між собою. Одним з таких є лічильне сортування (counting sort), яке працює шляхом підрахунку кількості появ кожного унікального значення в наборі даних. Цей алгоритм ефективний для сортування цілих чисел у невеликому діапазоні і має складність $O(n)$, що робить його надзвичайно швидким на практиці. Проте, його використання обмежене специфічними умовами: він потребує значної кількості додаткової пам'яті, якщо діапазон можливих значень великий, і не підходить для сортування дійсних чисел або рядків.

Ще одним незрівняним алгоритмом є сортування за розрядами (radix sort), яке сортує числа за окремими цифрами або розрядами. Воно спочатку сортує числа за найменш значущим розрядом, потім за наступним, і т. д., доки всі розряди не будуть оброблені. Це дозволяє алгоритму досягати складності $O(n)$, що робить його ефективним для великих наборів даних, особливо у

випадках, коли числа мають обмежену кількість розрядів. Як і лічильне сортування, цей алгоритм застосовується лише в специфічних випадках і потребує значної кількості пам.

1.2. Класифікація алгоритмів за складністю

Алгоритми з $O(n^2)$

Алгоритми сортування з часовою складністю $O(n^2)$ включають в себе одні з найпростіших і найзрозуміліших алгоритмів, таких як сортування бульбашкою, вибором і вставками. Ці алгоритми, хоча й мають невисоку продуктивність на великих масивах даних, часто застосовуються для навчальних цілей або в задачах, де кількість елементів є відносно невеликою. Незважаючи на однакову асимптотичну складність, кожен з цих алгоритмів має свої унікальні характеристики, які роблять їх корисними в певних ситуаціях [2].

Сортування бульбашкою – один з найпростіших і найменш ефективних алгоритмів, оскільки кожен елемент масиву послідовно порівнюється зі своїм наступним сусідом, і якщо вони стоять у неправильному порядку, то міняються місцями. Цей процес повторюється багаторазово, доки весь масив не буде впорядкований. Основною перевагою сортування бульбашкою є його простота: його дуже легко зрозуміти і реалізувати, і він не потребує додаткової пам'яті, оскільки працює на місці. Однак його недоліком є те, що кількість операцій порівняння і обміну зростає квадратично з кількістю елементів. Для кожної ітерації алгоритм проходить по всьому масиву, порівнюючи сусідні елементи, тому при великій кількості елементів цей процес стає дуже повільним. У найгіршому випадку, коли масив відсортований у зворотному порядку, алгоритм виконує близько n^2 операцій, що робить його непридатним для сортування великих наборів даних. Втім, сортування бульбашкою може бути ефективним на невеликих масивах або коли масив вже майже відсортований.

Сортування вибором діє за іншим принципом, хоча його складність також $O(n^2)$. Алгоритм починається з пошуку найменшого (або найбільшого) елемента у невідсортованій частині масиву і його переміщення на початок. Потім алгоритм повторює цей процес для решти елементів масиву, зменшуючи область пошуку після кожної ітерації. У підсумку, після $n-1$ кроку весь масив стає відсортованим. Основною перевагою сортування вибором є те, що воно виконує мінімальну кількість обмінів елементів, оскільки кожен елемент переміщується лише один раз на своє остаточне місце. Однак недоліком є те, що кількість порівнянь між елементами залишається квадратичною, навіть якщо масив майже відсортований, що робить його повільним для великих наборів даних. Сортування вибором є корисним у ситуаціях, коли кількість обмінів є критичною, а також у випадках, коли необхідно економити пам'ять, оскільки цей алгоритм не вимагає додаткового простору для зберігання проміжних результатів.

Сортування вставками дещо більш оптимізоване у порівнянні з попередніми алгоритмами. Воно працює шляхом побудови відсортованої частини масиву шляхом послідовного вставлення елементів з невідсортованої частини на відповідні місця у впорядкованій частині. Алгоритм проходить через кожен елемент масиву, порівнюючи його з елементами в уже відсортованій частині, і вставляє його на відповідну позицію. Як і в попередніх алгоритмах, у найгіршому випадку його складність становить $O(n^2)$, проте сортування вставками є значно швидшим на частково відсортованих або невеликих масивах, оскільки воно є адаптивним і виконує менше порівнянь у таких ситуаціях. Це робить його ефективним для впорядкування невеликих або майже відсортованих наборів даних, де кількість переміщень і порівнянь буде мінімальною.

Сортування вставками також має важливу властивість стабільності, що означає, що однакові елементи зберігають свій початковий порядок після сортування. Це може бути важливим у випадках, коли дані мають кілька критеріїв сортування, і другорядний порядок має значення. Наприклад, якщо

необхідно спочатку відсортувати список за одним критерієм (наприклад, за датою), а потім за іншим (наприклад, за іменем), стабільний алгоритм сортування забезпечить правильне впорядкування даних без порушення попередньої сортування.

Незважаючи на те, що всі ці алгоритми мають однакову асимптотичну складність $O(n^2)$, кожен з них підходить для різних типів задач. Алгоритми сортування бульбашкою, вибором і вставками добре підходять для невеликих наборів даних або для випадків, коли простота реалізації є ключовим фактором. Вони також можуть бути корисними для навчання, оскільки демонструють базові принципи сортування і дозволяють зрозуміти, як можна впорядковувати елементи через порівняння і обміни.

Однак ці алгоритми не рекомендуються для великих масивів даних через їх низьку ефективність. У реальних додатках, де обробка великих наборів даних є критично важливою, більш складні і ефективні алгоритми, такі як швидке сортування або сортування злиттям, мають набагато кращу продуктивність. Попри це знання і розуміння алгоритмів з $O(n^2)$ все ще є необхідними, оскільки вони можуть бути оптимізовані для певних випадків або використані в поєднанні з іншими методами для вирішення специфічних завдань.

Алгоритми з $O(n \log n)$

Алгоритми сортування з часовою складністю $O(n \log n)$ є найефективнішими для великих наборів даних і широко застосовуються в сучасних системах. До таких алгоритмів належать швидке сортування, сортування злиттям і пірамідальне сортування. Кожен з цих алгоритмів використовує різні підходи для оптимізації процесу впорядкування, що робить їх придатними для широкого спектра задач. Основною перевагою цих алгоритмів є їх продуктивність: навіть при великій кількості елементів вони забезпечують високий рівень ефективності, особливо в порівнянні з алгоритмами з $O(n^2)$.

Швидке сортування (Quicksort) є одним з найвідоміших і найефективніших алгоритмів для впорядкування масивів. Воно базується на принципі розділяй і володарюй: алгоритм вибирає опорний елемент (півот) і розділяє масив на дві частини – одну, що містить елементи, менші за півот, і іншу, де елементи більші. Після цього кожна частина сортується рекурсивно, доки масив не стане впорядкованим. У середньому швидке сортування має складність $O(n \log n)$, що робить його надзвичайно ефективним для великих наборів даних. Проте, у найгіршому випадку, коли вибір опорного елемента є невдалим (наприклад, якщо масив уже відсортований або майже відсортований), складність алгоритму може зрости до $O(n^2)$. Однак такі ситуації можна уникнути завдяки вдосконаленим технікам вибору півота, таким як вибір середнього з трьох елементів або випадковий вибір.

Однією з головних причин популярності швидкого сортування є те, що воно є сортуванням на місці, тобто не потребує додаткової пам'яті для зберігання проміжних результатів. Це важливо для застосувань, де пам'ять є обмеженим ресурсом. Крім того, швидке сортування легко модифікувати для різних потреб, наприклад, для роботи з великими файлами або потоками даних. Незважаючи на те, що швидке сортування не є стабільним (тобто елементи з однаковими значеннями можуть змінювати свій порядок), його гнучкість і висока продуктивність роблять його одним з найбільш широко використовуваних алгоритмів у різних галузях, від баз даних до обробки текстових файлів.

Сортування злиттям (Merge Sort) – це ще один потужний алгоритм, який також використовує принцип розділяй і володарюй. Проте його підхід відрізняється від швидкого сортування. Спочатку масив ділиться на дві частини рівного розміру, кожна з яких рекурсивно сортується. Після цього дві впорядковані частини зливаються в один відсортований масив. Операція злиття виконується за лінійний час, а розбиття масиву на частини забезпечує загальну складність алгоритму $O(n \log n)$, яка залишається сталою навіть у найгіршому випадку. Це робить сортування злиттям особливо корисним у

ситуаціях, коли необхідно забезпечити стабільну продуктивність незалежно від початкового стану даних.

Сортування злиттям - стабільний алгоритм, що означає, що однакові елементи зберігають свій початковий порядок після сортування. Це є важливою властивістю в багатьох практичних додатках, де стабільність впорядкування має значення, наприклад, коли дані сортуються за кількома критеріями. Однак одним з недоліків сортування злиттям є те, що воно вимагає додаткової пам'яті для зберігання тимчасових масивів під час процесу злиття, що може бути проблемою для обмежених систем. Незважаючи на це, сортування злиттям є дуже ефективним і надійним алгоритмом для великих наборів даних, особливо коли стабільність є критичною вимогою.

Пірамідальне сортування (Heapsort) також забезпечує складність $O(n \log n)$ і має унікальну особливість у порівнянні з іншими алгоритмами. Воно базується на використанні структури даних, відомої як купа або піраміда, яка є частковим деревом, де кожен батьківський елемент більше або дорівнює своїм нащадкам (для максимальних куп) або менше або дорівнює (для мінімальних куп). Алгоритм починається з того, що перетворює невпорядкований масив у максимальну купу, де найбільший елемент розташований на вершині. Після цього цей елемент обмінюється з останнім елементом масиву, виключається з купи, а решта масиву перебудовується, щоб зберегти властивість купи. Цей процес повторюється, доки весь масив не буде впорядкованим.

Пірамідальне сортування - стабільне в тому сенсі, що його складність залишається $O(n \log n)$ незалежно від початкового порядку елементів. Однією з його основних переваг є те, що воно не вимагає додаткової пам'яті, оскільки сортування виконується без створення нових масивів. Це робить його придатним для задач, де пам'ять є обмеженим ресурсом. Пірамідальне сортування також відзначається тим, що має хорошу продуктивність на практиці, оскільки його поведінка залишається сталою навіть у найгіршому випадку, на відміну від швидкого сортування. Проте його реалізація

складніша, ніж у швидкого сортування, і на практиці пірамідальне сортування використовують рідше через дещо повільнішу роботу в середньому.

Ці три алгоритми, хоч і мають однакову асимптотичну складність $O(n \log n)$, суттєво відрізняються за своїми характеристиками та оптимальністю в різних ситуаціях. Швидке сортування є найшвидшим у середньому і використовується найчастіше, але воно може бути нестабільним в найгірших випадках. Сортування злиттям забезпечує стабільну продуктивність і стабільність сортування, але потребує додаткової пам'яті, що може бути недоліком. Пірамідальне сортування, навпаки, не вимагає додаткової пам'яті і завжди гарантує $O(n \log n)$ незалежно від початкового стану даних, хоча й може бути повільнішим у середньому.

Кожен з цих алгоритмів знаходить своє застосування в реальних задачах. Швидке сортування залишається незамінним для сортування в системах з обмеженою пам'яттю, де неважливо, чи є масив стабільним. Сортування злиттям застосовується там, де необхідна стабільність і додаткова пам'ять не є проблемою, наприклад у системах баз даних. Пірамідальне сортування використовується для ситуацій, коли стабільність часу роботи є критичною вимогою, особливо у випадках, коли необхідна висока надійність і обмежений доступ до оперативної пам'яті.

Алгоритми з $O(n)$

Алгоритми сортування з часовою складністю $O(n)$ включають у себе ефективні методи для сортування, коли можна скористатися специфічними особливостями даних, щоб обійти стандартні порівняльні методи, які мають нижню межу $O(n \log n)$. Такі алгоритми, як лічильне сортування, сортування по розрядах і сортування за допомогою комірок (bucket sort), досягають лінійної складності, коли виконуються певні умови, наприклад, обмежений діапазон значень або можливість порівняння часткових елементів. Хоча ці методи не є універсальними, вони забезпечують високу ефективність у задачах, де їх можна застосувати.

Лічильне сортування (Counting Sort) - один з найефективніших алгоритмів, коли йдеться про сортування цілих чисел з обмеженого діапазону. Основна ідея цього алгоритму полягає в тому, щоб підрахувати кількість кожного елемента в масиві і потім використовувати ці підрахунки для побудови відсортованого масиву. Лічильне сортування працює наступним чином: спочатку створюється допоміжний масив, де кожен індекс відповідає можливому значенню елементів в оригінальному масиві, а значення цього індексу зберігає кількість появ цього числа в масиві. Потім, використовуючи накопичені значення цього масиву, відбувається побудова відсортованого результату.

Основною перевагою лічильного сортування є його швидкість – воно сортує масив за час $O(n)$, де n – це кількість елементів, а k – діапазон значень елементів. Проте це сортування підходить лише тоді, коли діапазон значень є відносно малим у порівнянні з кількістю елементів. Наприклад, якщо ми маємо великий масив чисел, але всі вони лежать у діапазоні від 1 до 100, лічильне сортування буде надзвичайно ефективним. Однак якщо діапазон значень є занадто великим (наприклад, від 1 до 10^9), то цей алгоритм стане неефективним, оскільки для кожного значення необхідно буде виділяти місце в допоміжному масиві. Також лічильне сортування вимагає додаткової пам'яті для зберігання цього допоміжного масиву, що може бути проблемою для систем з обмеженими ресурсами.

Окрім цього, лічильне сортування – дуже стабільний алгоритм, тобто порядок елементів з однаковими значеннями зберігається після сортування, що робить його корисним для задач, де важлива стабільність впорядкування. Лічильне сортування також можна адаптувати для сортування символів, дат або інших категорійних даних з обмеженою кількістю можливих значень.

Сортування по розрядах (Radix Sort) – ще один лінійний алгоритм, який часто використовується для сортування цілих чисел або інших об'єктів, які можна розбити на розряди (наприклад, символи в рядках). Принцип роботи сортування по розрядах полягає в тому, щоб сортувати елементи не за їх

повним значенням, а за кожним окремим розрядом (цифрою або символом) починаючи з найменшого розряду і закінчуючи найбільшим. Наприклад, при сортуванні чисел спочатку можна сортувати їх за одиничними цифрами, потім за десятковими, сотенними і т. д.

Для реалізації сортування по розрядах зазвичай використовують стабільний алгоритм сортування, наприклад, лічильне сортування або сортування за допомогою комірок. У результаті кожного кроку елементи впорядковуються за певним розрядом, і на наступному етапі алгоритм повторюється для наступного розряду. Важливо зазначити, що сортування по розрядах працює ефективно тільки тоді, коли кількість розрядів або символів є відносно невеликою, інакше алгоритм може втратити свою лінійну складність.

Основною перевагою сортування по розрядах є те, що воно може ефективно сортувати великі набори даних, особливо тоді, коли числа або символи мають обмежену кількість розрядів. Наприклад, алгоритм часто використовується для сортування чисел у фіксованій системі числення або рядків з однаковою довжиною. У такому випадку його складність залишається лінійною, що робить його привабливим варіантом для багатьох застосувань, де традиційні алгоритми на основі порівняння будуть менш ефективними. Важливо також зазначити, що сортування по розрядах є стабільним, що робить його корисним для задач, де збереження порядку однакових елементів є важливим.

Сортування за допомогою комірок (Bucket Sort) - використовується для сортування даних, що рівномірно розподілені по діапазону. Основна ідея алгоритму полягає в тому, щоб розділити масив на кілька "комірок" (bucket), кожна з яких містить певний діапазон значень. Після цього кожен елемент масиву поміщається в одну з цих комірок відповідно до свого значення. Потім кожна комірка сортується окремо за допомогою іншого алгоритму сортування, наприклад, вставками або швидким сортуванням. На завершальному етапі відсортовані комірки зливаються в один відсортований масив.

Сортування за допомогою комірок працює особливо ефективно, коли елементи рівномірно розподілені по всьому діапазону, оскільки в цьому випадку кожна комірка містить приблизно однакову кількість елементів, і процес сортування всередині кожної комірки буде швидким. У таких умовах загальна складність алгоритму наближається до $O(n)$. Однак, якщо розподіл елементів є нерівномірним (наприклад, більшість елементів потрапляє в одну комірку), ефективність алгоритму може суттєво знизитися.

Сортування за допомогою комірок також може бути корисним в задачах, де відомо, що дані мають певний діапазон значень і їх розподіл є передбачуваним. Алгоритм також можна адаптувати для роботи з різними типами даних, не лише числами. Важливим аспектом сортування за допомогою комірок є правильний вибір кількості та діапазонів для комірок, оскільки від цього залежить ефективність алгоритму. Якщо комірки вибрані невдало, алгоритм може перетворитися на звичайне сортування з $O(n^2)$, оскільки всі елементи потраплять в одну або кілька комірок.

Отже, лінійні алгоритми сортування, такі як лічильне сортування, сортування по розрядах і сортування за допомогою комірок, є надзвичайно ефективними в певних умовах. Вони використовують нестандартні підходи для впорядкування даних, що дозволяє їм досягати високої продуктивності при сортуванні великих масивів. Однак важливо враховувати специфіку даних, з якими працює алгоритм, оскільки лінійна складність може бути досягнута тільки за певних умов, таких як обмежений діапазон значень або рівномірний розподіл даних.

1.3. Аналіз теоретичної складності

Тимчасова складність

Тимчасова складність дозволяє оцінити, як швидко алгоритм буде працювати в залежності від обсягу вхідних даних. У комп'ютерних науках для оцінки тимчасової складності часто використовується нотація "О-велике" (Big-O), яка описує асимптотичну поведінку алгоритму. Це означає, що ми

фокусуємося на тому, як змінюється час виконання алгоритму, коли розмір вхідного масиву зростає [3].

Тимчасова складність базується на кількості базових операцій, які алгоритм виконує у процесі своєї роботи. Цими базовими операціями можуть бути, наприклад, порівняння елементів, обміни значень чи інші прості дії, які мають фіксовану вартість. Нотація Big-O дозволяє ігнорувати константи та малозначущі доданки, що дає можливість зосередитися на найважливіших аспектах, зокрема на тому, як алгоритм буде масштабуватися при збільшенні кількості даних.

Розглядаючи різні типи тимчасової складності, можна виділити кілька основних категорій. Наприклад, $O(1)$ — це найкращий можливий варіант, коли час виконання алгоритму не залежить від обсягу вхідних даних. Однак на практиці таку складність можна спостерігати лише для простих операцій, таких як доступ до конкретного елемента масиву. Жоден реальний алгоритм сортування не може похвалитися загальною складністю $O(1)$.

Ще однією категорією є $O(\log n)$. Ця складність притаманна алгоритмам, які працюють за принципом поділу даних на частини, де кожна частина обробляється окремо. Хоча така складність не є кінцевою для алгоритмів сортування, вона може бути частиною більш складних оцінок, таких як у швидкому сортуванні або пірамідальному сортуванні, де масив ділиться на менші частини, що й зумовлює логарифмічну складність.

Далі, складність $O(n)$ означає, що час виконання алгоритму пропорційний кількості елементів. Це свідчить про те, що кожен елемент повинен бути оброблений один раз. Алгоритми сортування з лінійною складністю можливі, але тільки в специфічних умовах, коли можна уникнути порівнянь. Лінійне сортування, сортування по розрядах і сортування за допомогою комірок досягають лінійної складності, якщо виконуються певні умови, наприклад, коли є обмежений діапазон значень.

Складність $O(n \log n)$ - найбільш ефективна для порівняльних алгоритмів сортування. Вона характерна для таких алгоритмів, як швидке

сортування, сортування злиттям та пірамідальне сортування. У цих алгоритмах масив спочатку ділиться на менші частини, а потім кожна частина сортується. Така комбінація обробки даних забезпечує хорошу продуктивність для більшості реальних наборів даних.

Тимчасова складність $O(n^2)$ є квадратичною і вважається менш ефективною для великих наборів даних, адже час виконання алгоритму зростає пропорційно квадрату кількості елементів. Алгоритми з такою складністю зазвичай виконують багато порівнянь або обмінів, що робить їх непридатними для масштабних задач. Сортування бульбашкою, сортування вибором і сортування вставками є прикладами алгоритмів з $O(n^2)$. Хоча вони можуть бути корисними для малих наборів даних, їх продуктивність залишається обмеженою.

Аналіз тимчасової складності також включає в себе різні сценарії роботи алгоритму. Найгірший випадок характеризує ситуацію, коли вхідні дані є найбільш несприятливими. Наприклад, у швидкому сортуванні найгірший випадок трапляється тоді, коли кожен вибір опорного елемента є невдалим. У такій ситуації часові витрати можуть досягти $O(n^2)$. Найгірший випадок є важливим, оскільки він допомагає оцінити надійність алгоритму в критичних системах, де невдалий сценарій може суттєво вплинути на продуктивність.

Середній випадок враховує всі можливі вхідні дані та обчислює середній час роботи алгоритму. У швидкому сортуванні, наприклад, середній випадок відповідає $O(n \log n)$, оскільки у більшості ситуацій вибір опорного елемента є достатньо вдалим для рівномірного поділу масиву. Найкращий випадок є оптимістичним сценарієм, коли вхідні дані вже частково або повністю відсортовані. У цьому випадку для сортування вставками найкращий випадок настає тоді, коли масив вже відсортований, що дозволяє алгоритму працювати за лінійний час $O(n)$.

Отже, аналіз тимчасової складності - це важливий критерій для вибору алгоритмів сортування. Він дозволяє передбачити, як алгоритм працюватиме на великих даних, забезпечує розуміння його переваг та недоліків у різних

сценаріях, а також допомагає прийняти обґрунтоване рішення для конкретного завдання з урахуванням характеристик вхідних даних і вимог до системи. Розуміння тимчасової складності алгоритму стає ключовим для розробників, які прагнуть до оптимізації продуктивності програмного забезпечення, адже воно дозволяє ефективно вирішувати задачі, з якими вони стикаються.

Просторова складність

Просторова складність визначає, скільки пам'яті або обсягу пам'яті потребує алгоритм у процесі виконання, в залежності від розміру вхідних даних. Аналіз просторової складності є критично важливим, особливо при роботі з великими наборами даних, оскільки обмеження пам'яті можуть суттєво впливати на продуктивність програми та її здатність обробляти інформацію.

При оцінці просторової складності використовують аналогічну нотацію "О-велике" (Big-O), яка описує, як змінюється обсяг пам'яті, необхідний для виконання алгоритму, у залежності від розміру вхідних даних. Це дозволяє оцінити, наскільки ефективно алгоритм використовує доступну пам'ять.

Просторова складність зазвичай складається з двох основних компонентів:

Перша компонента — це постійні витрати пам'яті, які не залежать від розміру вхідних даних. Це можуть бути змінні, які використовуються в алгоритмі, константні значення, а також фіксовані структури даних, наприклад, масиви або списки, які мають визначений розмір. Ці витрати залишаються незмінними незалежно від обсягу даних, з якими працює алгоритм.

Друга компонента — це змінні витрати пам'яті, які залежать від розміру вхідних даних. Це може включати динамічно виділену пам'ять, яка змінюється під час виконання алгоритму в залежності від обсягу даних. Наприклад, якщо алгоритм використовує рекурсію або створює додаткові структури даних для

обробки масиву, обсяг пам'яті, який він потребує, зростатиме пропорційно до розміру вхідних даних.

При аналізі просторової складності важливо також розуміти, що не всі алгоритми потребують однакової кількості пам'яті. Деякі з них можуть працювати в умовах обмеженої пам'яті, у той час як інші можуть вимагати значних обсягів ресурсів. Наприклад, алгоритми, які використовують масиви, можуть потребувати додаткової пам'яті для зберігання допоміжних структур даних. Це, наприклад, стосується алгоритмів сортування злиттям, які створюють нові масиви для зберігання проміжних результатів, що може призвести до збільшення просторової складності до $O(n)$.

У той же час, деякі алгоритми, такі як швидке сортування, можуть працювати в режимі "in-place", тобто вони сортують елементи без створення додаткових масивів, і їхня просторовість залишається $O(\log n)$, оскільки вони використовують пам'ять тільки для зберігання рекурсивних викликів. Це робить швидке сортування дуже ефективним у контексті використання пам'яті.

Просторова складність також включає в себе неявні витрати пам'яті, які можуть виникати внаслідок реалізації алгоритму, наприклад, внаслідок зберігання стекових кадрів під час рекурсивних викликів. У випадку рекурсивних алгоритмів, таких як деякі методи сортування, кожен виклик функції може споживати пам'ять, яка накопичується в стеку. Це особливо важливо враховувати, коли розмір вхідних даних є великим, адже глибина рекурсії може призвести до перевантаження стеку і, в результаті, до аварійного завершення програми.

Загалом, розуміння просторової складності необхідне для ефективного використання пам'яті при розробці програмного забезпечення. Вона допомагає виявити потенційні проблеми з ресурсами, що можуть виникнути внаслідок великих обсягів даних або неефективних алгоритмів. Розробники повинні зважати на те, наскільки ефективно алгоритми використовують пам'ять, щоб забезпечити стабільну і надійну роботу своїх програм.

Отже, просторову складність слід розглядати поряд із тимчасовою складністю для отримання повної картини про ефективність алгоритмів. Знання про те, як алгоритми використовують пам'ять, допоможе розробникам обирати найкращі рішення для конкретних задач, враховуючи обмеження системи і характеристики вхідних даних. Цей аналіз дозволяє створювати оптимізовані та ефективні програми, що працюють з великими наборами даних.

Порівняння складності різних алгоритмів

Порівняння складності різних алгоритмів є ключовим аспектом при виборі найбільш ефективного рішення для певних задач. У випадку алгоритмів сортування важливо оцінювати як їхню тимчасову, так і просторову складність, оскільки ці два аспекти можуть впливати на загальну продуктивність програми. Алгоритми сортування відрізняються за підходами, структурою та результатами, і тому їх варто розглядати в контексті конкретних умов, в яких вони будуть використовуватися.

Розглядаючи алгоритми з квадратичною складністю $O(n^2)$, можна відзначити, що вони підходять для невеликих масивів, але не є ефективними для великих обсягів даних. Наприклад, якщо у вас є масив з декількох десятків або сотень елементів, ці алгоритми можуть продемонструвати задовільні результати. Але коли розмір масиву починає зростати до тисяч чи мільйонів елементів, час виконання швидко зростає, що робить їх непридатними для використання в реальних задачах[4].

На противагу їм, алгоритми з логарифмічною складністю $O(n \log n)$, такі як швидке сортування, сортування злиттям і пірамідалне сортування, демонструють значно кращу продуктивність. Ці алгоритми використовують методи розділення та злиття, що дозволяє їм ефективно працювати навіть з великими масивами. Наприклад, швидке сортування, яке в середньому має складність $O(n \log n)$, може швидко обробляти великі дані, якщо опорний елемент вибирається правильно. Однак у найгіршому випадку, коли вибір

елемента невдалий, його складність може сягати $O(n^2)$, що потрібно враховувати при виборі алгоритму.

Додатково, алгоритми з лінійною складністю $O(n)$ можуть бути надзвичайно ефективними, якщо умови їх застосування задовольняють вимоги, такі як обмежений діапазон значень або специфічні формати даних. Лічильне сортування, сортування по розрядах і сортування за допомогою комірок є прикладами алгоритмів, які можуть виконуватися за $O(n)$, але вимагають певних передумов. Вони не є порівняльними алгоритмами, а це означає, що їхня ефективність вища за порівняльні алгоритми в специфічних умовах.

Крім того, важливо враховувати просторову складність алгоритмів. Алгоритми з $O(n^2)$, як правило, вимагають менше додаткової пам'яті, оскільки зазвичай реалізуються "in-place". У той час як алгоритми злиття потребують $O(n)$ для зберігання проміжних масивів, що може бути критично важливим у системах з обмеженими ресурсами. Швидке сортування, якщо реалізоване в "in-place" режимі, займає $O(\log n)$ пам'яті, що є значною перевагою в контексті використання ресурсів.

При порівнянні складності різних алгоритмів також варто звернути увагу на їхню стабільність. Стабільні алгоритми зберігають порядок елементів з однаковими значеннями після сортування, що є важливим фактором у багатьох додатках. Наприклад, сортування злиттям є стабільним, в той час як швидке сортування зазвичай не є стабільним, хоча його можна модифікувати для досягнення стабільності. У випадках, коли важливий порядок елементів з однаковими значеннями, вибір стабільного алгоритму стає критично важливим.

Отже, при виборі алгоритму сортування варто розглядати не лише його теоретичну складність, але й реальні умови застосування. Різні алгоритми мають свої переваги і недоліки, і їхня ефективність може суттєво змінюватися в залежності від обсягу даних, їх структури та характеристик. Важливо також враховувати обмеження системи, адже алгоритм, який показує хороші

результати на тестових даних, може виявитися неефективним у реальних умовах через обмеження пам'яті чи час виконання.

Розуміння і порівняння складності різних алгоритмів дозволяє розробникам обирати найбільш підходящі рішення для своїх задач, враховуючи як вимоги до продуктивності, так і обмеження ресурсів. Це знання є основою для створення оптимізованих і надійних програмних продуктів, що відповідають сучасним вимогам в обробці даних.

1.4. Аналіз продуктивності алгоритмів сортування

Аналіз продуктивності алгоритмів сортування передбачає розгляд їхньої поведінки в трьох основних сценаріях: середній, найгірший та найкращий випадки. Це дозволяє отримати повне уявлення про ефективність алгоритму в різних умовах і передбачити, як він буде працювати на практиці.

Середній, найгірший і найкращий сценарії

У середньому випадку більшість алгоритмів демонструють свою стандартну поведінку. Це означає, що дані подаються у випадковому порядку, без жодних специфічних властивостей, які могли б сприяти або заважати швидкому сортуванню. Наприклад, швидке сортування в середньому працює з тимчасовою складністю $O(n \log n)$, оскільки в кожному кроці масив ділиться на дві рівні частини, що дає стабільний розподіл навантаження на кожен виклик функції.

У найгіршому випадку продуктивність алгоритму знижується через певні характеристики вхідних даних, що спричиняють найменш ефективне виконання алгоритму. Наприклад, у швидкому сортуванні найгірший випадок виникає, коли опорний елемент кожного разу обирається найменш вдало, наприклад, коли масив вже відсортований або коли всі елементи однакові. У цьому випадку складність підвищується до $O(n^2)$, оскільки кожен розподіл є дуже нерівномірним і фактично перетворюється на лінійний розподіл, що призводить до зростання кількості порівнянь.

У найкращому випадку алгоритм працює з найвищою ефективністю, коли структура вхідних даних ідеально підходить для його виконання. Для швидкого сортування це може бути ситуація, коли масив вже частково відсортований або коли кожен розподіл є рівномірним. У таких випадках продуктивність досягає оптимального значення $O(n \log n)$, оскільки кількість порівнянь та переміщень мінімальна, і кожен поділ масиву виконується максимально рівномірно.

Кількість порівнянь і обмінів

Кількість порівнянь і обмінів є одними з найважливіших показників, що дозволяють оцінити ефективність алгоритмів сортування. Кожен алгоритм виконує певні операції порівняння між елементами та обмінює їх місцями, коли необхідно досягти правильної послідовності. Тому, коли говоримо про кількість порівнянь і обмінів, ми фактично оцінюємо, наскільки добре алгоритм справляється зі своєю роботою при різних обставинах, і як багато операцій потрібно для того, щоб відсортувати масив елементів.

Порівняння — це фундаментальний крок у більшості класичних алгоритмів сортування. Саме за допомогою порівняння алгоритм визначає, чи правильно розташовані два елементи відносно один одного або їх потрібно переставити. Для простих алгоритмів, таких як сортування бульбашкою, вибором чи вставками, кожен елемент порівнюється з кожним іншим на різних етапах, що призводить до великої кількості порівнянь, особливо на великих наборах даних. В таких алгоритмах кількість порівнянь зазвичай відповідає $O(n^2)$, що робить їх малоефективними для великих масивів. Наприклад, у бульбашковому сортуванні кожен елемент послідовно порівнюється з наступним, і якщо він більше, то вони міняються місцями. Такий підхід призводить до того, що кількість порівнянь стрімко зростає з кількістю елементів.

Обміни, або переміщення елементів, відбуваються кожного разу, коли алгоритм виявляє, що два елементи знаходяться не у правильному порядку, і

їх потрібно переставити. У простих алгоритмах кількість обмінів часто порівняна з кількістю порівнянь, оскільки кожен невірно розташований елемент потребує перестановки. У бульбашковому сортуванні, наприклад, кожне порівняння, яке виявляє, що елементи не на своїх місцях, призводить до обміну. Це призводить до великої кількості обмінів, що знижує продуктивність алгоритму, особливо коли масив має значну кількість елементів, і їх необхідно часто переміщати для досягнення відсортованого стану.

Однак існують алгоритми, де кількість обмінів є значно меншою порівняно з кількістю порівнянь. Наприклад, у швидкому сортуванні кількість порівнянь у середньому випадку відповідає $O(n \log n)$, але кількість обмінів залежить від вибору опорного елемента та структури даних. Якщо опорний елемент обирається вдало, обміни відбуваються ефективніше, оскільки масив ділиться на дві приблизно рівні частини, що дозволяє швидко знайти правильні місця для елементів. Проте у найгіршому випадку, коли опорний елемент обирається невдало (наприклад, коли масив вже відсортований або майже відсортований), кількість обмінів може суттєво зрости, і продуктивність алгоритму погіршується до $O(n^2)$.

У сортуванні злиттям кількість порівнянь залишається $O(n \log n)$ у всіх випадках, але цей алгоритм потребує додаткової пам'яті для зберігання тимчасових масивів, що означає, що обмінів як таких може бути менше, оскільки елементи не переміщуються безпосередньо у вихідному масиві, а копіюються у нові масиви. Це робить сортування злиттям більш стабільним щодо кількості обмінів, але воно потребує більше пам'яті для виконання.

Існують також алгоритми, які мінімізують кількість порівнянь та взагалі не виконують обмінів у звичному розумінні, такі як лічильне сортування, сортування по розрядах та сортування за допомогою комірок (bucket sort). Ці алгоритми працюють на основі підрахунку або класифікації елементів, що дозволяє уникнути прямого порівняння кожної пари елементів. Наприклад, у лічильному сортуванні алгоритм підраховує кількість кожного елемента та

розміщує їх на відповідні позиції в кінцевому масиві без порівнянь між елементами. Це дозволяє досягати тимчасової складності $O(n)$, але такий підхід працює лише для певних типів даних (наприклад, цілі числа в обмеженому діапазоні).

Вплив розміру даних на ефективність

Розмір даних безпосередньо впливає на ефективність будь-якого алгоритму сортування, оскільки збільшення обсягу інформації, яку потрібно обробити, веде до зростання кількості операцій, які алгоритм повинен виконати для досягнення кінцевого результату. Цей вплив проявляється як у кількості порівнянь і обмінів, так і в споживанні пам'яті та часу на виконання алгоритму[5].

Коли дані невеликі за розміром, різниця між алгоритмами може бути не такою помітною. Просте сортування, наприклад бульбашкове, може виконуватися досить швидко навіть на невеликих масивах, попри те, що його тимчасова складність — $O(n^2)$. При цьому сортування бульбашкою буде легко реалізувати і не потребує додаткових ресурсів, як-от додаткова пам'ять, що робить його достатнім для малих наборів даних. У таких випадках навіть неефективні алгоритми можуть справлятися з поставленим завданням за прийнятний час.

Однак, коли розмір даних починає рости, ефективність алгоритму стає критично важливою. Алгоритми з квадратичною складністю, як-от сортування вибором або вставками, починають помітно відставати. Це відбувається тому, що їхня кількість операцій зростає експоненціально з розміром масиву. Наприклад, для масиву розміром 100 елементів сортування бульбашкою потребує приблизно 10 000 операцій, а для масиву з 1000 елементів — вже близько 1 000 000 операцій. Зростання обсягу даних суттєво збільшує час, необхідний для виконання сортування, і цей процес стає повільним, малопридатним для реальних додатків, де дані часто мають великий розмір.

На великих наборах даних більше підходять алгоритми з логарифмічною складністю, як-от швидке сортування, сортування злиттям та

пірамідальне сортування. Їхня складність $O(n \log n)$ дозволяє ефективніше обробляти дані, оскільки кількість операцій зростає значно повільніше порівняно з квадратичними алгоритмами. Важливо розуміти, що для великих масивів навіть незначне покращення у складності може привести до величезної різниці в продуктивності. Наприклад, швидке сортування для масиву з 1000 елементів потребує лише близько 10 000 операцій, у той час як сортування бульбашкою — мільйони.

Окрім тимчасової складності, розмір даних впливає і на просторові вимоги алгоритмів. Деякі алгоритми сортування, як-от сортування злиттям, потребують додаткової пам'яті для тимчасових масивів, які використовуються під час сортування. Тому при зростанні обсягу даних пам'ять може стати вузьким місцем для виконання такого алгоритму, особливо в системах з обмеженими ресурсами. У таких випадках сортування, яке використовує менше пам'яті, як-от швидке сортування, може бути кращим вибором, попри можливість найгіршого випадку $O(n^2)$.

Крім того, велика кількість даних може впливати на продуктивність через такі фактори, як кешування або обмеження апаратних ресурсів. Більші масиви даних важче помістити в кеш процесора, що може спричинити більші затримки через роботу з оперативною пам'яттю або навіть із дисковими накопичувачами. Алгоритми, які працюють з локальністю даних і ефективно використовують кеш, можуть мати перевагу при роботі з великими наборами даних.

Таким чином, вплив розміру даних на ефективність алгоритмів сортування стає вирішальним фактором при виборі конкретного методу. Якщо при роботі з малими обсягами даних різниця в ефективності може бути мінімальною, то при збільшенні розміру масиву неефективні алгоритми можуть зазнати значних проблем, знижуючи продуктивність і затримуючи процес сортування. Тому для великих наборів даних важливо вибирати алгоритми, здатні обробляти їх у розумний термін і з мінімальними витратами ресурсів.

1.5. Переваги та недоліки основних алгоритмів сортування

Кожен алгоритм має свої сильні і слабкі сторони, які можуть впливати на його продуктивність в залежності від контексту використання. Основними критеріями для порівняння є швидкість, стабільність, вимоги до пам'яті та адаптивність.

Порівняння за швидкістю, стабільністю, вимогами до пам'яті та адаптивністю.

Порівняння алгоритмів сортування є важливим аспектом при виборі найбільш підходящого методу для конкретного завдання. Кожен алгоритм має свої сильні і слабкі сторони, які можуть впливати на його продуктивність в залежності від контексту використання. Основними критеріями для порівняння є швидкість, стабільність, вимоги до пам'яті та адаптивність.

Швидкість алгоритму — один з найважливіших критеріїв, оскільки він безпосередньо впливає на загальний час виконання програми. Різні алгоритми мають різну тимчасову складність, яка зазвичай вимірюється в найгіршому, найкращому та середньому випадках. Наприклад, швидке сортування та сортування злиттям зазвичай мають складність $O(n \log n)$, що робить їх ефективними для великих обсягів даних. В той же час, алгоритми з квадратичною складністю, такі як сортування бульбашкою, вибором або вставками, значно відстають в швидкості на великих масивах, оскільки їхня продуктивність погіршується з ростом розміру даних. У випадках, коли дані вже частково відсортовані, адаптивні алгоритми, такі як сортування вставками, можуть виявитися швидшими, оскільки вони використовують цю особливість для зменшення кількості операцій[6].

Стабільність алгоритму також є важливим фактором, особливо в ситуаціях, коли потрібно зберегти відносний порядок еквівалентних елементів. Стабільні алгоритми сортування гарантують, що якщо два елементи мають однакові значення, їх порядок у вихідному масиві залишиться незмінним. Наприклад, сортування злиттям та сортування вставками є стабільними, тоді як швидке сортування за замовчуванням не є стабільним. У

багатьох випадках, особливо при роботі з об'єктами, де важливими є не лише значення, а й інші атрибути, вибір стабільного алгоритму може бути критично важливим.

Вимоги до пам'яті також варіюються між різними алгоритмами. Деякі алгоритми, як-от сортування злиттям, потребують додаткової пам'яті для тимчасових масивів, що може стати вузьким місцем при роботі з великими наборами даних. У той же час, алгоритми, такі як швидке сортування, можуть бути реалізовані з використанням тільки постійної кількості додаткової пам'яті ($O(\log n)$ для рекурсивних викликів), що робить їх більш пам'яттезберігаючими. Сортування бульбашкою та вставками взагалі не вимагають додаткової пам'яті, оскільки вони виконують операції на місці, але це досягається за рахунок їхньої низької продуктивності на великих даних.

Адаптивність алгоритму — це його здатність ефективно працювати з частково впорядкованими даними. Адаптивні алгоритми скорочують кількість порівнянь і обмінів, якщо дані вже в певному ступені відсортовані. Наприклад, алгоритм сортування вставками адаптивний, оскільки в найкращому випадку, коли дані вже упорядковані, його складність зменшується до $O(n)$. Сортування злиттям також може бути адаптивним у певних реалізаціях, де алгоритм виявляє вже відсортовані ділянки і працює з ними більш ефективно.

У результаті, при виборі алгоритму сортування важливо враховувати всі ці фактори. Швидкість може бути критично важливою для великих обсягів даних, але стабільність і вимоги до пам'яті також можуть мати вирішальне значення в залежності від специфіки завдання. Адаптивність алгоритму може значно підвищити його ефективність, якщо дані не є абсолютно випадковими. Тому при виборі методу сортування доцільно зважати на характер даних, вимоги до продуктивності та специфічні умови виконання програми.

1.6. Модифікації класичних алгоритмів сортування

Паралельне сортування

Паралельне сортування є однією з найбільш ефективних модифікацій класичних алгоритмів сортування, яка націлена на прискорення процесу сортування шляхом розподілу роботи між декількома процесорами або потоками. У традиційних алгоритмах сортування вся робота виконується послідовно, що означає, що лише один процесор обробляє дані в кожен момент часу. Це стає проблемою, коли розмір набору даних зростає, оскільки навіть найбільш ефективні алгоритми з $O(n \log n)$ складністю можуть почати потребувати занадто багато часу для виконання[7].

Застосування паралельного підходу дозволяє розбити велику задачу сортування на кілька менших підзадач, які можуть виконуватися одночасно на різних процесорах або ядрах. Це значно підвищує швидкість виконання, оскільки кожен процесор обробляє свою частину даних незалежно від інших, а після завершення сортування частин відбувається об'єднання результатів.

Одним з найпоширеніших підходів до паралельного сортування є модифікація алгоритму сортування злиттям. Алгоритм сортування злиттям вже на своєму базовому рівні передбачає поділ масиву на менші підмасиви для подальшого злиття, що робить його ідеальним кандидатом для паралельної обробки. При паралельному сортуванні злиттям початковий масив може бути розділений на кілька частин, кожна з яких сортується окремо на своєму процесорі або ядрі. Після цього відбувається паралельне злиття відсортованих підмасивів, що дозволяє зберегти ефективність $O(n \log n)$, але з додатковим прискоренням завдяки одночасній роботі кількох потоків.

Швидке сортування також може бути модифіковане для паралельної обробки. В базовому алгоритмі швидкого сортування масив ділиться на дві частини відносно опорного елемента, після чого кожна частина обробляється окремо. У паралельній версії ці дві частини можна сортувати на різних процесорах, що дозволяє виконувати сортування одночасно в кількох частинах масиву. Таким чином, підвищується продуктивність, оскільки два

процесори працюють одночасно над різними частинами масиву, а після їх сортування об'єднують результати.

Іншою важливою модифікацією є використання паралельного сортування у гібридних підходах. Наприклад, алгоритм Timsort, який використовує поєднання сортування вставками та сортування злиттям, може бути вдосконалений за допомогою паралельного оброблення. Гібридні алгоритми, поєднуючи переваги кількох методів, можуть забезпечити ефективність на різних типах даних, а з додаванням паралелізму — ще й високу швидкість обробки великих наборів.

Використання паралельних алгоритмів не обмежується багатоядерними процесорами. У сучасних умовах паралельне сортування також може виконуватися на графічних процесорах (GPU). Завдяки їхній архітектурі, що дозволяє одночасно виконувати тисячі операцій, GPU є ідеальним середовищем для виконання паралельного сортування великих наборів даних. Паралельні версії алгоритмів, таких як Bitonic Sort або Radix Sort, використовують можливості GPU для прискорення сортування в масових обчисленнях.

Однак, як і з будь-якими модифікаціями, паралельне сортування має свої виклики. Один з них — це потреба в ефективному розподілі роботи між процесорами. Якщо задача поділена нерівномірно або якщо є затримки в обміні даними між процесорами, це може призвести до ситуації, коли одні процесори простоюють, очікуючи результатів від інших. Це знижує загальну ефективність алгоритму і може навіть зробити його повільнішим, ніж послідовне виконання. Тому важливо ретельно планувати стратегії розподілу даних і синхронізації роботи процесорів.

Іншою проблемою є обмеження апаратних ресурсів. Якщо наявна система має лише декілька процесорів або ядер, то паралелізм не дасть значного приросту в продуктивності. У таких випадках може бути доцільніше використовувати більш прості та менш ресурсомісткі алгоритми.

Алгоритми сортування для специфічних задач

Алгоритми сортування для специфічних задач відіграють важливу роль, оскільки класичні методи не завжди оптимальні для певних типів даних або завдань. Задачі, такі як сортування частково відсортованих даних або сортування рядків і символів, вимагають спеціальних підходів, що дозволяють максимально ефективно використовувати властивості цих даних для підвищення продуктивності[8]

Почнемо з задачі сортування частково відсортованих даних. Це ситуація, коли масив або список, який ми маємо сортувати, вже частково впорядкований. У таких випадках деякі з класичних алгоритмів сортування можуть виявитися надмірними, оскільки виконують більше операцій, ніж це потрібно. Для таких задач ідеально підходить сортування вставками. Цей алгоритм дуже ефективний, коли дані вже майже відсортовані, оскільки кількість операцій, необхідних для впорядкування елементів, є мінімальною. У найкращому випадку, коли масив майже повністю відсортований, його складність $O(n)$, що робить сортування вставками одним з найшвидших алгоритмів для таких специфічних випадків.

Ще один приклад — алгоритм Timsort, який розроблений спеціально для роботи з реальними даними, що часто мають частково відсортовану структуру. Цей алгоритм використовує гібридну стратегію, що поєднує сортування вставками та сортування злиттям. Він розпізнає вже відсортовані ділянки в масиві і обробляє їх окремо, що робить його дуже ефективним для випадків, коли дані вже мають певний ступінь впорядкованості. Timsort застосовується у багатьох мовах програмування як стандартний алгоритм сортування, включаючи Python і Java.

Інша специфічна задача — сортування рядків і символів. Це значно відрізняється від сортування числових даних, оскільки тут ключову роль відіграють лексикографічні порівняння символів і рядків. Кожен символ у рядку має свій код, і сортування рядків проводиться на основі послідовного порівняння символів із початку кожного рядка. Найчастіше для таких задач

використовують модифікації алгоритмів швидкого сортування або сортування злиттям, оскільки вони добре підходять для ефективного порівняння елементів у рядках.

Однак для роботи з рядками існують спеціалізовані алгоритми, які використовують певні властивості символів. Наприклад, сортування по розрядах (radix sort) може бути використане для сортування рядків, коли всі рядки мають однакову довжину або можна використовувати нульові символи для вирівнювання довжини. Радиксне сортування працює за принципом обробки символів по чергово — спочатку сортуються рядки за останнім символом, потім за передостаннім і так далі, поки не будуть враховані всі символи в рядку. Цей метод є дуже ефективним для задач, де дані мають обмежений набір символів (наприклад, ASCII або цифри).

Ще один важливий метод — сортування за допомогою префіксного дерева (trie). Цей підхід застосовується для ефективного сортування великих обсягів текстової інформації, наприклад, словників або наборів слів. Trie-дерево дозволяє зберігати рядки так, щоб однакові префікси ділили одну гілку дерева, що значно економить час на сортування і пошук. Це особливо корисно, коли потрібно обробити величезну кількість рядків або текстів.

Для задач сортування символів унікальним є алгоритм лічильного сортування (counting sort). Цей алгоритм дуже ефективний, коли ми маємо справу з обмеженим набором символів. Наприклад, якщо потрібно відсортувати текстовий файл за частотою появи символів, counting sort дозволяє дуже швидко підрахувати кількість кожного символа і потім розмістити їх у відсортованому порядку за частотою.

1.7. Висновки до розділу 1

У першому розділі роботи було розглянуто ключові аспекти класифікації та аналізу алгоритмів сортування, що є критично важливими для розуміння ефективності обробки даних у різних контекстах. Основні поняття сортування включають в себе визначення терміна "сортування", а також опис

різних типів даних, які можуть підлягати цьому процесу. Важливим є розмежування алгоритмів на порівняльні та непорівняльні, що дозволяє з'ясувати, як дані можуть бути впорядковані і які алгоритми підходять для різних сценаріїв.

Подальша класифікація алгоритмів за складністю дозволяє нам зрозуміти їхню продуктивність у контексті обсягу даних. Алгоритми з квадратичною складністю, такі як сортування бульбашкою, вибором та вставками, демонструють помітно гіршу продуктивність на великих масивах даних у порівнянні з алгоритмами, що мають складність $O(n \log n)$, такими як швидке сортування і сортування злиттям. Окрім того, алгоритми з лінійною складністю, такі як лічильне сортування та сортування по розрядах, показують найкращі результати, але їх застосування обмежене специфічними умовами.

Аналіз теоретичної складності алгоритмів допомагає визначити, які з них є більш ефективними в різних сценаріях. Тимчасова складність вказує на кількість операцій, які алгоритм виконує, в той час як просторова складність відображає вимоги до пам'яті. У порівнянні складності різних алгоритмів можна побачити, як різні підходи впливають на загальну продуктивність.

Продуктивність алгоритмів сортування була проаналізована з точки зору середнього, найгіршого і найкращого сценаріїв, а також кількості порівнянь і обмінів. Вплив розміру даних на ефективність алгоритму також був розглянутий, оскільки він може суттєво змінювати результати.

РОЗДІЛ 2

ОПТИМІЗАЦІЯ АЛГОРИТМІВ СОРТУВАННЯ ДЛЯ ВЕЛИКИХ НАБОРІВ ДАНИХ

2.1. Аналіз проблем сортування великих наборів даних

Обмеження пам'яті

Коли йдеться про великі набори даних, обмеження пам'яті можуть стати серйозним бар'єром. Багато класичних алгоритмів сортування, такі як сортування злиттям, потребують значної кількості додаткової пам'яті для тимчасових масивів, що може бути критичним у ситуаціях, коли обсяги даних перевищують наявну оперативну пам'ять. Наприклад, якщо ми намагаємося відсортувати величезний масив, розмір якого перевищує доступну пам'ять, алгоритми, що вимагають додаткової пам'яті, можуть призвести до значних затримок через часті операції з введенням/виведенням, оскільки система буде змушена звертатися до повільніших зовнішніх сховищ, таких як жорсткі диски.

Для подолання цих обмежень, розроблено декілька підходів. Один з них — це алгоритм зовнішнього сортування, який спеціально призначений для роботи з великими наборами даних, що не поміщаються у пам'ять. Зовнішнє сортування розбиває великий набір даних на менші частини, які можуть бути відсортовані окремо, а потім об'єднані в один відсортований масив. Це дозволяє значно зменшити вимоги до пам'яті під час виконання сортування. Алгоритми зовнішнього сортування, такі як сортування злиттям, в основному базуються на принципі "зрозуміти все, перш ніж з'єднати", що забезпечує ефективне управління пам'яттю та ресурсам .

Кафедра КІТ				ДНП ДУ КАІ 24 16 03 000 ПЗ						
	<i>ПІБ</i>			РОЗДІЛ 2. ОПТИМІЗАЦІЯ АЛГОРИТМІВ СОРТУВАННЯ ДЛЯ ВЕЛИКИХ НАБОРІВ ДАНИХ	<i>Літ.</i>	<i>Аркуш</i>	<i>Архів</i>			
<i>Розроб.</i>	Грицак Б. Я.						42	22		
<i>Керівник</i>	Сидоренко В. М.				М-122-23-1-ТП					
<i>Н. Контр.</i>	Толстікова О.В.									

Інший підхід полягає в використанні алгоритмів, які потребують менше додаткової пам'яті. Наприклад, сортування бульбашкою або вставками працюють на місці і не потребують додаткових масивів, однак вони значно програють у швидкості на великих наборах даних через свою квадратичну складність. У таких випадках важливо зважати на баланс між використанням пам'яті та швидкістю виконання алгоритму[9].

Крім того, важливу роль у вирішенні проблем з пам'яттю відіграє вибір структур даних. Наприклад, використання деревоподібних структур даних або графів може допомогти зменшити обсяги, які потрібно зберігати у пам'яті одночасно. Взаємодія між алгоритмами сортування і структурою даних, що використовується, може суттєво вплинути на загальну продуктивність програми.

Узагальнюючи, проблеми, пов'язані з обмеженнями пам'яті, є критично важливими при розробці алгоритмів сортування для великих наборів даних. Розуміння того, як пам'ять впливає на виконання алгоритмів, може допомогти програмістам вибрати найбільш ефективні методи та оптимізувати їх реалізації, враховуючи специфіку даних, з якими вони працюють. Вибір правильного підходу до сортування в умовах обмежень пам'яті є ключовим етапом у досягненні високої продуктивності системи, що працює з великими обсягами інформації.

Проблеми з продуктивністю на великих наборах

Проблеми з продуктивністю на великих наборах даних можуть виникати з різних причин і суттєво впливають на ефективність алгоритмів сортування. Коли мова йде про обробку великих обсягів інформації, важливо враховувати не лише алгоритмічну складність, а й специфіку даних, апаратне забезпечення та середовище виконання.

Однією з основних причин проблем з продуктивністю є квадратична складність деяких класичних алгоритмів сортування, таких як сортування бульбашкою або вставками. Коли обсяг даних збільшується, час виконання цих алгоритмів зростає в геометричній прогресії, що робить їх непридатними

для великих наборів. Навіть при середніх розмірах масивів ці алгоритми можуть демонструвати значні затримки в роботі. У таких випадках доцільно використовувати більш ефективні алгоритми, такі як швидке сортування або сортування злиттям, які мають кращу асимптотичну складність ($O(n \log n)$).

Суттєво вплинути на продуктивність може і частота операцій з введенням/виведенням (I/O). Коли обсяги даних виходять за межі оперативної пам'яті, системі доводиться часто звертатися до повільних зовнішніх накопичувачів, що призводить до значних затримок. Алгоритми зовнішнього сортування, які спеціально розроблені для роботи з великими наборами даних, можуть допомогти вирішити цю проблему, але вони також вимагають ретельного управління ресурсами та пам'яттю, щоб оптимізувати час виконання.

Ще однією проблемою може стати ненадійність апаратного забезпечення. Під час обробки великих обсягів даних, особливо в умовах обмежених ресурсів, можуть виникати збої, які ускладнюють виконання алгоритмів сортування. Такі збої можуть бути спричинені перевантаженням процесора, недостатньою оперативною пам'яттю або проблемами з системою зберігання даних. Результатом цього може бути втрата даних або некоректне завершення програми. Для вирішення цих проблем можуть використовуватися паралельні обчислення, які дозволяють розподілити навантаження на кілька ядер процесора, що призводить до зменшення часу виконання завдань.

Використання неефективних структур даних може суттєво знизити продуктивність, особливо в операціях, які передбачають часті вставки, видалення або пошук. Вибір оптимальних структур даних, таких як хеш-таблиці або дерева, може підвищити загальну продуктивність алгоритмів сортування.

Адаптивність алгоритмів також може впливати на продуктивність. Деякі алгоритми показують кращі результати на частково відсортованих даних, що дозволяє зменшити кількість порівнянь і обмінів. Наприклад, сортування вставками може бути дуже ефективним у таких випадках, оскільки його

складність в найкращому випадку зменшується до $O(n)$. Проте, якщо дані повністю випадкові, алгоритми з кращими асимптотичними характеристиками, як-от швидке сортування, можуть забезпечити кращу продуктивність.

Вплив дискових операцій і I/O на швидкість обробки

Вплив дискових операцій і введення/виведення (I/O) на швидкість обробки даних - одна з найбільш критичних тем, коли ми розглядаємо алгоритми сортування, особливо для великих наборів даних. Це питання не лише стосується теоретичних аспектів алгоритмів, а й практичних викликів, з якими стикаються розробники програмного забезпечення та інженери з обробки даних.

Коли ми говоримо про обробку великих обсягів даних, важливо усвідомлювати, що оперативна пам'ять — це лише частина загальної картини. Дані, які не вміщуються у пам'яті, починають "втікати" на диск, і саме тут починається справжня гра. Диски, навіть SSD, мають велику латентність у порівнянні з оперативною пам'яттю. Це означає, що кожен раз, коли алгоритм сортування змушений звертатися до диска для читання або запису даних, відбувається затримка, яка може суттєво знизити загальну продуктивність програми[11].

Час, що витрачається на дискові операції, може бути в кілька разів більшим, ніж час, необхідний для виконання математичних операцій, які фактично здійснює алгоритм. Якщо алгоритм сортування потребує численних звернень до диска, навіть найшвидші алгоритми можуть перетворитися на тривалі і неефективні процедури. Це особливо критично, якщо дані мають бути прочитані з кількох різних місць на диску, оскільки переміщення головки диска може зайняти додатковий час.

Для оптимізації цих процесів важливо використовувати стратегії, які зменшують кількість дискових операцій. Один із способів полягає в групуванні даних для зменшення частоти звернень до диска. Алгоритми зовнішнього сортування, наприклад, проектують обробку даних так, щоб

звести до мінімуму кількість звернень до диска, обробляючи частини даних поетапно. Замість того, щоб зчитувати дані з диска по одному елементу, ці алгоритми намагаються зчитати великі блоки, обробляти їх у пам'яті, а потім записати назад на диск[10].

Крім того, паралелізація може бути корисною в цьому контексті. Якщо система має кілька ядер, розподіл завдань між ними дозволяє одночасно виконувати кілька операцій введення/виведення, що значно прискорює загальний процес обробки. Це означає, що алгоритми можуть максимально ефективно використовувати ресурси, що є особливо важливим для великих наборів даних.

Не можна також забувати про структуру даних, яка використовується для зберігання даних перед їх обробкою. Наприклад, використання ефективних структур даних, таких як дерева або хеш-таблиці, може дозволити зменшити кількість необхідних дискових операцій, оскільки вони дозволяють швидше знаходити й організувати дані.

2.2. Підходи до оптимізації

Паралельне та розподілене сортування

Паралельне та розподілене сортування — це потужні методи оптимізації, які дозволяють справлятися з величезними обсягами даних, використовуючи ресурси кількох обчислювальних одиниць. У сучасних умовах, коли дані зростають з небаченими темпами, традиційні алгоритми, які виконуються на одному процесорі, часто не здатні впоратися з таким навантаженням. Тому важливо вміти ефективно розподіляти завдання, щоб підвищити продуктивність.

Паралельне сортування передбачає одночасну обробку різних частин масиву на кількох процесорах або ядрах. Наприклад, алгоритм сортування злиттям може бути модифікований так, щоб розділити масив на кілька менших частин, які потім обробляються окремо. Кожна частина сортується в пам'яті, а

потім відсортовані підмасиви об'єднуються в один великий масив. Цей підхід дозволяє значно скоротити час виконання, оскільки обробка відбувається паралельно.

Розподілене сортування, в свою чергу, дозволяє розподіляти дані не лише між ядрами одного процесора, а й між кількома фізичними машинами. Тут важливу роль відіграє модель MapReduce, яка була розроблена Google для обробки великих обсягів інформації. В основі MapReduce лежать два ключові етапи: "Map", на якому дані обробляються і фільтруються, і "Reduce", де результати агрегуються в один підсумковий масив. Кожен етап виконується паралельно на різних вузлах, що дозволяє максимально ефективно використовувати ресурси.

Hadoop — це відкрите програмне забезпечення, яке реалізує концепцію MapReduce і забезпечує зберігання даних на розподілених системах. Це дозволяє користувачам працювати з великими наборами даних, розподіляючи їх між багатьма комп'ютерами. Завдяки Hadoop можна обробляти величезні обсяги інформації, зберігаючи дані на недорогих серверах. Наприклад, завдяки цій платформі компанії можуть обробляти дані, які не вміщуються у пам'ять одного комп'ютера, просто розподіляючи їх між кількома машинами.

Однак розподілене сортування має свої виклики. Управління даними стає більш складним, адже необхідно стежити за узгодженням результатів, які можуть бути у різних станах через затримки в мережі. Конфлікти між процесами також можуть призвести до проблем з коректністю даних, тому важливо мати належні механізми управління.

Паралельне та розподілене сортування відкриває нові можливості для обробки великих наборів даних. Вони забезпечують гнучкість, масштабованість та можливість використання сучасних обчислювальних потужностей, що дозволяє швидко реагувати на зростаючі вимоги до обробки даних у різних сферах. Це особливо актуально для бізнесу, науки, медицини та інших галузей, де дані становлять величезну цінність.

Мультиядерні процесори та їх вплив на продуктивність алгоритмів

Мультиядерні процесори стали стандартом для сучасних комп'ютерних систем, що призвело до значних змін у способах розробки та оптимізації алгоритмів. Вони дозволяють виконувати кілька операцій одночасно, що суттєво підвищує продуктивність, особливо для обчислювально інтенсивних завдань, таких як сортування великих наборів даних. Використання мультипроцесорності дозволяє ефективно розподіляти навантаження та знижувати час виконання алгоритмів.

Основна перевага мультиядерних процесорів полягає в їхній здатності виконувати паралельні обчислення. Коли алгоритми адаптуються до роботи на кількох ядрах, кожне ядро може обробляти свою частину даних, що дозволяє зменшити загальний час обробки. Це особливо важливо для алгоритмів сортування, де дані можуть бути розділені на менші частини, які обробляються незалежно. Наприклад, у випадку сортування злиттям, масив можна розділити на кілька підмасивів, які будуть сортуватися одночасно на різних ядрах, а потім ці підмасиви об'єднуються.

При цьому важливо враховувати, що не всі алгоритми можуть бути ефективно паралелізовані. Деякі з них мають серйозні залежності між етапами обробки, що ускладнює їх адаптацію до паралельного виконання. Наприклад, алгоритми, що використовують багато серійних кроків, можуть не вигравати від використання мультиядерності, оскільки частина обчислень залежить від результатів попередніх етапів.

Однак технології, що дозволяють ефективно використовувати мультиядерні процесори, продовжують розвиватися. Багато сучасних алгоритмів, такі як QuickSort і MergeSort, вже адаптовані для роботи в паралельному режимі, що дає змогу значно покращити їхню продуктивність. Паралелізація також може бути реалізована за допомогою бібліотек, таких як OpenMP або MPI, які надають програмістам можливості для ефективного управління паралельними обчисленнями.

Важливим аспектом використання мультиядерних процесорів є управління ресурсами, оскільки паралельні обчислення можуть призвести до конкурентних умов. Коли кілька ядер намагаються отримати доступ до одних і тих же ресурсів, таких як пам'ять або введення/виведення, це може призвести до затримок і зниження продуктивності. Тому ефективні стратегії управління пам'яттю і синхронізації процесів стають критично важливими для досягнення оптимальної продуктивності.

GPU-сортування

GPU-сортування використовує графічні процесори для прискорення обробки даних, включаючи алгоритми сортування. Графічні процесори, або GPU, мають архітектуру, оптимізовану для виконання великої кількості одночасних обчислень. Це робить їх ідеальними для завдань, що вимагають паралельних обчислень, таких як обробка великих наборів даних. Використання GPU для сортування даних дозволяє досягти значного приросту продуктивності в порівнянні з традиційними процесорами.

Однією з основних переваг GPU-сортування є можливість обробляти тисячі або навіть мільйони потоків одночасно. Це означає, що алгоритми, які можуть бути адаптовані для паралельного виконання, можуть суттєво виграти від використання графічних процесорів. Наприклад, алгоритми, такі як битове сортування або сортування злиттям, можуть бути реалізовані на GPU, що дозволяє одночасно обробляти величезні обсяги даних.

Одним із найбільш відомих підходів до GPU-сортування є алгоритм "Bitonic Sort", який використовує ідеї, засновані на битонних послідовностях. Цей алгоритм має високу ступінь паралелізму, що робить його особливо підходящим для виконання на графічних процесорах. Інший популярний алгоритм — "Radix Sort", який також можна адаптувати для роботи на GPU, завдяки своїй простій структурі, що дозволяє легко розподілити обчислення між різними потоками.

Використання GPU для сортування вимагає специфічних знань і навичок програмування, оскільки розробники повинні розуміти, як працює

архітектура графічних процесорів і як оптимізувати код для досягнення максимальної продуктивності. Для цього існують бібліотеки, такі як CUDA (Compute Unified Device Architecture) від NVIDIA, яка надає програмістам інструменти для роботи з GPU. За допомогою CUDA можна створювати програми, які використовують обчислювальні потужності графічних процесорів для виконання паралельних обчислень, включаючи сортування[12].

Однак, як і у випадку з будь-якими технологіями, GPU-сортування має свої обмеження. Не всі алгоритми можуть бути ефективно адаптовані для графічних процесорів, і в деяких випадках обробка даних на GPU може виявитися менш ефективною, ніж на CPU, особливо якщо дані не підходять для паралельної обробки. Додатково, переміщення даних між пам'яттю графічного процесора та основною пам'яттю комп'ютера може бути дорогим за часом, тому важливо оптимізувати цей процес.

Таким чином, GPU-сортування - це потужний метод, який може значно прискорити процес обробки великих наборів даних, якщо алгоритми правильно адаптовані для роботи на графічних процесорах. Завдяки своїм можливостям для паралельних обчислень, GPU стають дедалі більш популярними у сферах, де потрібно обробляти великі обсяги інформації з високою швидкістю.

2.3. Ефективне використання кеш-пам'яті

Кеш-обізнані та кеш-ефективні алгоритми

Ефективне використання кеш-пам'яті має критичне значення для підвищення продуктивності алгоритмів, оскільки доступ до кешу значно швидший, ніж до основної пам'яті. Кеш-пам'ять розташована близько до процесора і призначена для зберігання даних, до яких часто звертаються. Коли дані знаходяться в кеші, час доступу до них суттєво зменшується, що може позитивно вплинути на швидкість виконання алгоритмів, зокрема тих, що працюють з великими наборами даних.

Кеш-обізнані алгоритми враховують структуру кешу при проектуванні своїх операцій. Вони намагаються зменшити кількість кеш-промахів, що виникають, коли процесор намагається звернутися до даних, яких немає в кеші. Один із способів досягнення цього — оптимізація порядку обробки даних, щоб дані, які часто використовуються разом, були розміщені в сусідніх адресах пам'яті. Це дозволяє алгоритмам завантажувати цілі блоки даних у кеш за один раз, зменшуючи потребу у повторних зверненнях до повільнішої основної пам'яті[13].

Кеш-ефективні алгоритми йдуть ще далі, розробляючи стратегії, які максимально використовують кеш-пам'ять під час виконання. Наприклад, алгоритми сортування, які реалізують методи, що обробляють дані блочно, можуть бути більш ефективними в контексті кешу. Вони обробляють дані не по одному елементу, а блоками, зменшуючи кількість звернень до пам'яті і використовуючи кеш для зберігання найбільш актуальних даних. Це дозволяє уникнути кеш-промахів і знижує затримки.

Загалом, підхід до розробки кеш-обізнаних і кеш-ефективних алгоритмів вимагає детального розуміння архітектури системи, на якій ці алгоритми будуть виконуватися. Наприклад, різні процесори мають різні розміри кешу, а також різні стратегії кешування, що може вплинути на вибір алгоритму. Розробники можуть використовувати інструменти профілювання для аналізу кеш-ефективності алгоритмів і виявлення вузьких місць, пов'язаних із кешем.

Кеш-пам'ять зазвичай розділяється на кілька рівнів: L1, L2 та L3, кожен з яких має різні розміри та швидкості. Алгоритми можуть бути оптимізовані для використання конкретного рівня кешу, щоб зменшити затримки. Наприклад, часті звернення до даних, які перебувають у кеші L1, забезпечують найшвидший доступ, тому алгоритми можуть бути спроектовані так, щоб зосереджуватися на максимально можливому використанні кешу L1.

Ефективне використання кеш-пам'яті безпосередньо впливає на загальну продуктивність системи. Алгоритми, які усвідомлюють кеш, можуть

значно знизити затримки та підвищити швидкість виконання, що є особливо важливим у випадках, коли обробка великих наборів даних вимагає максимальної ефективності.

Алгоритми, оптимізовані для багаторівневої ієрархії пам'яті

Алгоритми, оптимізовані для багаторівневої ієрархії пам'яті, враховують різні рівні пам'яті, які використовуються в сучасних комп'ютерах, щоб забезпечити максимально ефективно виконання завдань. Багаторівнева ієрархія пам'яті зазвичай складається з декількох рівнів, включаючи регістри, кеш-пам'ять (L1, L2, L3) і основну пам'ять. Кожен з цих рівнів має свої характеристики швидкості та обсягу, тому алгоритми повинні бути спроектовані так, щоб мінімізувати затримки доступу до даних.

Першим етапом в оптимізації алгоритмів є розуміння принципів роботи кешу. Алгоритми, які використовують локальність посилань, можуть значно зменшити кількість кеш-промахів. Локальність посилань поділяється на дві категорії: тимчасова локальність, що передбачає повторне використання даних у найближчому часі, і просторову локальність, яка передбачає, що, якщо один елемент масиву був доступний, то, ймовірно, сусідні елементи також будуть доступні найближчим часом.

Підходи до оптимізації алгоритмів для багаторівневої ієрархії пам'яті:

1. Блочне оброблення: блочні алгоритми розбивають дані на менші частини або блоки, що дозволяє максимально використовувати кеш. Цей підхід є особливо корисним для алгоритмів матричних обчислень, таких як множення матриць. Замість того, щоб працювати з усіма даними одночасно, алгоритм обробляє їх блоками, що знижує кількість звернень до повільнішої основної пам'яті.
2. Паралелізація: алгоритми можуть бути оптимізовані для багатоядерних архітектур, де кілька ядер виконують частини одного завдання одночасно. З використанням паралельних обчислень можна одночасно обробляти кілька блоків даних, що зменшує загальний час виконання.

3. Оптимізація доступу до пам'яті: порядок, в якому дані обробляються, може вплинути на продуктивність. Наприклад, послідовний доступ до масиву є більш ефективним, ніж випадковий. Це дозволяє кешу зберігати найбільш актуальні дані, що зменшує затримки.
4. Адаптивні структури даних: вибір відповідної структури даних також може вплинути на продуктивність. Наприклад, дерева, хеш-таблиці або масиви можуть бути більш оптимальними в залежності від типу операцій, які виконуються. Використання адаптивних структур даних може зменшити кількість звернень до пам'яті та покращити загальну швидкість виконання.
5. Зменшення розміру робочого набору: зменшення кількості даних, що обробляються одночасно, може призвести до зменшення кількості кеш-промахів. Алгоритми можуть бути спроектовані так, щоб працювати з меншими підмножинами даних, зберігаючи інші дані у пам'яті, коли вони не потрібні.
6. Гібридні методи: комбінація декількох стратегій, таких як блочне оброблення та адаптивні структури даних, може призвести до значного підвищення продуктивності. Гібридні алгоритми можуть бути налаштовані для конкретних задач, максимально використовуючи архітектуру пам'яті.

2.4. Зовнішні алгоритми сортування

Алгоритми, що працюють з даними, що не поміщаються в оперативну пам'ять

Алгоритми, що працюють з даними, які не поміщаються в оперативну пам'ять, відомі як зовнішні алгоритми сортування (external sort). Вони використовуються в ситуаціях, коли обсяг даних значно перевищує доступну оперативну пам'ять комп'ютера, і вимагають спеціальних стратегій для ефективної обробки та зберігання інформації на зовнішніх носіях, таких як жорсткі диски або SSD. Основна мета цих алгоритмів — мінімізувати кількість

операцій з читання і запису даних, оскільки доступ до дискової пам'яті є набагато повільнішим, ніж доступ до оперативної пам'яті.

Одним з найбільш поширених методів зовнішнього сортування є зовнішнє злиття. Цей метод складається з кількох етапів, які дозволяють ефективно обробляти великі обсяги даних. Спочатку великі дані розбиваються на менші частини, які можуть поміститися в оперативній пам'яті. Ці частини сортуються за допомогою одного з внутрішніх алгоритмів сортування, таких як сортування вставками або бульбашкою. Після цього відсортовані частини зберігаються на зовнішньому носії.

Коли всі частини відсортовані, наступний крок полягає в їх об'єднанні. Для цього використовується метод злиття, який полягає в тому, що алгоритм читає перші елементи з кожної відсортованої частини, порівнює їх і записує найменший з них у фінальний відсортований файл. Цей процес продовжується до тих пір, поки всі елементи не будуть перенесені в фінальний вихідний файл. Зовнішнє злиття є дуже ефективним, оскільки воно використовує невелику кількість оперативної пам'яті для зберігання лише кількох елементів, що забезпечує швидкий доступ до даних.

Іншим підходом є зовнішнє сортування з використанням дерев. Цей метод дозволяє зберігати дані в структурі, схожій на бінарне дерево, і використовувати її для оптимізації процесу сортування. Деревоподібні структури даних дозволяють зберігати відсортовані елементи в пам'яті та поступово зливати їх, що підвищує ефективність доступу до даних на зовнішніх носіях.

Крім того, алгоритми, що працюють з зовнішніми даними, повинні враховувати особливості файлової системи та організації даних на диску. Наприклад, важливо мінімізувати фрагментацію, що може виникати в процесі запису великих обсягів даних. Алгоритми можуть бути налаштовані для роботи з даними, які організовані у спеціальні блоки або файлові структури, щоб зменшити затримки під час читання та запису.

Алгоритм зовнішнього злиття

Алгоритм зовнішнього злиття, або External Merge Sort, є однією з найпопулярніших стратегій для сортування великих обсягів даних, які не поміщаються в оперативну пам'ять. Він ґрунтується на принципах внутрішнього сортування та ефективного злиття відсортованих частин. Цей алгоритм особливо корисний у ситуаціях, коли потрібно обробляти дані, що зберігаються на зовнішніх носіях, таких як жорсткі диски або SSD, оскільки він максимально знижує кількість операцій введення/виведення [14].

Процес зовнішнього злиття можна поділити на кілька основних етапів. Перший етап передбачає розбивку великого набору даних на менші частини, які можуть бути завантажені в оперативну пам'ять. Ці частини називаються "блоками". Розмір кожного блоку визначається залежно від доступної оперативної пам'яті.

Після цього алгоритм здійснює наступні кроки:

1. Кожен блок завантажується в оперативну пам'ять, де застосовується один з внутрішніх алгоритмів сортування, таких як сортування вставками або швидке сортування. Після сортування блоки записуються назад на диск у відсортованому вигляді. Цей етап дозволяє перетворити необроблені дані на кілька відсортованих частин.
2. Після того як всі блоки були відсортовані, алгоритм переходить до злиття. На цьому етапі алгоритм відкриває всі відсортовані блоки одночасно та використовує метод, схожий на злиття в алгоритмі злиття (Merge). Алгоритм зчитує перші елементи з кожного блоку, порівнює їх і записує найменший у фінальний відсортований вихідний файл. Цей процес повторюється до тих пір, поки всі елементи з усіх блоків не будуть перенесені у фінальний файл.
3. Після того як всі повні блоки оброблені, алгоритм може залишити деякі елементи, які не потрапили в остаточний файл. Якщо блоки не рівномірно розподілені, можуть залишитися незавершені записи, які також потрібно злиття.

Оскільки доступ до зовнішньої пам'яті є повільнішим, алгоритм оптимізує кількість операцій читання і запису. Зазвичай зовнішнє злиття використовує підходи до зменшення фрагментації та оптимізації доступу до пам'яті, що робить його дуже ефективним для того щоб працювати з обмеженнями пам'яті

Ефективність алгоритму зовнішнього злиття полягає в кількості зчитуваних та записуваних блоків, а також розміром блоку. Алгоритм вважається дуже оптимальним для роботи з великими наборами даних, оскільки він забезпечує високу швидкість та ефективність. Зовнішнє злиття використовується в багатьох сучасних системах, таких як бази даних та великі системи обробки даних, де необхідно обробляти великі обсяги інформації без значних затримок.

Каскадне злиття

Каскадне злиття, або *multiway merge*, є вдосконаленою формою злиття, яка використовується в алгоритмах зовнішнього сортування для обробки великої кількості відсортованих файлів. Це метод, що дозволяє зливати декілька відсортованих списків одночасно, що значно підвищує ефективність злиття, особливо коли мова йде про велику кількість відсортованих блоків.

Процес каскадного злиття починається з того, що відкриваються всі відсортовані блоки, і алгоритм зчитує з них перші елементи. Замість того, щоб обробляти лише пару блоків, каскадне злиття дозволяє обробляти, скажімо, 4 або 8 блоків одночасно, що підвищує швидкість злиття. Коли алгоритм зчитує перші елементи з усіх блоків, він порівнює їх і визначає найменший з них, який потім записується у фінальний вихідний файл. Процес повторюється, поки всі елементи не будуть оброблені.

Цей підхід має кілька переваг. По-перше, він знижує кількість операцій зчитування та запису, оскільки одночасно обробляються кілька блоків, що зменшує час доступу до дискової пам'яті. По-друге, завдяки одночасному зчитуванню даних з кількох блоків, каскадне злиття може використовувати локальність посилань, що призводить до підвищення продуктивності.

Алгоритм зможе залишити більше даних у кеші, що також позитивно вплине на швидкість виконання.

Проте каскадне злиття потребує більшої кількості оперативної пам'яті, оскільки одночасно необхідно зберігати елементи з усіх блоків, що обробляються. Кількість блоків, що можуть бути оброблені одночасно, залежить від обсягу доступної оперативної пам'яті. Якщо пам'ять обмежена, алгоритм може потребувати зменшення кількості оброблюваних блоків, що, в свою чергу, може знизити його ефективність.

Каскадне злиття активно використовується в системах, де обробляються великі обсяги даних, наприклад, у БД і в системах обробки даних, оскільки дозволяє значно зменшити час, витрачений на злиття відсортованих блоків. Цей метод робить злиття більш гнучким і швидким, що в умовах сучасних вимог до обробки даних є надзвичайно важливим.

2.5. Оптимізація розподілу ресурсів

Балансування використання процесора, пам'яті та диска.

Оптимізація розподілу ресурсів, особливо балансування використання процесора, пам'яті та диска, має критичне значення для забезпечення високої продуктивності алгоритмів сортування, особливо в умовах обробки великих обсягів даних. Коли мова йде про сортування даних, правильне управління ресурсами може істотно вплинути на швидкість і ефективність виконання.

Процесор, пам'ять і диск виконують різні функції в системі. Процесор відповідає за обчислення, пам'ять зберігає дані, а диск служить для довготривалого зберігання інформації. Коли алгоритми сортування активні, важливо, щоб усі ці ресурси працювали узгоджено і не створювали "вузьких місць". Наприклад, якщо процесор обробляє дані дуже швидко, а система не може встигнути записати результати на диск, це призведе до того, що процесор буде змушений чекати, що знизить загальну продуктивність [15].

Балансування використання ресурсів починається з аналізу алгоритмів, які можуть бути виконані паралельно. У цьому випадку розподіл роботи між

кількома процесорами чи ядрами дозволяє зменшити навантаження на окремі компоненти системи. Наприклад, при використанні паралельного сортування, коли дані діляться на частини, кожна частина обробляється окремим потоком. Це допомагає забезпечити максимальне використання процесорного часу і зменшує час виконання в цілому.

Пам'ять також потребує оптимізації. У разі недостатньої оперативної пам'яті алгоритми можуть використовувати сторінкову пам'ять або обробляти дані частинами, що створює додаткові затримки. Зокрема, при зовнішньому сортуванні важливо мати належний розмір буферів для зберігання проміжних результатів. Занадто малий буфер може призвести до частих операцій запису та читання, в той час як занадто великий може призвести до перевантаження пам'яті.

Диск також має бути оптимізований для роботи з даними. Швидкість доступу до даних на диску може значно вплинути на швидкість сортування. Використання SSD замість традиційних жорстких дисків може суттєво підвищити продуктивність. Крім того, належна організація даних на диску, наприклад, у вигляді блоків або сегментів, може зменшити фрагментацію та прискорити доступ до даних[16].

Загалом, балансування використання процесора, пам'яті та диска вимагає комплексного підходу до розробки алгоритмів, які здатні оптимально використовувати всі доступні ресурси. Це передбачає не лише ефективну реалізацію алгоритмів, але й розуміння архітектури апаратного забезпечення, на якому вони працюють. Правильне управління ресурсами дозволяє суттєво підвищити ефективність обробки даних, забезпечуючи швидке та ефективне сортування, навіть для великих наборів інформації.

Технології управління потоками

Потоки дозволяють програмам виконувати кілька завдань одночасно, що значно підвищує ефективність використання ресурсів системи. У контексті сортування даних управління потоками включає в себе паралельне виконання, асинхронне програмування та оптимізацію обміну даними між потоками.

Коли говорять про паралельне виконання, мається на увазі, що одна програма може розподілити свою роботу на кілька потоків, кожен з яких може виконувати певну частину завдання одночасно. Наприклад, у алгоритмах, які реалізують паралельне сортування, дані діляться на кілька частин, і кожна частина обробляється окремим потоком. Це дозволяє значно зменшити час виконання, оскільки кілька обчислювальних ядер можуть працювати одночасно [17].

Асинхронне програмування також відіграє важливу роль в управлінні потоками. Замість того, щоб блокувати виконання програми, коли одне завдання потребує часу для завершення (наприклад, операції введення/виведення), асинхронний підхід дозволяє програмі продовжувати виконувати інші завдання. Це означає, що під час очікування результатів операцій зчитування або запису даних, програма може виконувати додаткові обчислення, що підвищує загальну продуктивність.

Коли кілька потоків працюють над одними й тими ж даними або взаємодіють один з одним, необхідно забезпечити належний механізм для синхронізації доступу до цих даних. Це може включати використання блокувань, семафорів або інших механізмів синхронізації, які дозволяють уникнути конфліктів при доступі до спільних ресурсів. Однак надмірна синхронізація може призвести до затримок, тому важливо знайти баланс між безпекою доступу до даних і продуктивністю.

Управління потоками також вимагає належного планування розподілу ресурсів. Кожен потік споживає частину системних ресурсів, і при великій кількості потоків може виникнути конкуренція за ці ресурси. Технології, які дозволяють динамічно регулювати кількість активних потоків залежно від навантаження на систему, можуть суттєво підвищити ефективність виконання.

Усвідомлення особливостей управління потоками дозволяє розробникам створювати більш продуктивні алгоритми сортування, які можуть працювати з великими наборами даних ефективно та швидко. Ці технології стають все більш актуальними в умовах, коли обсяги даних

зростають, а вимоги до їх обробки стають дедалі складнішими. Застосування технологій управління потоками сприяє досягненню кращої продуктивності і знижує затримки при виконанні алгоритмів сортування.

2.6. Оцінка продуктивності та експерименти

Тестування продуктивності різних оптимізованих алгоритмів на великих наборах даних

Дослідження тестування продуктивності різних оптимізованих алгоритмів на великих наборах даних дозволяє зрозуміти, як різні алгоритми поведуться при обробці великих обсягів інформації, і які з них є найбільш ефективними в конкретних умовах. Важливість такого тестування зростає в умовах, коли дані стають дедалі більшими і складнішими, а вимоги до швидкості обробки залишаються високими[18].

Важливо визначити критерії, за якими будуть оцінюватися алгоритми. Це може включати час виконання, використання пам'яті, кількість операцій вводу/виводу та інші фактори, які впливають на загальну продуктивність. Також необхідно забезпечити однакові умови для тестування, щоб результати були порівнянними. Наприклад, всі алгоритми повинні бути протестовані на однакових наборах даних, що мають схожі характеристики: обсяги, рівень відсортованості, розподіл значень тощо.

Після підготовки умов для тестування важливо створити набори даних, які б відображали реальні сценарії використання. Це можуть бути як випадкові дані, так і частково відсортовані набори, які часто зустрічаються в практиці. Для тестування алгоритмів можна використовувати різні методи генерації даних, щоб охопити різноманітні випадки використання. Чим різноманітніші дані, тим краще можна оцінити, як алгоритми справляються з реальними умовами.

Після виконання тестів результати аналізуються для визначення найкращих алгоритмів за кожним із критеріїв. Це може включати порівняння швидкості виконання, де можна виявити, який алгоритм найшвидший для

даних, що підлягають сортуванню. Важливо також звернути увагу на використання пам'яті, оскільки алгоритм, який швидко виконує сортування, може виявитися непридатним, якщо споживає надмірну кількість ресурсів.

Крім того, тестування повинно включати сценарії з різними обсягами даних. Наприклад, результати можуть змінюватися в залежності від того, чи обробляються мільйони, мільярди або навіть трильйони записів. Це дозволяє виявити, як алгоритми масштабуються з ростом обсягу даних і чи зберігають вони свою ефективність.

Зібрані дані можуть бути представлені у вигляді графіків та таблиць, що полегшить їх сприйняття. Візуалізація результатів може допомогти зрозуміти, як алгоритми поведуться в різних умовах. Це також може бути корисно для прийняття рішень щодо вибору найбільш ефективного алгоритму для конкретного завдання або для подальших оптимізацій.

Тестування продуктивності — це не просто формальність, а важливий процес, який може суттєво вплинути на вибір алгоритму в умовах, коли дані продовжують зростати, а вимоги до швидкості обробки залишаються на високому рівні.

Оцінка впливу паралельної обробки та розподілених обчислень

Коли мова йде про обробку великих наборів даних, паралельні та розподілені системи здатні забезпечити значний приріст швидкості завдяки одночасному виконанню кількох завдань. Це дозволяє не тільки зменшити час виконання алгоритмів, але й ефективно використовувати ресурси системи.

Паралельна обробка передбачає використання кількох процесорів або ядер для виконання одного і того ж завдання, що розподіляє навантаження та зменшує час, необхідний для обробки даних. В алгоритмах сортування, таких як швидке сортування або сортування злиттям, дані можуть бути розділені на частини, які обробляються окремими потоками. Це дозволяє зменшити загальний час виконання, оскільки кілька ядер можуть виконувати обчислення одночасно. Однак важливо врахувати, що не всі алгоритми легко піддаються

паралелізації. Тому важливим аспектом є аналіз, як конкретний алгоритм може бути адаптований для паралельного виконання.

Розподілені обчислення, в свою чергу, включають обробку даних на кількох машинах або серверах, що дозволяє масштабувати систему для роботи з надзвичайно великими наборами даних. У таких системах дані можуть бути зберігані на кількох вузлах, а обчислення виконуються паралельно на цих вузлах. Технології, такі як Hadoop та MapReduce, стали стандартом для реалізації розподілених обчислень. Вони дозволяють ефективно обробляти великі дані, зберігаючи їх на різних машинах і виконуючи обчислення паралельно. При цьому важливо звернути увагу на комунікаційні витрати, оскільки передача даних між вузлами може створювати затримки, які знижують загальну продуктивність.

Одним із ключових факторів оцінки впливу паралельної обробки та розподілених обчислень є так званий "speedup" — це відношення часу виконання алгоритму на одному процесорі до часу, витраченого на виконання того ж алгоритму з використанням багатьох процесорів. Чим вищий показник "speedup", тим краще алгоритм справляється з паралельною обробкою. Однак цей показник може варіюватися в залежності від архітектури системи, структури даних і способу реалізації алгоритму.

Крім того, важливо врахувати "паралельну ефективність", яка визначає, наскільки ефективно використовується доступна обчислювальна потужність. Це відношення часу, що витрачається на обчислення в паралельному режимі, до теоретично максимального часу, який можна було б досягти при ідеальному розподілі завантаження..

Оцінка впливу цих підходів потребує не лише тестування продуктивності, але й глибокого аналізу ресурсів, які використовуються під час виконання алгоритмів. Важливо виявити можливі "вузькі місця", які можуть негативно вплинути на загальну продуктивність. Наприклад, затримки в обміні даними, недостатнє використання пам'яті або процесорних ресурсів можуть суттєво знизити ефективність паралельних та розподілених обчислень.

Наприклад, ми маємо великий набір даних, наприклад, 1 мільярд записів, які потрібно відсортувати. Використання традиційних алгоритмів сортування може виявитися недостатнім через обмеження в пам'яті та час виконання. Замість цього ми можемо використовувати паралельну обробку на кількох машинах, застосовуючи MapReduce для розподіленого сортування.

На першому етапі, відомому як "Map", дані розбиваються на менші частини, які потім обробляються паралельно на різних вузлах. Кожен вузол виконує локальне сортування своїх частин, застосовуючи, наприклад, алгоритм злиття. Таким чином, кожен вузол отримує відсортовані підмножини даних.

На другому етапі, "Reduce", відсортовані підмножини з кожного вузла зливаються в єдиний відсортований набір. Цей етап також може бути реалізований паралельно, оскільки ми можемо використовувати кілька потоків для злиття частин.

Переваги такої паралельної обробки очевидні. По-перше, завдяки розподілу завантаження на кілька машин, час виконання алгоритму значно зменшується. Для набору даних у 1 мільярд записів, який зазвичай потребує декількох годин для виконання на одному процесорі, використання паралельної обробки може зменшити цей час до кількох хвилин, залежно від кількості доступних ресурсів.

Однак слід врахувати, що паралельна обробка та розподілені обчислення також можуть мати свої недоліки. Наприклад, надмірна комунікація між вузлами може призвести до затримок. Якщо дані, що передаються між вузлами, є значними, це може вплинути на загальний час виконання. Тому важливо оптимізувати не лише алгоритм, але й саму архітектуру системи.

2.7. Висновки до розділу 2

Серед ефективних підходів до оптимізації варто відзначити паралельне та розподілене сортування, які, завдяки своїй здатності розподіляти навантаження на кілька обчислювальних одиниць, демонструють значні

переваги у порівнянні з традиційними методами. Впровадження технологій, таких як MapReduce і Hadoop, дозволяє не лише швидше виконувати обчислення, але й ефективно масштабувати систему в умовах зростаючих обсягів даних.

Також важливим аспектом є використання мультитядерних процесорів і графічних процесорів (GPU) для підвищення швидкості обробки. Це підкреслює значення адаптації алгоритмів до новітніх технологій, що відкриває нові можливості для їхньої продуктивності.

Ефективне використання кеш-пам'яті та алгоритмів, оптимізованих для багаторівневої ієрархії пам'яті, дозволяє знизити затримки при доступі до даних, що суттєво підвищує швидкість обробки. Зовнішні алгоритми сортування, такі як зовнішнє злиття і каскадне злиття, демонструють свою важливість при роботі з великими обсягами даних, що не поміщаються в оперативну пам'ять.

Проведене тестування продуктивності різних оптимізованих алгоритмів і оцінка впливу паралельної обробки та розподілених обчислень підтверджують необхідність постійного вдосконалення алгоритмів для забезпечення їхньої ефективності в умовах зростаючих обсягів даних. Отримані результати вказують на перспективи подальшого розвитку в цій галузі, що дозволить розробляти ще більш продуктивні рішення для вирішення сучасних викликів у обробці даних.

РОЗДІЛ 3

РЕАЛІЗАЦІЯ ОПТИМІЗАЦІЇ АЛГОРИТМУ

3.1. Огляд алгоритму для оптимізації

MapReduce дозволяє обробляти величезні обсяги даних, розподіляючи їх на окремі частини, які можуть оброблятися паралельно. Це забезпечує високий рівень масштабованості, який необхідний для оптимізації роботи з великими обсягами даних, де традиційні алгоритми сортування можуть бути занадто повільними або ресурсозатратними[19].

В умовах розподілених систем, таких як кластери або хмарні сервіси, MapReduce, використовуючи розподілену обробку даних, дає можливість розділити задачу на багато етапів і обробляти їх паралельно, знижуючи затримки та підвищуючи загальну продуктивність системи.

Одним з основних викликів при роботі з великими наборами даних є обмеження на обсяг оперативної пам'яті. MapReduce добре справляється з цією проблемою, оскільки він дозволяє працювати з даними, що зберігаються на дисках, і мінімізує кількість операцій вводу-виводу (I/O), що може бути значним обмеженням для традиційних алгоритмів сортування.

Модель MapReduce дозволяє оптимізувати використання ресурсів, таких як пам'ять і процесор, шляхом розподілу задач між різними вузлами в кластері. Кожен етап (Map і Reduce) можна налаштувати так, щоб він максимально ефективно використовував ресурси системи, забезпечуючи балансування навантаження та зменшення часу на виконання операцій.

MapReduce добре підтримуються в різних популярних фреймворках для обробки великих даних, таких як Hadoop та Apache Spark. Ці фреймворки дозволяють значно спростити реалізацію і обслуговування алгоритмів, а також

Кафедра КІТ				ДНП ДУ КАІ 24 16 03 000 ПЗ							
	<i>ПІБ</i>			РОЗДІЛ 3. РЕАЛІЗАЦІЯ ОПТИМІЗАЦІЇ АЛГОРИТМУ			<i>Літ.</i>	<i>Аркуш</i>	<i>Аркушів</i>		
<i>Розроб.</i>	Грицак Б. Я.								65	27	
<i>Керівник</i>	Сидоренко В. М.						М-122-23-1-ТП				
<i>Н. Контр.</i>	Толстікова О.В.										

забезпечують широкий спектр інструментів для роботи з даними, таких як зберігання, обробка та аналіз, що робить MapReduce оптимальним вибором для реалізації практичних задач.

Хоча MapReduce - це потужний інструмент для розподілених обчислень, програмування в ньому є відносно простим. Програмісти можуть зосередитися на визначенні функцій `map` і `reduce`, не турбуючись про низькорівневі деталі управління розподіленими обчисленнями або зберіганням даних. Це спрощує процес реалізації і тестування алгоритмів.

3.2. Основні концепції MapReduce

1. Map Stage (Етап мапування)

Фаза мапування — це перший етап обробки даних у MapReduce, де вхідні дані (наприклад, великі набори текстових даних або чисел) обробляються за допомогою функції Map. На цьому етапі дані розділяються на менші частини, які можуть оброблятися паралельно на різних вузлах. Кожна частина даних передається функції Map, яка виконує певну операцію, таку як фільтрація або перетворення даних[20].

Функція Map приймає пару (ключ, значення) на вході та видає нову пару (ключ, значення) на виході. Тобто, Map функція бере кожен елемент даних і перетворює його в іншу пару ключ-значення, де ключ може бути використаний для подальшої обробки. Після цього всі пари ключ-значення сортуються та передаються до наступного етапу — Reduce (рис.3.1).

```
def map_function(key, value):
    words: object = value.split() # Розділяємо текст на слова
    for word in words:
        # Для кожного слова створюємо пару (слово, 1)
        yield (word, 1)
```

Рис. 3.1. Приклад функції Map

У цьому прикладі функція Map приймає пару (ключ, значення) — це може бути рядок тексту, і розбиває текст на окремі слова, видаючи пару (слово,

1) для кожного слова. Тобто, кожне слово буде асоційовано з числом 1, що є частиною задачі підрахунку слів у тексті.

2. Shuffle and Sort

Після виконання функції Map, усі пари ключ-значення передаються до стадії сортування та переміщення (Shuffle and Sort). На цьому етапі всі однакові ключі групуються разом, що дозволяє подальшу обробку за одним ключем. Цей етап критичний для коректного виконання функції Reduce, оскільки саме на основі ключів буде виконана агрегація результатів.

3. Reduce Stage (Етап зведення)

Фаза зведення — це останній етап, де функція Reduce обробляє пари ключ-значення, отримані після етапу мапування та сортування. Функція Reduce приймає всі значення для одного ключа (всі пари, які мають однаковий ключ) і виконує операцію агрегації. Наприклад, вона може підсумовувати значення, підраховувати їх кількість, знаходити максимальні або мінімальні значення, тощо [21].

Функція Reduce може приймати список значень для кожного ключа та обробляти їх, щоб отримати кінцевий результат. Наприклад, при підрахунку слів кожна пара (слово, 1) з усіх маперів буде передана на функцію Reduce, де вони будуть сумовані для отримання підсумкового рахунку (рис.3.2).

```
def reduce_function(key, values):  
    total = sum(values) # Сумуємо всі значення для одного ключа  
    return (key, total)
```

Рис. 3.2. Приклад функції Reduce

Функція Reduce отримує список значень (кількість одиниць для кожного слова) для кожного ключа (слова), підсумовує їх і повертає пару (слово, кількість).

Принципи роботи Map Reduce

- Розподіл роботи- дані, що потребують обробки, розбиваються на частини, і кожна частина передається на окремий вузол у кластері.

- Map - кожен вузол обробляє свої частини даних, застосовуючи функцію Map. Це призводить до створення проміжних результатів у вигляді пар ключ-значення.
- Shuffle and Sort - після виконання функції Map дані передаються на етап Shuffle, де пари з однаковим ключем збираються разом і сортуються.
- Reduce - для кожної групи ключів функція Reduce обробляє їх значення, зводячи результат до фінальної пари (ключ, значення).

3.3. Основні недоліки Map Reduce та можливі вирішення

1. Обмежена швидкість через I/O-операції

MapReduce значною мірою залежить від зчитування та запису даних на диск під час переходу між стадіями Map та Reduce. Це може призводити до значних затримок через високе навантаження на дискову систему.

Щоб вирішити цю проблему, можна застосувати кілька підходів, що дозволяють зменшити навантаження на диск і зменшити кількість I/O-операцій.

Використання кешування в оперативній пам'яті (In-memory caching)

Одним із способів вирішення проблеми є збереження проміжних результатів в оперативній пам'яті замість запису їх на диск. Це дозволяє скоротити час, витрачений на I/O-операції, і значно пришвидшити процес обробки даних.

Рішення за допомогою Apache Spark: Apache Spark дозволяє зберігати проміжні результати в пам'яті, що дає змогу обробляти великі обсяги даних без необхідності частих записів на диск (рис. 3.3).

```

from pyspark import SparkContext

# Створення контексту Spark
sc = SparkContext("local", "In-memory Cache Example")

# Завантаження даних
data = sc.textFile("input.txt")

# Преобразування даних (наприклад, підрахунок слів)
words = data.flatMap(lambda line: line.split())

# Кешування результатів в пам'яті
words.cache()

# Виконання підрахунку слів
word_counts = words.countByValue()

for word, count in word_counts.items():
    print(f"{word}: {count}")
sc.stop()

```

Рис. 3.3. Приклад використання Spark для кешування даних у пам'яті під час обробки

Ми використовуємо метод `cache()` для збереження даних у пам'яті. Це дозволяє уникнути додаткових операцій запису на диск між етапами Map і Reduce, що значно прискорює обробку.

Використання великих буферів для зменшення кількості записів на диск

Інший підхід полягає в тому, щоб зменшити кількість записів на диск шляхом буферизації результатів. Замість того, щоб писати дані після кожного етапу, можна обробляти та зберігати їх в великих блоках, щоб мінімізувати кількість дискових операцій.

Наприклад, якщо ми працюємо з великими текстовими файлами, можна зберігати результати етапів в пам'яті, поки не накопиться певна кількість даних, після чого записувати на диск одним великим блоком (рис. 3.4).

```

from pyspark import SparkContext

sc = SparkContext("local", "Buffered Output Example")

# Завантаження даних
data = sc.textFile("input.txt")

# Перетворення даних
processed_data = data.map(lambda x: x.upper())

# Замість запису на диск після кожного етапу, записуємо
дані після того, як вони оброблені
processed_data.saveAsTextFile("output_folder")

sc.stop()

```

Рис. 3.4. Приклад з використанням Hadoop

Використання технології Apache Flink **або** Apache Kafka для обробки поточкових даних

Якщо дані надходять потоком або обробка має бути більш реалістичною в режимі реального часу, можна використовувати технології, такі як Apache Flink або Apache Kafka. Вони дозволяють значно зменшити затримки завдяки обробці даних в пам'яті або за допомогою малих частин, що надходять безпосередньо до споживача.

Використання алгоритмів сортування, що мінімізують дискові операції

При роботі з великими наборами даних важливо правильно організувати сортування, щоб зменшити кількість записів і зчитувань з диска. Можна використовувати алгоритми, що оптимізують цей процес, такі як external merge sort, де дані сортуються по частинах, і лише після цього з'єднуються.

Паралельне оброблення з використанням кількох процесів чи потоків.

Інший підхід — це використання паралельних процесів для обробки даних. Замість того, щоб обробляти дані по черзі, можна створити кілька паралельних процесів, які працюють одночасно, зменшуючи час на I/O-операції.

У Python можна використати бібліотеки для паралельного виконання, такі як multiprocessing або concurrent.futures (рис.3.5).

```

from multiprocessing import Pool

def map_function(data):
    return data.upper()

if __name__ == '__main__':
    data = ["hello", "world", "map", "reduce"]

    with Pool(processes=4) as pool:
        result = pool.map(map_function, data)

    print(result)

```

Рис. 3.5. Приклад паралельного виконання

2. Неоптимальне використання ресурсів

MapReduce часто обробляє всі дані, навіть якщо частина з них не потрібна для отримання кінцевого результату. Це може призводити до неефективного використання процесора та пам'яті.

Фільтрація даних на етапі Map

Замість обробки всіх даних, можна додати логіку фільтрації, щоб виключати непотрібну інформацію ще на етапі Map (рис. 3.6.). Це дозволяє зменшити обсяг даних, що передається на наступні етапи.

```

from pyspark import SparkContext
sc = SparkContext("local", "Filter Example")

# Завантаження даних
data = sc.parallelize([
    ("key1", 10),
    ("key2", 20),
    ("key3", 30),
    ("key4", 40)])

# Фільтрація даних: залишаємо лише ті, де значення > 20
filtered_data = data.filter(lambda x: x[1] > 20)

# Обробка після фільтрації
result = filtered_data.mapValues(lambda x: x * 2).collect()
print(result) # Виведе [('key3', 60), ('key4', 80)]
sc.stop()

```

Рис. 3.6. Приклад фільтрації

Попередня агрегація (Combiner)

Використання Combiner дозволяє об'єднати проміжні результати на стадії Map, зменшуючи обсяг даних, що передаються на стадію Reduce (рис 3.7).

```
public class CombinerExample extends Reducer<Text,
IntWritable, Text, IntWritable> {
    @Override
    protected void reduce(Text key, Iterable<IntWritable>
values, Context context)
        throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable value : values) {
            sum += value.get();
        }
        context.write(key, new IntWritable(sum));
    }
}
```

Рис. 3.7. Combiner з використанням Hadoop MapReduce (Java)

Інтелектуальне розділення даних (Partitioning)

Правильне розділення даних між вузлами дозволяє уникнути надмірного завантаження окремих вузлів, забезпечуючи рівномірне використання ресурсів а також дозволяє розподілити обчислення між вузлами і уникнути їх перевантаження (рис. 3.8).

```
data = sc.parallelize(range(1000), numSlices=4) #
Розділення на 4 вузли
partitioned_data = data.map(lambda x: (x % 4,
x)).partitionBy(4) # Користувацьке розділення
```

Рис. 3.8. Приклад розділення в PySpark

Уникнення повторної обробки даних (Data Skipping)

Якщо вхідні дані мають інформацію, яка не змінюється, їх можна пропустити або зберігати в проміжному вигляді для подальшого використання.


```

from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("Data Skipping
Example").getOrCreate()

# Завантаження даних
data = spark.read.csv("data.csv", header=True,
inferSchema=True)

# Фільтрація без необхідності обробляти весь набір даних
filtered_data = data.filter(data["column_name"] > 50)

filtered_data.show()

```

Рис. 3.9. Приклад із використанням Spark SQL

Тут ми уникаємо обробки рядків, які не відповідають умові.

Оптимізація блокового розміру даних

Встановлення оптимального розміру блоків дозволяє зменшити накладні витрати на розподіл і передачу даних між вузлами.

У файлі `hdfs-site.xml` можна змінити розмір блоку (рис. 3.10). Оптимальний розмір блоку залежить від типу і розміру оброблюваних даних.

```

<property>
  <name>dfs.blocksize</name>
  <value>134217728</value> <!-- 128 MB -->
</property>

```

Рис. 3.10. Приклад для Hadoop

Розподілена обробка із використанням DataFrames

Spark DataFrame API дозволяє автоматично оптимізувати обробку завдяки Catalyst Optimizer, який виконує розумну фільтрацію й скорочує обсяг даних для обробки (рис. 3.11).

```

from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("Optimized DataFrame
Example").getOrCreate()

# Завантаження даних
df = spark.read.json("large_data.json")

# Фільтрація та агрегація
optimized_df = df.filter(df["value"] >
100).groupBy("key").sum("value")

optimized_df.show()

```

Рис. 3.11. Catalyst Optimizer

Зменшення кількості ключів на етапі Shuffle

Зменшення кількості унікальних ключів дозволяє скоротити час обміну даними між вузлами під час операції shuffle (рис. 3.12).

```

data = sc.parallelize([
    ("key1", 1), ("key1", 2),
    ("key2", 3), ("key2", 4)
])

# Використання попередньої агрегації для зменшення ключів
reduced_data = data.reduceByKey(lambda a, b: a +
b).collect()

print(reduced_data) # Виведе [('key1', 3), ('key2', 7)]

```

Рис. 3.12. Приклад у PySpark

Використання спеціалізованих обчислювальних платформ (Spark, Flink)

Платформи, такі як Apache Spark, автоматично зменшують накладні витрати завдяки оптимізаціям на рівні оператора, таких як pipeline execution.

Приклад pipeline execution на рис. 3.13.

```

result = (sc.textFile("data.txt")
    .map(lambda x: x.split(","))
    .filter(lambda x: int(x[1]) > 10)
    .reduceByKey(lambda a, b: a + b))

result.saveAsTextFile("output_folder")

```

Рис. 3.13. Приклад пайплайну у Spark

3. Обмежена підтримка складних обчислень

MapReduce добре підходить для простих завдань, таких як сортування або фільтрація, але зі складними обчисленнями, що потребують багатоетапного процесу, можуть виникати труднощі.

Ось кілька методів вирішення цієї проблеми:

Створення багатоступневих завдань (Chained Jobs)

Одним із рішень є створення ланцюжка завдань, де результати одного MapReduce завдання передаються як вхідні дані для наступного. Хоча це збільшує час обробки через додаткові етапи, але дозволяє розбити складне завдання на менші та легше керовані частини і крок за кроком вирішувати складніші завдання (рис. 3.14).

```
from pyspark import SparkContext

sc = SparkContext("local", "Chained Jobs Example")

# Перший етап - фільтрація даних
data = sc.textFile("input_data.txt")
filtered_data = data.filter(lambda line: "important" in
line)

# Другий етап - підрахунок частоти
words = filtered_data.flatMap(lambda line: line.split(" "))
word_counts = words.map(lambda word: (word,
1)).reduceByKey(lambda a, b: a + b)

# Збереження результату
word_counts.saveAsTextFile("output_data")

sc.stop()
```

Рис. 3.14. Chained Jobs

Використання DAG (Directed Acyclic Graph) підходу

Платформи на кшталт Apache Spark замінюють традиційний підхід MapReduce на більш потужний DAG (орієнтований ациклічний граф). Це дозволяє планувати виконання завдань більш ефективно, враховуючи залежності між етапами обчислень (рис. 3.15).

```

from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("DAG
Example").getOrCreate()

# Завантаження даних
df = spark.read.csv("data.csv", header=True,
inferSchema=True)

# Виконання складних обчислень у кілька етапів
result = (df.filter(df["value"] > 100)
          .groupBy("category")
          .agg({"value": "avg"}))
          .orderBy("avg(value)", ascending=False))

result.show()

```

Рис. 3.15. DAG на Spark з використанням DataFrame API

У цьому прикладі Spark автоматично керує залежностями між завданнями та оптимізує обчислення.

Розбиття завдань на мікросервіси або окремі компоненти

Цей метод особливо добре підходить для оптимізації обробки даних, оскільки розділення великих і складних завдань на дрібні компоненти допомагає краще використовувати ресурси, підвищити ефективність та зменшити складність коду. В основі цього підходу лежить ідея модульності, коли кожен компонент системи відповідає лише за свою окрему частину роботи, і всі компоненти взаємодіють між собою за допомогою стандартизованих протоколів.

При розбитті складного завдання на окремі компоненти, завдання стає набором незалежних сервісів, які можуть працювати паралельно. Це означає, що кожен сервіс має власну чітко визначену роль і може працювати незалежно від інших, взаємодіючи лише через чітко визначені інтерфейси, такі як REST API, gRPC або повідомлення через черги, наприклад, RabbitMQ або Apache Kafka.

Переваги підходу:

1. Модульність. Можливість розділяти складну логіку на окремі сервіси, кожен з яких легко тестувати, оновлювати та масштабувати.
2. Паралелізм. Оскільки сервіси працюють незалежно, їх можна виконувати паралельно, що значно підвищує швидкість обробки великих наборів даних.
3. Масштабованість. Кожен сервіс можна масштабувати окремо, залежно від навантаження.
4. Гнучкість. Можливість використовувати різні технології та мови програмування для різних компонентів системи.

Для демонстрації використаємо Python і Flask, а для зберігання даних Redis.

Сервіс на порту 5001 отримує сирі дані та фільтрує їх. Результати зберігаються у Redis (рис. 3.16).

```
from flask import Flask, request
import json
import redis

app = Flask(__name__)
r = redis.StrictRedis(host='localhost', port=6379, db=0)

@app.route('/filter', methods=['POST'])
def filter_data():
    data = request.json
    filtered_data = [entry for entry in data if
entry['status'] == '200'] # Фільтруємо тільки успішні
запити
    r.set('filtered_data', json.dumps(filtered_data))
    return {'status': 'Data filtered successfully'}, 200

if __name__ == '__main__':
    app.run(port=5001)
```

Рис. 3.16. Сервіс 1: Первинна обробка даних (фільтрація)

Сервіс на порту 5002 отримує відфільтровані дані з Redis, обчислює статистику та зберігає їх назад у Redis (рис. 3.17).

```
from flask import Flask
import json
import redis
from collections import Counter

app = Flask(__name__)
r = redis.StrictRedis(host='localhost', port=6379, db=0)

@app.route('/compute', methods=['GET'])
def compute_statistics():
    data = json.loads(r.get('filtered_data'))
    stats = Counter(entry['url'] for entry in data) #
    Підрахунок частоти відвідувань за URL
    r.set('statistics', json.dumps(stats))
    return {'status': 'Statistics computed successfully'},
    200

if __name__ == '__main__':
    app.run(port=5002)
```

Рис. 3.17. Сервіс 2: Обчислення статистики

Сервіс на порту 5003 отримує обчислені статистичні дані та формує підсумковий звіт. (рис.3.18)

```
from flask import Flask
import json
import redis

app = Flask(__name__)
r = redis.StrictRedis(host='localhost', port=6379, db=0)

@app.route('/aggregate', methods=['GET'])
def aggregate_results():
    stats = json.loads(r.get('statistics'))
    # Уявімо, що тут відбувається формування звіту
    report = {'total_urls': len(stats), 'top_urls':
sorted(stats.items(), key=lambda x: x[1],
reverse=True)[:5]}
    return report, 200

if __name__ == '__main__':
    app.run(port=5003)
```

Рис. 3.18. Сервіс 3: Агрегація результатів

Комунікація між сервісами відбувається через Redis або інші системи черг. Це дозволяє кожному сервісу бути незалежним та автономним, а також спрощує їх масштабування та оновлення. Якщо, наприклад, обчислення статистики вимагають більше потужності, сервіс можна просто запустити на кількох машинах паралельно.

Недоліки підходу:

- Мікросервіси додають складності в управлінні, моніторингу та відлагодженні системи.
- Часта взаємодія між сервісами може створити додаткове навантаження на мережу.
- Для управління мікросервісами потрібні додаткові інструменти на зразок Kubernetes чи Docker Compose.

Перехід на інші парадигми обробки даних

У випадках, коли MapReduce не справляється зі складними обчисленнями, доцільно перейти на інші парадигми, які надають більшу гнучкість у роботі з даними. Наприклад:

- Apache Flink — система обробки потоків, що підтримує складні аналітичні обчислення з низькою затримкою.
- Dask — Python-бібліотека для обробки великих даних, що дозволяє обчислення з використанням паралельних задач і побудови DAG
- Presto або Apache Drill — інструменти для SQL-запитів по великих наборах даних, які можуть бути більш ефективними для складних аналітичних задач.

```

import dask.dataframe as dd

# Завантаження великого CSV файлу з використанням Dask
df = dd.read_csv("large_data.csv")

# Виконання складних обчислень: фільтрація, групування та
агрегація
result = df[df["value"] >
100].groupby("category").mean().compute()

print(result)

```

Рис. 3.19. Приклад з використанням Dask

Dask дозволяє обробляти великі дані та підтримує складні обчислення з оптимізацією ресурсів (рис.3.19).

Використання спеціалізованих бібліотек для складних обчислень

У деяких випадках доцільно використовувати спеціалізовані бібліотеки, які призначені для конкретних типів завдань, наприклад, MLib у Spark для машинного навчання (рис.3.20). Це дозволяє уникнути написання складної логіки обробки та спростити процес.

```

from pyspark.ml.feature import Tokenizer
from pyspark.ml.classification import LogisticRegression
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("MLlib
Example").getOrCreate()

# Завантаження даних
data = spark.read.csv("data.csv", header=True,
inferSchema=True)

# Токенізація тексту (попередня обробка даних)
tokenizer = Tokenizer(inputCol="text", outputCol="words")
wordsData = tokenizer.transform(data)

# Побудова та навчання моделі
lr = LogisticRegression(maxIter=10, regParam=0.01)
model = lr.fit(wordsData)

# Прогнозування на нових даних
predictions = model.transform(wordsData)
predictions.show()

```

Рис. 3.20. Приклад використання MLib у Spark

4. Високе навантаження на вузли-майстри

MapReduce має центральну точку управління — вузол-майстер, який керує розподілом завдань та моніторингом процесів. Це може призводити до вузьких місць у системі, особливо при роботі з великими кластерами.

Розподіл ролі майстра

Один із підходів до вирішення проблеми — це розподіл обов'язків між кількома вузлами-майстрами, щоб кожен з них керував певною частиною обчислень. Це дозволяє розвантажити основний вузол-майстер, оскільки робоче навантаження буде розподілене між кількома вузлами. Цей підхід використовується в деяких системах, таких як Apache YARN, де роль майстра розподіляється між декількома компонентами.

Якщо розглядати варіант з Hadoop і YARN, тут вузол-майстер розподіляє завдання між різними Node Managers (рис. 3.21):

```
# Налаштування файлу yarn-site.xml для активації YARN
<property>
  <name>yarn.resourcemanager.address</name>
  <value>localhost:8032</value>
</property>
<property>
  <name>yarn.nodemanager.aux-services</name>
  <value>mapreduce_shuffle</value>
</property>
```

Рис. 3.21. Розподіл обов'язків

Це дозволяє зменшити навантаження на один центральний вузол за рахунок розподілу ресурсів між декількома вузлами в кластері.

Використання динамічних планувальників завдань

Замість того, щоб один вузол-майстер обробляв усі завдання, можна використовувати динамічні планувальники завдань, які автоматично розподіляють завдання між різними вузлами в кластері на основі поточного навантаження (рис. 3.22):. Динамічний планувальник може відстежувати стан кожного вузла і призначати завдання вузлам з найменшим навантаженням, таким чином зменшуючи ймовірність перевантаження майстра.

```

from pyspark import SparkConf, SparkContext

# Налаштування конфігурації Spark для динамічного
планування ресурсів
conf = SparkConf() \
    .setAppName("Dynamic Scheduling Example") \
    .setMaster("yarn") \
    .set("spark.dynamicAllocation.enabled", "true") \
    .set("spark.dynamicAllocation.minExecutors", "1") \
    .set("spark.dynamicAllocation.maxExecutors", "10")

sc = SparkContext(conf=conf)

# Виконання завдання з динамічним розподілом ресурсів
rdd = sc.textFile("hdfs://data/bigdatafile.txt")
result = rdd.flatMap(lambda line: line.split()) \
    .map(lambda word: (word, 1)) \
    .reduceByKey(lambda a, b: a + b)
result.saveAsTextFile("hdfs://output/wordcount")

sc.stop()

```

Рис. 3.22. Приклад реалізації динамічного планування в Spark

Динамічне планування допомагає більш ефективно використовувати ресурси кластера, автоматично додаючи або зменшуючи кількість виконавців в залежності від навантаження.

Використання систем моніторингу та автоматичного масштабування

Моніторинг кластера з використанням інструментів, таких як Prometheus або Grafana, дозволяє відслідковувати навантаження на майстра в реальному часі та налаштовувати автоматичне масштабування. Це дозволяє додавати нові вузли або змінювати конфігурацію кластера залежно від потреб (рис. 3.23):.

```

# Приклад налаштування моніторингу Spark з Prometheus
spark.executor.extraJavaOptions=-
javaagent:/path/to/jmx_prometheus_javaagent.jar=7071:/path/
to/prometheus-config.yaml
spark.driver.extraJavaOptions=-
javaagent:/path/to/jmx_prometheus_javaagent.jar=7072:/path/
to/prometheus-config.yaml

```

Рис. 3.23. Приклад моніторингу кластера з Prometheus

В результаті можна отримати деталізовану статистику роботи кластеру та розподілу навантаження між вузлами, що допоможе ефективніше керувати ресурсами.

Впровадження відмовостійкості майстра

Висока доступність та відмовостійкість вузла-майстра може бути досягнута за рахунок використання кількох резервних майстрів, які автоматично переймають управління у випадку збою основного вузла. У Hadoop можна налаштувати "high availability" для вузла-майстра (рис. 3.24):

```
<!-- hdfs-site.xml -->
<property>
  <name>dfs.nameservices</name>
  <value>mycluster</value>
</property>
<property>
  <name>dfs.ha.namenodes.mycluster</name>
  <value>nn1,nn2</value>
</property>
<property>
  <name>dfs.namenode.rpc-address.mycluster.nn1</name>
  <value>node1.example.com:8020</value>
</property>
<property>
  <name>dfs.namenode.rpc-address.mycluster.nn2</name>
  <value>node2.example.com:8020</value>
</property>
```

Рис. 3.24. Налаштування високої доступності

Налаштування високої доступності забезпечує автоматичне перемикавання на резервний вузол у випадку відмови основного.

Зменшення кількості зворотних зв'язків до вузла-майстра

Якщо в MapReduce є багато завдань, що потребують частих зворотних зв'язків до вузла-майстра, це може викликати перевантаження. Замість цього можна використовувати механізми кешування проміжних даних на локальних вузлах або розподіляти обробку на рівні самих вузлів без необхідності постійного звернення до майстра.

5. Проблеми зі збоями

Збої вузлів у кластері — це неминуча реальність у розподілених обчисленнях, таких як MapReduce. Проблеми можуть виникати через апаратні збої, мережеві проблеми або збої програмного забезпечення, і це може призводити до втрати даних або переривання обчислювальних завдань.

Реплікація даних

Один з основних підходів для забезпечення відмовостійкості — це реплікація даних. Дані, які зберігаються в кластері, копіюються на кілька вузлів. Якщо один вузол виходить з ладу, дані все ще доступні на інших вузлах. Це особливо важливо для ключових даних, необхідних для обчислень.

В Hadoop за замовчуванням дані реплікуються на три вузли. Це налаштовується через файл конфігурації `hdfs-site.xml` (рис. 3.25):

```
<!-- hdfs-site.xml -->
<property>
  <name>dfs.replication</name>
  <value>3</value>
</property>
```

Рис. 3.25. Приклад реплікації в Hadoop

Завдяки цьому, якщо один вузол виходить з ладу, інші вузли зберігають копії даних, що дозволяє продовжувати обчислення без переривання.

Автоматичне перезапуск завдань

Багато платформ MapReduce, таких як Hadoop, мають вбудовані механізми для автоматичного перезапуску завдань у разі збою. Якщо під час виконання обчислення вузол виходить з ладу, завдання може бути перенесене на інший доступний вузол і запущене знову.

Приклад механізму перезапуску завдань:

Hadoop автоматично перезапускає завдання у випадку невдачі. Кількість спроб можна налаштувати у файлі конфігурації `mapred-site.xml` (рис. 3.26):

```

<!-- mapred-site.xml -->
<property>
  <name>mapreduce.map.maxattempts</name>
  <value>4</value>
</property>
<property>
  <name>mapreduce.reduce.maxattempts</name>
  <value>4</value>
</property>

```

Рис. 3.26. Механізм перезапуску завдань

Це означає, що якщо завдання Map або Reduce не вдалося виконати, система автоматично спробує його виконати ще три рази.

Зберігання проміжних результатів на стійких носіях

Зберігання проміжних результатів на стійких носіях, таких як HDFS або інші розподілені файлові системи, допомагає уникнути втрати даних у разі збою. Навіть якщо вузол виходить з ладу, проміжні результати будуть доступні після його відновлення.

Приклад збереження проміжних даних в HDFS:

У Hadoop усі проміжні результати записуються в HDFS, що забезпечує їхню стійкість до збоїв. Це дозволяє зберігати всі результати обчислень на стійкому носії, доступному для подальшого використання навіть після відмови вузла (рис. 3.27).

```

# Збереження проміжних результатів у HDFS
hadoop jar myjob.jar MyMapReduceJob -D
mapreduce.output.fileoutputformat.outputdir=hdfs://output_d
ir

```

Рис. 3.27. Код для збереження результатів

Регулярне збереження контрольних точок (Checkpoints)

Регулярне створення контрольних точок (checkpoints) під час виконання завдання дозволяє зберігати стан обчислень на певному етапі. Якщо відбувається збій, обчислення можна продовжити з останньої контрольної точки, а не починати з нуля.

Spark підтримує механізм контрольних точок для зберігання проміжних станів обчислень. Воно дозволяє відновити обчислення з конкретного етапу після збою (рис. 3.28).

```
from pyspark import SparkConf, SparkContext

conf = SparkConf().setAppName("Checkpoint Example")
sc = SparkContext(conf=conf)

# Вказуємо директорію для контрольних точок
sc.setCheckpointDir("hdfs://checkpoints")

# Створення RDD та встановлення контрольної точки
rdd = sc.textFile("hdfs://data/bigdatafile.txt")
rdd.checkpoint()

# Виконання обчислень
result = rdd.flatMap(lambda line: line.split()) \
              .map(lambda word: (word, 1)) \
              .reduceByKey(lambda a, b: a + b)

result.saveAsTextFile("hdfs://output/result")

sc.stop()
```

Рис. 3.28. Використання контрольних точок у Spark

Моніторинг і проактивне управління ресурсами

Регулярний моніторинг стану вузлів і проактивне управління ресурсами можуть допомогти уникнути несподіваних збоїв. Наприклад, використання інструментів моніторингу, таких як Apache Ambari або Cloudera Manager, дозволяє відслідковувати навантаження на вузли та вчасно додавати ресурси, якщо система наближається до критичного стану.

Apache Ambari дозволяє налаштовувати моніторинг стану вузлів кластера, попередження про можливі збої та автоматичні заходи для виправлення ситуації. У випадку досягнення критичних навантажень на вузол-майстер або робочий вузол система може автоматично додати нові ресурси або перевести обчислення на резервні вузли.

Використання більш стійких платформ

Деякі платформи, такі як Apache Spark, мають додаткові механізми відновлення після збоїв, що дозволяють більш ефективно обробляти випадки відмов вузлів. Spark, наприклад, використовує Directed Acyclic Graph (DAG) для управління обчисленнями, що дозволяє легше відновлювати обчислення після збоїв (рис.3.29).

Spark може автоматично повторно виконати завдання після збою, використовуючи DAG для зберігання проміжного стану:

```
# Spark автоматично відновлює завдання завдяки використанню DAG
rdd = sc.textFile("hdfs://data/bigdatafile.txt")

# Виконання обчислень, навіть якщо відбувається збій,
завдання повторюється
result = rdd.flatMap(lambda line: line.split()) \
             .map(lambda word: (word, 1)) \
             .reduceByKey(lambda a, b: a + b)

result.saveAsTextFile("hdfs://output/wordcount")
```

Рис. 3.29. Відновлення обчислень у Spark

6. Складність налаштування та оптимізації

MapReduce вимагає глибокого розуміння параметрів налаштування для досягнення оптимальної продуктивності. Це може створювати труднощі для новачків і навіть для досвідчених користувачів, які не мають спеціальних знань про налаштування кластера.

Використання попередньо налаштованих платформ

Складність налаштування можна зменшити, використовуючи попередньо налаштовані платформи для роботи з MapReduce, такі як Amazon EMR (Elastic MapReduce) або Google Dataproc. Ці платформи пропонують готові до використання кластери з оптимальними налаштуваннями, що допомагає зменшити необхідність ручної оптимізації.

На Amazon EMR можна створити кластер через веб-інтерфейс без глибокого знання конфігурацій. Використання таких сервісів значно спрощує процес налаштування, знижуючи бар'єр для входу (рис. 3.30).

```
# Створення кластера EMR за допомогою AWS CLI
aws emr create-cluster --name "MyCluster" --use-default-
roles \
  --release-label emr-5.32.0 \
  --instance-count 3 \
  --instance-type m5.xlarge \
  --applications Name=Hadoop Name=Hive Name=Spark
```

Рис. 3.30. Налаштування Amazon EMR

Налаштування параметрів MapReduce

Щоб досягти кращої продуктивності, необхідно налаштувати параметри MapReduce відповідно до специфіки завдання. Це включає вибір оптимального розміру блоків даних, кількості вузлів, що виконують Map і Reduce завдання, і обсягу пам'яті, яку потрібно виділити.

У файлі конфігурації `mapred-site.xml` можна вказати оптимальні параметри:

```
<!-- mapred-site.xml -->
<property>
  <name>mapreduce.task.io.sort.mb</name>
  <value>256</value> <!-- Пам'ять для сортування -->
</property>
<property>
  <name>mapreduce.task.io.sort.factor</name>
  <value>100</value> <!-- Кількість файлів, які
об'єднуються одночасно під час злиття -->
</property>
<property>
  <name>mapreduce.job.reduces</name>
  <value>10</value> <!-- Кількість Reduce задач -->
</property>
```

Рис. 3.31. Приклад налаштування параметрів у Hadoop

Моніторинг і тюнінг ресурсів кластера

Моніторинг ресурсів, таких як процесор, пам'ять і дисковий простір, допомагає виявити "вузькі місця" в MapReduce завданнях. Інструменти

моніторингу, такі як Apache Ambari, Ganglia або Cloudera Manager, дозволяють відстежувати метрики системи та вносити корективи в налаштування.

Приклад використання моніторингу з Apache Ambari:

Apache Ambari надає інформацію про використання ресурсів у реальному часі та допомагає виявити проблеми, які впливають на продуктивність:

- Моніторинг навантаження на CPU, пам'ять, I/O.
- Виявлення вузьких місць у зберіганні та обробці даних.
- Налаштування параметрів на основі отриманих даних для покращення продуктивності.

Навчання та використання автоматизованих інструментів

Складність налаштування можна зменшити, навчившись ефективно використовувати автоматизовані інструменти для оптимізації, такі як Apache Tuning Tool або Job History Server. Ці інструменти допомагають ідентифікувати проблемні місця та пропонують варіанти налаштування параметрів.

Приклад використання автоматизованих інструментів:

- Використання Job History Server в Hadoop для аналізу продуктивності завдань.
- Налаштування параметрів на основі рекомендацій Apache Tuning Tool, що аналізує історію виконання завдань.

7. Затримки при обміні даними

Передача великих обсягів даних між етапами Map та Reduce може призводити до затримок, особливо у великих розподілених кластерах.

Паралельна передача даних (Shuffle та Sort оптимізація)

В етапі Shuffle дані передаються між вузлами для підготовки до Reduce стадії. Щоб зменшити затримки, можна використовувати паралельну передачу даних, щоб ефективніше використовувати мережеві ресурси та прискорити передачу.

Оптимізація етапу Shuffle та Sort. Налаштування параметрів Hadoop може допомогти оптимізувати етап Shuffle (рис. 3.32).

```
<!-- mapred-site.xml -->
<property>
  <name>mapreduce.reduce.shuffle.parallelcopies</name>
  <value>20</value> <!-- Кількість паралельних потоків для
Shuffle -->
</property>
<property>
  <name>mapreduce.task.io.sort.mb</name>
  <value>512</value> <!-- Розмір пам'яті для сортування -->
</property>
<property>
  <name>mapreduce.reduce.input.buffer.percent</name>
  <value>0.7</value> <!-- Розмір буфера для проміжних даних -->
</property>
```

Рис. 3.32. Паралельна передача даних

Використання комбінаторів та агрегаторів

У деяких випадках можна використовувати агрегатори, щоб зменшити кількість проміжних даних. Це дозволяє виконувати агрегацію даних безпосередньо на вузлах, на яких виконуються Map завдання.

Якщо завдання потребує підрахунку кількості записів, можна використати `reduceByKey`, щоб виконати агрегацію даних на стадії Map (рис. 3.33). Це дозволяє зменшити кількість даних, які необхідно передати на етапі Shuffle.

```
rdd = sc.textFile("hdfs://data/largefile.txt")

# Виконання агрегації на стадії Map
result = rdd.flatMap(lambda line: line.split()) \
              .map(lambda word: (word, 1)) \
              .reduceByKey(lambda a, b: a + b)

result.saveAsTextFile("hdfs://output/wordcount")
```

Рис. 3.33. Приклад `reduceByKey`

Уникнення передачі великих обсягів даних за допомогою фільтрації

Якщо частина даних не потрібна для аналізу, варто використовувати фільтрування на стадії Map, щоб мінімізувати обсяги даних, які потрібно передати до наступних етапів (рис. 3.34).

```
rdd = sc.textFile("hdfs://data/largefile.txt")

# Фільтрація непотрібних даних на стадії Map
filtered_rdd = rdd.filter(lambda line: "relevant_keyword" in
line)

# Подальша обробка
result = filtered_rdd.flatMap(lambda line: line.split()) \
    .map(lambda word: (word, 1)) \
    .reduceByKey(lambda a, b: a + b)

result.saveAsTextFile("hdfs://output/filtered_wordcount")
```

Рис. 3.34. Приклад фільтрації даних на стадії Map

Використання фільтрації допомагає зменшити обсяги даних, які потрібно обробляти та передавати.

Компресія даних під час передачі

Компресія даних перед передачею може суттєво зменшити обсяги даних, що передаються, та прискорити процес. Використання стиснення даних дозволяє ефективніше використовувати мережеву пропускну здатність.

```
<!-- mapred-site.xml -->
<property>
  <name>mapreduce.map.output.compress</name>
  <value>>true</value> <!-- Увімкнути стиснення результатів Map
-->
</property>
<property>
  <name>mapreduce.map.output.compress.codec</name>
  <value>org.apache.hadoop.io.compress.SnappyCodec</value> <!--
Використати Snappy для стиснення -->
</property>
```

Рис. 3.35. Приклад використання стиснення в Hadoop

3.4. Висновки до розділу 3

Впровадження покращень в MapReduce, значно підвищили ефективність обробки великих обсягів даних, зменшуючи затримки та оптимізуючи використання обчислювальних ресурсів. За допомогою таких підходів, як використання Combiner, вдалося суттєво знизити кількість проміжних даних,

що передаються між вузлами, а паралельне виконання завдань у стадії Shuffle зменшило загальний час обробки. Оптимізація процесу передачі даних за допомогою кешування в оперативній пам'яті та використання стиснення сприяли зниженню навантаження на мережу і дискові ресурси. Розбиття складних обчислень на окремі компоненти та використання мікросервісів дозволили ізолювати найскладніші елементи процесу, що підвищило гнучкість і адаптивність системи. Вдалося також мінімізувати втрати від можливих збоїв завдяки застосуванню резервування даних та механізмів повторного виконання завдань.

Впроваджені зміни значно знизили витрати на обчислення та підвищили стабільність роботи системи під великим навантаженням. Важливо зазначити, що ці покращення не лише підвищили продуктивність MapReduce, але й зробили його більш надійним та зручним для використання. Попри початкову складність налаштування та оптимізації, застосовані методи показали свою ефективність у різних сценаріях, що підтвердило їхню важливість для оптимізації роботи з великими обсягами даних. Загальний ефект від впроваджених покращень виявився значним, що свідчить про актуальність та необхідність продовження досліджень у цьому напрямі.

ВИСНОВКИ

У процесі роботи було проведено детальне дослідження, аналіз та експериментальна перевірка ефективності різних підходів до оптимізації алгоритмів сортування, зокрема, у середовищі великих даних. Було виявлено, що класичні алгоритми сортування мають свої обмеження, коли йдеться про обробку великих наборів даних, що зумовлює необхідність використання спеціальних методів оптимізації, таких як паралельна обробка, розподілені обчислення та ефективне управління пам'яттю.

Аналіз алгоритмів сортування показав, що незважаючи на те, що класичні підходи, такі як швидке сортування чи сортування злиттям, мають високу ефективність для середніх розмірів даних, їх продуктивність суттєво знижується при масштабуванні до великих обсягів даних. Важливою частиною дослідження стало порівняння алгоритмів з різною теоретичною та фактичною складністю, де враховувалася не лише тимчасова складність, але й використання оперативної пам'яті та вплив дискових операцій.

Однією з ключових частин роботи було вивчення підходів до оптимізації сортування великих наборів даних. Використання паралельного та розподіленого обчислення дозволило значно зменшити час обробки великих наборів завдяки одночасній обробці кількох частин даних на різних вузлах. Зокрема, технологія MapReduce була обрана як основний інструмент для реалізації розподілених обчислень. Вона продемонструвала свою ефективність для сортування великих обсягів даних завдяки простоті реалізації та можливості масштабування. Зокрема, було показано, як за допомогою MapReduce можна реалізувати сортування даних у великому кластері, розподіливши завдання між різними вузлами та об'єднуючи результати обробки.

У ході експериментів з використанням MapReduce вдалося визначити ряд проблем, таких як обмежена швидкість через I/O-операції, неоптимальне використання ресурсів та високе навантаження на вузли-майстри. Кожна з цих проблем потребувала відповідного методу вирішення, який був детально

розглянутий та протестований в роботі. Одним із рішень стало використання кешування в оперативній пам'яті для зменшення кількості дискових операцій та покращення продуктивності обробки даних.

Результати практичної частини роботи підтвердили ефективність розглянутих оптимізаційних підходів. Використання таких методів, як розбиття завдань на мікросервіси, введення проміжного кешування, ефективне управління ресурсами та балансування навантаження, дозволили підвищити швидкість та стабільність сортування великих наборів даних. При цьому особливу увагу було приділено вирішенню проблем підтримки масштабованості та толерантності до збоїв.

На основі проведених досліджень та експериментів можна зробити висновок, що оптимізація алгоритмів сортування для великих наборів даних потребує комплексного підходу, де враховуються особливості архітектури апаратного забезпечення, обсяги даних, можливості розподілених обчислень та обмеження пам'яті. Запропоновані рішення можуть бути корисними для розробників, які працюють з великими даними, та аналітиків, що шукають ефективні підходи до обробки інформації у великих масштабах.

Важливо також зазначити, що хоча MapReduce продемонструвала свою ефективність для ряду задач, для більш складних обчислень може знадобитися комбінування з іншими інструментами, такими як Apache Spark або більш сучасні платформи. Отримані результати також відкривають перспективи для подальших досліджень в області оптимізації алгоритмів сортування, зокрема в напрямку використання штучного інтелекту для автоматичної адаптації алгоритмів під конкретні завдання та структуру даних.

Оптимізація алгоритмів сортування для великих наборів даних — це не просто вибір правильного алгоритму, а системний підхід, який включає аналіз даних, вибір оптимальної архітектури обчислень, правильне використання ресурсів та гнучке налаштування параметрів під конкретні потреби.

СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. Apache Hadoop GIT HUB — URL: <https://github.com/apache/hadoop/tree/trunk/hadoop-mapreduce-project> (дата звернення: 30.08.2024)
2. АЛГОРИТМИ ТА СТРУКТУРИ ДАНИХ— URL: <https://ela.kpi.ua/server/api/core/bitstreams/0db974f9-16fa-459c-9f19-fab0021222ed/content> (дата звернення: 03.09.2024)
3. Алгоритми сортування в теорії та на практиці— URL: [https://javarush.com/ua/groups/posts/uk.1997.algoritmi-sortuvannja-v-teor-ta-na-practic](https://javarush.com/ua/groups/posts/uk.1997.algoritmi-sortuvannja-v-teor-ta-na-praktic) (дата звернення: 03.09.2024)
4. Sorting algorithms: Slower to faster— URL: <https://builtin.com/machine-learning/fastest-sorting-algorithm> (дата звернення: 04.09.2024)
5. Sorting algorithms — URL: <https://www.geeksforgeeks.org/sorting-algorithms/> (дата звернення: 04.09.2024)
6. Understanding sorting algorithms — URL: <https://medium.com/@noransaber685/understanding-sorting-algorithms-5575d52f5a18> (дата звернення: 04.09.2024)
7. Т. Б. Мартинюк Б. І. Круківський / КЛАСИФІКАЦІЙНИЙ АНАЛІЗ МЕТОДІВ СОРТУВАННЯ 2023 р.— URL: https://www.researchgate.net/publication/372672191_Classification_Analysis_of_Sorting_Methods (дата звернення: 06.09.2024)
8. All types of sorting algorithms — URL: <https://www.wscubetech.com/resources/dsa/sorting-algorithms> (дата звернення: 07.09.2024)
9. Інформаційне забезпечення роботехнічних систем — URL: <https://ela.kpi.ua/server/api/core/bitstreams/92cc8880-a066-40b9-ad77-9e576ac8baf0/content> (дата звернення: 07.09.2024)
10. Сучасні ПЕОМ; Апаратне та програмне забезпечення— URL:

<http://lib.kart.edu.ua/bitstream/123456789/6840/1/%D0%9A%D0%BE%D0%BD%D1%81%D0%BF%D0%B5%D0%BA%D1%82%20%D0%BB%D0%B5%D0%BA%D1%86%D1%96%D0%B9.pdf> (дата звернення: 09.09.2024)

11. Операційні системи — URL:

https://er.chdtu.edu.ua/bitstream/ChSTU/1041/1/%D0%9E%D0%9F%D0%95%D0%A0%D0%90%D0%A6%D0%86%D0%99%D0%9D%D0%86%20%D0%A1%D0%98%D0%A1%D0%A2%D0%95%D0%9C%D0%98_%D0%BD%D0%B0%D0%B2%D1%87.%D0%BF%D0%BE%D1%81..pdf (дата звернення: 09.09.2024)

12. Sorting algorithms animatio — URL:

<https://www.toptal.com/developers/sorting-algorithms> (дата звернення: 12.09.2024)

13. [Introduction to Algorithms” / Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest and Clifford Stein / 2018 р. (дата звернення: 13.09.2024)

14. [Електронний ресурс] — URL:

<https://www.oreilly.com/library/view/algorithms-in-a/9781491912973/ch04.html> (дата звернення: 13.09.2024)

15. [Sort and Search Algorithms / Yu Zhang and Mathias Funk / 2020 р.

16. Sorting algorithms — URL:

<https://www.geeksforgeeks.org/analysis-of-different-sorting-techniques/>
(дата звернення: 14.09.2024)

17. Sorting algorithms compared— URL:

https://adacomputerscience.org/concepts/sort_sorting_compared (дата звернення: 14.09.2024)

18. Sorting algorithms: an efficiency comparison— URL:

<https://medium.com/@lazaro.exe/sorting-algorithms-an-efficiency-comparison-c0f82fcd5db5> (дата звернення: 20.09.2024)

19. [Електронний ресурс] — URL:

<https://builtin.com/machine-learning/fastest-sorting-algorithm>
(дата звернення: 21.09.2024)

20. Sorting methods — URL:

<https://www.cprogramming.com/tutorial/computersciencetheory/sortcomp.html> (дата звернення: 22.09.2024)

21. Sorting algorithms comparison — URL:

<https://medium.com/@tssovi/comparison-of-sorting-algorithms-298fdf037c8f> (дата звернення: 26.09.2024)

22. Положення про кваліфікаційні роботи здобувачів вищої освіти Національного авіаційного університету. СМЯ НАУ П 03.01(10) – 03 – 2024.