

Иванченко Диана Олеговна, магистрант,
ФГАОУ ВО "МГТУ "СТАНКИН"

Гаврилов Андрей Геннадьевич, к.т.н., доцент,
ФГАОУ ВО "МГТУ "СТАНКИН"

МЕТОДИКА ВОССТАНОВЛЕНИЯ ИСТОРИЧЕСКОЙ ДОСТОВЕРНОСТИ ДАННЫХ В HADOOP НА ОСНОВЕ МЕХАНИЗМОВ СМЕЩЕНИЙ И «СЫРЫХ» СЛОЕВ ХРАНЕНИЯ

Аннотация. В статье предлагается методика обеспечения целостности и исторической достоверности данных при их передаче из операционных CRM-систем в аналитические хранилища на базе Hadoop. Основное внимание уделено стратегии восстановления информации (Recovery Point Objective) в случае программных или логических сбоев. Описаны механизмы использования технических смещений (offsets) в брокере сообщений Apache Kafka и организация «сырого» слоя хранения (landing zone) в Hadoop как инструментов гарантированной минимизации потерь данных.

Ключевые слова: Big Data, Hadoop, Apache Kafka, смещения (offsets), «сырой» слой данных, отказоустойчивость, историческая достоверность, идиempotentность.

В условиях современной цифровой экономики операционные CRM-системы становятся источником колоссальных объемов транзакционных данных, необходимых для глубокой аналитики. Однако архитектурная изолированность реляционных СУБД, на которых базируются CRM, создает барьеры для обработки исторических массивов информации. Создание сквозного конвейера данных для их консолидации в масштабируемой среде Hadoop требует не только высокой пропускной способности, но и бескомпромиссной отказоустойчивости. Любой технический сбой на этапе интеграции может привести к потере «исторической правды», что сделает невозможным построение достоверных предиктивных моделей.

Архитектурное обоснование механизмов восстановления

Проектируемая архитектура интеграции базируется на парадигме Event-Driven Architecture (EDA), где связующим звеном выступает распределенный брокер сообщений Apache Kafka. Выбор данного инструмента обусловлен его способностью хранить сообщения в неизменяемом логе транзакций (commit log), что позволяет изолировать контексты функционирования CRM-системы и Hadoop-кластера.

Фундаментальным инструментом восстановления данных в этой схеме выступает механизм технических смещений (offsets). Смещение – это уникальный порядковый номер, присваиваемый каждому сообщению внутри топика Kafka. В рамках разработанной методики система-потребитель (Hadoop) фиксирует последний успешно обработанный offset непосредственно в своей базе данных.

Методика восстановления на основе смещений Kafka

Основной риск при передаче данных связан с ситуацией, когда сообщение было получено из брокера, но не было успешно записано в хранилище из-за сетевого сбоя или аварии на стороне потребителя. Традиционные системы, ориентированные на дату внутри бизнес-сообщения, в таких случаях могут терять данные из-за временных задержек в источнике.

Предложенная методика реализует следующий алгоритм восстановления:

1. Мониторинг прогресса: Система-потребитель ориентируется не на бизнес-время события, а на физическую позицию в логе Kafka.



2. Политика хранения (Retention Policy): В брокере настраивается окно хранения сообщений N дней (N подбирается в зависимости от ресурсов системы). Этого времени должно быть достаточно для обнаружения и устранения последствий масштабных аварий.

3. Откат смещения: В случае обнаружения пробелов в данных ИС-потребитель инициирует повторную выгрузку путем «отката» технического смещения на последний сохраненный индекс.

Благодаря неизменяемости лога Kafka такой подход гарантирует, что данные будут восстановлены в той же последовательности, в которой они поступили от CRM-системы, обеспечивая соблюдение принципа идемпотентности.

Организация «сырого» слоя хранения в Hadoop

Вторым уровнем защиты исторической достоверности выступает выделенный «сырой» слой хранения (landing zone) внутри экосистемы Hadoop. В отличие от распределенного брокера, который оперирует исходными потоками сообщений, данный слой в Hadoop представляет собой копию уже обработанных данных из основного хранилища.

Согласно методике, срок хранения информации в этом слое ограничен выбранным сроком (срок зависит от объемов системы), что создает временную репликацию данных с заданным «сроком жизни». Наличие такой копии делает систему-приемник независимой от оперативных журналов брокера сообщений. Если в процессе агрегации данных в итоговые аналитические таблицы будет обнаружена логическая ошибка трансформации, информация может быть повторно извлечена из «сырого» слоя и переобработана. Это критически важно для минимизации нагрузки на операционные базы данных CRM-системы, так как восстановление происходит исключительно внутренними ресурсами Hadoop.

Обеспечение уникальности и хронологии

Для исключения дублирования информации при повторных выгрузках используется специальный технический атрибут времени. Данный параметр фиксирует точное время физической загрузки записи в систему-источник и выступает единственным уникальным ключом в историческом хранилище. Использование этого ключа в сочетании с меткой времени фактической обработки позволяет корректно восстанавливать историю изменений объектов даже при асинхронном поступлении данных.

Заключение

Разработанная методика восстановления данных, сочетающая механизмы управления смещениями в Apache Kafka и создание временной реплики обработанных данных в Hadoop, формирует отказоустойчивый фундамент для работы с Большими данными. Предложенный подход позволяет гарантировать сохранность транзакционной информации в течение определенного срока после возникновения сбоя, обеспечивая полную прозрачность аналитических отчетов. Практическая значимость работы заключается в возможности создания надежных интеграционных конвейеров, способных восстанавливать данные без прямого воздействия на производительность операционных CRM-систем.

Список литературы:

1. Галигузова Е. В., Илларионова Ю. Е. Сравнение реляционных и нереляционных СУБД // Символ науки. 2023. №1-2.
2. Селезнёв А. И., Селезнёв И. Л. Особенности организации конвейера данных с использованием брокера сообщений Apache Kafka в системах обработки данных // Молодой ученый. 2025. № 48 (599). С. 15-19.
3. Соломонов А. А. Оптимизация ETL-процессов для больших данных // Вестник науки. 2024. №9 (78).



4. Розен Н. Б. ETL-системы и базы данных. – 2023.

5. Пантелеева А. И. Интеграция больших данных и облачных платформ для анализа влияния экономической политики на финансовые рынки // Научный журнал. 2024.

