

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

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

АРХИТЕКТУРНЫЕ ПРЕИМУЩЕСТВА ИСПОЛЬЗОВАНИЯ АРАСНЕ КАФКА ДЛЯ ОБЕСПЕЧЕНИЯ ОТКАЗОУСТОЙЧИВОСТИ В ИНТЕГРАЦИОННЫХ КОНТУРАХ CRM-НАДООР

Аннотация. В статье исследуются методы построения отказоустойчивого интеграционного решения для передачи данных из операционных CRM-систем в аналитические платформы на базе Hadoop. Проанализирована роль распределенного брокера сообщений Apache Kafka как ключевого элемента обеспечения технологической надежности. Описаны механизмы слабой связности, гарантированной доставки и использования технических смещений для восстановления данных после сбоев. Особое внимание уделено алгоритмам восстановления данных при сбоях на основе «сырго» слоя хранения в Hadoop и механизмов управления смещениями (offsets) в брокере сообщений.

Ключевые слова: Big Data, Apache Kafka, Hadoop, CRM-система, отказоустойчивость, интеграция данных, брокер сообщений.

Современные системы обработки данных сталкиваются с необходимостью управления взрывообразно растущими информационными потоками. Операционные CRM-системы, построенные на реляционных СУБД, обеспечивают высокую целостность транзакций, однако их архитектурная изолированность создает барьеры для глубокого анализа больших массивов исторической информации. Для решения задач консолидации данных в масштабируемых хранилищах, таких как Hadoop, критически важным становится обеспечение отказоустойчивости интеграционного решения: сбой в любом узле не должен приводить к потере данных или деградации производительности CRM-системы.

Архитектурная роль Apache Kafka и принципы слабой связности

Проектируемая архитектура интеграционного взаимодействия базируется на принципах парадигмы Event-Driven Architecture (EDA) и концепции слабой связности (loose coupling). Распределенный брокер сообщений Apache Kafka выступает в роли надежного буфера, разделяющего производителей данных (CRM) и потребителей (Hadoop). Это позволяет изолировать контексты функционирования систем: сообщения накапливаются в брокере до момента их готовности к обработке приемником, что исключает риск отказа из-за перегрузки.

Для обеспечения высокой пропускной способности и горизонтального масштабирования системы используется механизм партиционирования топиков, позволяющий распределять нагрузку между узлами кластера. Такой подход в сочетании с распределенной природой Hadoop, обеспечивающей отказоустойчивость за счет программной избыточности и репликации данных, формирует стабильный фундамент для обработки терабайтов информации.

Механизмы обеспечения технологической надежности

Отказоустойчивость в интеграционном контуре достигается за счет сочетания нескольких механизмов:

1. **Персистентность и неизменяемость.** Каждое сообщение в Kafka записывается в неизменяемый лог транзакций на дисковое пространство и имеет возможность реплицирования, что гарантирует сохранность данных даже при выходе из строя отдельных серверов.



2. **Гарантированная доставка.** На стороне CRM-системы реализован алгоритм повторных попыток отправки в случае сетевых сбоев. Поддержка режимов «как минимум один раз» (at-least-once) и «строго один раз» (exactly-once) исключает риск потери или дублирования транзакций.

3. **Управление техническими смещениями (Offsets).** Использование смещений позволяет потребителю точно фиксировать позицию в потоке. В случае аварии система инициирует повторную загрузку данных путем «отката» смещения на нужный порядковый номер, используя окно хранения сообщений от 1 до N дней (N подбирается в зависимости от ресурсов Kafka).

4. **Многоуровневая верификация.** Отказоустойчивость логики обеспечивается двухэтапным контролем: технической валидацией формата JSON на уровне брокера и семантическим анализом бизнес-атрибутов на стороне потребителя.

Многоуровневая верификация и механизмы восстановления данных

Отказоустойчивость логики обеспечивается двухэтапным контролем: технической валидацией формата (JSON) на уровне брокера и семантическим анализом бизнес-атрибутов на стороне потребителя.

Фундаментальным аспектом методики является стратегия восстановления информации при критических сбоях. В Hadoop организован выделенный «сырой» слой хранения, где данные, которые были уже обработаны на стороне ИС-потребителя, хранятся в течение одного месяца. Наличие этого слоя позволяет повторно прогрузить данные в основную таблицу при логических ошибках в основных аналитических таблицах без обращения к системе-источнику.

Однако, в сценарии, когда данные не были успешно записаны даже в «сырой» слой Hadoop (например, из-за аварии на стороне потребителя), применяется механизм технических смещений (offsets). Система-потребитель фиксирует последний успешно обработанный порядковый номер сообщения в топике Kafka. Благодаря настроенной политике хранения (retention policy) в определенное количество дней, ИС-потребитель может инициировать повторную выгрузку путем «отката» смещения на последний сохраненный в базе данных индекс. Это гарантирует восстановление потерянных данных с соблюдением принципа идемпотентности – ни одно сообщение не будет пропущено или обработано дважды.

Заключение Использование Apache Kafka как промежуточного слоя между CRM-системой и платформой Hadoop формирует отказоустойчивую среду, способную эффективно распределять нагрузки и гарантировать сохранность данных. Сочетание механизмов распределенного буферирования, управления смещениями и многоуровневой верификации минимизирует риски потери информации и исключает негативное влияние интеграции на оперативную деятельность предприятия. Практическая значимость работы заключается в возможности масштабирования данного подхода на любые высоконагруженные системы предприятия для построения надежных аналитических репозиториях Больших данных.

Список литературы:

1. 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.

2. Соломонов А. А. Оптимизация ETL-процессов для больших данных // Вестник науки. 2024. № 9 (78).

3. Галигузова Е. В., Илларионова Ю. Е. Сравнение реляционных и нереляционных СУБД // Символ науки. 2023. № 1-2.

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

