

**Московский авиационный институт
(национальный исследовательский университет)**

**Факультет информационных технологий и прикладной
математики**

Кафедра вычислительной математики и программирования

Лабораторная работа №7 по курсу «Информационный поиск»

Студент: В. М. Филиппов
Преподаватель: А. А. Кухтичев
Группа: М8О-410Б
Дата:
Оценка:
Подпись:

Москва, 2025

Лабораторная работа №7 «Сжатие»

В этом задании необходимо применить алгоритмы сжатия к индексу. Исследовать изменения в размерах частей индекса, влияние на скорость индексации и поиска. В отчёте нужно указать:

- Выбранный метод сжатия. Привести побитовую схему хранения данных в индексе.
- Описать причины, по которым был выбран именно этот метод сжатия. Влияние сжатия на размер и скорость прохождения индексу.
- Обосновать, почему поиск после внедрения сжатия работает корректно. Как производилось тестирование?

1 Описание

При выполнении данной лабораторной работы я реализовал два метода сжатия

1. Delta encoding
2. Variable-Length Quantity (VLQ) / Varint

1 Delta encoding

В инвертированных индексах ID документов всегда идут по возрастанию (например: 100, 105, 110...). Если хранить их напрямую, они занимают много места. Если хранить разницы (100, 5, 5...), значения становятся маленькими. В сочетании со вторым методом сильно экономит нам память при хранении индекса.

2 Varint

Суть метода заключается в том, чтобы использовать 7 бит каждого байта для хранения данных, а старший бит для указания, есть ли еще продолжение у этого числа.

- Если старший бит равен 1, значит, следующий байт тоже относится к этому числу.
- Если старший бит равен 0, значит, этот байт последний.

При таком кодировании небольшие числа будут занимать меньше места. Также кодирование и декодирование происходит очень быстро (всего несколько битовых сдвигов). Однако если числа всегда большие (например, больше 2^{28}), такой формат может занять даже больше места (5 байт вместо 4), так как на каждые 7 бит данных добавляется 1 служебный бит.

2 Исходный код

Теперь при дампинге или чтении в функциях присутствует флаг zip, который указывает на то, каким образом нам нужно читать или записывать данные.

Реализованы три служебные функции writeVarInt, getVarIntSize, readVarInt. Первая нужна, чтобы записать число как varint, второе чтобы понять, сколько байт это число будет занимать, а третья, чтобы прочитать varint число из памяти. Вторая функция необходима для того, чтобы мы могли посчитать корректно смещение списков вхождения.

Корпус документов из 305135 страниц с размером текста 2661.05 Мб сдамплен за 14.1937 секунд, а итоговый размер получился 376500489 байт (385.9 МБ), что в 3,16 раза меньше, чем несжатый индекс. Это очень существенная экономия памяти. Кроме того, внедрение сжатия увеличило скорость доступа к данным, поскольку кэширование стало более эффективным, в оперативную память не загружается сразу очень много чисел при чтении. Кроме того, сокращаются операции работы с диском). Прочитать с устройства 386 МБ быстрее, чем 1.2 ГБ, даже с учетом времени на распаковку процессором.

Но есть и минусы, например

- Распаковка Varint требует арифметических операций и проверок битов, что даёт дополнительную нагрузку на CPU.
- Невозможность произвольного доступа к документу в списке вхождений.

1 Примеры запросов

Кстати, очень заметно, что поиск выполняется немного дольше, если делать запросы "на холодную". То есть, когда блоки памяти не закэшированы процессором, поиск работает медленнее.

```
1 Enter query: (manchester & united) | messi
2 Enter query: (NBA | NFL | NHL) & champion
3 Tokens searching: [nba, nfl, nhl, champion, ]
4 Top 10 results
5 [1] https://en.wikipedia.org/wiki/
    Major_professional_sports_leagues_in_the_United_States_and_Canada TF-IDF: 87.9164
6 [2] https://en.wikipedia.org/wiki/List_of_sports_figures_considered_the_greatest TF-
    IDF: 78.4068
7 [3] https://en.wikipedia.org/wiki/List_of_Jews_in_sports TF-IDF: 78.1504
8 [4] https://en.wikipedia.org/wiki/Draft_blunder TF-IDF: 73.9464
9 [5] https://en.wikipedia.org/wiki/Draft_steal TF-IDF: 73.9464
10 [6] https://en.wikipedia.org/wiki/Draft_(sports) TF-IDF: 73.9464
11 [7] https://en.wikipedia.org/wiki/Sports_broadcasting_contracts_in_the_United_States
    TF-IDF: 69.7583
```

```

12 [8] https://en.wikipedia.org/wiki/History\_of\_the\_National\_Hockey\_League\_on\_television TF-IDF: 67.8845
13 [9] https://en.wikipedia.org/wiki/History\_of\_ESPN\_on\_ABC TF-IDF: 67.1572
14 [10] https://en.wikipedia.org/wiki/List\_of\_multi-sport\_athletes TF-IDF: 66.54
15 Query time: 0.0131194 sec
16 Number of results: 2949 items

1 Enter query: (NBA | NFL | NHL) & champion & !football
2 Tokens searching: [nba, nfl, nhl, champion, footbal, ]
3 Top 10 results
4 [1] https://en.wikipedia.org/wiki/NBA\_regular\_season\_records TF-IDF: 53.3538
5 [2] https://en.wikipedia.org/wiki/List\_of\_teams\_that\_have\_overcome\_380931\_series\_deficits TF-IDF: 46.7721
6 [3] https://en.wikipedia.org/wiki/Los\_Angeles\_Lakers TF-IDF: 46.1695
7 [4] https://en.wikipedia.org/wiki/Stanley\_Cup\_playoffs TF-IDF: 43.3573
8 [5] https://en.wikipedia.org/wiki/Atlanta\_Hawks TF-IDF: 42.7092
9 [6] https://en.wikipedia.org/wiki/Presidents%27\_Trophy TF-IDF: 42.547
10 [7] https://en.wikipedia.org/wiki/Octagon\_\(mixed\_martial\_arts\) TF-IDF: 40.8788
11 [8] https://en.wikipedia.org/wiki/Ultimate\_Fighting\_Championship TF-IDF: 40.8788
12 [9] https://en.wikipedia.org/wiki/Edmonton\_Oilers TF-IDF: 40.751
13 [10] https://en.wikipedia.org/wiki/List\_of\_Dancing\_with\_the\_Stars\_\(American\_TV\_series\)\_competitors TF-IDF: 40.6825
14 Query time: 0.00811675 sec
15 Number of results: 1036 items

1 Enter query: (premier league) & !(championship) & !(league 1)
2 Tokens searching: [premi, leagu, championship, leagu, 1, ]
3 Top 10 results
4 [1] theguardian.com/football/2025/dec/17/premier-league-teams-africa-cup-of-nations-sunderland-morocco-chelsea-arsenal-aston-villa TF-IDF: 27.3792
5 [2] theguardian.com/football/2023/nov/17/everton-deducted-10-points-premier-league-guilty-financial-fair-play-breach TF-IDF: 26.753
6 [3] theguardian.com/football/2024/mar/11/premier-league-faces-backlash-after-failing-to-agree-financial-deal-with-epl TF-IDF: 26.4991
7 [4] https://en.wikipedia.org/wiki/List\_of\_New\_York\_metropolitan\_area\_sports\_teams TF-IDF: 26.2424
8 [5] theguardian.com/football/premier-league-2012-13/2013/may/20/all TF-IDF: 25.9398
9 [6] theguardian.com/football/2017/apr/05/premier-league-elite-player-performance-plan TF-IDF: 25.7756
10 [7] theguardian.com/football/2025/aug/18/premier-league-kick-it-out-funding-deal-reduced-concerns-charity-income TF-IDF: 25.7756
11 [8] theguardian.com/football/2023/dec/21/european-super-league-boost-court-of-justice-ruling-uefa-fifa-eu-law TF-IDF: 25.3414
12 [9] theguardian.com/football/2018/apr/19/santi-cazorla-arsenal-career-not-over-says-hopeful-arsene-wenger TF-IDF: 24.9507
13 [10] theguardian.com/football/blog/2025/jun/10/ange-postecoglou-tottenham-sacked-europa-league-post-truth-philosophy TF-IDF: 24.9428
14 Query time: 0.0247088 sec
15 Number of results: 711 items

```

```
1 | Tokens searching: [a, th, ]
2 | Top 10 results
3 | [1] https://en.wikipedia.org/wiki/Wikipedia\_talk:Manual\_of\_Style/Biography/2021\_archive TF-IDF: 0.224592
4 | [2] https://en.wikipedia.org/wiki/Tennis\_performance\_timeline\_comparison\_\(women\)\_1884%E2%80%931977 TF-IDF: 0.221983
5 | [3] https://en.wikipedia.org/wiki/Glossary\_of\_baseball\_terms TF-IDF: 0.21911
6 | [4] https://en.wikipedia.org/wiki/Baseball\_Annie TF-IDF: 0.21911
7 | [5] https://en.wikipedia.org/wiki/Glossary\_of\_nautical\_terms\_\(A%E2%80%93L\) TF-IDF: 0.204274
8 | [6] https://en.wikipedia.org/wiki/Wikipedia\_talk:Manual\_of\_Style/Dates\_and\_numbers/Archive\_160 TF-IDF: 0.203136
9 | [7] https://en.wikipedia.org/wiki/List\_of\_Super\_Bowl\_commercials TF-IDF: 0.201375
10 | [8] https://en.wikipedia.org/wiki/Glossary\_of\_cue\_sports\_terms TF-IDF: 0.200915
11 | [9] https://en.wikipedia.org/wiki/Tennis\_performance\_timeline\_comparison\_\(men\) TF-IDF: 0.20053
12 | [10] https://en.wikipedia.org/wiki/List\_of\_attacks\_related\_to\_secondary\_schools TF-IDF: 0.200482
13 | Query time: 0.082748 sec
14 | Number of results: 305064 items
```

3 Выводы

Во время выполнения данной лабораторной работы я узнал про способы сжатия индексов Varint и Delta Encoding. Важно, что эти способы сжатия используются не только при компрессии поискового индекса, но и при разработке БД, а также например в формате protobuf. Под каждым из этих алгоритмов сжатия лежит простая идея, однако в этой простоте и находится элегантность и эффективность решения.

Список литературы

- [1] Маннинг, Рагхаван, Шютце *Введение в информационный поиск* — Издательский дом «Вильямс», 2011. Перевод с английского: доктор физ.-мат. наук Д. А. Клюшина — 528 с. (ISBN 978-5-8459-1623-4 (рус.))
- [2] Список использованных источников оформлять нужно по ГОСТ Р 7.05-2008