

Способы сжатия данных в PostgreSQL

Жилин Михаил

PGMeetup.DBA

 PostgresPro

PGMeetup.DBA

О себе:

Telegram: [@mizhka](https://t.me/mizhka)
E-mail: mizhka@gmail.com



PostgresPro

- Физтех
- 14 лет работы с системной производительностью
- 2 года в компании Postgres Professional
- Контрибьютор в open source проекты
- Счастливый муж и отец
- Волейбол, дайвинг, путешествия

PGMeetup.DBA

PostgresPro

План разговора

- Полезно ли сжатие для баз данных?
- Сжатие бывает разным
- Встроенные механизмы и не только
- Примеры и рекомендации

PGMeetup.DBA

PostgresPro

Полезно ли сжатие для баз данных?



PGMeetup.DBA

Полезно ли сжатие для баз данных?

Достоинства

- Экономия дискового пространства
- Экономия оперативной памяти
- Меньше нагрузка на диски

PGMeetup.DBA

Полезно ли сжатие для баз данных?

Достоинства

- Увеличение пропускной способности системы (больше tps)
- Улучшается время ответа (меньше мс)

Полезно ли сжатие для баз данных?

Недостатки

Overhead	Shared Object	Symbol
63.82%	postgres	[.] pglz_decompress
16.04%	libc-2.31.so	[.] __memcpy_avx_unaligned_erms
1.99%	postgres	[.] LWLockAttemptLock
0.89%	postgres	[.] hash_search_with_hash_value
0.89%	postgres	[.] LWLockRelease
0.88%	postgres	[.] _bt_compare

Overhead	Shared Object	Symbol
46.70%	liblz4.so.1.9.3	[.] LZ4_decompress_safe
5.81%	postgres	[.] LWLockAttemptLock
2.59%	postgres	[.] LWLockRelease
2.19%	postgres	[.] _bt_compare
2.14%	postgres	[.] hash_search_with_hash_value

PGLZ

LZ4

PGMeetup.DBA

PostgresPro

Сжатие бывает разным



PGMeetup.DBA

Сжатие бывает разным

Особенности

- Без потерь
- Максимально компактное
- За максимально короткое время
- Текстовые и бинарные данные
- А что такое сжатие?

PGMeetup.DBA

PostgresPro

Сжатие бывает разным

Зависимости

- 1, 2, 3, 4, 5, 6, 7, 8, 9...
- 1, 2, 3, 5, 7, 11, 13, 17, 19, 23...
- 53363, 40550, 33191...

PGMeetup.DBA

PostgresPro

Сжатие бывает разным

Теория информации

- Алгоритмическая теория информации
Андрей Николаевич Колмогоров (1960-ые)
 - <https://nautil.us/kolmogorov-complexity-and-our-search-for-meaning-237158/>
 - Минимальная длина программы, которая выводит необходимый набор символов (данных)
- Парадокс неизмеримости: **всегда может существовать алгоритм с меньшей длиной программы**

PGMeetup.DBA

Сжатие бывает разным

Энтропийное сжатие

- Шеннон-Фано (1948-1949)
 - Вероятность символов
 - Чем чаще символ, тем меньше бит на символ
- Хаффман (1954)
 - Бинарное дерево частот и префиксы
 - Оптимальный набор

Сжатие бывает разным

Сжатие по словарю

- Lempel-Ziv-Welch (1978-1984)
 - рекурсивная копия, повторяемость последовательностей символов
 - LZ1 (LZ77) / LZ2 (LZ78) / LZW
 - gif
- Deflate (1991) с добавлением энтропийного сжатия
 - png, zip, gzip and others

Сжатие бывает разным

Соревнование на скорость

- LZO (1996) - Markus F.X.J. Oberhumer
 - быстрая распаковка: read-only filesystems
 - <https://github.com/markus-oberhumer>
- LZ4 (2011) + Zstandard (2016) - Yann Collet
 - tradeoff и оптимизации под 32-/64-бит архитектуры
 - <https://github.com/Cyan4973>

PGMeetup.DBA

Сжатие бывает разным

Современные реалии

- Смесь из энтропийного сжатия и сжатия по словарю
- Степень сжатия: от 1 до N (размер плавающего окна)
- Неэффективность на маленьких объемах
- LZ - словарь находится в сжатом поток

Сжатие бывает разным

PGLZ internals

00001810	c1	40	00	00	00	00	00	00	40	1f	00	00	b8	18	00	00	.@.....@.....	toast_compress_header→tc_i
00001820	00	42	75	6d	70	73	20	5b	65	00	76	65	6e	74	73	6f	.Bumps [e.ventso	
00001830	75	72	00	63	65	5d	28	68	74	74	70	00	73	3a	2f	2f	ur.ce](http.s://	
00001840	67	69	74	68	00	75	62	2e	63	6f	6d	2f	45	15	01	20	gith.ub.com/E..	ctl:15=10101, len=4, off=32
00001850	53	02	20	2f	08	2c	29	20	66	00	72	6f	6d	20	31	2e	S. /.,) f.rom 1.	
00001860	30	2e	80	37	20	74	6f	20	31	2e	01	02	00	0a	3c	64	0..7 to 1.....<d	
00001870	65	74	61	69	6c	00	73	3e	0a	3c	73	75	6d	6d	00	61	etail.s>.<summ.a	
00001880	72	79	3e	43	68	61	6e	80	67	65	6c	6f	67	3c	2f	05	ry>Chan.ge log</.	
00001890	13	00	0a	3c	70	3e	3c	65	6d	3e	05	03	57	64	03	4b	...<p>..Wd.K	
000018a0	3c	61	20	68	72	10	65	66	3d	22	0f	85	18	2f	62	6c	<a href=".../bl	
000018b0	00	6f	62	2f	6d	61	73	74	65	00	72	2f	48	49	53	54	.ob/maste.r/HIST	

PGMeetup.DBA

Сжатие бывает разным

LZ4 vs PGLZ

- Нет побитовых операций
- Копирование по несколько байт вместо по одному

Сжатие бывает разным

LZ4 vs PGLZ

Overhead	Shared Object	Symbol
63.82%	postgres	[.] pglz_decompress
16.04%	libc-2.31.so	[.] __memcpy_avx_unaligned_erms
1.99%	postgres	[.] LWLockAttemptLock
0.89%	postgres	[.] hash_search_with_hash_value
0.89%	postgres	[.] LWLockRelease
0.88%	postgres	[.] _bt_compare

Overhead	Shared Object	Symbol
46.70%	liblz4.so.1.9.3	[.] LZ4_decompress_safe
5.81%	postgres	[.] LWLockAttemptLock
2.59%	postgres	[.] LWLockRelease
2.19%	postgres	[.] _bt_compare
2.14%	postgres	[.] hash_search_with_hash_value

PGLZ

LZ4

PGMeetup.DBA

Сжатие бывает разным

zstd vs LZ4

- Добавлено энтропийное сжатие (Huffman + FSE)
- Многопоточный
- Отрицательные уровни сжатия (-7 - самый быстрый)
- Колоссальные размеры окон (до 128MiB)

PGMeetup.DBA

PostgresPro

Сжатие бывает разным

Benchmarks

Алгоритм	Описание	Сжатие	Распаковка	Степень
deflate	добрый старый	5-100MiBps	10-200MiBps	2.8
lzo	быстрая распаковка	8MiBps	850MiBps	2.8
lz4	самый быстрый	780MiBps	4500MiBps	2.1
zstd	лучший баланс	480MiBps	1200MiBps	2.8

PGMeetup.DBA

PostgresPro

Встроенные механизмы сжатия данных



PGMeetup.DBA

Встроенное сжатие: inline и TOAST

Для строк/значений
больше чем 2килобайта

- короткие в таблице
- длинные в TOAST
- PGLZ
- LZ4 (с 14ой версии)

PGMeetup.DBA

Встроенное сжатие: Index Key
Dedup

Группировка TID-ов в
Posting List

- PostgreSQL 12+
- PostgresPro 10+

PGMeetup.DBA

Встроенное сжатие: WAL Full Page Image

Сжатие блоков в WAL
файлах

- PGLZ
- LZ4 (с 15ой версии)

PGMeetup.DBA

Дополнительно...

Колоночное хранилище

- GreenPlum & ZedStore (fork)
 - Citus Columnar & cstore_fdw (extension)
 - Append-only optimizations
 - lz4, zstd, zlib, rle
- Нет индексного сжатия
 - Ряд ограничений

PGMeetup.DBA

Дополнительно...

Сжатие на уровне OS

- OpenZFS
- Прозрачно для базы данных
- lz4, zstd, snapshot-ы

- Пределы масштабирования
- Требуется опыта в настройках

Может быть завтра...

- Быстрого TOAST-а
- Дедупликации значений в таблицах
- Сжатие ключей в индексах

Блочное сжатие: PostgresPro CFS

- Покрывает все типы данных (индексы, таблицы)
- Прозрачно для пользователей
- Простота конфигурации
 - задаётся на уровне tablespaces-ов
- lz4, zstd, zlib, pglz, levels....
- Входит в PgPro Enterprise 9.6+

PGMeetup.DBA

PostgresPro

Типовые сценарии использования сжатия



PGMeetup.DBA

Типовые сценарии использования сжатия

Современные базы
небольшого размера

- < 500GiB
- < 500tps
- Достаточное количество RAM для кэширования
- SSD/NVMe

Встроенного сжатия, скорее всего, достаточно

PGMeetup.DBA

PostgresPro

Типовые сценарии использования сжатия

Современные базы
большого размера

- > 1TiB

CFS или ZFS позволит сохранить терабайты

PGMeetup.DBA

PostgresPro

Типовые сценарии использования сжатия

Медленные диски и/или мало
RAM

- HDD / NAS
- Всего 2-4GiB RAM

CFS или ZFS могут значительно ускорить работу PostgreSQL

PGMeetup.DBA

PostgresPro

Типовые сценарии использования сжатия

Файлохранилка (PDF,
photos)

- Хранится в TOAST-е
- Медленно, но надёжно
- По копии на каждой реплике и в бекапах

**Выносите хранение сжатых файлов из PostgreSQL
на внешнее хранилище а-ля S3**

PGMeetup.DBA

Типовые сценарии использования сжатия

Аналитические базы
данных, витрины

- Колонки без индексов
- Низкоэнтропийные данные

Попробуйте колоночное хранилище

PGMeetup.DBA

Типовые сценарии использования сжатия

Зашифрованные данные

- Энтропия
- Нет встроенного шифрования
- Не включайте сжатие TOAST для зашифрованных данных

Сначала сжимать, потом шифровать

PGMeetup.DBA

Итого:

Use compression, Luke!

- Вопросы?

PGMeetup.DBA

СПАСИБО!

postgrespro.ru