

Онлайн-дедупликация для баз данных

Лянхун Сюй
Университет Карнеги-Меллона
lianghon@andrew.cmu.edu

Эндрю Павло
Университет Карнеги-Меллона
pavlo@cs.cmu.edu

Судипта Сенгупта
Исследования Майкрософт
sudipta@microsoft.com

Грегори Р. Гангер
Университет Карнеги-Меллон
ganger@ece.cmu.edu

АБСТРАКТ

dbDedup — это схема дедупликации на основе подбоя для систем управления базами данных в режиме онлайн (СУБД). Помимо сжатия на уровне блоков отдельных страниц баз данных или сообщений журнала операций (oplog), используя уменьшение размера данных, хранящихся на диске, по сравнению с тем, что обеспечивается традиционными схемами сжатия, и (2) уменьшение объема данных, передаваемых по сети для служб репликации. Читая о нашей работе, мы выделили dbDedup в распределенную СУБД NoSQL и проанализировали ее свойства с использованием чужих реальных наборов данных. Наши результаты показывают, что dbDedup обеспечивает до 37-кратного сокращения размера хранилища и трафика репликации баз данных самостоятельно, но до 61-кратного сокращения в сочетании с блочным сжатием СУБД. dbDedup обеспечивает оба преимущества с незначительным влиянием на пропускную способность СУБД или задержку клиента (средние конечные).

1. ВВЕДЕНИЕ

Темп роста данных превышает стоимость оборудования. Сжатие баз данных является одним из решений этой проблемы. Для хранения баз данных, в дополнение к экономии места, сжатие помогает сократить количество дисковых операций ввода-вывода и повысить производительность, поскольку запрашиваемые данные помещаются на меньшее количество страниц. Для распределенных баз данных, реплицированных по географическим регионам, также существует потребность в сокращении объема переданных данных, используемого для синхронизации реплик. Наиболее широко используемый подход к сокращению данных в операционных СУБД — это сжатие на уровне блоков [30, 37, 46, 43, 3, 16]. Такие СУБД используют для поддержки полнотекстовых приложений, которые выполняют простые запросы для небольшого количества операций записи за раз (в отличие от выполнения сложных запросов, которые сканируют большие сегменты баз данных). Хотя сжатие на уровне блоков является простым и эффективным, оно не решает проблему избыточности между блоками, следовательно, остается значительная возможность для улучшения для многих приложений (например, из-за управления версиями на уровне приложений в вики или илчестного копирования записей на дисках объектных). Дедупликация (dedup) стала популярной в системах резервного копирования для устранения дублирующегося контента по

Разрешение на создание цифровых или печатных копий всей или части работ для личного или учебного использования предоставляется бесплатно при условии, что копии не будут сделаны или распространены с целью получения прибыли или коммерческой выгоды и что копии будут содержать это уведомление и полную ссылку на первоисточник. Авторские права на компоненты этой работы принадлежат другим лицам, помимо ACM, должны соблюдаться. Реферирование с указанием источника разрешено. Для копирования иным способом или повторной публикации, размещенной на серверах или распространяемой по электронной почте, требуется предварительное разрешение и/или плата. Запросите разрешения по адресу permissions@acm.org.

SIGMOD'17, 14–19 мая 2017 г., Роли, Северная Каролина, США ©
2017 ACM. ISBN 978-1-4503-4197-1/17.04.05... 15,00 долларов США
DOI: <http://dx.doi.org/10.1145/3035918.3035938>

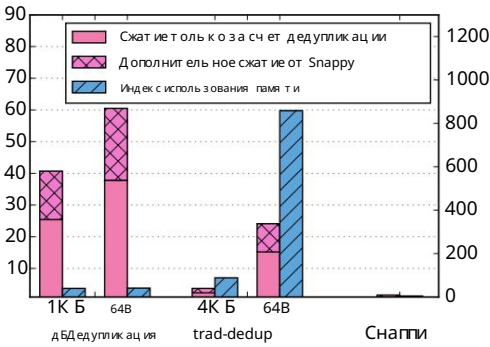


Рисунок 1: Коэффициент сжатия и использование памяти и индекс для данных Википедии, хранящихся в различных конфигурациях MongoDB: с dbDedup (размер фрагмента 1 КБ и 64 Б), с традиционной дедупликацией (4 КБ и 64 Б) и с Snap (сжатие на уровне блоков). dbDedup обеспечивает более высокую степень сжатия и меньшие накладные расходы на память с индексом, чем традиционная дедупликация. Snap обеспечивает такое же сжатие 1,6х для данных после дедупликации или исходных данных.

весь корпус данных, частота доступа гораздо выше, чем коэффициент сжатия. Резервный поток делится на фрагменты и в качестве идентификатора каждого фрагмента используется устойчивый к коллизиям хэш (например, SHA-1). Система дедупликации поддерживает глобальный индекс всех хэшей и использует его для обнаружения дубликатов. Дедупликация хорошо работает как для основных, так и для резервных наборов данных хранения, которые состоят из больших файлов, которые редко изменяются (если изменяются, то изменения редки).

К сожалению, традиционные схемы дедупликации на основе фрагментов не подходят для операционных СУБД, где приложения выполняют запросы на обновление, изменяющие отдельные записи. Количество дублирующихся данных в отдельных записях, скорее всего, незначительно. Но большие размеры фрагментов (например, 4–8 КБ) являются нормой, что обычно бежит огромных индексов в памяти и больше операций чтения с диска.

В этой статье представлено dbDedup, облегченная схема для систем баз данных в режиме онлайн, которая использует дедупликацию на основе сходства [65] для сжатия отдельных записей. Вместо индексации каждого хэша фрагмента, dbDedup выбирает небольшой поднабор хэшей фрагментов для каждой новой записи баз данных, а затем использует этот образец для идентификации похожих записей в базе данных. Затем он использует дельта-сжатие на уровне байтов для двух записей, чтобы уменьшить как используемое онлайн-хранилище, так и пропускную способность удаленной репликации. dbDedup обеспечивает более высокие коэффициенты сжатия с меньшими накладными расходами на память, чем дедупликация на основе фрагментов, и хорошо сочетается со сжатием на уровне блоков, как показано на рис. 1.

Мы представляем и описываем несколько методов для достижения этого эффекта. Прежде всего, мы представляем новое двукратное кодирование для эффективной передаточной скорости новых записей (прямое кодирование) в удаленные реплики, сохраняя новые записи с закодированными формами выбранных исходных записей (обратное кодирование). В результате нет де-

код т ребует ся для распрост раненного случ ая доста па к самой послед ней записи в ц епоч ке к одирования (например к послед ней версии Википедии).

Ч тобы из беж ать наклады расходов на произ вод ите ль ность при обновлении исходных записей, мы так же вводим кэш обрат ной записи с пот еря ми, наст роенный на макс имиз ац ию коэффи циент а сжат ия , избег ая приэ том к онк уренц ии ввода-вывода. Наш подход так же использ ует новую техн ику, называемую од ированием ск ач ко в, к от орая минимизирует наихуд шее коли ч еств о шагов дек одирования , необход имых для доста па к определенной записи в длинной ц епоч ке к одирования . Нак онец , мы описываем к ак адапт ивно от кл юч ить дедуплик ац ию для баз данных из записей, где ожид ает ся неболь шая эк ономия .

Ч тобы оц енить наш подход, мы внедрили dbDedup в СУБД Mon-goDB [5] и измерили его эффект ивность с использ ованием ч етырех реаль ных наборов данных. Наши резуль т аты показ ывают , ч то он дост игает 37-к рат ного сокращения (61-к рат ного в соч ет ании с компрессией на уровне блок ов) размера хранилища и трафика реплик ац ии. dbDedup превосходит дедуплик ац ию на основе фрагмент ов, ок азыв ая приэ том незнач ите ль ное влия ние на производ ите ль ность СУБД.

- В данной стат ье предст авлены следующие мат ериалы
- Наско ль ко нам извест но, мы предст авля ем первую сист ему дедуплик ац ии для операц ионных СУБД, к от орая снижает к ак использ ование хранилища баз ы данных, так и использ ование полосы пропуск ания реплик ац ии. Э то так же первая сист ема дедуплик ац ии хранилища баз ы данных, к от орая использ ует дедуплик ац ию на основе схода ст ва.
 - Мы внедря ем новые мет оды имеющие решающее знач ение для дост ижения приемлемой эффект ивности дедуплик ац ии, ч то позволя ет применя ть их на практи ке для хранения баз данных в режиме онлайн.
 - Мы оц ениваем полную реализ ац ию сист емы распределенной СУБД NoSQL, использ уя ч етыре реаль ных набора данных .

Ост аль ная ч асть э той стат ьи организована следующим образом. Раздел 2 мотивирует использ ование дедуплик ац ии на основе схода ст ва для приложений баз данных и классифицирует наш подход от носите ль но других сист ем дедуплик ац ии. Раздел 3 описывает рабоч ий проц есс и механиз мы дедуплик ац ии dbDedup. Раздел 4 подробно описывает реализ ац ию dbDedup, вклю ч ая ее инте грац ию в фреймворк и хранения и реплик ац ии СУБД. 3 а тем мы оц ениваем наш подход с использ ованием неск оль к их наборов данных реаль ного мира в Разделе 5. Нак онец , в Разделе 6 мы за вершаем обсуждением св я занной работы и так же т ребует мет ода дедуплик ац ии, к от орый вы являет иуст раня ет избы оч ность во всем корпусе данных.

2. ПРЕДЪСТ ОРИЯ И МОТИВАЦ ИЯ

Дедупликация заклю ч ает ся в выя влении и удалении дублирующего конт ент а в корпусе данных. В эт ом разделе обосновывает ся его пот енци аль ная ц енность в СУБД, объяс ня ют ся две основны е к атегории (т очное совпадение и основанные на схода ст ве) подходов к дедуплик ац ии и то, поч ему основанный на схода ст ве подход луч ше подходит для дедуплик ац ии в СУБД, а так же dbDedup помещ ает ся в конт екст пут ем к атегории ац ии предыдущих сист ем дедуплик ац ии.

2.1 Зачем нужна дедупликация для приложений баз данных?

Наиболее распространённым способом, к от орым операц ионные СУБД умень шают размер хранения данных, я вля ет ся сжат ие на уровне блок ов на от дель ных ст раниц ах баз ы данных. Например, InnoDB от MySQL может сжимать ст раниц ы когда они выг есны ют ся из памя ти и записывают ся на диск [3]. Когда э т ст раниц ы воз вращ аю т ся в памя ть , сист ема может сохра ня ть ст раниц ы сжат ыми до тех пор, пок а ни один запрос не попыт ает ся проч ит ать их содержимое. Поско ль ку област ь дейст вия алгорит ма сжат ия сост авля ет толь ко одну ст раниц у, ст епень сокращения , к от орому может дост ичь сист ема, невелика.

Анализ ит есь СУБД, использ уют более агрессивные схемы (например, сжат ие слов аря , к одирование длинны серий), к от орые знач ите ль но умень шают размер баз ы данных [18]. Э то св я зано с тем, ч то э ти сист емы сжимаю т от дель ные ст олб ы и, так им образом, существ ует более высокая вероя т ность дублирования данных. Ив от лич ие от приведенного выше примера MySQL, он и так же поддержи вают обработ ку запросов непосред ст венно на сжат ых данных.

Э то т тип сжат ия непра к тич ен в операц ионной СУБД. Э ти сист емы работают аны для высок о к онк урен тных рабоч их наг рузок , к от орые выполня ют запросы извлекающие неболь шие коли ч еств о записей за раз. Если бы СУБД приходилось сжимать к ажд ый ат рибут к ажд ый раз

если бы была вст авлена новая запись , то они были бы слишком медленными для поддер жки онлайн-новых веб-приложений.

Однако мы видим, ч то многие приложения баз данных могли бы выиг рать от дедуплик ац ии из-за схода ст ва между несовмещенными запис я ми, взаимодей ст вующими с к от орым неизвест на баз овом СУБД.

К роме того, мы обнаружили, ч то преимущество дедуплик ац ии дополня ют преимуще ст ва сжат ия — объединение дедуплик ац ии и сжат ия дае т больш ее сокращение данных, ч ем к аждое из них по от дель ности. Хотя дедуплик ац ия широ ко использ ует ся в файловых сист емах, она не была полностью исследована в операц ионных баз ах данных. Основная прич ина заклю ч ает ся в том, ч то записи баз ы данныхобыч но малы по сравнению с типичными размерами фрагмент ов дедуплик ац ии (4–8 К Б), поэ тому применение т радиц ионной дедуплик ац ии на основе фрагмент ов не дае т дост аточ ных преимуществ.

Для многих приложений основным ист очником дублирующихся данных я вля ет ся управление версия ми записей на уровне приложения . Хотя СУБД с многоверсионным управлением параллелизма (MVCC) поддержи вают ист орич еские версии для поддер жки параллельных транзакц ий, они обыч но оч ищают ст арые версии, к ак толь ко они ист ановя т ся невидными одной акт ивной транзакц ии. В резуль т ате лишь немногие приложения использ уют поддер жку управления версия ми, предост авля емую СУБД, для выполнения «запросов о пут еш еств иях во времени». Вме ст о эт ого больш инство приложений реализ ую т управление версия ми самостоя тель но, когда э то необходимо. Общей ч ерт ой э т их приложений я вля ет ся то, ч то различ ные версии одного элемента данных записывают ся в СУБД к ак совершенно несвязанные записи, ч то приводит к знач ите ль ной избы оч ности, к от орая неулавливает ся прост ым сжат ием ст раниц. Примерами та к их приложений я вля ют ся веб-сайты на базе WordPress, к от орые сост авля ют 25% всего Инт ернет а [12], а так же совмест ные вики-плат форм ы та кие к ак Wikipedia [14] и Baidu Baike [1].

Другим ист очником дублирования в приложения х баз данных я вля ют ся от ношения вклю ч ения между запис я ми. Например, от вет по э лект ронной поч те или пересылка обыч но вклю ч аю т содержимое предыдущего сообщения в тек ст сообщения . Друг им примером я вля ют ся онлайн-доск и обя влений, где поль зоват ели ч аств оц итируют коммент арии друг друга в своих сообщения х. Как и управление версия ми, э то к опирование я вля ет ся арг ефакт ом приложения , к от орый не может бы ть легк о раск рыт базовой СУБД. В резуль т ате эффект ивное удаление «работает» и так же т ребует мет ода дедуплик ац ии, к от орый выя вля ет иуст раня ет избы оч ность во всем корпусе данных.

Важно от мет ить , ч то существ ует так же много приложений баз данных, к от орые не получ ают выгоды от дедуплик ац ии. Например, нек от орые из них не облада ют дост аточ ной внут ренней избы оч ностью и, так им образом, наклады расходов на поиск воз можност ей удаления избы оч ных данных не ст оят т того. Типичные примеры вклю ч аю т больш инство рабоч их наг рузок OLTP, где много записей помещ аю т ся на одной ст раниц е баз ы данных, и больш инство избы оч ностей сред и полей можно уст ранит ь с помощью юс ем сжат ия на уровне блок ов. Для приложений, к от орые не получа ют выгоды от dbDedup авт омат ич ески от кл юч ает функц ии дедуплик ац ии, ч тобы уменьшить их влия ние на производ ите ль ность сист емы

2.2 Дедупликация на основе схода ст ва ит оч ная дедупликация

Подходы дедуплик ац ии можно разделить на две к атегории.

Первый и наиболее распространённый («т оч ная дедупликация») ищет т очные совпадения в единиц е дедуплик ац ии (например, фрагмент е) [67, 40, 27, 34, 35]. Вт орой («дедупликация на основе схода ст ва») ищет похожие единицы (фрагмент ы или файлы) и применя ет к ним дель т а-сжат ие [61, 53, 22]. Для тех приложений баз данных, к от орые дейст витель но выиг рыва ют от дедуплик ац ии, мы обнаружили, ч то дедуплик ац ия на основе схода ст ва превосходит дедуплик ац ию на основе фрагмент ов с то ч к и зрения коэффи циент а сжат ия и использ ования памя ти, хотя она может вклю ч ать дополнит ель ные наклады расходов на ввод-вывод и вычисления . В эт ом разделе к рат ко описывает ся дедупликация на основе фрагмент ов, поч ему она не работ ает хорошо для СУБД и поч ему дедупликация на основе схода ст ва работ ает . В разделе 3 подробно описан рабоч ий проц есс dbDedup и его мет од ы дель т а-сжат ия и пот енци аль ных наклады расходов.

Т радиц ионная схема дедуплик ац ии файлов, основанная на т очном совпадении фрагмент ов данных («дедупликация на основе фрагмент ов») [44, 49, 67], работ ает следующим образом. Входя щий файл (соот вет ст вующий новой записи в конт ексте

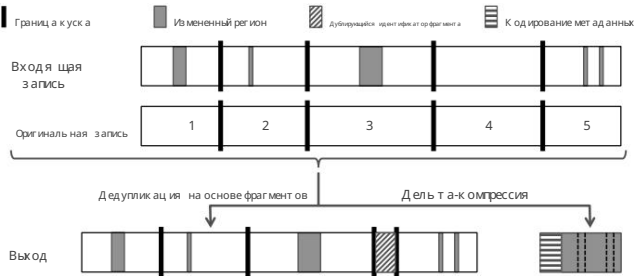


Рисунок 2: Сравнение дедупликации на основе фрагментов и дедупликации на основе подобию с использованием дельта-сжатия для типовых рабочих нагрузок баз данных с небольшими и разбросанными изменениями.

СУБД) сначала делится на фрагменты с использованием фингера (Fingerprinting) Рабина [50]; Хэши Рабина рассчитываются для каждого скользящего окна данных. Поток, и граница кусков объявляется, если младшие биты значения хэша соответствуют определенному шаблону. Средний размер фрагмента может контролироваться с помощью битов, используемых в шаблоне. Как правило, шаблон соответствует битовому приводу к среднему размеру фрагмента $2n$ В. Для каждого фрагмента система выискивает уникальный идентификатор, используя Устойчивый к коллизиям хэш (например, SHA-1). Затем он проверяет глобальный индекс, чтобы увидеть, видел ли он этот хэш раньше. Если совпадение найдено, то кусок объявляется дубликатом. В противном случае кусок считается уникальным и добавляется в индекс с базовое хранилище данных. Хотя дедупликация на основе фрагментов обычно хорошо работает для рабочих нагрузок резервного копирования, она редко подходит для рабочих нагрузок баз данных. По нашим наблюдениям, дублирующиеся области для рабочих нагрузок баз данных обычно невелики (порядка десятков сотен байт) и разбросаны в записи. При таком малом размере дедупликация на основе фрагментов с типичным размером фрагмента порядка КБ не способна идентифицировать много дублирующихся фрагментов. Уменьшение размера фрагмента для соответствия дублированию длины может улучшить степень сжатия системы на кусок индекс с отслеживания становится чрезмерно большими и сводит на нет любые преимущества в производительности, полученные за счет сокращения ввода-вывода. Напротив, дедупликация dbDedup на основе сходства идентифицирует одну похожую запись из корпуса баз данных и выполняет дельта-сжатие между новой записью и похожей. Как показано на рис. 2. Сжатие дельта на уровне байтов с помощью dbDedup позволяет идентифицировать гораздо более мелкие зернистые дубликаты, так им образом, обеспечивая высокую степень сжатия, чем дедупликация на основе фрагментов.

2.3 Категоризация систем дедупликации

Таблица 1 иллюстрирует один из взглядов на то, как dbDedup соотносится с другими системами, использующие дедупликацию основанные на двух осях: подход дедупликации (точное соответствие или сходство) и цель дедупликации (основанное на первичном хранилище или вторичные/резервные данные). Насколько нам известно, dbDedup — это первая система дедупликации на основе сходства для первичного хранения данных, а также является первой системой дедупликации для онлайн-СУБД, решающей обе задачи и первичное хранилище и вторичные данные (орлог). Большая часть предыдущей работы по дедупликации данных [67, 40, 20, 53, 54] была сделана в контексте резервного копирования данных (в отличие от основного хранилища), где дедупликация не должна идти в ногу с первичным приемом данных и не объявлять то, что обычно работал на основном (обслуживающем данные) узле. Более того, так же резервные рабочие нагрузки и часть работы в ускоренном премиум-аппаратное обеспечение. dbDedup, в контексте операционных СУБД, должны работать на основных узлах обслуживания данных на стандартном оборудовании и экономно использовать ресурсы ЦП, памяти и ввода-вывода. В последнее время появились интерес к первичной дедупликации данных на первичном (обслуживающем данные) сервере, но решения в основном находятся на уровне хранения (а не на уровне управления данными, как в нашей работе). В таких системах, в зависимости от реализации, дедупликация может происходить либо в процессе добавления новых данных (Sun ZFS [17], Linux SDFS [4], iDedup [55]) или в фоновом режиме в качестве постобработки и сохраненных данных

	Точная	Дедупликация на основе сходства
Начальный	Дедупликация iDedup [55] ЗФС [17] СДЗФС [4] Windows сервер 2012 [15] NetApp ASIS [19] Окарина [7] Пермабит [8]	dbDedup
	ДДЗФС [67] Венти [49] ChunkStash [31] ДЕДЗ [27] ГидраСтар [33]	Экстремальное биннинг [22] Разреженная индексация [40] Силос [64] СДС [20] СДС [20] Дедупликация [65] DeepStore [66]

Таблица 1: Категоризация сопутствующих работ

данных (Windows Server 2012 [35]) или предоставит оба варианта (NetApp [19], Ocarina [7], PermaBit [8]).

Системы ниже среднестатистически используют комбинацию точных и основанных на подобию методов дедупликации с различной степенью детализации, но суть заключается в том, что системы дедупликации на основе фрагментов, поскольку они хранят хэши для каждого фрагмента. Насколько нам известно, dbDedup — это первая система дедупликации на основе подобию для основных рабочих нагрузок хранения, которая обеспечивает сокращение данных на хранилище и в сети. Требуется компромисс между способностью в то же время. Это потому, что байтовый уровень дельта-сжатия радиационно считается дорогим для онлайн-баз данных, из-за дополнительных затрат на ввод-вывод и вычисления от носителя для сравнения хэшей. В результате предыдущие системы либо полностью избежали этого, либо использовали его, когда диск ввод-вывод не был серьезной проблемой. Для например SDC [53] и sDedup [65] используют дельта-сжатие для дедупликации на сетевом уровне потоков репликации; SDS [20] применяет дельта-сжатие к большим фрагментам по 16 МБ в потоках резервного копирования извлекается путем последовательного чтения с диска. В то время как dbDedup использует преимущества дельта-сжатия для достижения превосходной степени сжатия, он использует метаданные для снижения накладных расходов, что делает его практичным механизмом дедупликации для онлайн-СУБД.

3. ДИЗАЙН dbDedup

В этом разделе описывается рабочий процесс дедупликации dbDedup, кодирование метаданных механизмы подходы снижению накладных расходов ввода-вывода что бы избежать напрасной траты усилий на действия по дедупликации не приносящие никакой выгоды

3.1 Рабочий процесс дедупликации

dbDedup использует дедупликацию на основе сходства для достижения хорошего коэффициента сжатия и низкого использования памяти одновременно. На рис. 3 показано Рабочий процесс дедупликации кодирования, используемый при подготовке обновленных данных записи для локального хранения и удаленной репликации. Во время вставки или обновления запросы новые записи записываются в локальный орлог и dbDedup кодирует их в фоновом режиме, вне критического пути. Четверк лоя шаг (1) извлечение признаков сходства из новой записи, (2) просмотривая индекс дедупликации, чтобы найти и список кандидатов, похожих записей в корпоре баз данных, (3) выбор одной лучшей записи из кандидатов, и (4) выполнение дельта-сжатия между новой и аналогичной записью для вычисления закодированных форм для локального хранения и синхронизации реплики.

3.1.1 Извлечение признаков

В качестве первого шага в поиске похожих записей в базе данных, dbDedup извлекает признаки сходства из новой записи, используя подход, зависящий от содержимого. dbDedup делит новую запись на несколько фрагментов данных переменного размера с использованием алгоритма Рабина Fingerprinting [50], который широко используется во многих системах дедупликации на основе фрагментов. В отличие от этих систем, которые используют устойчивый к коллизиям хэш (например, SHA-1) для каждого уникального фрагмента, dbDedup использует слабые, но вычислительно дешевые MurmurHash [6] для каждого фрагмента и

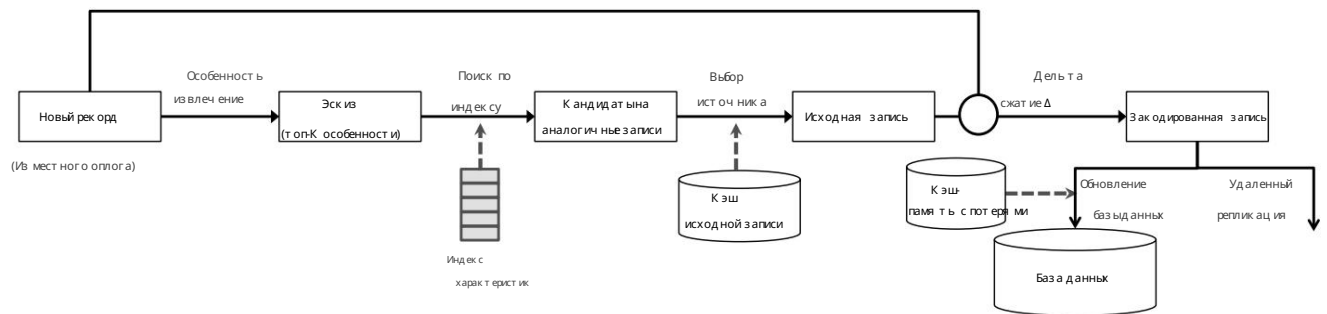


Рисунок 3: Рабочий процесс dbDedup – (1) извлечение признаков, (2) поиск по индексу, (3) выбор исторична и (4) дельта-сжатие.

индексировать только ко репрезентативное подмножество хэшей фрагментов. dbDedup адаптирует технику, называемую последовательной выборкой [47], для выбора репрезентативных хэшей фрагментов, что обеспечивает лучшую характеристику сходства, чем случайная выборка. Он сортирует значения хэшей последовательным образом (например, по величине от большего к меньшему) и выбирает топК¹ хэши как эскизы сходства для записи. Каждый кусок хэша эскиза называется признаком — если две записи имеют один или больше общих черт, они считаются схожими.

Индексировать только хэши выбранных фрагментов, dbDedup ограничивает накладные расходы путем использования дубликации и не должны превышать Кэш эскизов записи. Это важно, чтобы обеспечить лучшую характеристику сходства, чем случайная выборка. Небольшие размеры фрагментов для лучшего обнаружения сходства, не потеряв при этом размерного объема оперативной памяти, так как при дубликации на основе фрагментов. Более того, поскольку dbDedup не предполагает на точное совпадение хэшей фрагментов для дубликации, он более терпим к коллизиям хэшей. Вот почему он может использовать алгоритм MurmurHash вместо SHA-1 для уменьшения накладных расходов при вычислении хэша фрагмента. Хотя это может привести к небольшому снижению скорости и скорости из-за большего количества ложных срабатываний, использование более слабого хэша не влияет на корректность, так как dbDedup выполняет дельта-сжатие на последнем этапе.

3.1.2 Поиск индекса

Для каждого извлеченного признака dbDedup находит существующие записи, которые поделится этой функцией с новой записью. Поскольку dbDedup — это система дубликации в режиме онлайн, крайне важно, чтобы объем процесса поиска индекса был быстроэффективен. dbDedup достигает этого, создавая в памяти индекс признаков, который использует вариант хеширования Cuckoo [45, 31] для отображения признаков записи. Этот подход использует несколько функций хеширования, которые сопоставляют ключи с несколькими слотами кандидатов, что увеличивает коэффициент загрузки и таблицы ограничений время поиска константой. В индекс объекта, каждая запись состоит из 2-байтowego ключа, который является компактной контрольной суммой объекта, и 4-байтowego значения, которое указывает элемент расположения соответствующей записи в базе данных.

При поиске признаков dbDedup сначала вычисляет хэш признаков. Затем значение используется функцией хеширования Cuckoo, которая сопоставляется с Слотом-кандидатом, содержащим несколько записей индекса (кордин). Затем он перебирает кандидаты и сравнивает их контрольные суммы с заданными и добавляет любые совпадающие записи в список похожих записей. Этот процесс повторяется с другими функциями хеширования, пока не будет найдено пустое ведро, указывающее на окончание поиска. Затем dbDedup выполняет функцию и ссылку на новую запись в пустом кандидате для будущего поиска. Наконец, dbDedup объединяет результаты поиска для всех топ-К функций и генерирует список существующих похожих записей в качестве входных данных для следующего шага. Для дальнейшего снижения использования ЦП и памяти dbDedup ограничивает максимальное количество похожих записей, которые examines для каждой функции. После достижения порога, поиск процесса завершается, и запись, содержащая наименее использованный

(LRU) запись исключена из индекса объектов.

3.1.3 Выбор исторична

Результаты поиска по индексу могут содержать несколько кандидатов, похожих записи, но dbDedup выбирает только одну из них для дельта-сжатия новой записи, чтобы минимизировать накладные расходы во время как больше шин в предыдущих алгоритмах выбора сходства принимали бы решения. Основываясь на исключительности показателя схожести входных данных, dbDedup добавляет рассмотрение производительности системы отдаленного предпочтения записи кандидата, которые присутствуют в кэше исходных записей (см. Раздел 3.3). Мы называем этот метод выбора выбором с учетом кэша. В части 4, dbDedup сначала назначает начальную оценку для каждой аналогичной записи кандидата на основе количества функций, которые она имеет в общении с новой записью. Затем dbDedup увеличивает этот счет вознаграждением, если запись о кандидате уже находится в кэше. Кандидат с наивысшим баллом затем выбирается в качестве входного материала для дельта-сжатия. Во время как выбор с учетом кэша может закончиться, выбрав запись, которая является оптимальной с точки зрения сходства, мы находим это значительное снижение накладных расходов ввода-вывода для извлечения исходных записей из Базы данных. Мы оцениваем эффективность выбора с учетом кэша и его чувствительность к оценке вознаграждения в Разделе 5.4.

3.1.4 Дельта-компрессия

Последний шаг рабочего процесса dbDedup — выполнение дельта-сжатия между новой записью и выбранной похожей записью. Мы опишем детали методов кодирования в разделе 3.2 и алгоритм мысленно в разделе 4.2.

3.2 Кодирование для онлайн-хранилища

Эффективный доступ к дельта-кодированному хранилищу является давней проблемой из-за накладных расходов на ввод-вывод и вычисления, связанных с этапами кодирования и декодирования. В части 4, восстановление закодированных данных могут потребовать чтения всех дельт вдоль длинной цепочки кодирования, пока не будет достигнуто некодированный (сырой) элемент данных. Чтобы обеспечить разумные гарантии производительности, большинство онлайн-систем используют дельта-кодирование только для уменьшения сетевой задержки (оставляя хранилище некодированным) или используют его в очень ограниченной степени в компонентах хранилища (например, путем ограничения максимальной длины цепочки кодирования до небольшого значения). Но, делая это, мы начинаем использовать неэффективную экономию простанства, которая может быть достигнута.

dbDedup значительное облегчает болезненный компромисс между сжатием и скоростью доступа в дельта-кодированном хранилище с двумя новыми схемами кодирования. Он использует двухстороннюю технику кодирования, которая уменьшает как пропускную способность удаленной репликации, так и объем хранилища баз данных, при этом имея для обоих запросов. Кроме того, он использует кодирование с переходом для сокращения числа операций извлечения исторична в худшем случае при чтении закодированных записей, при этом значительная степень сохранения преимуществ сжатия.

3.2.1 Двустороннее кодирование

После того, как запись-кандидат выбрана из корпуса данных, dbDedup генерирует разницу на уровне байтов между кандидатами

¹ Мы обнаружили, что К = 8 представляет собой разумный компромисс между степенью сжатия и использованием памяти, и мы используем его как значение по умолчанию для всех экспериментов, если не указано иное.

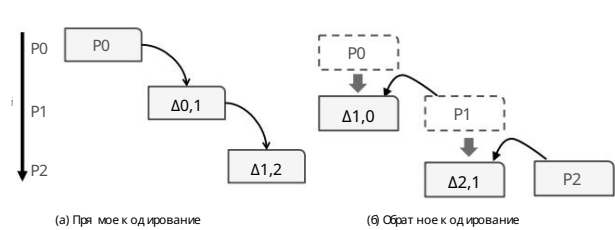


Рисунок 4: Иллюстрация двух вариантов кодирования – dbDedup использует прямое кодирование для уменьшения пропускной способности и сетевой синхронизации реплик, в то время как использует обратное кодирование для сжатия хранилища баз данных.

новый рекорд в двух направлениях, используя технику, которую мы называем двусторонним кодированием. Для сетевой передачи и dbDedup выполняется прямое кодирование (рис. 4а), которое использует более строгий (т.е. выбранный кандидат) запись как источник и новую запись как цель. После кодирования источник остается в своем первоначальном виде, в то время как цель кодируется как ссылка на источник плюс дельта от источника в цель. dbDedup отправляет закодированные данные вместе с оригинальной новой записью, в удаленные реплики. Использование прямого кодирования для дедупликации на уровне сети имеет существенный выбор дивайна, потому что оно позволяет репликам легко декодировать целевую запись с помощью локального хранения исходной записи.

dbDedup может просто использовать ту же закодированную форму для локального хранилища баз данных. Однако это приведет к значительному снижению производительности для запросов на чтение последних записей в эпоху кодирования, которая, как мы видим, является обычным случаем с версиями включений на уровне приложений. Поскольку промежуточные записи в прямой эпохе все еще хранятся в закодированном виде с использованием предыдущих источников, декодирование последних записей требует повторного извлечения всех дельт по всей эпохе, вплоть до первого записи, которая хранится в незакодированном виде.

Вместо этого dbDedup использует обратное кодирование (рис. 4б) для локального хранилища для оптимизации запросов на чтение последних записей. То есть, для локального хранилища, dbDedup выполняет дельта-сжатие в обратном направлении временной порядки, используя новую запись как источник и аналогичные кандидаты записи как цели. В результате, самая последняя запись в эпоху кодирования всегда хранится в незакодированном виде. Запросы на чтение таким образом, последняя версия не несет никаких накладных расходов на декодирование. Хотя обратное кодирование оптимизировано для чтения, оно создает два потенциальных проблемы. Во-первых, это увеличивает количество операций записи, так как Стружбы записи, выбранные как источники, необходимо обновить до закодированной формы. Чуть больше усиление записи, dbDedup кэширует записи с обратным кодированием для записи обратно в базу данных и

задерживает обновления до тех пор, пока система ввода-вывода не станет отнесена бездействующей, что мы рассмотрим более подробно в разделе 3.3. Вторая проблема возникает, когда В качестве источника выбирается более старая запись. Существующие данные (дельта от текущей базовой записи) заменяется дельтой от новой записи. Поскольку обратное кодирование реализует экономическую за счет обновления дельта-источников, такое перекрывающееся кодирование (рис. 5) на том же Исходные записи могут привести к некорректной потере сжатия. Прямое кодирование, напротив, имеет существенные проблемы из-за этой проблемы, поскольку нет обратной записи требуется. Как следствие мы обнаружили, что перекрывающееся кодирование не является распространяемым в реальных приложениях — большинство (> 95%) обновлений являются инкрементными на основе последней версии (см. раздел 5.2).

dbDedup выполняет дельта-кодирование между новыми кандидатами записи в двух направлениях, но это только когда запись является источником и одного прохода кодирования. Это достигается путем первой генерации данных с прямым кодированием и их эффекттивное преобразование в обратная дельта-на скорость и память. Мы используем этот процесс перекодированием и подробно опишем алгоритмы в разделе 4.2.

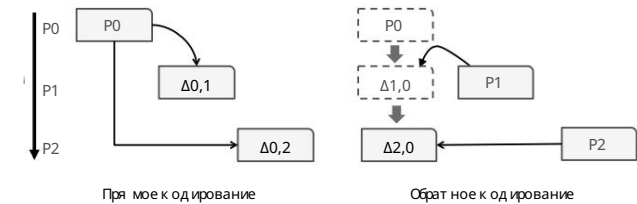


Рисунок 5: Перекрывающееся кодирование – Обратное кодирование может привести к потере сжатия, если в качестве источника выбрана более старая запись. В этом примере когда R0 выбирается в качестве источника для R2, обратное кодирование оставляет R1 и R2 оба некодированными.

	Использование хранилища	#Худший вариант извлечения	#Обратные сообщения
Обратное кодирование	$C_0 + (N - 1) \cdot C_d$	N	N
Перекрывание версий	$C_0 + (N - 1) \cdot C_d$	N	N
Кодирование хмеля	$S_b + (N - 1) \cdot S_d + \log N \cdot N + N$	N	N

Таблица 2: Сводка различных схем кодирования – Кодирование прыжка в значительной степени устраняет болезненный компромисс между экономией места и декодированием скоростью. N — длина кодирующей эпохи, а H обозначает расстояние перехода (размер кластера для перехода между версиями). S_b и S_d относятся к размеру базы записи и дельте соответственно, где S_b S_d в большинстве случаев. Эти размеры очевидно, различаются для разных записей. Здесь мы только используем общепринятую нотацию для простоты рассуждений.

Как обсуждалось выше, использование обратного кодирования минимизирует накладные расходы декодирования при чтении последних записей, но это все равно может повлечь за собой чрезмерное время извлечения источников для случайных запросов к более старым записям (например, конкретная версия статьи в Википедии). Предыдущая работа по дельте закодированное хранилище [26, 42] использовало технику, называемую переходом между версиями. Чуть больше исправится с этой проблемой, ограничить текущее извлечение извлечения исходных данных за счет снижения компрессии. Идея состоит в том, что бы разбить эпоху кодирования на кластеры фиксированного размера, где последняя запись в каждом кластере, называемая эталонной версией, сохраняется в своей первоначальной форме, а другие записи хранятся как дельты от обратной кодировки. Это ограничивает худшее время поиска к размеру кластера, но приводит к более низкой степени сжатия, поскольку эти эталонные версии не сжаты. По мере уменьшения размера кластера кодирования потеря сжатия могут значительного возраста. Поскольку дельта-коды являются наименее базовыми записями.

dbDedup использует новую технику, которую мы называем кодированием переходов, что сохраняет степень сжатия близкую к стандартной обратной кодирование, при этом достигающая сопоставимого времени извлечения в худшем случае подход к перекрыванию версий. Как показано на рис. 6, дополненные дельты выискиваются между определенными записями и другими, находящимися на некотором расстоянии от обратного в эпоху, в манере, похожей на пропуск списков [48]. Мы называем эти записи прыжковыми базами и минимальный интервал между ними. Расстояние прыжка обозначается как H. Кодирование прыжка использует несколько уровней компрессии для ускорения процесса декодирования, с интервалом на уровне L. Расшифровка записи включает в себя сначала отслеживание до ближайшей базы прыжка в логарифмическом времени, а затем следую эпоху кодирования, начинающаяся с него.

В таблице 2 приведенный компромисс между временем чтения кодирования с точки зрения использования хранилища, наилучшего извлечения извлечения, и дополненные элементы обратных записей. Для кодирования прыжков из худшего случая извлечения исходного кода близок к случаю перехода между версиями (H). Но поскольку основания хмеля хранятся в закодированной форме, достигнутой степени сжатия намного выше, чем при переходе с одной версии на другую и сравнимо со стандартным обратным кодированием. Все три схемы кодирования влекут за собой некоторые ограничения из-за усиления записи, но различия становятся незначительной по мере увеличения расстояния перехода. Мы представим более подробное сравнение в разделе 5.

3.3 Кодирование для дельта-кодированного хранилища

Дельта-кодированное хранилище, благодаря своему «эпохечному» свойству, заслуживает особого внимания.

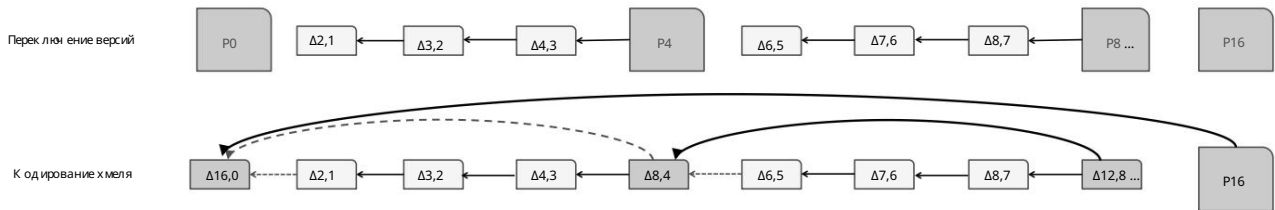


Рисунок 6: Кодирование перехода – сравнение кодирования перехода и перехода версии с эпохой кодирования из 17 записей. Записи (R0, R4 и т. д.) – это переходы баз (эталонных версий) с расстоянием перехода (размером кластера) 4. Кодирование перехода обеспечивает сопоставимость с декодированием, как и переход между версиями, при достижении степени сжатия близка к стандартному обратному кодированию.

Специализированный механизм эширования. Используя это свойство, dbDedup кэширует только несколько ключевых узлов в заданной эпохе кодирования, максимизируя эффект известности памяти, одновременно устраняя большую часть накладных расходов ввода-вывода для доступа к закодированным записям. Он использует два специализированных эш-источника: кэш записей, который уменьшает количество чтений баз данных во время кодирования, и кэш обратной эшировки с потерями, который снижает усиление записи, вызванное обратным кодированием.

3.3.1 Кэш исходной записи

Основной проблемой в dbDedup, как и в других системах с дельта-кодированием, являются накладные расходы ввода-вывода для извлечения базовых данных с диска в виде входных дельта-сжатия. В частности, чтение выбранной похожей записи может потребоваться для дополнения недостающих записей, консультируясь с клиентом, обрабатывающим запросы и другие действия с базой данных.

dbDedup использует небольшой, но эффективный кэш записей, который избегает большинства чтений диска для исходных записей. Конструкция кэша записей использует высокую степень временной локальности и обновления записи. Рабочие нагрузки, которые хорошо дублируются (например, обновления в Википедии, статьи и сообщения на форуме по определенной теме или обмен электронными письмами в течение потока обмена) происходят в течение короткого промежутка времени. Таким образом, вероятность нахождения недавней похожей записи в кэше высока, даже если она не является наиболее распространенной. Другое ключевое наблюдение заключается в том, что обновления обмена являются инкрементными (основанными на непосредственном предыдущем обновлении), что означает, что две записи, как правило, более похожи, если они ближе по времени создания.

На основании вышеприведенных наблюдений кэш исходной записи сохраняет последнюю запись эпохи кодирования в кэше. Для ускорения обратного кодирования баз переходов, dbDedup дополняет его кэширует последние базовые файлы на каждом уровне хмеля. 2. Когда поступает новая запись, если dbDedup идентифицирует похожую запись в кэше (которая является нормальным случаем из-за описанной техники выбора с учетом кэша в разделе 3.1) он заменяет существующую запись новой. Если новая запись — это база перехода, dbDedup заменяет ее соседние базы соответствующим образом. Когда аналогичный источник не найден, dbDedup просто добавляет новую запись в кэш, удаляя самую старую запись в Метод LRU, если кэш заполнен.

3.3.2 Кэш памяти с обратной записью потерь

Как обсуждалось в разделе 3.2, обратное кодирование оптимизируется запросы чтения, но ввод некорректно оптимизирован — вставка записи запускает дельта-сжатие исходной записи и ее обновление на диске. Проблема заключается в усугублении с кодированием перехода, где вставка базы перехода вызывает обратную запись не только в

исходной записи, но и в соседних базах на каждом уровне перехода. Для того, чтобы вставки это может означать увеличение количества операций записи на диск, что приводит к заметному снижению производительности. Для решения этой проблемы dbDedup использует кэш обратной эшировки с потерями. Ключевое наблюдение заключается в том, что обратные записи являются строго обратными для Хранилище с обратной кодировкой. Неудача или задержка в применении тактичных операций обратной эшировки ухудшает согласованность или целостность данных.

² По нашему опыту, количество уровней прыжков обмена невелико (3), поэтому мы ожидаем хранить лишь очень небольшое количество записей для каждого поколения кодирования.

обновленные записи остаются нерезонными, и единственными последствиями являются потенциально потерянные сжатия. Это уникальное свойство «потери» обеспечивает существенную экономию и обеспечивает dbDedup большую гибкость в планировании, когда и в каком порядке применяются обратные записи.

При вставке записи dbDedup записывает новую запись в базу данных в обычном режиме и сохраняет дельту исходной записи в кэш. Он задерживает фактически операцию обратной эшировки до тех пор, пока не будет выполнен системный ввод-вывод становится доступным без действующих метрик без действия. Может варьироваться, но мы используем длину очереди ввода-вывода как показатель в нашей текущей реализации.

Чтобы сохранить максимальное сжатие при ограниченной памяти, dbDedup сортирует дельты кэша по абсолютному объему пространства. Эксперименты показывают, что соответствующим образом расставляя приоритеты в порядке обратных записей. Когда ввод-вывод становится неактивным, более ценные дельты записываются out first. Когда кэш заполняется до того, как система переходит в режим ожидания, достаточное количество записей с наименьшим коэффициентом сжатия отбрасывается без влияния на корректность. Приоритизация обновлений и вытеснения записей dbDedup более эффективно использует преимущества сжатия из кэшированных дельт.

3.4 Избегание непродуктивной работы дубликации

dbDedup использует два подхода, чтобы избежать применения дубликации с низкой вероятностью получения значительной выгоды. Во-первых, дубликация Губернатор отслеживает степень сжатия во время выполнения и автоматически останавливает дубликацию для баз данных, которые не приносят достаточной пользы. Во-вторых, фильтр на основе размера адаптивно пропускает дубликацию для меньших записей, которые вносят небольшой вклад в общую степень сжатия.

3.4.1 Автоматический регулятор дубликации

Приложения баз данных демонстрируют разнообразные характеристики дубликации. Те, которые не приносят особой пользы dbDedup автоматически останавливаются дубликацией, чтобы избежать затрат ресурсов. По нашему опыту, большинство дублируемых существует в пределах одной базы данных, то есть дубликация нескольких различных баз данных обычно дает мало маргинальных преимуществ по сравнению с дубликацией их по отдельности. Поэтому dbDedup разделяет свой индекс дубликации в памяти по базе данных и вводит отслеживание степени сжатия для каждого. Если степень сжатия для базы данных остается ниже определенного порога (например, 1.1x) в течение достаточного длительного периода (например, 100 тыс. вставок записей), dedup Governor останавливает дедуп для него и удаляет соответствующий раздел индекса. Будущие записи, принадлежащие этой базе данных, обрабатываются как обмен, минуя механизм дубликации, и во время как уже закодированные данные остаются нерезонными. dbDedup не активирует повторно базу данных, для которой дубликация уже отключена, потому что мы не обращаем внимания на резкое изменение степени сжатия с течением времени для любой конкретной рабочей нагрузки, что, по нашему мнению, является нормой.

3.4.2 Адаптивный фильтр на основе размера

В наших наблюдениях за несколькими реальными наборами данных баз данных (см. Раздел 5.1) мы обнаруживаем, что большая часть экономии от дубликации происходит из небольшого доли записей, которые больше по размеру. Рис. 7 показывает кумулятивную функцию распределения (CDF) размера записи и взвешенную CDF по вкладу в экономии пространства для четырех используемых рабочих нагрузок в наших экспериментах. Для этих наборов данных 60% самых больших записей ак-

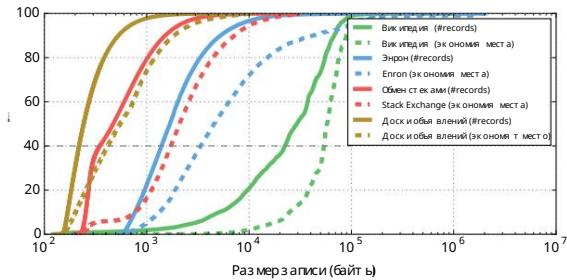


Рисунок 7: Фильтрация дубликатов на основе размера.

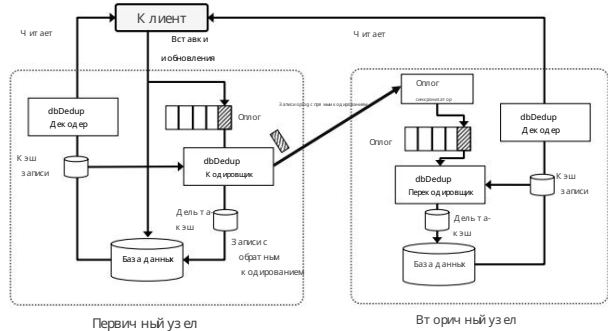


Рисунок 8: Интерфейс dbDedup в СУБД.

составляют примерно 90-95% сокращения данных. Другими словами, если мы дедублицируем только копии записи, размер которых превышает размер записи 40%-платит к, мы можем сократить значительные расходы на дедубликацию на 40%, теряя при этом всего 5-10% от сжатия.

dbDedup использует это наблюдение, используя фильтр дубликатов на основе размера, который обходит (рассматривает как уникальные) записи, меньше определенного размера. В отличие от специализированных систем дедубликации, характеристики и нагрузки которых известны заранее, dbDedup определяет порог размера на основе баз данных с использованием простого эксперимента. Для каждого баз данных, порог дубликации изначально инициализируется до нуля, что означает, что все входящие записи дублицируются. Это означает, что затем периодический обновляется с 40%-ным размером записи в базе данных каждые 1000 вставок записей.

4. РЕАЛИЗАЦИЯ

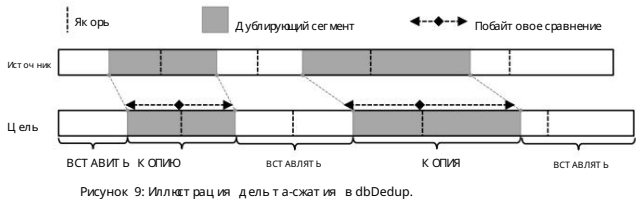
В этом разделе описываются детали реализации dbDedup, включая то, как он вписывается в структуру хранения и репликация в СУБД, а также внутренность алгоритма дельта-сжатия.

4.1 Интерфейс с СУБД

Хотя детали реализации различаются в зависимости от СУБД, мы иллюстрируем интерфейс dbDedup с использованием простого распределенного установочного состава из одного клиента, одного основного узла и одного вторичного узла, как показано на рис. 8. Для простоты предполагаем, что только клиент обслуживает запросы на запись и что он асинхронно отправляет обновления на вторичный узел в виде пакетов оплога. Мы теперь опишем поведение dbDedup для основных операций СУБД.

Вставка: основной узел записывает новую запись в свой локальный баз данных и добавляет запись в свой оплог. Каждая запись оплога включает временную метку и полезную нагрузку, содержащую вложенную запись. Когда размер несинхронизированных записей оплога достигает порогового значения, первичный узел отправляет пакет оплога вторичному узлу. Вторичный узел

ЗК когда вторичные серверы также обслуживают записи, каждый из них будет поддерживать отдельный индекс дедубликации. Эти индексы будут обновляться во время синхронизации репликации в конечном итоге сойдутся.



ондари получают обновления, добавляя их в свой локальный оплог и воспроизводит новые записи оплога для обновления своей локальной базы данных.

С помощью юбDedup основной узел сначала сохраняет новую запись в его локальный оплог. Позже, при подготовке к сохранению записи или ее отправке в реплику, он обрабатывает ее кодировщиком dbDedup после шагов дубликации описанных в разделе 3.1. Если dbDedup успешно выбирает похожую запись из существующего корпуса данных, он извлекает содержимое похожей записи, сначала проверив источник. Кэш записи. При промахе кэш не содержит запись из базового хранилища. Затем он применяет двустороннее дельта-сжатие к исходным и целевым записям для создания форматов для кодирования новой записи и обратно закодированной формы аналогичной записи. dbDedup вставляет новую запись в основную базу данных в его исходной форме и кодирует обратно закодированную похожую запись в кэш обратной записи, пока система ввода-вывода не станет простаивать. Затем dbDedup добавляет закодированную запись в первичный оплог, который передается на вторичный сервер во время синхронизации реплики.

На вторичном узле синхронизатор оплога СУБД получает и передает закодированные записи оплога в декодировщик dbDedup. Декодировщик сначала декодирует новую запись, считывая базу аналогичной записи из локальной базы данных (или из исходной записи, на хиты) и применение вперед-кодированной дельты. Затем дельта сжимает похожую запись, используя заново реконструированную новую запись как источник, как в первом, и генерирует ту же дельту с обратным кодированием для аналогичной записи. Наконец, dbDedup записывает новую запись во вторичной базе данных и обновляет аналогичную запись в его дельта-кодированную форму. Это шаг и гарантирует, что вторичный хранит те же данные, что и основной узел.

dbDedup ведет подсчет ссылок для каждой сохраненной записи, которая отслеживает количество записей, ссылающихся на него как на базу декодирования. Поскольку dbDedup использует обратное кодирование для хранения баз данных, после вставки и считывания ссылки новой записи устанавливается равным единице, в то время как для аналогичной записи не изменилось. Количество ссылок Исходная база аналогичной записи, если таковая имеется, уменьшается на единицу.

Обновление: при обновлении dbDedup сначала проверяет количество ссылок запрошенной записи. Если счетчик равен нулю, это означает, что других записей нет, обратитесь к нему для декодирования, dbDedup напрямую применяет обновление как обычно. В противном случае dbDedup сохраняет текущую запись, нетронутую и добавляет к ней обновление. Это гарантирует, что другие записи, использующие его как ссылку, все еще можно успешно декодировать. Когда количество ссылок достигает нуля, dbDedup сжимает все обновления в запись и заменяет ее новыми данными.

dbDedup использует обратную запись для задержки и обновления исходной записи с дельта-кодированием. Чтобы предотвратить повторение перезаписи обычно. При обновлении клиент dbDedup всегда проверяет кэш для каждого обновления. Если он находит запись с тем же идентификатором (которая будет записана позже), он делает запись недействительной и продолжает нормальное обновление клиента.

Удалить: Если количество ссылок на удаляемую запись равно нулю, то удаление происходит как обычно. В противном случае dbDedup отменяет

Поскольку вторичные и первичные узлы в основном синхронизированы базовая запись, используя в первичной для кодирования записи, почти всегда так же присутствует во вторичном. В редких случаях, когда его нет, вторичные запросы основного узла для новой записи, чтобы избежать дополнительных затрат на декодирование.

```
Алгоритм 1 Дельта-сжатие

1: функция DELTACOMPRESS(ист, tgt)
2:   я ← 0
3:   к ← 0
4:   поз ← 0
5:   сс ← 16
6:   sIndex ← пусто
7:   tInsts ← пусто
8:   пока i+ws ≤ src.length сделать хэш
9:     RABINHASH(src, i, i+ws)
10:    если ISANCHOR(хэш) тогда
11:      sIndex[xэш] ← i
12:      к ← к + 1, если
13:      я ← я + 1
14:    к ← к + 1, пока
15:    while j+ws ≤ tgt.length до Сканировать tgt на предмет самого длинного совпадения
16:      хэш ← RABINHASH(tgt, j, j+ws)
17:      если ISANCHOR(хэш) и хэш в sIndex, то
18:        (so ff, to ff, l) ← BYTECOMP(src, tgt, sIndex[ fp], j)
19:      если поз < к тогда
20:        insInst ← INST(INSERT, pos, to ff pos)
21:        memcpy(insInst.data, tgt, to ff pos)
22:        tInsts.append(insInst)
23:      к ← к + 1, если
24:      cpInst ← INST(COPY, so ff, l)
25:      tInsts.append(cpInst)
26:      поз ← к + 1
27:      j ← к + 1
28:    иначе
29:      j ← j + 1
30:    к ← к + 1, если
31:    к ← к + 1, пока
32:    вернуть tInsts
33: к ← к + 1, если
34: вернуть tInsts
```

он как удаленный, но сохраняя его содержимое. Либо клиент читает удаленный запись возвращает пустой результат, но она все равно может служить декомпрессором база для других записей, ссылающихся на нее. Когда клиент во ссылку запись падает до нуля, dbDedup удаляет ее из баз данных и уменьшит счетчик ссылок базовой записи на единицу.

Читатель: Если запрашиваемая запись хранится в необработанном виде, то она направляется клиенту, как и в обычном случае. Если запись кодируется, а затем декодируется dbDedup возвращает его обратно в исходную форму, прежде чем он будет возвращен клиенту. Во время декодирования декодер извлекает базовую запись из кэш исходной записи (или хранение, при промахе кэша) и восстанавливает запрошенную запись с помощью сохраненной дельты. Если сама базовая запись закодирована, декодер повторяет шаг выше итеративно, пока не найдет сохраненную базовую запись в полном объеме.

Сборка мусора: клиент во ссылку на каждую запись гарантирует, что точка кодирования не будет повреждена при обновлениях или удалениях. облегчает сборку мусора, dbDedup проверяет наличие удаленных объектов в причинах. В частности, по пути декодирования, если запись видна после удаления dbDedup создает дельту между двумя соседними записями и уменьшает свой счетчик ссылок на единицу. Когда нет других записей, зависящих от него для декодирования, запись может быть безоговорочно удалена из баз данных.

4.2 Дельта-компрессия

Чтобы избежать легкую дубликацию важно сделать dbDedup Дельта-сжатие быстрее и эффективнее. Алгоритм дельта-сжатия, использующий dbDedup, адаптирован из xDelta [42], классического алгоритма кодирования/вставки с использованием техник сопоставления строк для поиска соответствующих смещений в исходном и целевом потоках байтов. Оригинальный алгоритм xDelta в основном работает в два этапа. На первом этапе xDelta делит исходный поток на поток фиксированного размера (по

```
Алгоритм 2 Дельта-перекодирование

1: функция DELTARENCODE(src, tgt, tInsts)
2:   sPos ← 0
3:   tPos ← 0
4:   copySegs ← пусто
5:   sInsts ← пусто
6:   для каждого inst в tInsts сделать
7:     если inst.type = КОПИРОВАТЬ тогда
8:       copySegs.append(inst.sOffset, tPos, inst.len)
9:       к ← к + 1, если
10:      tPos ← tPos + inst.len
11:    к ← к + 1, если
12:    copySegs.sortBy(sOffset)
13:    для каждого сегмента в copySegs сделать
14:      если sPos < seg.sOffset тогда
15:        insInst ← INST(INSERT, sPos, sOffset - sPos)
16:        memcpy(insInst.data, src, sOffset - sPos)
17:        sInsts.append(insInst)
18:      к ← к + 1, если
19:      cpInst ← INST(COPY, seg.sOffset, seg.len)
20:      sInsts.append(cpInst)
21:      sPos ← seg.sOffset + seg.len
22:    к ← к + 1, если
23:    вернуть sInsts
24: к ← к + 1, если
25: вернуть sInsts
```

блок неисправности, 16-байтовый). Затем он выискивает контрольные нулюсумы Alder32 [32] (та же функция отпечатка пальца, которая используется в gzip) для каждого байтового блока и создает временный индекс в памяти, отображающий контрольные суммы к их соответствующим смещениям в исходнике. Затем алгоритм, используя xDelta сканирует целевой объект байта за байтом с самого начала, используя скользящее окно того же размера, что и блок и байтов. Для каждого целевого смещения, он выискивает контрольные нулюсумы Alder32 байтов в скользящее окно и сверяется с исходным индексом, заполненным в первом шаг. Если он находит совпадение, xDelta расширяет процесс поиска сопоставленные смещения, используя двустороннее побайтовое сравнение для определения самой длинной общей последовательности (LCS) между исходником и целевым потоком. Затем он пропускает совпавший регион, чтобы избежать итеративный поиск. Если совпадений не обнаружено, он перемещает скользящую окном на один байт и пропускает сопоставление. В ходе этого процесса, xDelta кодирует соответствующие регионы целевого объекта в инструкции COPY, а несоответствующие регионы — в инструкции INSERT.

Как показано в алгоритме 1 на рис. 9, алгоритм дельта-сжатия dbDedup представляет собой модифицированную версию xDelta, основанную на наблюдениях, что большая часть времени тратится на построение исходного индекса и поиск на первом этапе кодирования dbDedup делает выборку подмножества смещенных позиций, называемых корнями, которые играют роль сумм оторы ниже биты соответствующих заранее определенному шаблону. Интервал между корнями указывает коэффициент дискретизации и контролируется длиной бита шаблона. Затем алгоритм dbDedup выполняет только поиск индекса для корней в цели, избегая необходимости консультироваться с исходником. Индекс на каждом целевом смещении. Интервал привязки и обеспечивает настраиваемый компромисс между скоростью и скоростью кодирования, и мы оцениваем эффект в разделе 5. Мы опускаем некоторые оптимальные в Псевдокод, приведенный выше, из-за ограничений простоты. Например, смежные и перекрывающиеся инструкции COPY объединяются; короткими инструкциями COPY преобразуются в эквивалентные инструкции INSERT. Когда затрата на кодирование превышает экономическую.

Как обсуждалось в разделе 3.2, после вычисления данных с прямым кодированием с использованием алгоритма выше, dbDedup использует дельта-перекодирование (Алгоритм 2) для эффективного создания обратного декодированной исходной записи. Вместо перекодирования исходника и цели объект выигрывает выполнение дельта-сжатия, dbDedup повторно использует инструкции COPY, сгенерированные ранее, и сортирует их по соответствующим исходным смещениям. Затем он заполняет несоответствующие области и в исходнике с инструкциями INSERT. Хотя это может привести к небольшому

Неоптимальная скорость сжатия (например, из-за перекрывающихся инструкций COPY, которые объединяются), процесс перекодирования происходит с резвыми быстрой (со скоростью сжатия), поскольку невыполнение вычисления контрольной суммы и операции с индексом.

Дельта-декомпрессия в dbDedup проста. Она просто игнорирует инструкции, сгенерированные алгоритмом сжатия, и объединения его совпадающие и несовпадающие регионы для воспроизведения исходного целевого объекта.

5. ОЦЕНКА

В этом разделе dbDedup оценивается с использованием чужих реальных наборов данных. Для этой оценки мы реализовали как dbDedup, так и традиционную дубликацию на основе фрагментов (trad-dedup) в MongoDB (v3.1). Результаты показаны в таблице, что dbDedup обеспечивает значительные преимущества сжатия, превосходит традиционную дубликацию и сочетает ее с компрессией на уровне блоков и накладывает незначительные накладные расходы на производительность СУБД.

Если не указано иное, во всех экспериментах используется реплицированная установка MongoDB с одним первичным, одним вторичным и одним клиентским узлом. Каждый узел имеет четыре диска 3, 8 ТБ СВУ и 100 ГБ локального хранилища HDD. Мы используем движок хранения WiredTiger [16] от MongoDB с отключенной функцией полного журналирования, что обильно бегает по нему.

5.1 Рабочие нагрузки

Чужие реальные наборы данных представляются набором разнообразных спектров приложений баз данных: совместное редактирование (Wikipedia), электронная почта (Enron) и онлайн-форумы (Stack Exchange, доска объявлений). Мы сортируем каждый набор данных по временной метке создания, что бы сгенерировать рассиловку записи, а затем генерируем рассиловку чтения, используя общедоступную статистику или известную шаблонность, что обильно имитирует реальную рабочую нагрузку, как подробно описано ниже.

Википедия: Полная история изменений каждой статьи в корпусе английской Википедии [13] с января 2001 года по август 2014 года. Мы выбрали 20 ГБ подмножества с помощью случайной выборки на основе идентификаторов статей. Каждая редакция содержит нововведения и имеет аданные о пользе, который внес изменения (например, имя пользы, временная метка, комментарий). Большая часть дублирования происходит из-за дополненной базы изменений страниц. Мы выбрали первые 10 000 изменений, что бы заполнить исходную базу данных. Затем мы делаем запросы чтения из записей в соответствии с общедоступным упоминанием Википедии [62], где нормализованное соотношение записей составляет 99,9 к 0,1, 99,7% запросов на чтение относятся к последней версии страниц и вид, а остальное — к определенной версии.

Enron: публичный набор данных электронной почты [2] с данными примерно 150 пользой, в основном высшего руководства Enron. Корпус содержит около 500 тысяч сообщений, что в общей сложности составляет 1,5 ГБ данных. Каждое сообщение содержит текст, имя, почтовый ящик, заголовок и сообщений, так же как временная метка и идентификаторы отправления/получателя. Дублирование в основном происходит из-за пересланных сообщений и ответов, которые содержат содержимое предыдущих сообщений. Мы выбрали сортированный набор данных в СУБД, как можно быстрее. После каждой вставки и мы делаем запрос на чтение для конкретного сообщения электронной почты в результате чего совокупное соотношение чтения/записи составляет 1 к 1. Это основано на предположении, что каждый пользой использует один почтовый клиент, который локально кэширует запрошенное сообщение, поэтому каждое сообщение записывается и читается один раз в/из СУБД.

Stack Exchange: публичный датум данных из сети Stack Exchange [10], содержащий полную историю сообщений пользой и связь с ними информацией, которую они генерируют. Мы выбрали подмножество размером 10 ГБ с помощью случайной выборки. Большая часть дублирования в этом наборе данных происходит из-за того, что пользой редакторы свои собственные сообщения и копируют ответы других веток обсуждения. Мы выбрали сообщения в СУБД как новые записи во временном порядке. Для каждого сообщения мы ищем его столько же раз, сколько было просмотров. Совокупное соотношение чтения/записи составляет 99,9 к 0,1.

Диск и сообщения: набор данных форума объемом 10 ГБ, содержащий сообщения пользой, сканированные рядом общедоступных досок сообщений на базе vBulletin [11], которые охватывают широкий спектр тем, связанных с данными между собой, так же как спорт, автомобили и животные. Каждое сообщение содержит название форума, идентификатор темы, идентификатор сообщения, идентификатор пользы, и текст сообщения, включая цитаты из других сообщений. Это набор данных так же содержит количество просмотров, количество, которые мы используем для генерации синтетических запросов на чтение. Дублирование в основном происходит из-за того, что пользой цитируют комментарии других. Что бы имитировать поведение пользой на форуме обсуждений, для каждой вставки и сообщения мы делаем определенное количество «прочтений» (просмотров), которые запрашивают все предыдущие сообщения в содержащем просмотре. Количество прочтений просмотра на вставку выводит путь деления общего количества просмотров просмотра на количество содержащихся в нем сообщений.

5.2 Коэффициент сжатия и индексная память

Сначала мы оцениваем коэффициент сжатия dbDedup и использование памяти и индексации сравниваемых с trad-dedup и Snappy [9], компрессором MongoDB на уровне блоков по умолчанию. Для каждого набора данных мы загружаем записи в СУБД, как можно быстрее и измерим полученные размеры хранилища, объем данных, переданных по сети, и использование памяти и индексации.

На рис. 10 показаны результаты для пяти конфигураций: (1) dbDedup с фрагментами по 1 КБ или 64 байта, (2) trad-dedup с фрагментами по 4 КБ или 64 байта и (3) Snappy. Розовая (левая) полоса показывает коэффициент сжатия хранилища, указывающий на вклад только дубликации и сжатия после дубликации. Коэффициент сжатия определяется как исходный размер данных, деленный на размер сжатых данных, поэтому значение, равное единице, означает, что сжатие не достигнуто. Синяя (правая) полоса показывает использование памяти и индекса. Небольшой эскалометр записей (32 МБ, используется как dbDedup, так и trad-dedup) и кэш обратной записи потеряны (8 МБ, используется только dbDedup) не показаны.

Наибольшие преимущества наблюдаются для Wikipedia (рис. 10а). При размере фрагмента 1 КБ dbDedup сокращает хранилище данных в 26 раз (в 41 раз в сочетании со Snappy), используя 36 МБ индексной памяти. Уменьшение размера фрагмента до 64 Б увеличивает коэффициент сжатия до 37 раз (61 раз), используя всего 45 МБ индексной памяти. Уменьшение размера фрагмента для dbDedup не сильно увеличивает использование индексной памяти, поскольку dbDedup индексирует не более 4 записей на запись, независимо от размера фрагмента. Напротив, в то время как коэффициент сжатия trad-dedup увеличивается с 2,3 раз (3,7 раз) до 15 раз (24 раз) при использовании размера фрагмента 64 Б вместо 4 КБ, его индексная память увеличивается с 80 МБ до 780 МБ, что делает его непригодным для операционных СУБД. Это связано с тем, что trad-dedup индексирует каждый уникальный хэш фрагмента, что приводит к почти линейному увеличению накладных расходов индекса по мере уменьшения размера фрагмента, а также с тем, что он должен использовать гораздо больше ключей индекса (20-байтовый хэш SHA-1 против 2-байтовой контрольной суммы), поскольку коллизии приведут к повреждению данных. Потребляя на 40% меньше памяти и индекса, dbDedup с размером фрагмента 64 Б достигает коэффициента сжатия в 16 раз выше, чем trad-dedup с его типичным размером фрагмента 4 КБ. Snappy сжимает набор данных всего в 1,6 раза, поскольку он не может ускорить дублирование, выванное версионированием на уровне приложения, но не требует памяти и индекса. Он обеспечивает такое же сжатие в 1,6 раза при применении к дублицированным данным.

Для других наборов данных абсолютные преимущества меньше, но основные наблюдения схожи: dbDedup обеспечивает более высокую степень сжатия при меньшем использовании памяти, чем trad-dedup, а преимущества Snappy (1,6–2,3х) дополняют дубликацию. Для набора данных Enron (рис. 10b) dbDedup сокращает хранилище на 3,0х (5,8х), что согласуется с результатами, полученными нами в ходе экспериментов с данными из облачного развертывания серверов Microsoft Exchange, содержащих 16 реальных данных электронной почты пользой. Два набора данных форума (рис. 10c и 10d) демонстрируют столько же дублирования.

5К сожалению мы не можем раскрыть подробности из-за ограничений конфиденциальности.

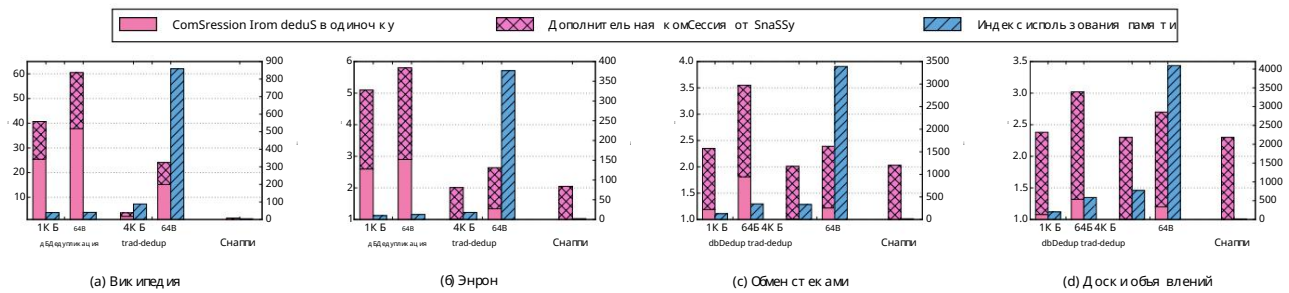


Рисунок 10: Коэффициент сжатия и индексная память – Коэффициент сжатия и использование индексной памяти для dbDedup (фрагменты по 1 К Б и 64 байт а), trad-dedup (4 К Б и 64 байт а) и Snappy. Верхняя часть каждой полоски представляет дополнительную преимущество сжатия после дедупликации.

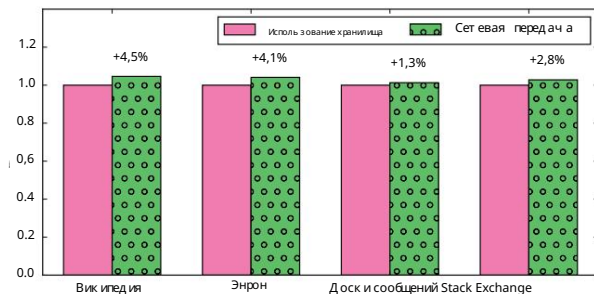


Рисунок 11: Экономия пропускной способности хранилища и сети — отнесенные к коэффициент сжатия, достигнутому с помощью dbDedup (с размером фрагмента 64 байт а) для локального хранилища и сетевой передачи для каждого из наборов данных, нормализованные по абсолютным коэффициентам сжатия хранилища, показанным на рис. 10 (для dbDedup с размером фрагмента 64 байт а).

tion как Wikipedia или набор данных электронной почты, потому что пользует ели не цитируют или не редактируют комментарии, а часть о, как редакция Wikipedia или пересылка и/от вет в электронной почте. Тем не менее, мы все равно наблюдаем, что dbDedup уменьшает хранилище в 1,3–1,8× (3–3,5×). Поскольку мы смогли проанализировать только последние сообщения в наборе данных Message Boards, коэффициент сжатия dbDedup является консервативным, не включая преимущество от дельта-сжатия редакций, пользующихся их собственными сообщениями x.6 B.

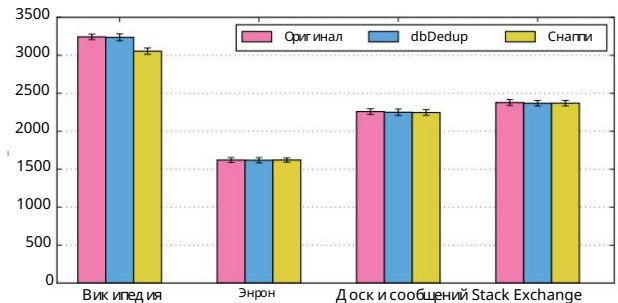
дополнение к использованию хранилища, dbDedup одновременно достигает значительного сжатия при передаче данных по сети с прямым кодированием. На рис. 11 показаны сжатие на уровне сети и как нормализованный результат по отношению к использованию хранилища (1,0 на оси Y для каждого набора данных). dbDedup достигает немного более низкого сжатия в хранилище баз данных, чем в данных, передаваемых по сети, в основном из-за перекрывающихся кодировок (раздел 3.2) и дельта-высечений из кэша обратной записи. Тем не менее, разница составляет менее 5% для всех наборов данных, поскольку к перекрывающимся кодировкам вст реч ается редко, а кэш обратной записи с пот еря ми использует приоритетное выделение.

5.3 Влияние на производительность во время выполнения. Это

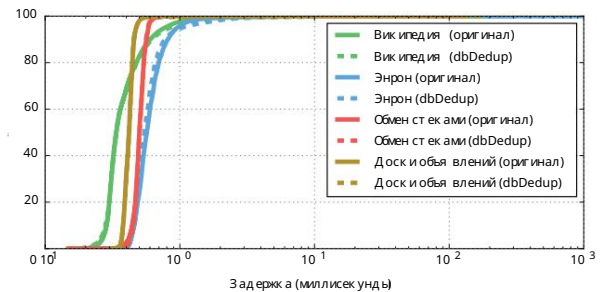
эксперимент предназначен для измерения влияния dbDedup на производительность СУБД. Мы сравниваем три конфигурации развертывания MongoDB: (1) Без сжатия («Исходная»), (2) dbDedup и (3) Snappy. Для каждой настройки мы запускаем эксперимент три раза для всех рабочих нагрузок и сообщаем среднее значение.

Пропускная способность: Рис. 12а показывает пропускную способность для четырех рабочих нагрузок. Мы видим, что dbDedup настраивает незначительные накладные расходы на пропускную способность. Snappy также немного ухудшает производительность для трех рабочих нагрузок, поскольку это быстрый и легкий встраиваемый компрессор. Исключением является Wikipedia, для которой использование Snappy приводит к снижению пропускной способности на 5%, поскольку куче больше записи Wikipedia

Мы обнаружили, что 0,15% сообщений редактируются как минимум один раз, а размер большинства от редактированных сообщений превышает средний размер сообщения.



(а) Пропускная способность



(б) Задержка

Рисунок 12: Влияние на производительность — из измерения производительности и задержки MongoDB во время выполнения для различных рабочих нагрузок и конфигураций.

не может поместиться на одной странице WiredTiger и требует дополнительной операции ввода-вывода.

Задержка: Рис. 12b показывает CDF задержки клиента. Для ясности мы показываем результат только для MongoDB с включенным dbDedup и без него. Опять же, мы видим, что dbDedup почти не влияет на производительность. Кривые распределения задержки с включенным dbDedup близки к соответствующим кривым без сжатия / дедупликации. Разница в задержке 99,9%-тилю и составляет менее 1% для всех рабочих нагрузок.

5.4 Эффекты эширования dbDedup

использует два специализированных кэша для минимизации накладных расходов ввода-вывода, связанных с чтением и обновлением исходных записей: кэш исходных записей (32 МБ) и кэш обратной записи с потерями (8 МБ). Теперь мы оцениваем эффект от этих кэшей.

Кэш исходной записи: Рис. 13а показывает влияние кэша исходной записи на коэффициент сжатия (левая ось Y) и процент извлечений исходной записи, требующих чтения СУБД (коэффициент промахов кэша; правая ось Y), с диапазоном значений оценок и вознаграждения для рабочей нагрузки и Wikipedia. Напомним, что dbDedup использует выборочное эширование записей с учетом кэша, назначая оценку вознаграждения кандидатам, которые присутствуют в кэше (см. Раздел 3.1.3).

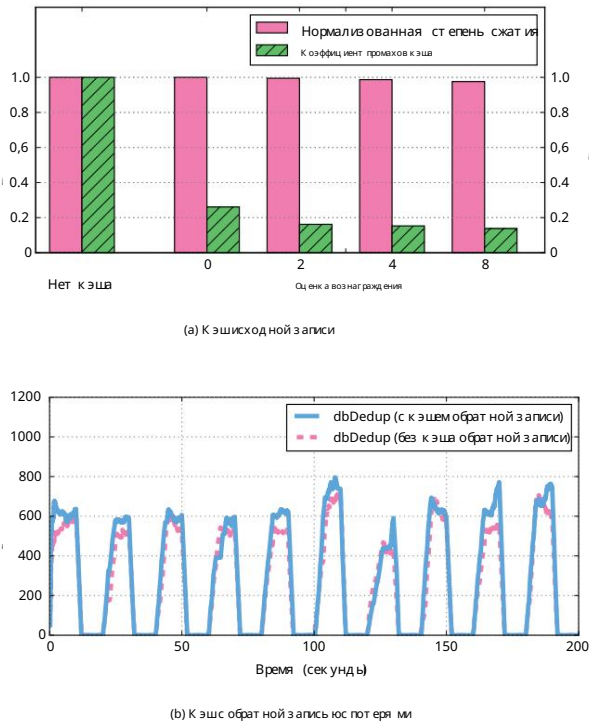


Рисунок 13: Эффект кэширования – измерения времени выполнения dbDedup механизмы зашифрования для рабочей нагрузки в Википедии

Когда кэшные используются (самые левые столбцы), каждое извлечение Исходная запись вызывает запрос на чтение. Даже без выбора с учетом оmissions (0 баллов вознаграждения) небольшой кэш исходной записи устаревает 74% эти запросов. Соотношение вознаграждения два (по умолчанию) метод выбора с учетом оmissions еще больше сокращает коэффициент промахов на 40% (до 16%), без заметного снижения степени сжатия. Далее увеличение оmissions вознаграждения незначительно не снижает коэффициент промахов кэша приводит к небольшому уменьшению степени сжатия, поскольку кэширование похоже Кандидаты больше вероятности будут выбраны в качестве исходных записей.

Кэш обратной записи с помощью dbDedup использует обратное кодирование что бы избежать дублирования причинения последних «версий» последовательности и обновления. Таким образом, дублирование новой записи включает в себя как запись полностью новая запись и замена исходной записи на дельта-кодированную данные. Дополнительная запись (замена) может привести к значительным проблемам с производительностью для интенсивных рабочих нагрузок ввода-вывода во время всплесков записи. Кэш обратной записи с помощью dbDedup снимает такие проблемы

Чтобы имитировать пиковую нагрузку с интенсивным вводом-выводом и периодами простоя, мы вставили данные Википедии на полную скорость и в течение 10 секунд и переходим в спящий режим. в течение 10 секунд, многократно. Рис. 13b показывает вставку MongoDB производит производительность с кэшем обратной записи и без него. Без кэша производительность СУБД заметно снижается в периоды отсутствия и из-за дополнительных записей в базу данных. Напротив, при использовании кэш обратной записи позволяет избежать замедления работы СУБД во время пиковых нагрузок, как показано на разнице между двумя линиями в различных точках времени (например, в секунды 0, 130, 170 и 190).

dbDedup имеет два основных настраиваемых параметра, помимо рассмотренных выше, которые влияют на компромисс между степенью сжатия и производительностью. Расстояние и интервал между кэшированием. Это подраздел который определяет эффект этих параметров и поясняет значения по умолчанию

5.5 Кодирование переходов

dbDedup использует кодирование переходов для уменьшения наихудшего случая извлечения

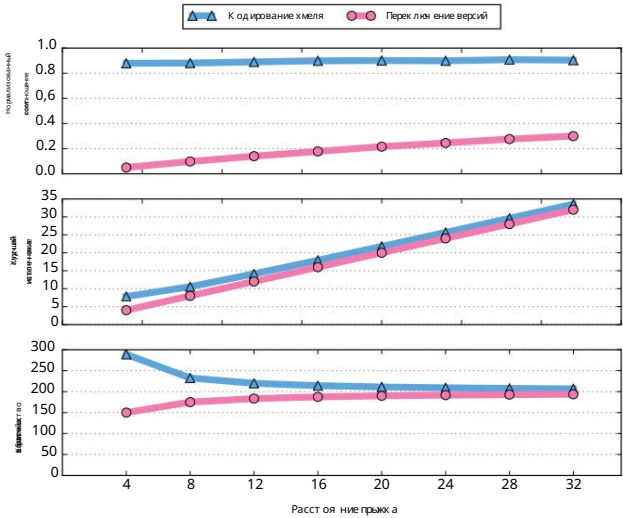


Рисунок 14: Кодирование скачков против перехода между версиями для Википедии рабочей нагрузки и умеренные расстояния перехода, кодирование перехода обеспечивает гораздо более высокую коэффициент сжатия с небольшим увеличением худшего случая извлечения источника и количество обратных записей.

раз, сохраняя преимущество сжатия. Чуть уменьшает его эффективность, мы также реализовали переход между версиями в MongoDB и сравнили две схемы кодирования.

На рис. 14 показаны результаты для трех методов в зависимости от прыжка. расстояние: коэффициент сжатия (нормализованный по отношению к стандартному обратному кодированию), наихудший источник извлечения источника (для кодирования Длина ценой к 200) и количество обратных записей. Переклечение версий приводит к значительному (на 60-90%) снижению степени сжатия, поскольку кэширование справочные версии хранятся в неактивном виде. Его степень сжатия улучшается по мере увеличения расстояния прыжка, поскольку меньше записей хранятся в неактивном виде. В отличие от этого, поскольку хмелевые основания хранятся как дельты кодирования скачков обеспечивает коэффициент сжатия в пределах 10% Полное обратное кодирование. Для кодирования прыжка коэффициент сжатия остается относительно стабильным по мере увеличения расстояния прыжка из-за наличия меньше, но менее схожие хмелевые базы

Число наихудших случаев получения источника для кодирования прыжка близок к версии jumping. С несколькими уровнями прыжков, отслеживание ближайшей базы хмеля занимает всего лишь логарифмическое время. По мере увеличения расстояния прыжка время декодирования доминирует прохождение обратных дельт между соседними базами прыжков. Нижний график показывает количество дополнительных обратных записей, необходимых в каждой схеме. В то время как кодирование прыжка вызывает больше обратных записей для небольших прыжков расстояние, обе схемы быстро приближаются к длине кодирующей цепи по мере увеличения расстояния прыжка. Эмпирически мы обнаруживаем, что прыжок расстояние 16 (по умолчанию) обеспечивает хороший компромисс между степенью сжатия и затратами на декодирование.

5.6 Оптимизация дельта-сжатия

dbDedup превосходит алгоритм Delta, сокращая накладные расходы вставки исходного индекса и поиска. Он вводит настраиваемый интервал привязки к отсрочке управления частотой выборки и смещение точек в исходном потоке байтов.

На рис. 15 показана степень сжатия (левая ось Y) и пропускная способность (правая ось Y) для различных значений интервала привязки для сравнения с xDelta, для рабочей нагрузки в Wikipedia. С интервалом привязки и 16 (размером по умолчанию в xDelta), dbDedup выполняет почти то же самое, что и xDelta. Скорость дельта-сжатия dbDedup улучшается с увеличением интервала привязки, поскольку оно уменьшает количество вставок и поисков в исходном индексе смещения.

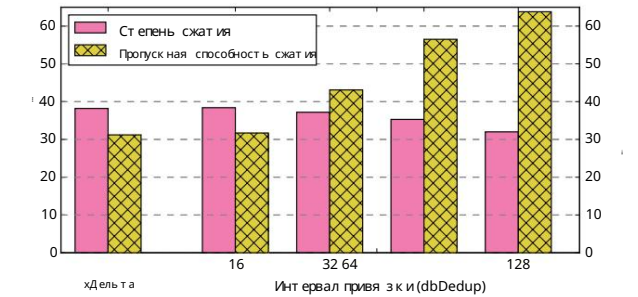


Рисунок 15: Оптимизация дельта-сжатия — сравнение оптимизированного варианта dbDedup и xDelta с использованием рабочей нагрузки из Wikipedia.

Коэффициент сжатия не уменьшается существенно, поскольку у dbDedup выполняется сравнение на уровне байтов в двух направлениях относительно сжатых данных. При интервале привязки 64 dbDedup превосходит xDelta на 80% с точки зрения пропускной способности сжатия, при этом теряется всего 7% коэффициента сжатия. Увеличение интервала привязки до 128 дополнительно увеличивает пропускную способность на 10%, приводит к потере 15% коэффициента сжатия. Мы используем 64 в качестве значения по умолчанию, чтобы обеспечить разумный баланс между коэффициентом сжатия и пропускной способностью.

6. ДОПОЛНИТЕЛЬНАЯ СВЯЗАННАЯ РАБОТА

Большая часть предыдущей работы по дедупликации обсуждается в разделе 2. Этот раздел обсуждается некоторая дополнительная связанная работа.

Сжатие баз данных: За последние несколько десятилетий было предложено несколько схем сжатия баз данных. Большинство операций СУБД, которые сжимают содержимое базы данных, используют сжатие на уровне страниц или блоков [30, 37, 46, 43, 3, 16]. Некоторые используют сжатие префиксов, которое ищет общие последовательности и в начале значимых полей для данного столбца по всем строкам на каждой странице.

Как и в случае с нашим подходом dbDedup, такое сжатие требует, чтобы СУБД распаровывала строки перед обработкой запроса.

В некоторых системах OLAP существуют схемы, которые позволяют СУБД обрабатывать данные в сжатом формате. Например, сжатие словаря заменяет повторяющиеся длинные значения домена короткими целыми значениями фиксированной длины. Этот подход обычно используется в хранилищах данных, ориентированных на столбцы [18, 36, 69, 51]. Эти системы обычно фокусируются на атрибутах, которые несут небольшую переменную и исследуют перекрестные аспекты значений, чтобы ускорить словарь управления размером [23]. Авторы [56] предлагают схему дельта-кодирования, где каждое значение в отсортированном столбце представлено дельтой от предыдущего значения. Хотя этот подход хорошо работает для числовых значений, он не подходит для строк.

Ни один из этих методов не обнаруживает и не устраняет избыточные данные с granularity, меньшей, чем одного поля, тем самым теряя потенциальные преимущества сжатия для многих приложений, которые изначально содержат такую избыточность. dbDedup, напротив, способен удалять гораздо более мелкие дубликаты с помощью дельта-сжатия на уровне байтов. В отличие от других схем строенного сжатия, dbDedup не находится на критическом пути записи для запросов и, следовательно, оказывает минимальное влияние на производительность СУБД во время выполнения. В дополнение к этому, поскольку dbDedup сжимает данные на уровне записи, он выполняет шаг дедупликации только один раз и использует закодированный результат как для хранения в базе данных, так и для передачи и по сети. Напротив, одна и та же запись будет сжата дважды (на странице баз данных и в пакете протокола) для схем сжатия страниц для доставки сжатия сжатия данных на обоих уровнях.

Дельта-кодирование: Было проведено больше исследований по предыдущим работам по методам дельта-кодирования, включая несколько универсальных алгоритмов, основанных на подходе Лемпела-Зива [68], таких как vcd-

iff [21], xDelta [42] и zdelta [60]. Специализированные схемы могут использоваться для определенных форматов данных (например, XML) для улучшения качества сжатия [28, 63, 39, 52]. Алгоритм дельта-сжатия, используемый в dbDedup, адаптирован из xDelta, связь с которым обсуждается в разделе 4.2.

Дельта-сжатие использовалось для сокращения сетевого трафика для потоков передатчиков и файлов и синхронизации. Большинство систем предполагают, что предыдущие версии одного и того же файла являются идентичными приложениям и дублирование существует только в средних версиях одного и того же файла [61, 57]. Исключением является TAPER [38], который сокращает сетевую нагрузку для синхронизации реплик файловой системы путем отправки дельта-кодированных файлов; он идентифицирует похожие файлы, выискивая коллизии в совпадающих битовых потоках, сгенерированных с помощью хэшей фрагментов файлов. dbDedup идентифицирует похожую запись из корпуса данных без необходимости приложения и, следовательно, является более общим подходом, чем большинство этих предыдущих систем.

Метод обратного кодирования, используемый в dbDedup, вдохновлен системами хранения данных с управлением версиями, такими как RCS [59] и XDFS [42]. Похожие методы использовались в системах контроля версий, таких как Git [41] и SVN [29], чтобы обеспечить возможность перемещения назад по истории коммитов. В отличие от этих систем, которые в основном поддерживают основную версию для всех файлов, dbDedup устаревает цепочку кодирования полностью на основе отклонений сходства между записями, таким образом, не требует поддержки на уровне системы для управления версиями. [54] использует дельта-кодирование для дедупликации резервного хранилища. Он использует прямое кодирование и допускает только коллекцию кодирования с максимальной длиной в два элемента. Подобно другим системам хранения с дельта-кодированием, которые используют зумп-переходы между версиями, он должен пожертвовать выигрышем в сжатии, чтобы ограничить наихудшие случаи извлечения для базовых данных.

Несколько названий, dbDedup — первая система, которая исследует различные формы кодирования для сжатия на уровне сетевых хранилищ и обеспечивает эффективные преобразования между ними. Благодаря новой схеме кодирования с переходом от dbDedup значительное облегчение болезненных компромиссов между степенью сжатия и издержками извлечения для дельта-кодированного хранилища. Кроме того, он вводит новые механизмы кодирования, специализированные для дельта-кодированного хранилища, значительное сокращение издержек ввода-вывода, при этом максимальное увеличение эффективности памяти. Все эти методы в совокупности делают онлайн-доступ к закодированному хранилищу практичным.

Обнаружение сходства: Предыдущие работы представили различные подходы к выискиванию эскизов (метрик сходства) для идентификации схожих элементов. Базовая техника идентификации признаков в объектах так же, как и в объектах, чтобы схожие объекты имели идентичные признаки, была впервые предложена Бродером [24, 25] в контексте веб-страниц. Несколько статей [58, 47, 20, 53, 65] предлагают методы выискивания эскизов для обнаружения сходства, которые являются надежными небольшими изменениями в данных. Подход к извлечению признаков, используемый в dbDedup, аналогичен подходу в DOT [47] и sDedup [65].

7. ЗАКЛЮЧЕНИЕ

dbDedup — это легкий механизм дедупликации на основе подбора для операций СУБД, который сокращает как использование хранилища, так и объем данных, переданных для удаленной репликации. Объединяя чистый нулевой индекс дельта-сжатия на уровне байтов, dbDedup достигает более высоких коэффициентов сжатия, чем сжатие на уровне блоков и дедупликация на основе фрагментов, при этом эффективно используя память. Он использует новые механизмы кодирования и кодирования, чтобы избежать значительных накладных расходов ввода-вывода, связанных с доступом к записи с дельта-кодированием. Экспериментальные результаты с участием реальных рабочих нагрузок показывают, что dbDedup способен достичь сокращения до 37 раз (в 61 раз при сочетании сжатия на уровне блоков) размера хранилища и трафика репликации, при этом накладывая незначительные накладные расходы на производительность СУБД.

Благодарности Мухомели бы благодарить Кейт Бостика за его руководство по внутренним компонентам WiredTiger. Мы благодарим анонимных рецензентов за полезные отзывы. Мы также благодарим членов и компаний консорциума PDL (включая Broadcom, Citadel, Dell EMC, Facebook, Google, Hewlett-Packard Labs, Hitachi, Huawei, Intel, Microsoft Research, NetApp, Oracle, Samsung, Seagate, Tintri, Two Sigma, Uber, Veritas и Western Digital) за их интерес, идеи, отзывы и поддержку. Это исследование было частично спонсировано Intel в рамках Научно-технического центра Intel по облачным вычислениям (ISTC-CC) и MongoDB Incorporated. Эксперименты были проведены благодаря щедрому пожертвованию оборудования от Intel и NetApp.

8. ССЫЛКИ

- [1] Байду Байе. <http://baidu.baidu.com/>.
- [2] Набор данных электронной почты Enron. <https://www.cs.cmu.edu/~enron/>.
- [3] Сжатие InnoDB. <http://dev.mysql.com/doc/refman/5.6/en/innodb-compression-internals.html>.
- [4] Linux SDFS. www.openedup.org.
- [5] MongoDB. <http://www.mongodb.org>.
- [6] MurmurHash. <https://sites.google.com/site/murmurhash>.
- [7] Ocarina Networks. www.ocarinanetworks.com.
- [8] Оптимизация данных Permabit. www.permabit.com.
- [9] Snapki. <http://google.github.io/snapki/>.
- [10] Архив данных Stack Exchange. <https://archive.org/details/stackexchange>.
- [11] Вбллетень. <https://www.vbulletin.com>.
- [12] W3Techs. <http://www.w3techs.com>.
- [13] Загрузки и Wikimedia. <https://dump.wikimedia.org>.
- [14] Википедия. <https://www.wikipedia.org/>.
- [15] Сервер хранения Windows. [http://technet.microsoft.com/en-us/library/gg232683\(WS.10\).aspx](http://technet.microsoft.com/en-us/library/gg232683(WS.10).aspx).
- [16] WiredTiger. <http://www.wiredtiger.com/>.
- [17] Дедупликация ZFS. blogs.oracle.com/bonwick/entry/zfs_dedup.
- [18] Д. Абади, С. Мэдден и М. Феррейра. Интеграция сжатия и выполнение в системах столбчатых баз данных. В SIGMOD, страницы 671–682, 2006.
- [19] C. Alvarez. Руководство по дедупликации NetApp для FAS и V-Series и ее внедрению 2010.
- [20] Л. Аронович, Р. Ашер, Э. Бахмат, Х. Бигнер, М. Хирш и ST Klein. Проектирование систем дедупликации на основе сходства. В SYSTOR, страница 6, 2009.
- [21] Дж. Бенгли и Д. Макилрой. Сжатие данных с использованием длинных общих строк. На конференции по сжатию данных, 1999. Труды DCC'99, страницы 287–295, 1999.
- [22] Д. Бхагват, К. Эшги, Д. Лонг и М. Лилибридж. Экстремально-биннинг: Масштабируемая параллельная дедупликация для резервного копирования файлов на основе фрагментов. В MASCOTS, страницы 1–9, 2009.
- [23] C. Binnig, S. Hildenbrand и F. Färber. Сжатие строк с сохранением порядка на основе словаря для хранилищ столбцов основной памяти. В SIGMOD, страницы 283–296, 2009.
- [24] Бродер А. Сходство и содержание документов: Сжатие и сложность последовательностей, 1997.
- [25] Бродер А. Выведение и фильтрация почтовых дубликатов документов 11-й ежегодный симпозиум по комбинаторному сопоставлению образов, 2000.
- [26] RC Burns и DD Long. Эффективное распределенное резервное копирование с дельта-сжатием. В Трудях пятого семинара по вводу-выводу в параллельных и распределенных системах, страницы 27–36,

1997.

- [27] А. Клементс, И. Ахмад, М. Вилайаннури и Дж. Ли. Децентрализованная дедупликация в кластерных файловых системах SAN. В USENIX ATC, 2009.
- [28] G. Cobena, S. Abiteboul и A. Marian. Обнаружение изменений в документах XML. В ICDE, страницы 41–52, 2002.
- [29] Б. Коллинз-Сасман, Б. Фицпатрик и М. Пилато. Версия контроля с помощью юсубверсии. 2004.
- [30] Кормак Г. В. Сжатие данных в системах баз данных. Сообщения ACM, 28(12):1336–1342, 1985.
- [31] Б. Дебнат, С. Сентуга и Дж. Ли. Chunkstash: Ускорение дедупликации встоего хранилища с использованием флешпамяти. На ежегодной технической конференции USENIX, 2010.
- [32] П. Дойч и Ж.-Л. Гейлли. Спецификация формата сжатых данных Zlib версии 3.3. Технический отчет, 1996.
- [33] К. Дубницкий, Л. Грив, Л. Хелдт, М. Качмарик, В. Килиан, П. Шельчак и Й. Шенксовски. Hydrastor: масштабируемое встраиваемое хранилище. В FAST, 2009.
- [34] К. Дубницкий, Л. Грив, Л. Хелдт, М. Качмарик, В. Килиан, П. Шельчак, Й. Шенксовски, К. Унгуряну и М. Велницкий. HYDRASrор: масштабируемое встраиваемое хранилище. В FAST, 2009.
- [35] А. Эль-Шими, Р. Качал, А.К. Адди, О.Дж. Ли и С. Сентуга. Первая дедупликация данных — крупномасштабное исследование и проектирование системы. На ежегодной технической конференции USENIX, 2012.
- [36] C. Харвиопулос, В. Лян, Д. Абади и С. Мэдден. Компромиссы производительности в базах данных, оптимизированных для чтения. В VLDB, страницы 487–498, 2006.
- [37] Б. Айери и Д. Уилхайт. Поддержка сжатия данных в базах данных. 1994.
- [38] Н. Джейн, М. Далин и Р. Тевари. Тарет: многоуровневый подход к ускорению избыточности при синхронизации и репликации. В FAST, 2005.
- [39] E. Leonardi и SS Bhowmick. Xanadue: система обнаружения изменений в XML-данных в реляционных базах данных, не поддерживающих рекурсивную структуру. В SIGMOD, страницы 1137–1140, 2007.
- [40] М. Лилибридж, К. Эшги, Д. Бхагват, В. Деолаликар, G. Trezise и P. Camble. Разреженное индексирование: крупномасштабная встраиваемая дедупликация с использованием выборки и локальности. В FAST, 2009.
- [41] Дж. Лелигер. Контроль версий с помощью югит: мощные инструменты для совместной разработки и программного обеспечения. 2009.
- [42] Дж. П. Макдональд. Поддержка файловой системы для дельта-сжатия. Магистерская диссертация, Калифорнийский университет, Беркли, 2000.
- [43] С. Мишра. Сжатие данных: стратегия, планирование емкости и лучшие практики. Техническая статья сайта SQL Server, 2009.
- [44] А. Мутитчарен, Б. Чен и Д. Мазьер. А сетевая файловая система с низкой пропускной способностью. В SOSR, 2001.
- [45] Р. Паг и Ф. Родлер. Улучшенное хеширование. Журнал алгоритмов, 51(2):122–144, 2004.
- [46] М. Поэссид, Потлов. Сжатие данных в Oracle. В VLDB, страницы 937–947, 2003.
- [47] Х. Пуча, Д. Г. Андерсен и М. Каминский. Эксплуатация сходства для загрузки нескольких объектов с использованием печатных файлов. В NSDI, 2007.
- [48] В. Пью и Список пропусков: вероятностная альтернатива сбалансированному дереву. В Workshop on Algorithms and Data Structures, страницы 437–449. Springer, 1989.
- [49] С. Куинлан и С. Дорвард. Venti: новый подход к архивному хранению. В FAST, 2002.
- [50] Рабин М. О. Дактилоскопия случайными полиномами.

- [51] В. Раман, Г. Атталури, Р. Барбер, Н. Чайнани, Д. Калмук, В. Куландай Сами, Дж. Леенстра, С. Лайтстун, С. Лю Г. М. Ломан и др. Db2 с использованием Blu: гораздо больше, чем просто хранилище столбцов. VLDB, 6(11):1080–1091, 2013.
- [52] С. Сакр. Методы сжатия XML: обзор и сравнение. Журнал компьютерных и системных наук, 75(5):303–322, 2009.
- [53] П. Шилан, М. Хуан, Г. Уоллес и В. Сюй. WAN-оптимизированная репликация резервных наборов данных с использованием потоковой дельта-компрессии. В FAST, 2012.
- [54] П. Шилан, Г. Уоллес, М. Хуан и В. Сюй. Дельта-сжатие и дедупликация хранилища с использованием потоковой локальности. USENIX Hot Storage, 2012.
- [55] К. Шринивасан, Т. Биссон, Г. Гудсон и К. Воргант. idedup: Учтивая задержку встроенная дедупликация данных для первичного хранилища. В FAST, 2012.
- [56] М. Стонфрейкер, DJ Абди, А. Баткин, Х. Чен, М. Черняк, М. Феррейра, Э. Лау, А. Лин, С. Мэдден, Э. О'Нил и др. C-store: СУБД, ориентированная на столбцы в VLDB, страницы 553–564, 2005 г.
- [57] Т. Суэлл и Н. Мемон. Алгоритмы дельта-сжатия и удаленной синхронизации файлов. Справочник по сжатию без потерь, 2002.
- [58] Д. Теодосиу, Ю. Гуревич, М. Манассе и Ж. Порка. Оптимизация репликаций файлов в сетях с ограниченной пропускной способностью с использованием удаленного дифференциального сжатия. Технический предсказатель. MSR-TR-2006-157, Microsoft Research, 2006.
- [59] WF Tichy. Rcs – система контроля версий. Программное обеспечение: практика и опыт, 15(7):637–654, 1985.
- [60] Д. Трендафилов, Н. Мемон и Т. Суэлл. delta: Эффективный Инструмент дельта-компрессии. Технический отчет TR-CIS-2002-02, Политехнический университет, 2002.
- [61] А. Тиджелл. Эффективные алгоритмы сортировки и синхронизации. В докторской диссертации, Австралийский национальный университет, 2000.
- [62] Г. Урданет, Г. Пьер и М. Ван Стин. Википедия: Анализ рабочей нагрузки и для децентрализованного хостинга. Компьютерные сети, 53(11):1830–1845, 2009.
- [63] Y. Wang, DJ DeWitt и J.-Y. Cai. X-diff: Эффективный алгоритм обнаружения изменений для XML-документов. В ICDE, страницы 519–530, 2003.
- [64] В. Ся, Х. Цзянь, Д. Фэн и Ю. Хуа. Сирос: А Схема дедупликации на основе сходства локальности и слияния издержек оперативной памяти и высокой пропускной способностью. На ежегодной технической конференции USENIX, 2011.
- [65] Л. Сюй, А. Павло, С. Сенгупа, Дж. Ли и Г. Р. Гангер. Сокращение пропускной способности репликаций для распределенных баз данных. В SoCC, страницы 222–235, 2015.
- [66] LL You, KT Pollack и DD Long. Глубокое хранилище: архитектура архивной системы хранения. В ICDE, страницы 804–815, 2005.
- [67] Б. Чжу, К. Ли и Р. Х. Паттерсон. Избегание диска узкое место в файловой системе дедупликации домена данных. В FAST, 2008.
- [68] Дж. Зив и А. Лемпель. Универсальный алгоритм для последовательного сжатия данных. Труды IEEE по теории информации, 23(3):337–343, 1977.
- [69] М. Жук овски, С. Хеман, Н. Нес и П. Бонц. Суперсжатие кэш RAM-CPU. В ICDE, страницы 59–59, 2006.