


Быстрый ETL для PostgreSQL

Арустамов Алексей, Loginom Company

Extract, Transform, Load

Разработчики баз данных с развитыми экосистемами предлагают инструменты для реализации ETL-процессов:

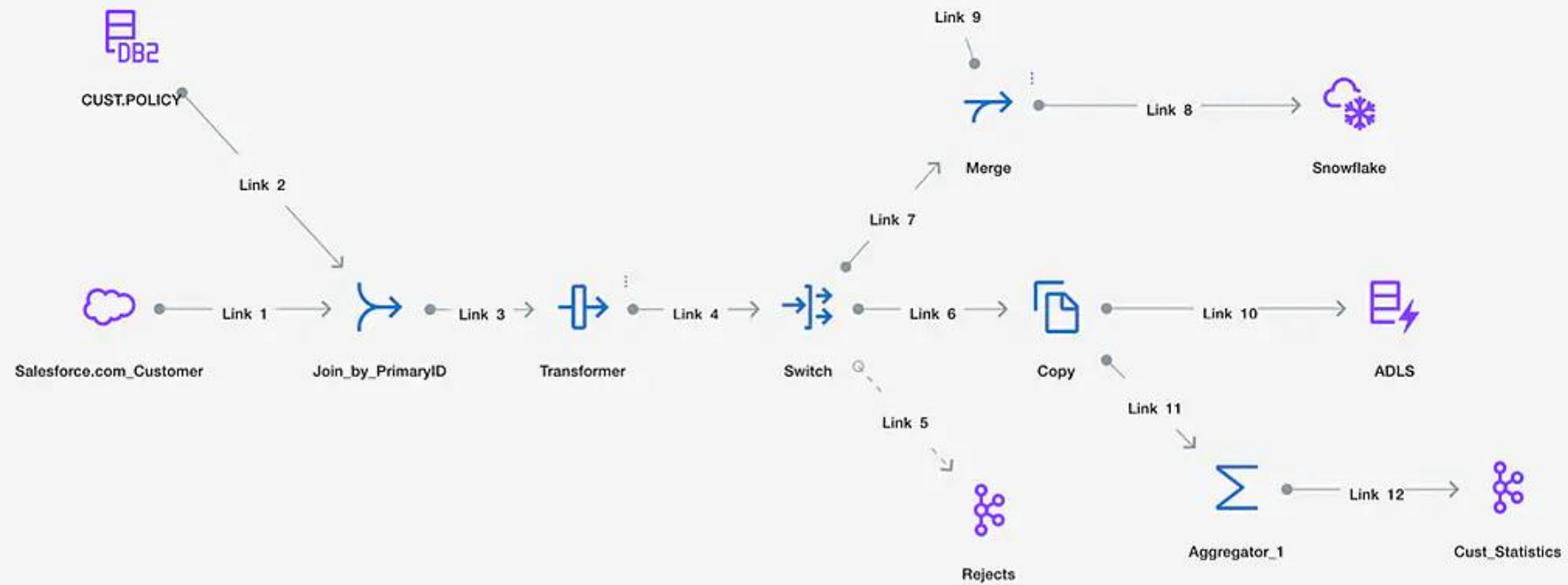
- Oracle Data Integrator
 - Microsoft SQL Server Integration Services
 - IBM DataStage
- 

Все ETL инструменты основаны на визуальном проектировании

Find palette nodes

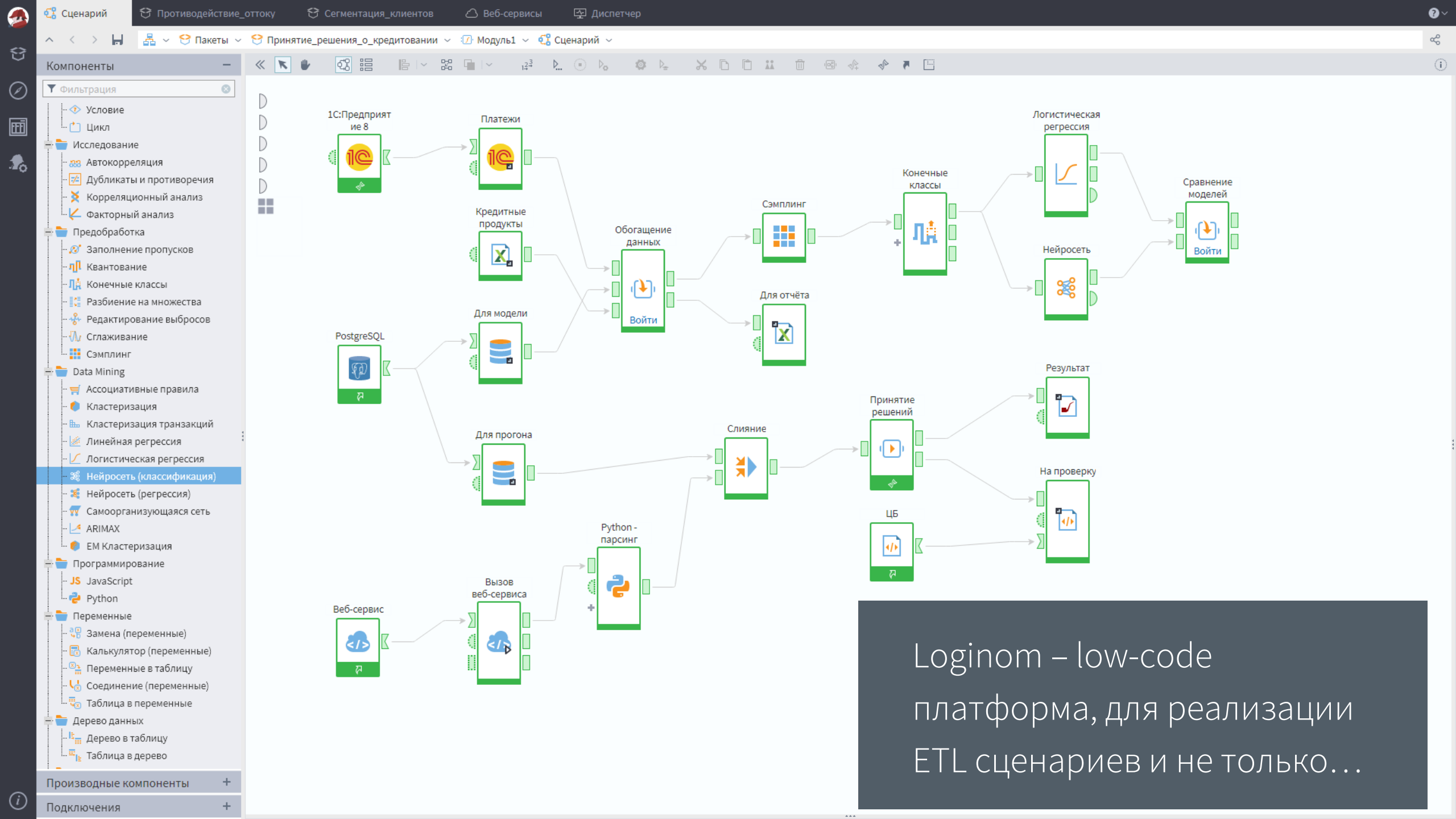
Save Compile Run

- Connectors
- Stages
 - Aggregator
 - Bloom Filter
 - Change Apply
 - Change Capture
 - Checksum
 - Column Export
 - Column Generator
 - Column Import
 - Compare
 - Compress
 - Copy
 - Decode
 - Difference
 - Encode
 - Expand
 - External Filter



Визуальный ETL

Реализация ETL-процесса, не требующая кодирования и (обязательного) знания SQL упрощает работу с базой данных, позволяя решать большинство задач пользователям без специальной подготовки. Self-service избавляет разработчиков от рутины.



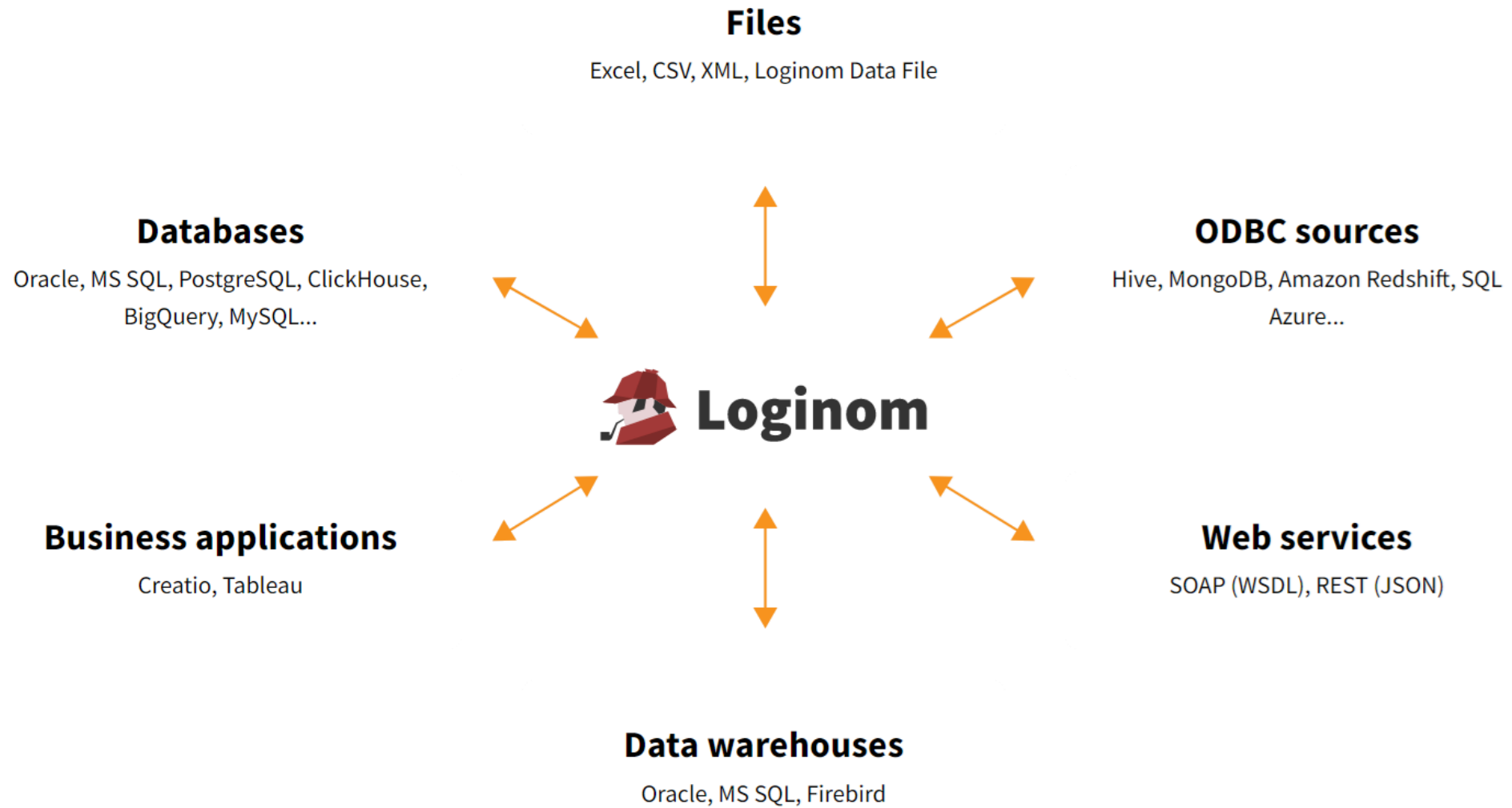
Loginom – low-code платформа, для реализации ETL сценариев и не только...



СУБД **Postgres Pro Standard** и **Postgres Pro Enterprise**

технологически полностью
совместимы с аналитической
low-code платформой **Loginom**

Извлечение и загрузка



Обработка и трансформация



Grouping



Cross Table



Union



Join



Row Filter



Correlation Analysis



Duplicates Detection



Eliminate Outliers



Sampling



Imputation



Clustering



Logistic Regression



Neural Network (Classification)



ARIMAX



XML Generation



Transaction Clustering




Smoothing



49+ components

Способы повышения скорости

1. Прямой доступ без драйверов
 2. Чтение/запись пачками
 3. Пул подключений
 4. Параллелизм
 5. In-memory обработка
 6. Ленивые вычисления
 7. Паттерн MapReduce
- 

Прямой доступ к PostgreSQL

1. Не требуется установка клиента
2. Поддержка СУБД версий от 8.0 до 14.4
3. Аутентификация
 1. До версии Logiном 6.5 – md5
 2. С версии Logiном 6.6 – md5 + scram-sha-256
4. Поддерживается Greenplum и Arenadata

Чтение/запись пачками

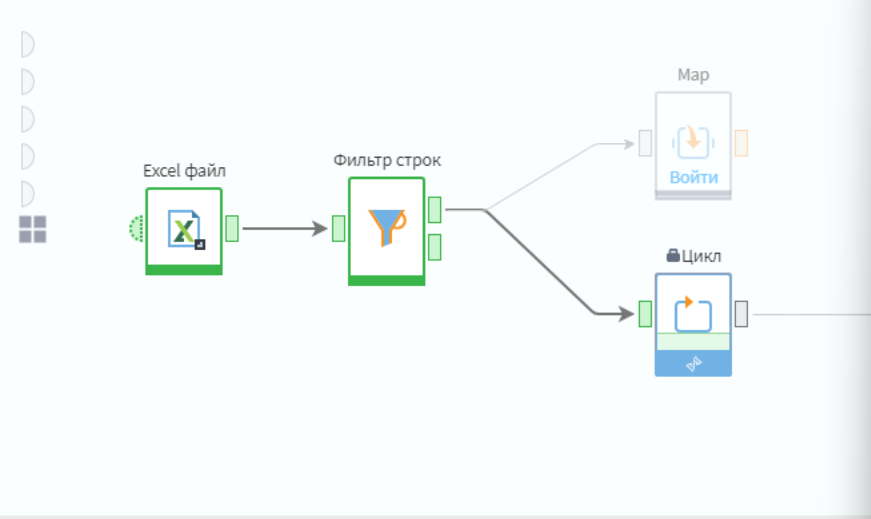
1. Чтение и запись производится пачками

Операция	Стандарт (сек.)	Пакетный (сек.)
Insert	346.7	1.69
Update	334.4	4.59
Delete	373.7	2.05

2. Размер пачки – 1000 строк

Пул подключений к PostgreSQL

1. Пул подключений уменьшает время затрачиваемое на повторное подключение. Подключение держится после завершения выполнения сценария, чтобы ускорить пакетное выполнение
2. При отсутствии подключения оно создается, по завершении выполнения запросов подключение возвращается в пул, при следующем обращении ищем есть ли в пуле подключение с такими же параметрами и используем его.
3. Уменьшается количество создаваемых подключений. Есть ограничение на количество одновременных подключений, по умолчанию 100. Если сценарий выполняется последовательно то одно подключение, если параллельно, то столько подключений, сколько параллельных потоков.



№	Процесс	%	Обработка	Ошибки
39	Активация узлов	27		
39.1	Цикл	27		
39.1.1	Итерация №0	100		
39.1.2	Итерация №1	100		
39.1.3	Итерация №2	100		
39.1.4	Итерация №3	25		
39.1.5	Итерация №4	22		
39.1.6	Итерация №5	22		
39.1.7	Итерация №6	22		
39.1.8	Итерация №7	22		
39.1.9	Итерация №8	22		
39....	Итерация №9	22		
39....	Итерация №10	22		
39....	Итерация №11	22		
39....	Итерация №12	16		
39....	Итерация №13	16		
39....	Итерация №14	16		
39....	Итерация №15	0		
39....	Итерация №16	0		
39....	Итерация №17	0		

Диспетчер задач

Файл Параметры Вид

Процессы Производительность Журнал приложений Автозагрузка Пользователи Подробности Службы

ЦП
87% 2,23 ГГц

Память
4,5/15,3 ГБ (29%)

Диск 0 (C:)
SSD
29%

Wi-Fi
Беспроводная сеть
0: 0 П: 0 кбит/с

Графический про
AMD Radeon(TM) Graph
6% (46 °C)

ЦП

AMD Ryzen 5 5500U with Radeon Graphics

% использования более 60 секунд

Использование: **87%**

Процессы: **222**

Время работы: **0:01:43:22**

Скорость: **2,23 ГГц**

Потоки: **2417**

Дескрипторы: **93742**

Базовая скорость: **2,10 ГГц**

Сокетов: **1**

Ядра: **6**

Логических процессоров: **12**

Виртуализация: **Включено**

Кэш L1: **384 КБ**

Кэш L2: **3,0 МБ**

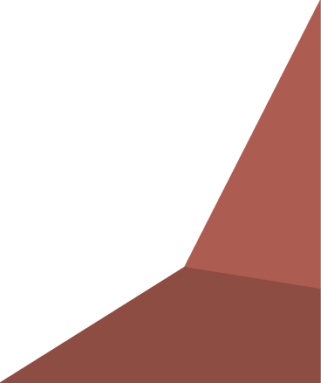
Кэш L3: **8,0 МБ**

Эффективная утилизация ресурсов за счет параллелизма

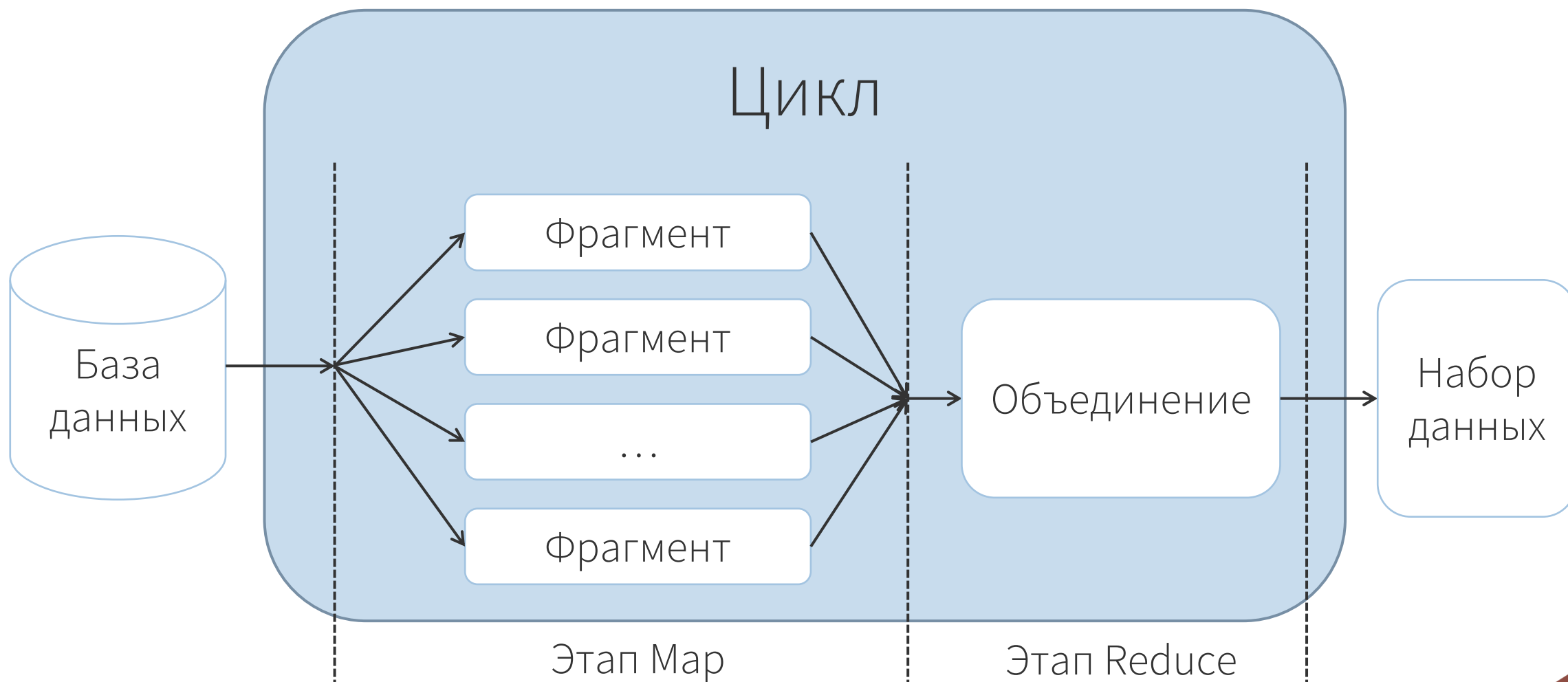
In-memory

1. Платформа пытается все данные держать в памяти
2. Оптимизировано хранение наборов в памяти, в частности, по умолчанию хранятся только уникальные данные
3. При необходимости можно задать точки кэширования данных

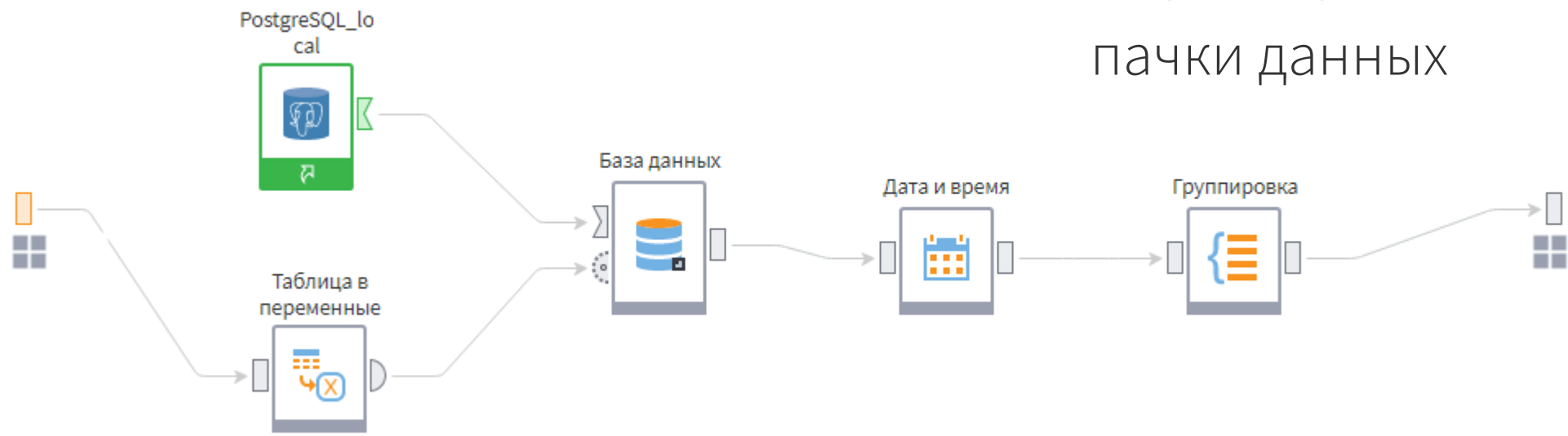
Ленивые вычисления

1. Вытягивать данные, а не толкать
 2. Не читать пока не будет запрошено
 3. Не считать пока не будет запрошено
 4. Не показывать пока не будет запрошено
- 

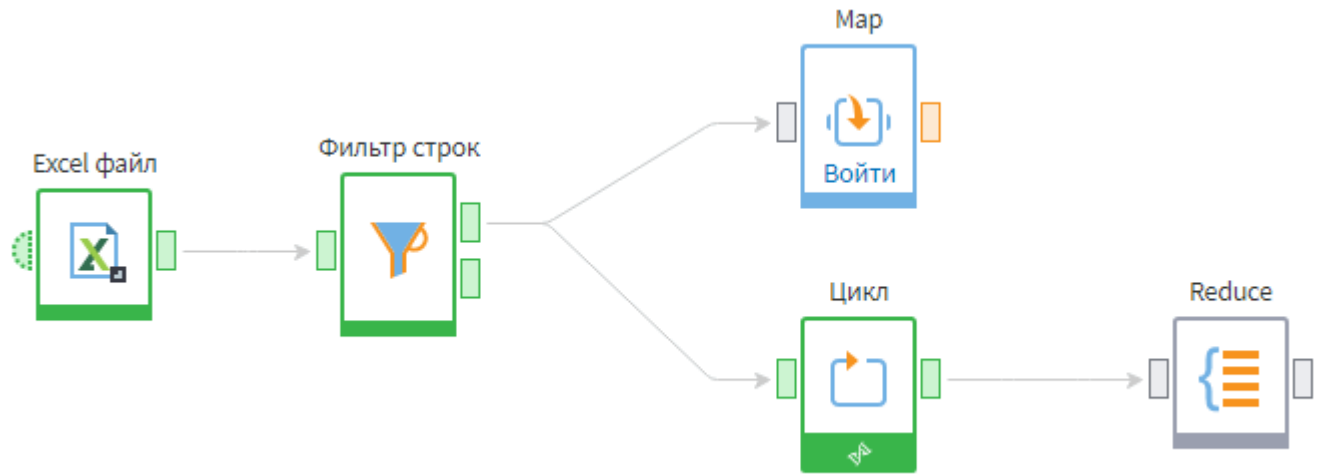
MapReduce в многопроцессорных системах



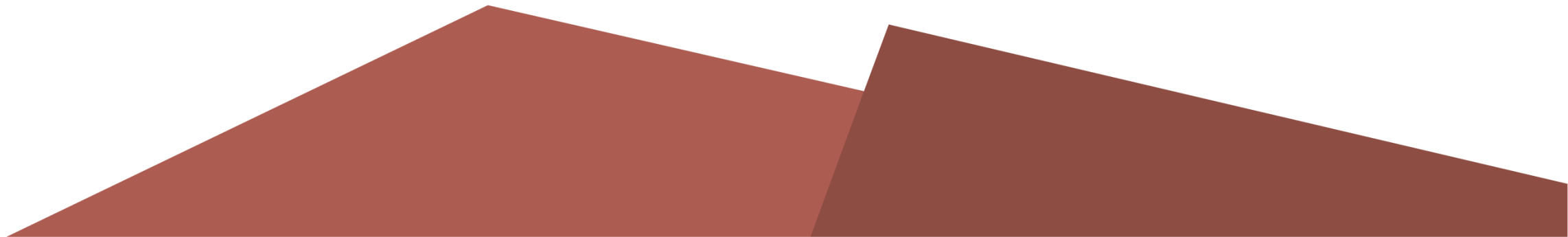
Map – обработка одной пачки данных

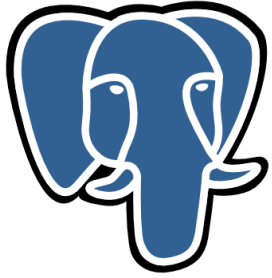


Выполнение в параллельном цикле и сбор итога



Демо реализации MapReduce



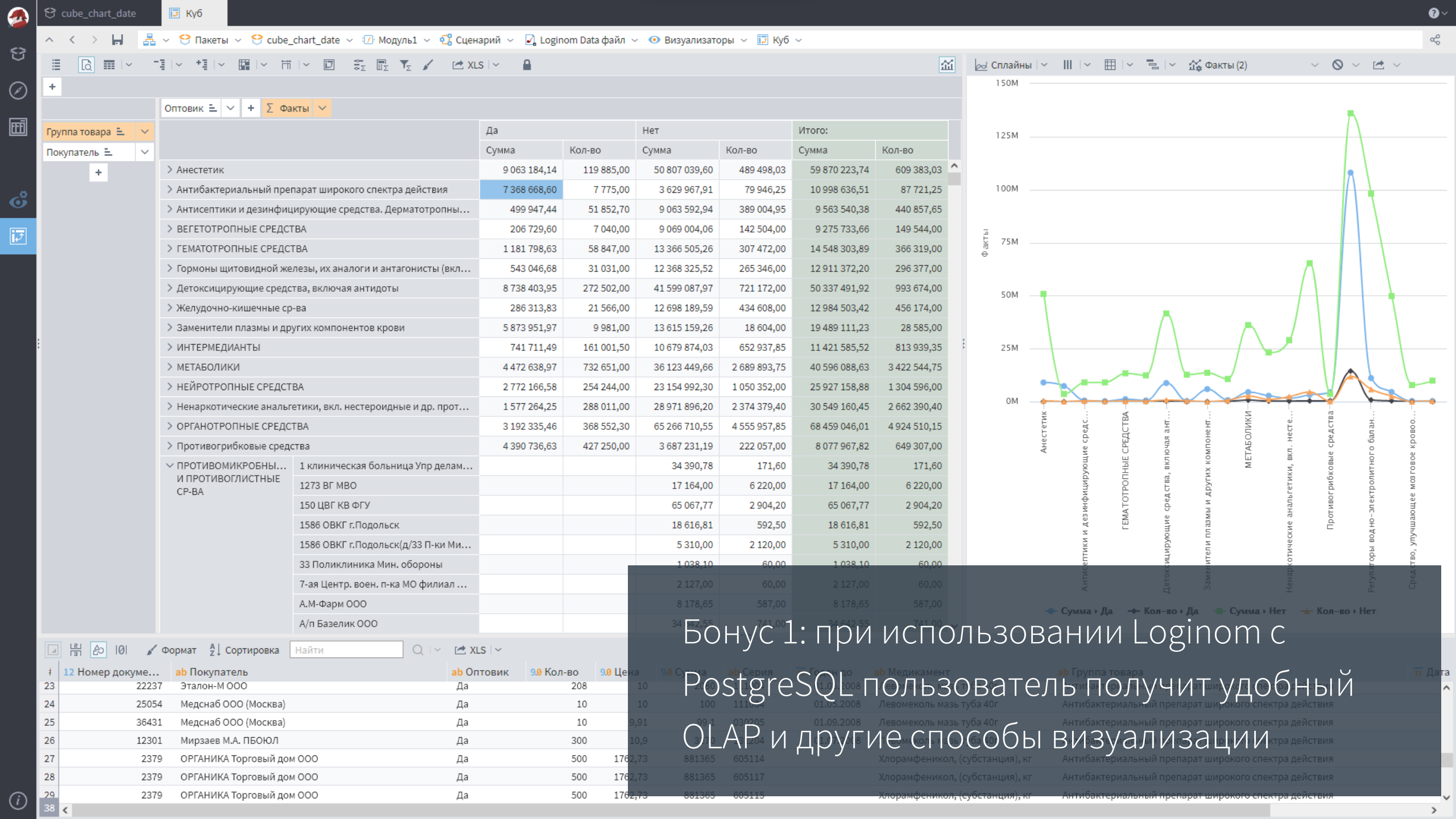


PostgreSQL



Loginom

Не только хранение данных
и обработка запросов, но и
сложный ETL + продвинутая
аналитика



Бонус 1: при использовании Loginom с PostgreSQL пользователь получит удобный OLAP и другие способы визуализации

Выбор диаграммы

- ROC-кривая
- PR-кривая
- Базовые показатели
- Диаграмма точности
- Диаграмма равновесия
- % распознанных событий
- Диаграмма роста
- Диаграмма отклика
- Диаграмма выигрыша

Кумулятивная

10 диапазонов

Множества

Обучающее Тестовое

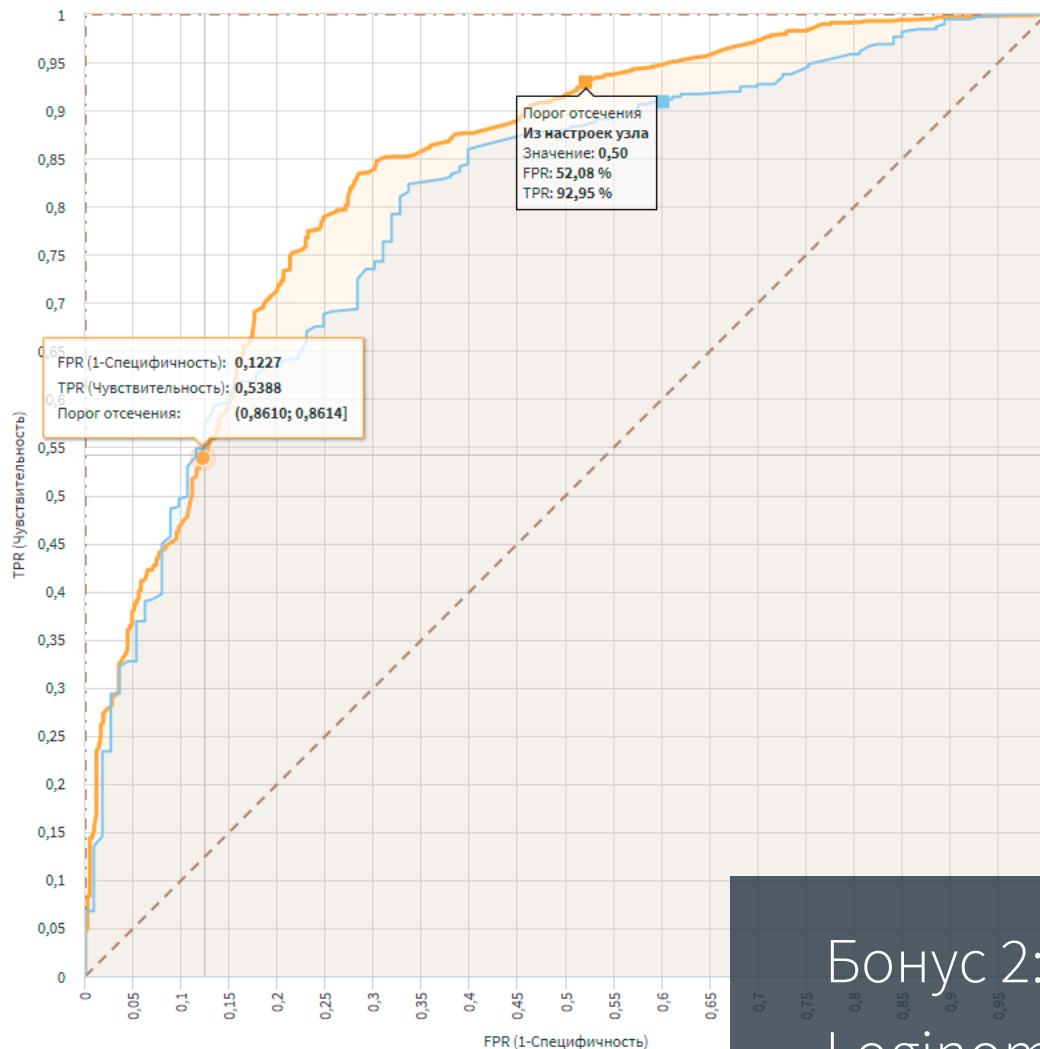
Порог отсекания

Из настроек узла

Значение порога:

ROC-кривая

Событие: Статус = хороший



- Обучающее множество
- Порог отсекания
- Базовая линия
- Идеальная линия
- Тестовое множество
- Порог отсекания

Оценки классификации

Показатель	Множества	
	Обучающее	Тестовое
Оценки классификатора		
AUC ROC	0,8362	0,8010
AUC PR	0,9424	0,9275
Коэффициент Джини	0,6724	0,6019
KS	55,2314	48,9691
Порог отсекания: Из настроек узла		
Значение	0,5000	0,5000
TPR (Чувствительность)	0,9295	0,9091
TNR (Специфичность)	0,4792	0,3982
FPR (1-Специфичность)	0,5208	0,6018
PPV	0,8658	0,8373
F1 Score	0,8965	0,8717
MCC	0,4604	0,3505

Матрицы ошибок

Классифицировано	Фактически		Итого
	Событие	Не-событие	
Обучающее			
Событие	1 561	432	1 676
Не-событие	1 451	225	317
Тестовое			
Событие	385	113	418
Не-событие	110	207	80

Распознано

Обучающее	83,19%
Тестовое	79,32%

Бонус 2: при использовании Loginom с PostgreSQL доступны методы машинного обучения

loginom.ru