

Качество данных в современном цифровом банке

Руслан Хохлов Департамент управления данными



Работа с данными





Система развития искусственного интеллекта.

Комплекс подходов, практик и процессов разработки и применения технологий ИИ для решения конкретных задач организации

Система управления данными.

Комплекс организационной структуры, процессов, практик, подходов и участников, обеспечивающий выполнение требований к данным.

Инфраструктура управления данными.

Включает в себя основные программно-аппаратные комплексы, осуществляющие обработку данных, а также системы, имеющие критическую важность для обеспечения доступности и качества данных.

Качество данных: почему это важно?



Ошибки в данных могут привести к выработке ошибочной стратегии

Как это может быть:

Ошибка в биржевых данных -> Ошибка моделей -> Убыточная сделка

Ошибка в метеорологических данных -> Прогноз хорошей погоды на выходных -> Пострадавшие на природе отдыхающие

Ошибка в клиентских данных -> Неверная атрибуция -> Финансирование терроризма

Ошибка в записи дня рождения тещи (свекрови)

-> Отсутствие поздравлений -> Испорченное настроение





Элементы системы автоматизированного контроля КД

Ввод информации (системы-источники) Консолидация данных (платформы и фабрики данных)

Предоставление данных потребителям (специализированные системы и витрины)



Элементы системы автоматизированного контроля КД

Ввод информации (системы-источники)

Консолидация данных (платформы и фабрики данных) Предоставление данных потребителям (специализированные системы и витрины)





Предотвращение ввода заведомо некорректной информации Повышение качества данных по накопленным массивам данных

Выявление специфичных аномалий



Элементы системы автоматизированного контроля КД

Ввод информации (системы-источники)

Консолидация данных (платформы и фабрики данных) Предоставление данных потребителям (специализированные системы и витрины)







Предотвращение ввода заведомо некорректной информации Повышение качества данных по накопленным массивам данных

Выявление специфичных аномалий







Форматно-логические контроли UI

Системы пост-контроля качества данных

Специализированные решения / Системы постконтроля КД



Проблематика стандартного подхода

Ввод информации (системы-источники)

Консолидация данных (платформы и фабрики данных) Предоставление данных потребителям (специализированные системы и витрины)







Форматно-логические контроли UI

Системы пост-контроля качества данных

Специализированные решения / Системы постконтроля КД







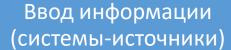
Значительный объем изменений на системах, сложности управления

Значительный объем контролей

Значительный объем контролей, специфичность задач



Поиск альтернативных решений



Консолидация данных (платформы и фабрики данных) Предоставление данных потребителям (специализированные системы и витрины)







Валидация данных на стороне системы КД. Источник только обрабатывает результаты валидации

Паттернализация контролей и автоматизированная генерация для однотипных данных на разных слоях

Self-service DQ – передача специализированных задач на сторону экспертов

Системы Офиса **CDO**в альтернативном подходе



Ввод информации (системы-источники) Консолидация данных (платформы и фабрики данных) Предоставление данных потребителям (специализированные системы и витрины)

Online / API

Валидация входящего вектора данных

BATCH

Выявление некорректных данных в массиве Поддержка и развитие пользовательского инструментария

Система управления качеством данных

Self-service DQ

Data Catalog

Обеспечение требуемыми метаданными



Паттернализация пост-контролей: проблематика

На различных слоях платформ консолидации имеются данные, относящиеся к одному и тому же атрибуту корпоративной модели данных.

Пример:

Реквизиты клиента ФЛ (источник CRM)

Реквизиты клиента ФЛ (источник АБС)

Реквизиты клиента ФЛ (источник CDI)

Реквизиты клиента ФЛ (источник CDI)

Реквизиты клиента ФЛ (витрина Риски)

Реквизиты клиента ФЛ (витрина вкладчики)

Реквизиты клиента ФЛ (витрина ...)

Физическая модель хранения разная на каждом источнике/слое, но правила контроля качества данных - одинаковые



Паттернализация пост-контролей: решение

Пост-контроли качества реализовать в терминах корпоративной модели данных, а реализацию для конкретной точки контроля генерировать автоматически на данных маппинга на физические модели (из каталога данных).

Основная информация Обла	сти контроля	Запрос	Даты Хар	арактеристика данных Триггеры проверок Периодичность
Ключ филиала:				Произвольный атрибут группировки результатов: ОСНОВНЫЕ СВЕДЕНИЯ.Статус клиента
Запрос выражение 1:	с выражение 1: ОСНОВНЫЕ СВЕДЕНИЯ,ФИО одной строкой			
Запрос выражение 2:	ОСНОВНЫЕ	СВЕДЕНИЯ.Идентифика		
Запрос выражение 3: ОСНОВНЫЕ СВ		СВЕДЕНИЯ.Унифициро	ванный идентификатор G	GUID
Запрос проверки (WHERE): *	regexp_like	("ОСНОВНЫЕ СВЕДЕНИЯ	Я.ФИО одной строкой", "[([0-9!#\$%&′*+/=?^`{[]~]') and "ОСНОВНЫЕ СВЕДЕНИЯ.ФИО одной строкой" IS NOT NULL
Основная и	нформация	Области контроля	3anpoc	Даты Характеристика данных Триггеры проверок Периодичность
Объект ко	нтроля	Область контроля	Входит в индикатор	Запрос
ОСНОВНЫ	Е СВЕДЕНИЯ	DO_PREP	N	SELECT CAST(decode(party_rk, -1, party_global_rk, party_rk) AS string) AS key_id, CAST(CAST(t_source_system_id AS string) AS string) AS source_system_id, NULL AS

Контроль входящих данных систем первичного ввода информации: проблематика



Пост-контроли качества данных позволяют выявлять и устранять значительное количество ошибок в данных, но только блокирование ввода заведомо некорректной информации может этот поток ошибок остановить.

- В крупных организациях данные вносятся во множество систем, включая вендорские решения. Возникает большая управленческая нагрузка на ведение таких изменений.
- Фронтальные системы это системы, связанные непосредственно с бизнесом, поэтому доработки, связанные с качеством данных имеют более низкий приоритет.
- Реализация контролей качества в системах нередко требует изменений в интеграции и совместных релизов.
- Требования к контролям качества данных достаточно часто изменяются.

Контроль входящих данных систем первичного

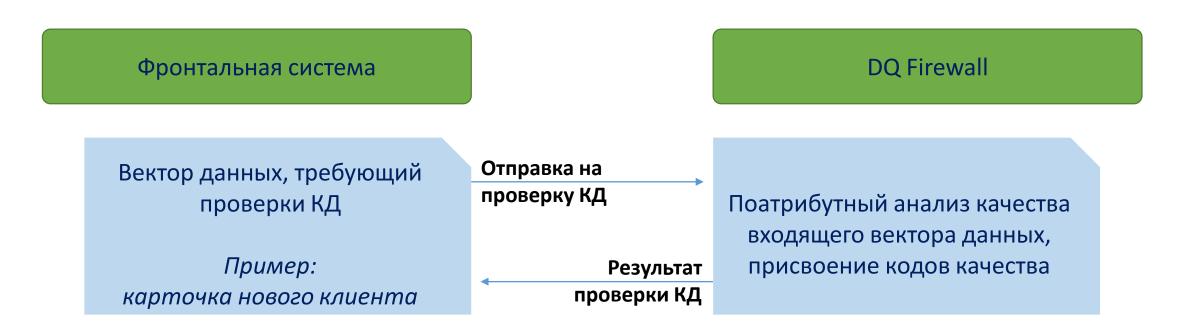


ввода информации: data quality firewall

Идея:

У нас есть библиотека пост-контролей в терминах корпоративной модели данных, которую мы развиваем и актуализируем.

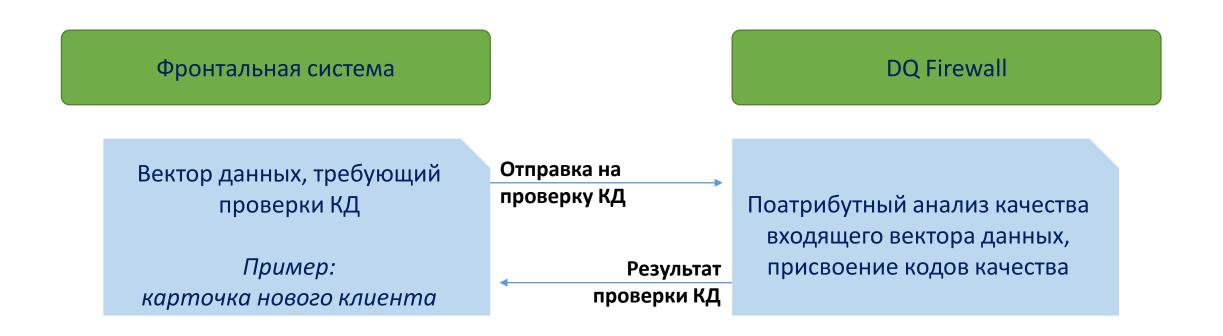
А если сделать сервис (API) для фронтальных систем по валидации данных в реальном времени по запросу? Например, как REST-сервис?



Data Quality Firewall: концепция

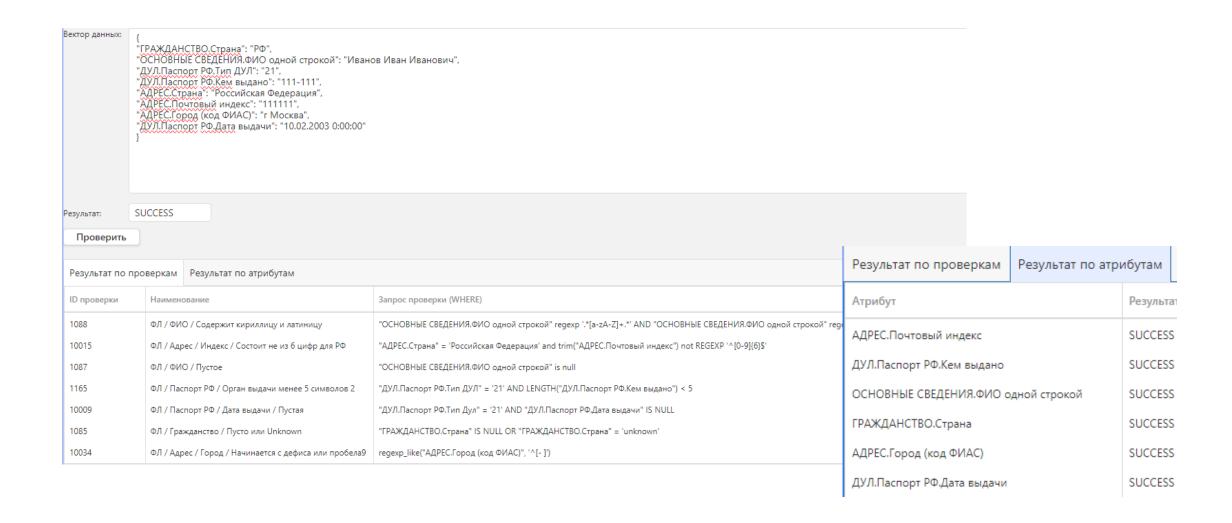


- 1. Синхронный сервис, получающий на вход вектор данных для валидации, возвращающий поатрибутный ответ в кодах качества.
- 2. Использует библиотеку контролей КД, разработанную для пост-контроля.
- 3. Для ускорения обслуживающие процессы (в т.ч. логирование) выведены в асинхронный режим и от основного процесса.



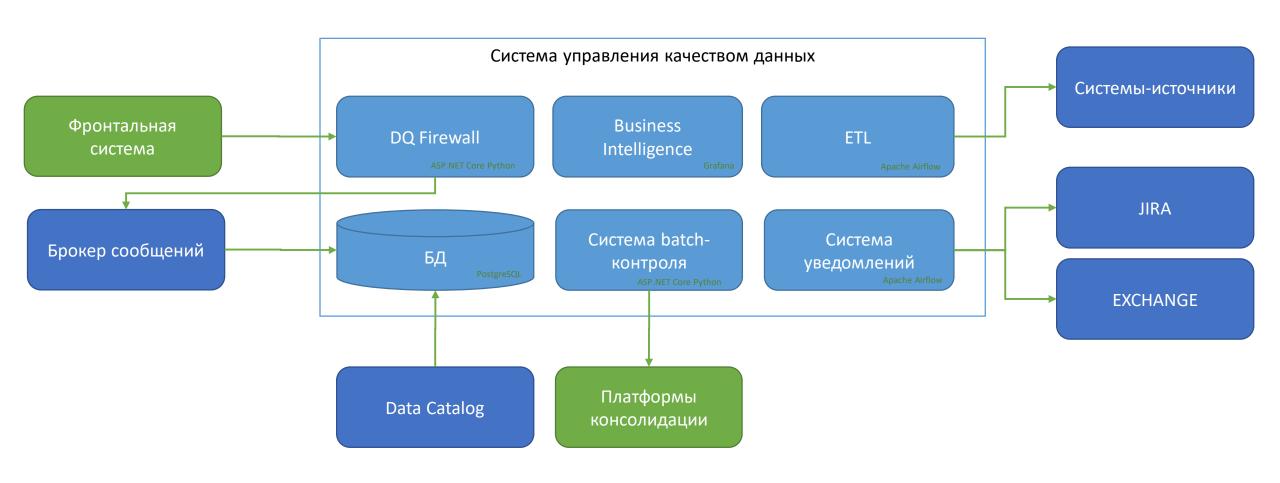


Data Quality Firewall: пример работы сервиса





Data Quality Firewall: архитектура





Пост-контроль в сервисной модели

Офис CDO

- Владеет и развивает библиотеку контролей качества данных в терминах корпоративной модели данных.
- По мере необходимости подключает новые объекты контроля КД и предоставляет на регулярной основе результаты контролей качества данных.
- Обрабатывает выявленные инциденты качества данных

Потребитель сервиса

- Использует результаты контроля качества данных для своих целей
- Учитывает факты возникновения инцидентов качества данных в своей работе



DQFirewall в сервисной модели

Офис CDO

- Владеет и развивает библиотеку контролей качества данных в терминах корпоративной модели данных.
- Развивает и сопровождает программный сервис (API) по валидации качества данных.
- Осуществляет валидацию входящего вектора данных и возвращает поатрибутный ответ о соответствии/несоответствии данных требованиям по качеству.

Потребитель сервиса

- Осуществляет подключение к программному сервису (API) своего бизнес-приложения.
- Программно направляет запрос на валидацию данных.
- Обрабатывает и выводит на экранные формы результаты работы сервиса.



Сервисная модель: сложности подхода

- 1. Основные потребители сервиса (DQ Firewall API) системы Mission Critical. Должен быть обеспечен аналогичный уровень отказоустойчивости.
- 2. Должно быть низкое время отклика такого сервиса.
- 3. Для batch-контроля должно быть обеспечено высокое качество маппинга в каталоге данных корпоративной модели данных на физические модели точек контроля.
- 4. Необходима готовность команды к большому объему интеграционных задач.



Сервисная модель: выгоды и плюсы

- 1. Радикальное сокращение time-to-market на изменения в контролях качества данных (достаточно внести изменения через интерфейс системы).
- 2. Снижение управленческих затрат за счет единой точки управления изменениями в контролях качества данных.
- 3. Единая методология контроля КД на всех узлах контроля. Используется сквозная библиотека контролей. Исключена рассинхронизация версий контролей при изменении происходит перегенерация исполняемого кода и обновление кэша для API.
- 4. Появление нового узла контроля типовая задача.
- 5. Возможность внедрения дополнительных инструментов контроля КД (в том числе интеграции с ГИС) прозрачно для потребителей сервисов.



Спасибо за внимание!

