



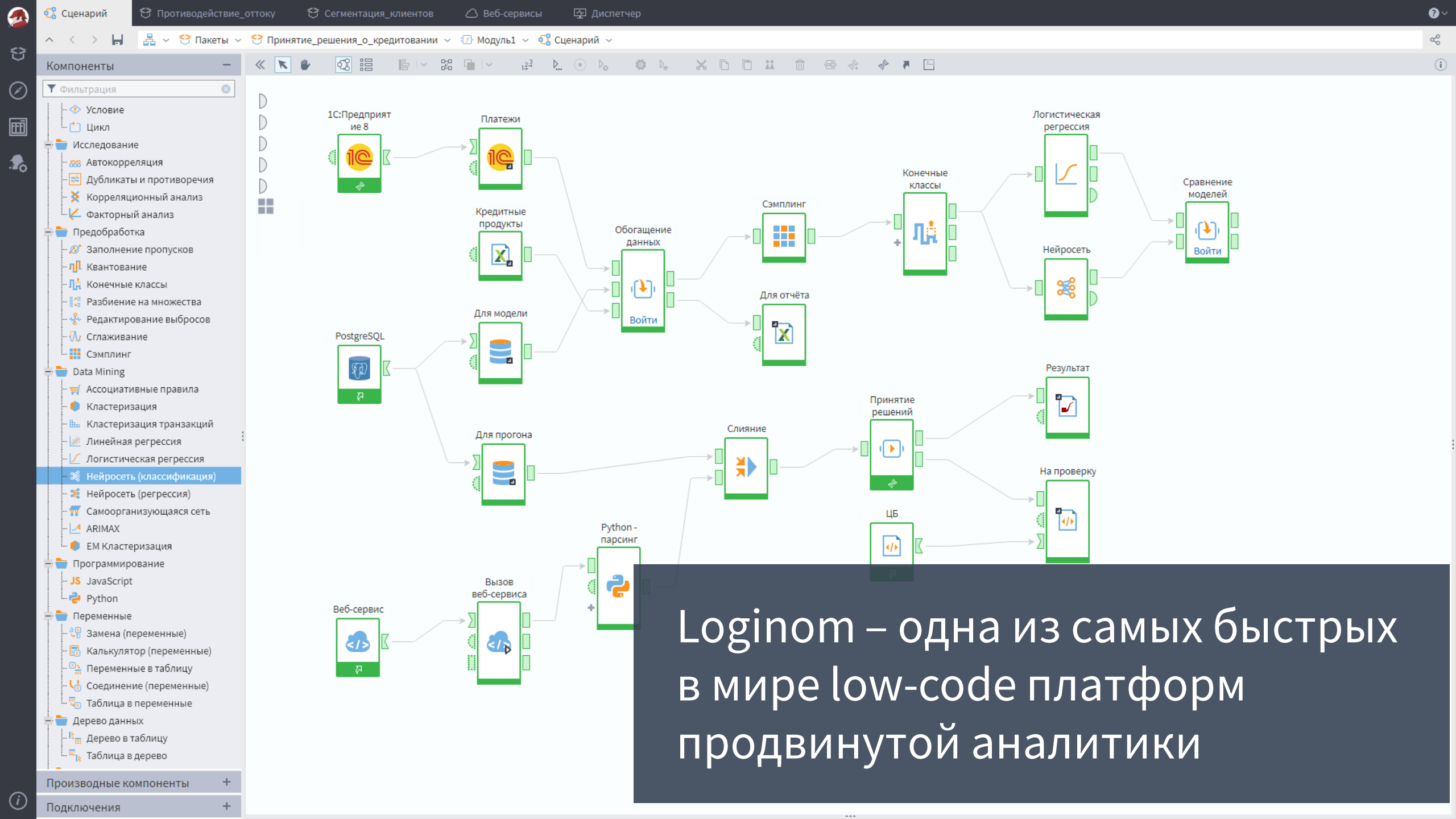
# **Loginom – low-code аналитика на предельной скорости**

Алексей Арустамов



Скорость – единственное свойство программ по которому есть консенсус: чем быстрее, тем лучше.

Высокая скорость должна обеспечиваться без усилий, «из коробки».



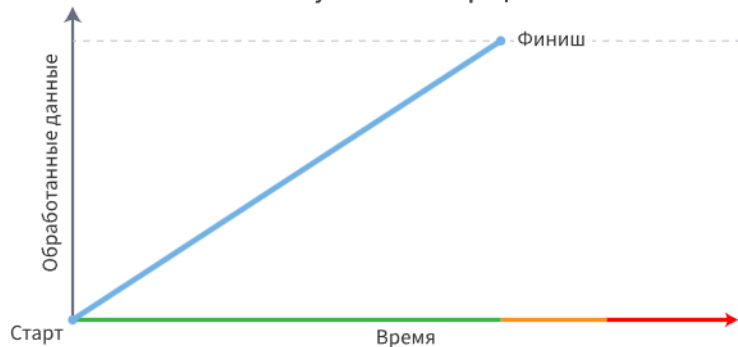
Факторы, влияющие на производительность:

1. Архитектура платформы
2. Алгоритмы обработки
3. Способы визуализации
4. Работа с процессорами, памятью, дисками
5. Работа с источниками/приемниками данных
6. Интеграция модулей Loginom между собой

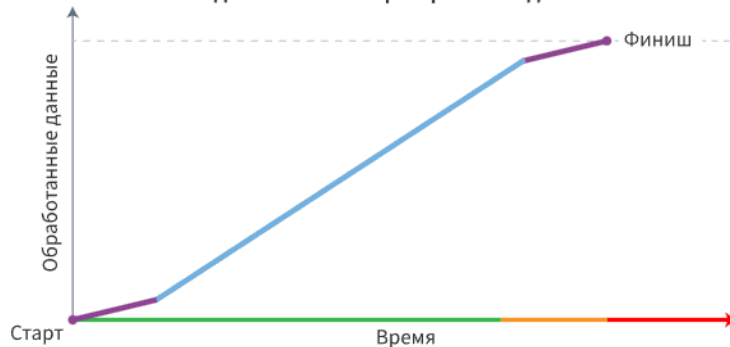
# Принципы оптимизации:

1. Оценивать производительность системы в целом, а не фрагмента
2. Улучшать то, что вызывается часто
3. Учитывать возможности железа
4. Искать первоисточник потерь

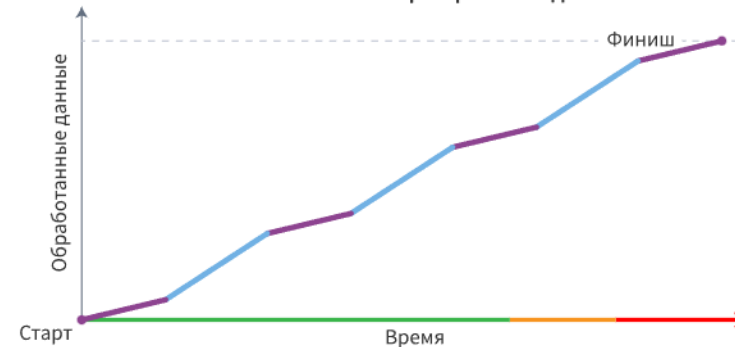
Отсутствие интеграций



Единый конвейер обработки данных



Составной конвейер обработки данных



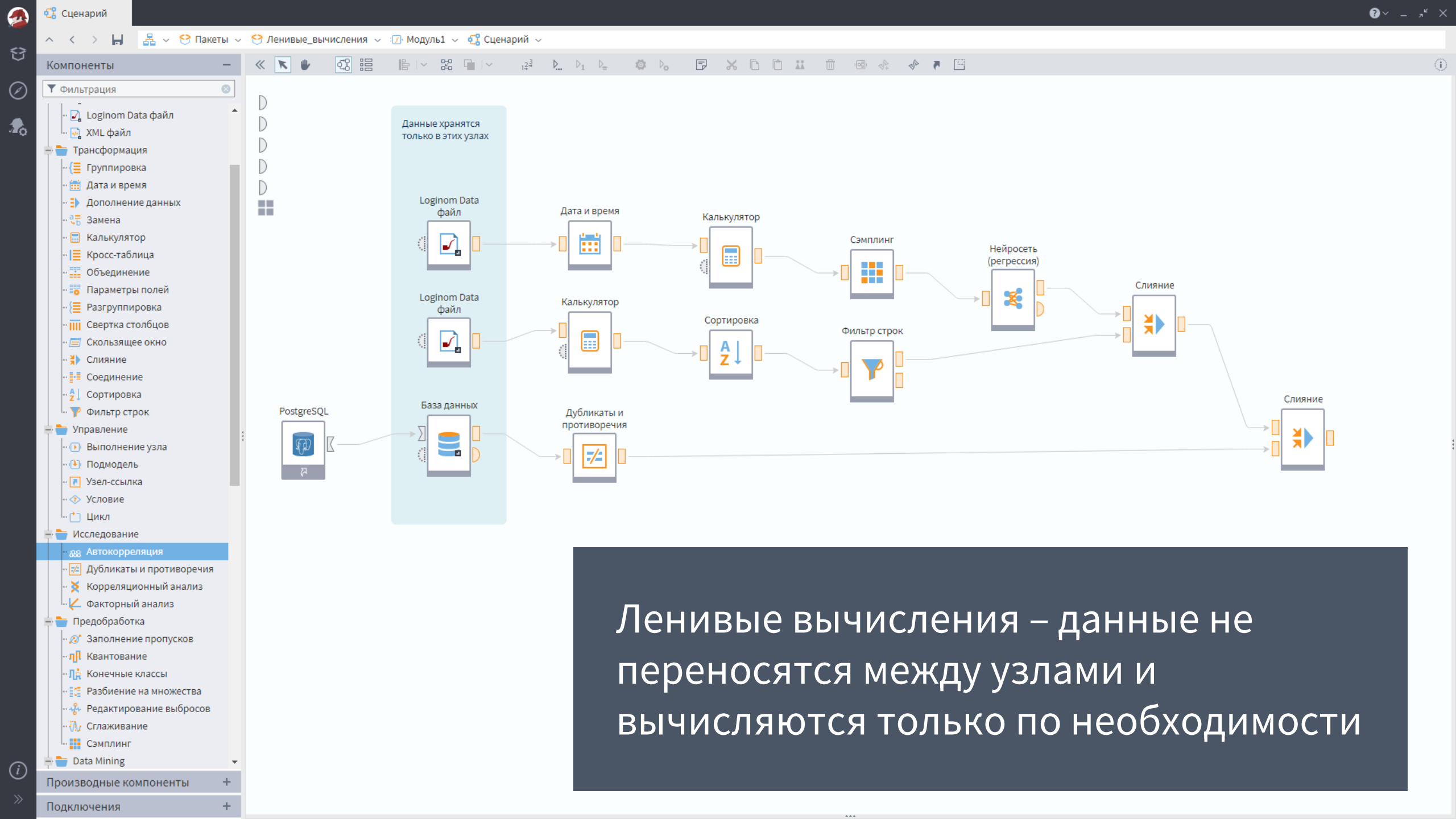
Чем больше в решении разнородных компонентов, тем больше точек интеграции и выше потери производительности

# Архитектуры платформы

Экономить ОЗУ, т.к. при работе с большими данными она заканчивается быстро:

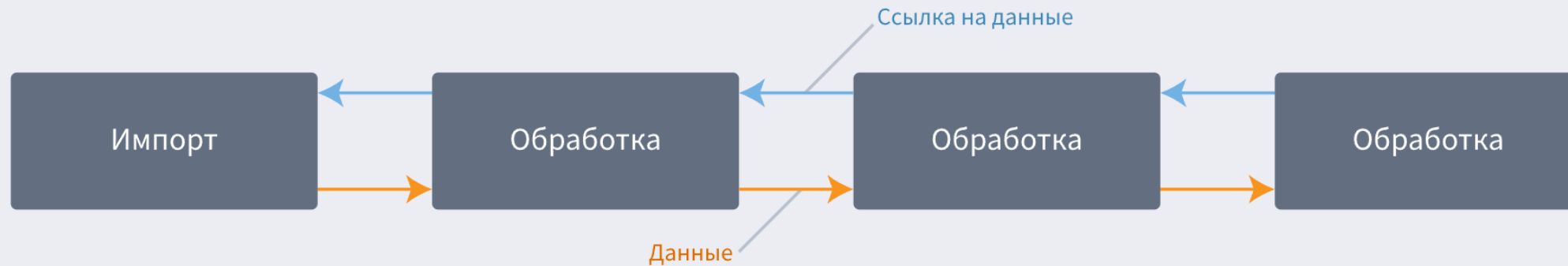
1. Работать с типизированными данными
2. Хранить в памяти уникальные значения
3. Использовать ленивые вычисления
4. ...



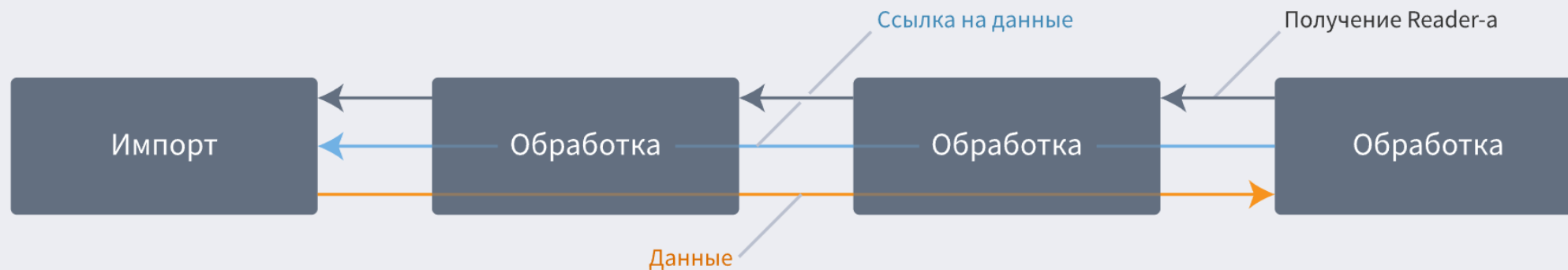


# Оптимизация ленивых вычислений

Было

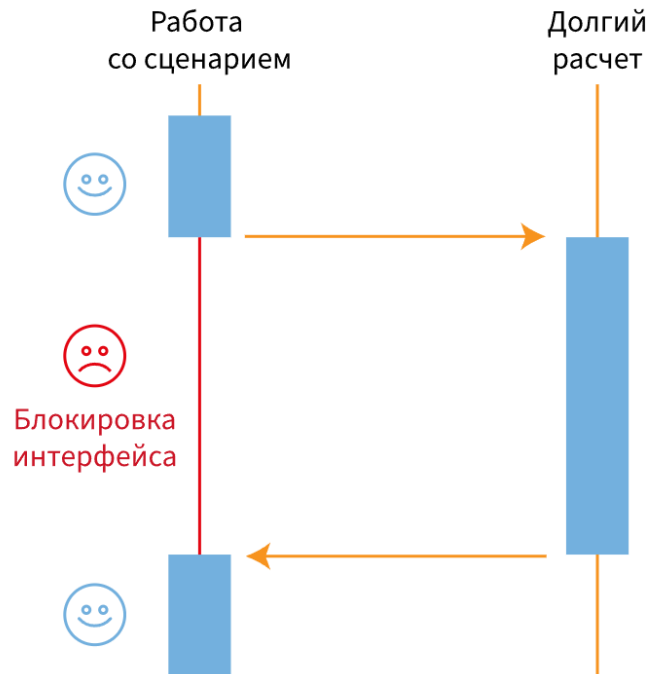


Стало

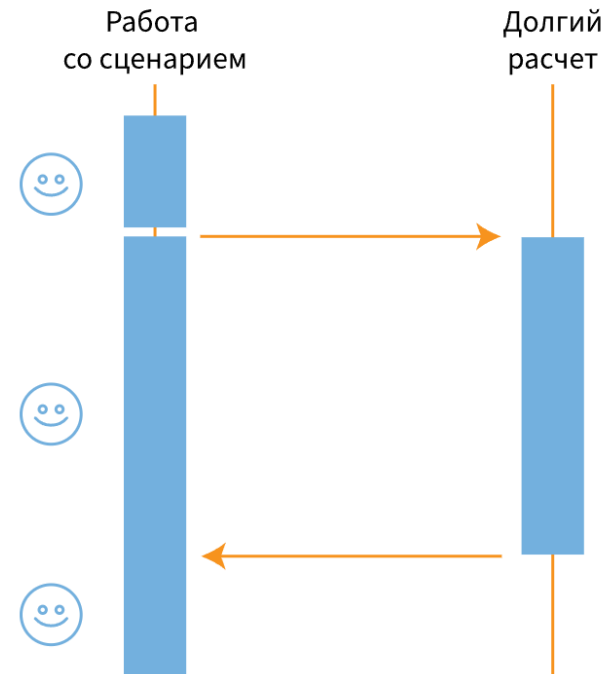




### Синхронный пользовательский интерфейс



### Асинхронный пользовательский интерфейс



Асинхронный пользовательский интерфейс не снижает время долгих расчетов, но не блокирует работу аналитика. За счет этого повышается отзывчивость системы и комфорт от работы.

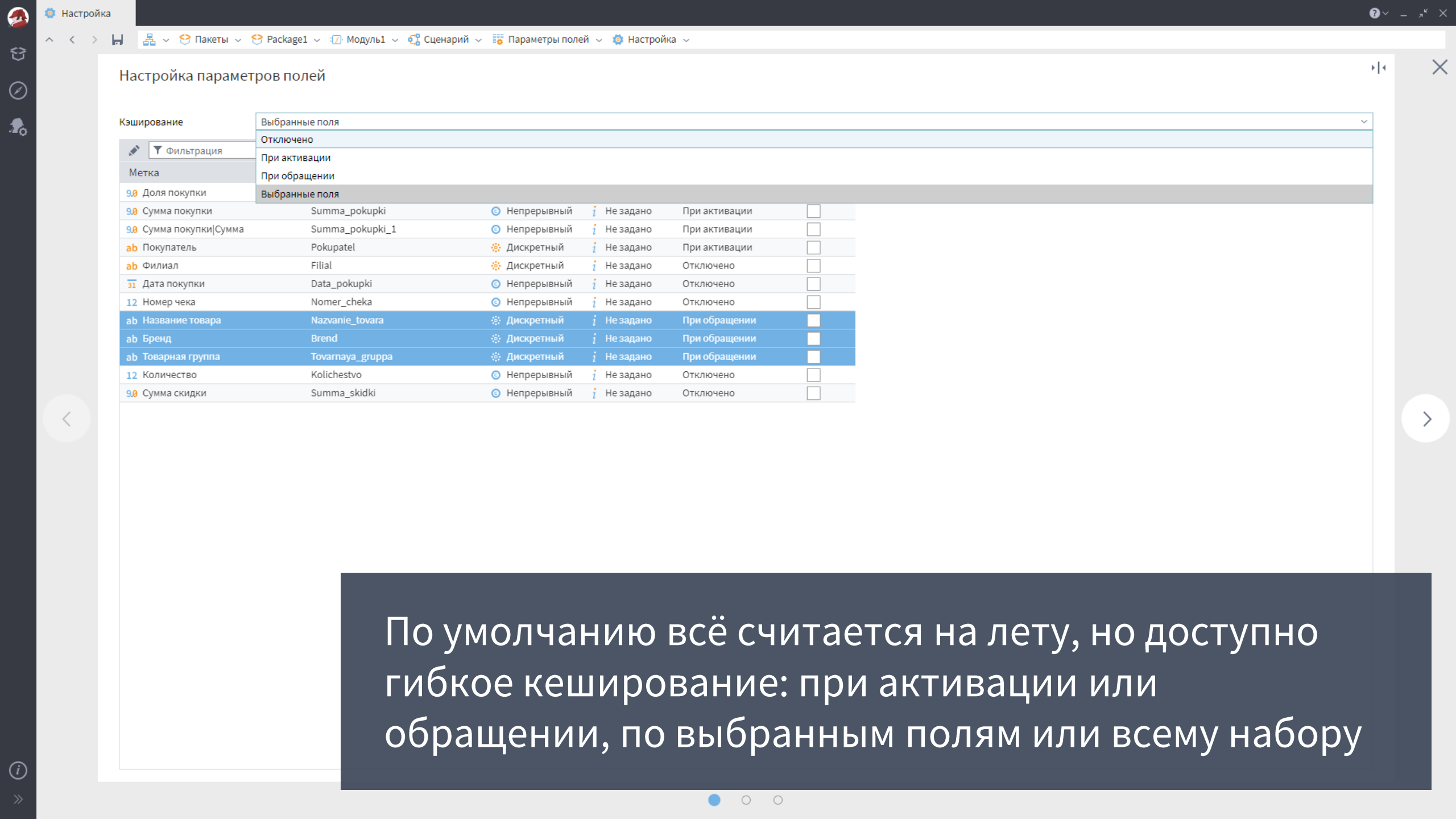
# Алгоритмы обработки

## Сторонние библиотеки

1. LIBLINEAR – одна из самых быстрых библиотек для построения моделей логистической регрессии
2. Intel MKL – одна из самых производительных библиотек для операций с матрицами и задач линейной алгебры
3. fast\_float – библиотека для быстрого преобразования строки в вещественное число
4. ALGLIB – оптимизированный расчет статистик, обучение нейросетей
5. ...

## Собственная реализация алгоритмов

1. Использование параллелизма при построении индексов (слиянии и дополнении данных)
2. Приблизительные вычисления при построении моделей машинного обучения
3. Подбор гиперпараметров нейросетей
4. Алгоритмы децимации данных для диаграмм
5. Многомерные расчеты в кубе
6. ...





# Визуализация

# Собственный RPC протокол обмена данными между клиентом и сервером:

1. Исключение всех промежуточных слоев
2. Обмен данными в бинарном виде для снижения объема передаваемой информации и минимизации потерь на сериализацию
3. Пакетная передача с объединением мелких запросов в крупные блоки
4. Асинхронный обмен

Прослойка между клиентом и сервером, учитывающая особенности визуализации проводит децимацию данных на сервере перед передачей клиенту.

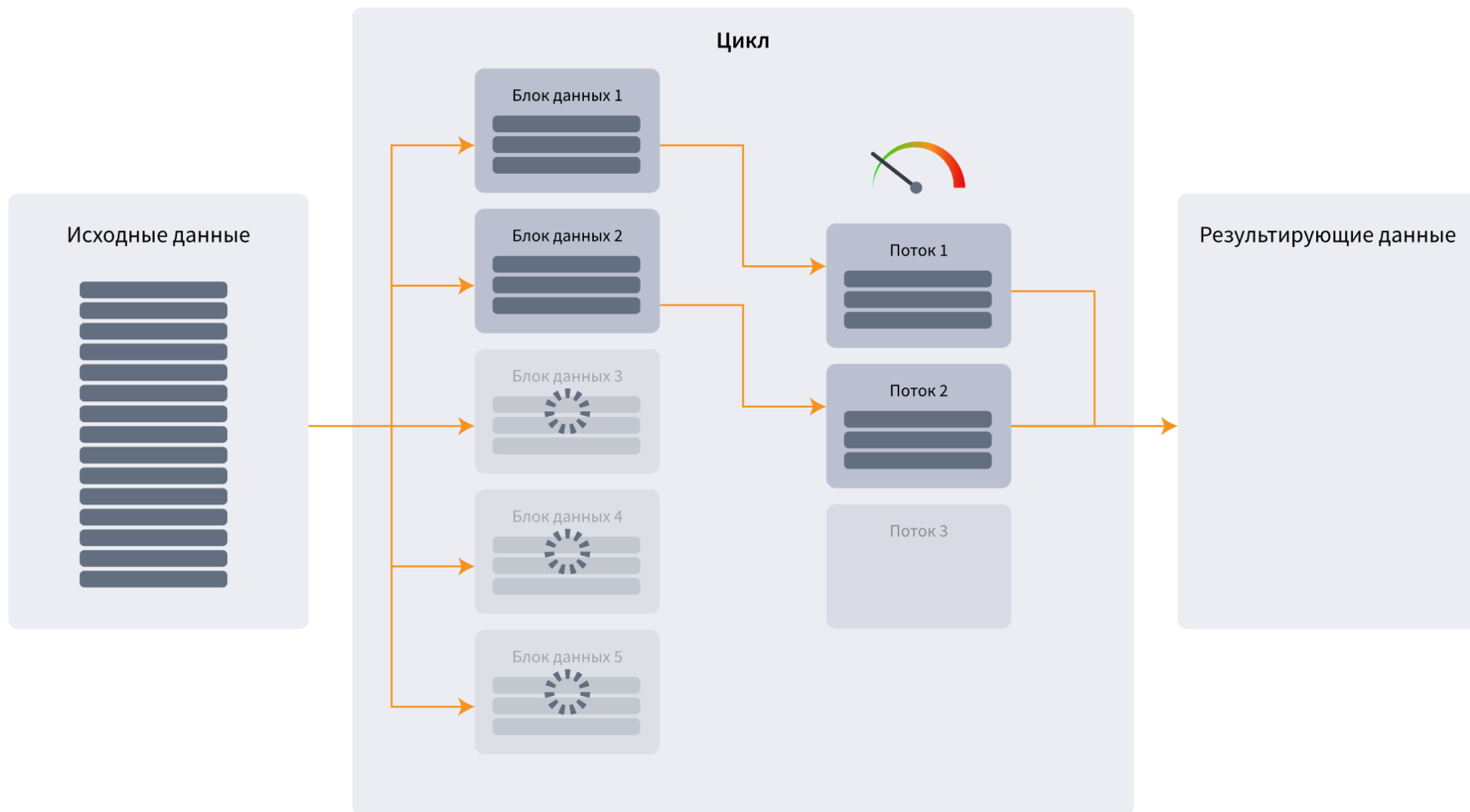


# Низкоуровневые операции

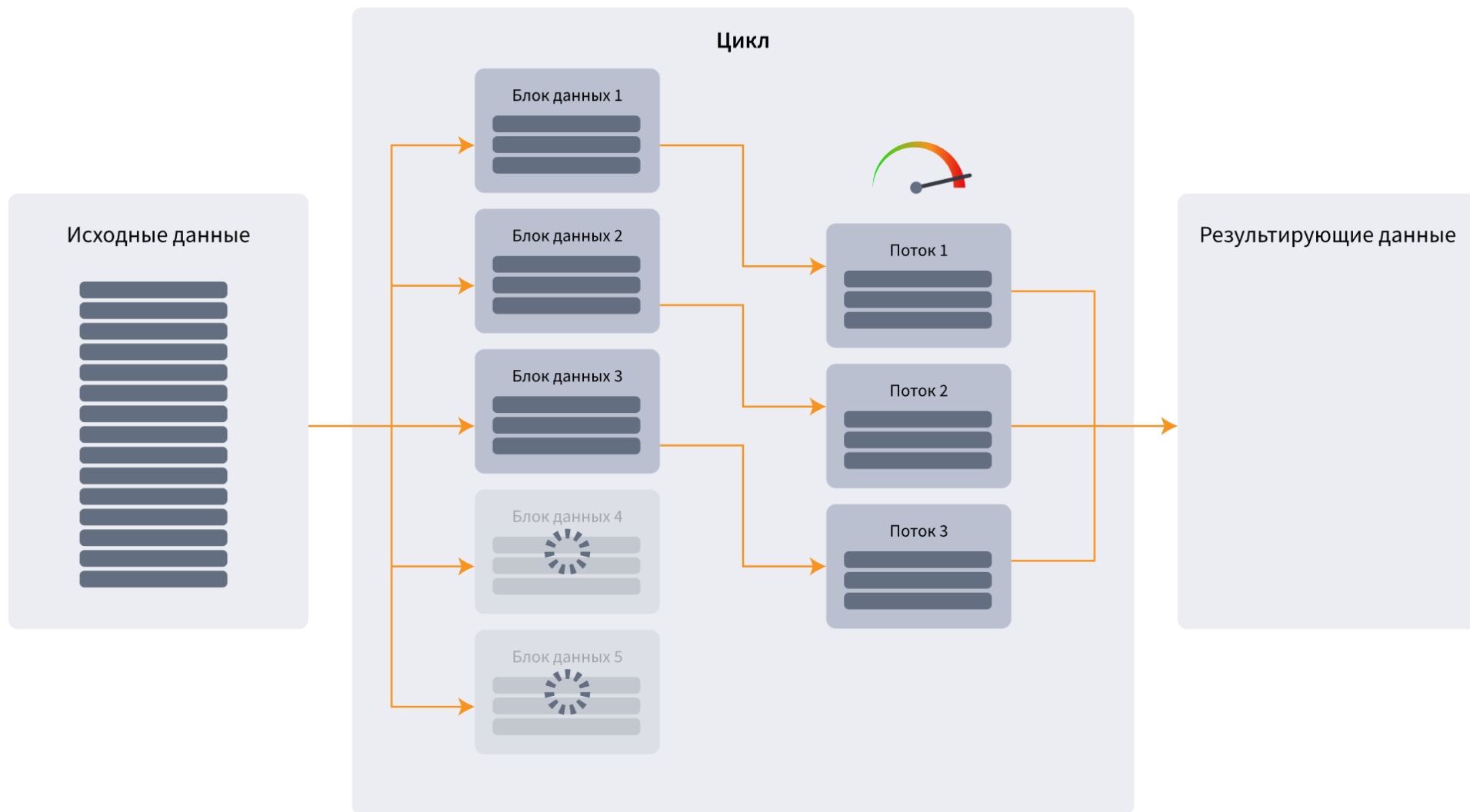
## Собственная реализация многопоточной обработки:

1. Формирование своего пула потоков
2. Мониторинг загрузки процессоров для добавления потока в случае не дозагрузки
3. Выполнение задач последовательно, если излишний параллелизм не выгоден
4. Освобождение потоков, если они не используются

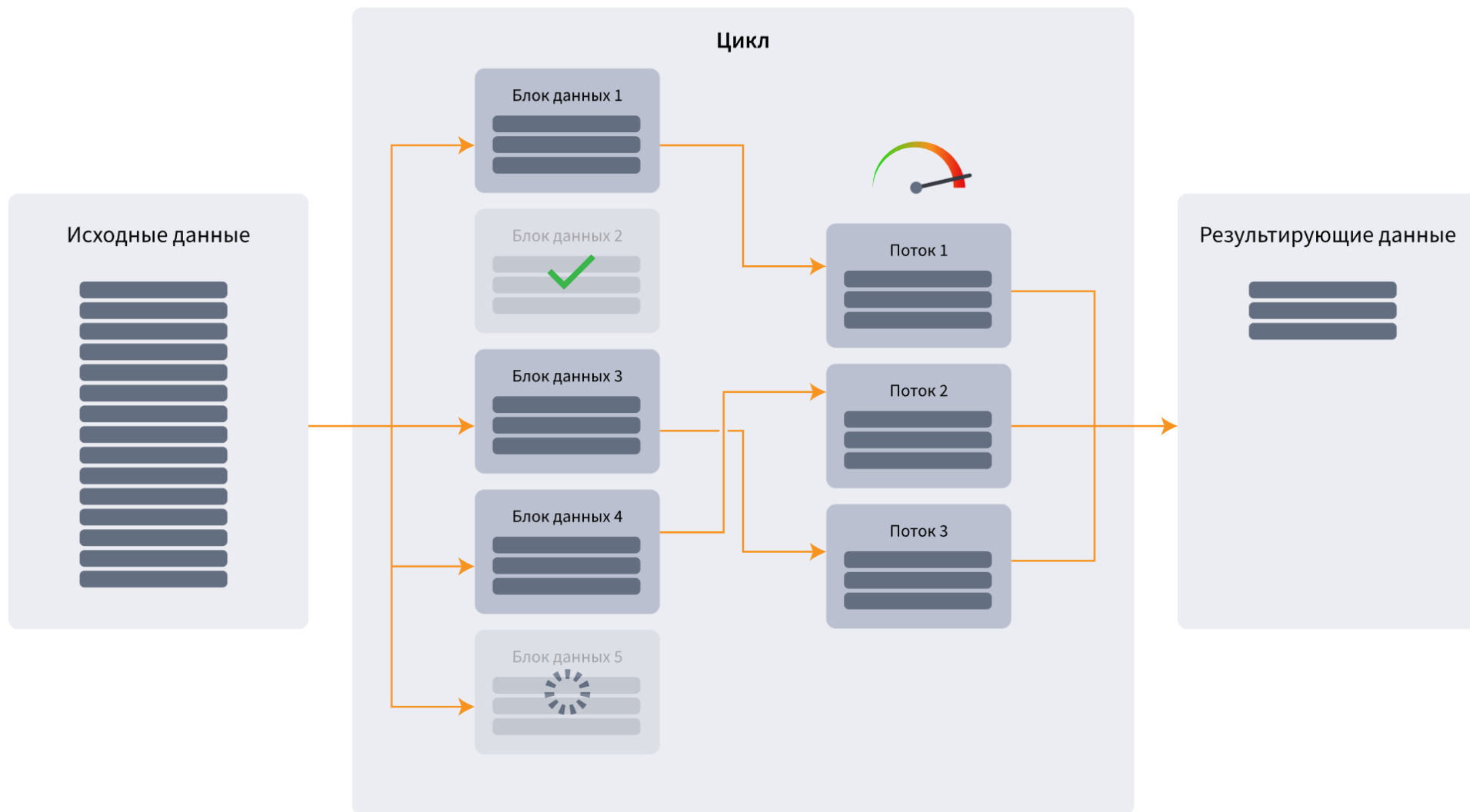
# Параллельный цикл в многопроцессорных системах (1)



## Параллельный цикл в многопроцессорных системах (2)

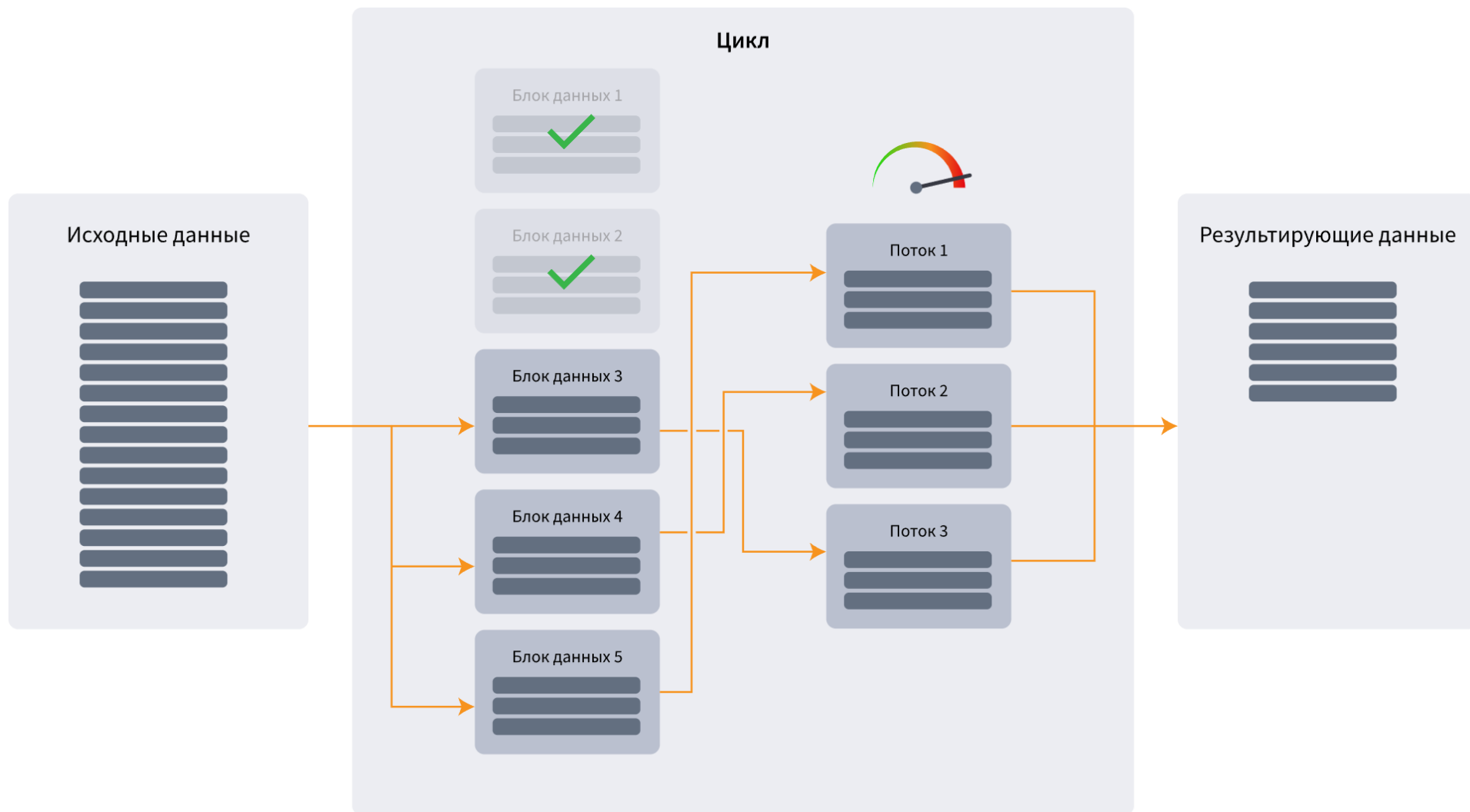


## Параллельный цикл в многопроцессорных системах (3)

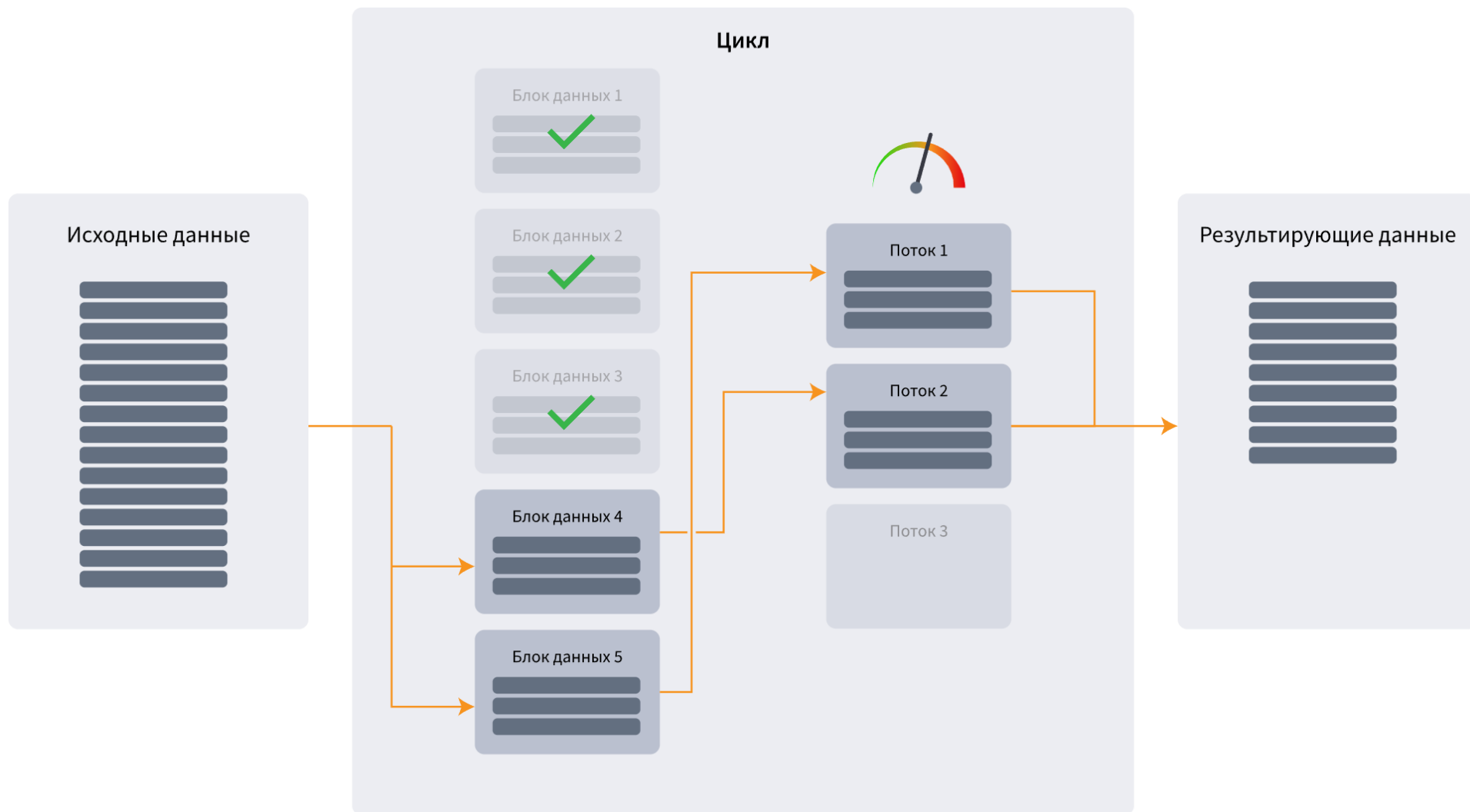




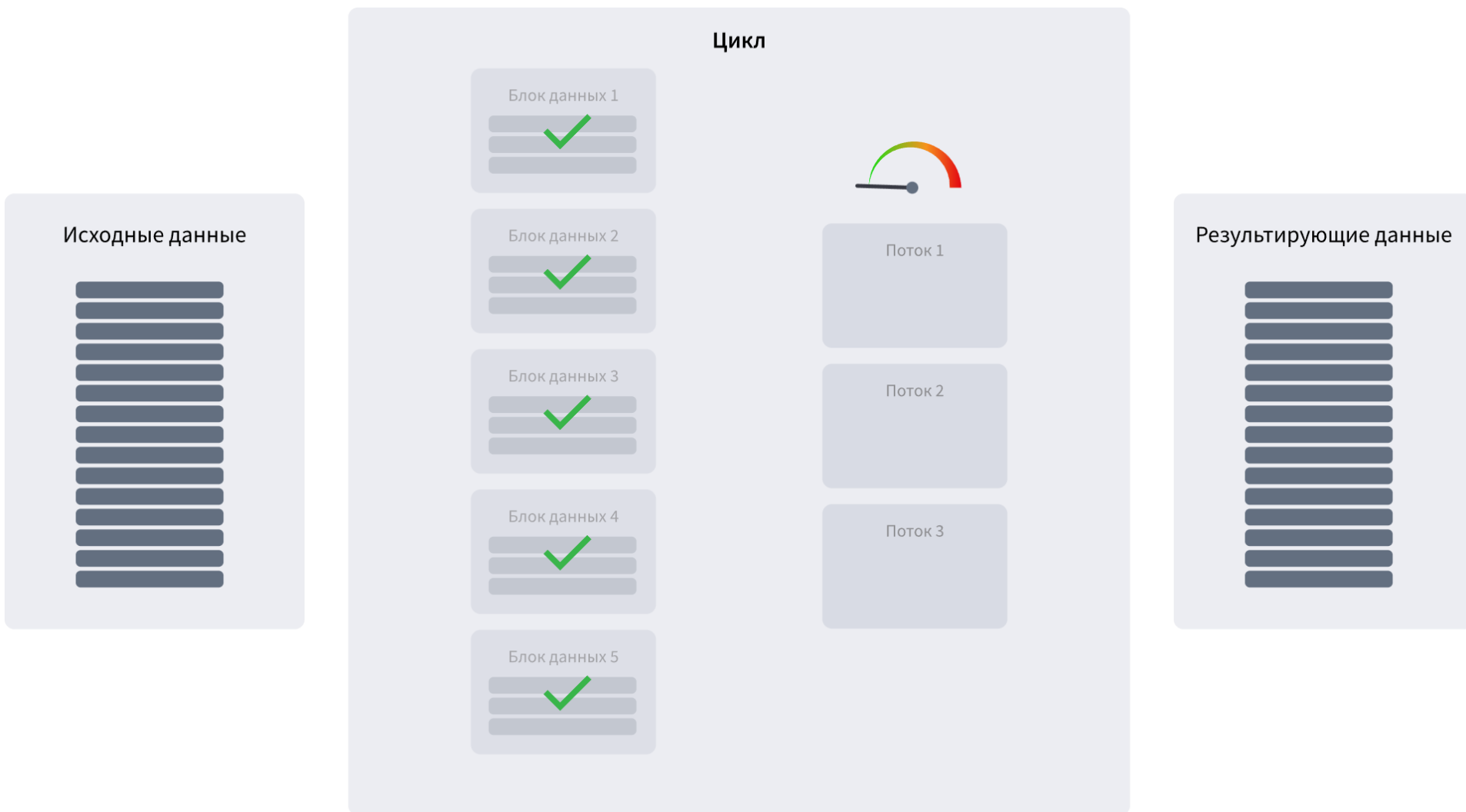
## Параллельный цикл в многопроцессорных системах (4)



## Параллельный цикл в многопроцессорных системах (5)



# Параллельный цикл в многопроцессорных системах (6)



## Эффективная работа со строками

Например, для быстрой фильтрации сравнение происходит при помощи специальной функции, исключающей атомарные операции обращения к общей памяти, блокирующие конвейер всех ядер процессора. Иначе в результате блокировок даже если ядер много фактически система работает в один поток.

Результат – на многоядерных серверах увеличение скорости фильтрации примерно на 20%.

# Источники/приемники данных

У всех СУБД запросы отправленные через один коннект выполняются последовательно. Для увеличения скорости требуется оптимизация:

1. При подключении к БД создается пул коннектов, за счет чего запросы выполняются параллельно
2. Для минимизации накладных расходов на коннект подключения удерживаются открытыми
3. Для экономии памяти на сервере можно настроить принудительное закрытие подключений

Оптимизация импорта из медленных источников (csv, xlsx...):

1. При чтении строковые данные кэшируются, кэш строк организован в виде хэш-таблицы.
2. Используется wuhash (самая быстрая и качественная хэш-функция, 2019 г.) — алгоритм, оптимизированный под современные процессоры.
3. Размер хеш-таблицы подобран так, чтобы он помещался в кэш-память и обеспечивал максимальную производительность.

Импорт замедлился, но уменьшился расход памяти и ускорилось выполнение последующих узлов в сценарии.

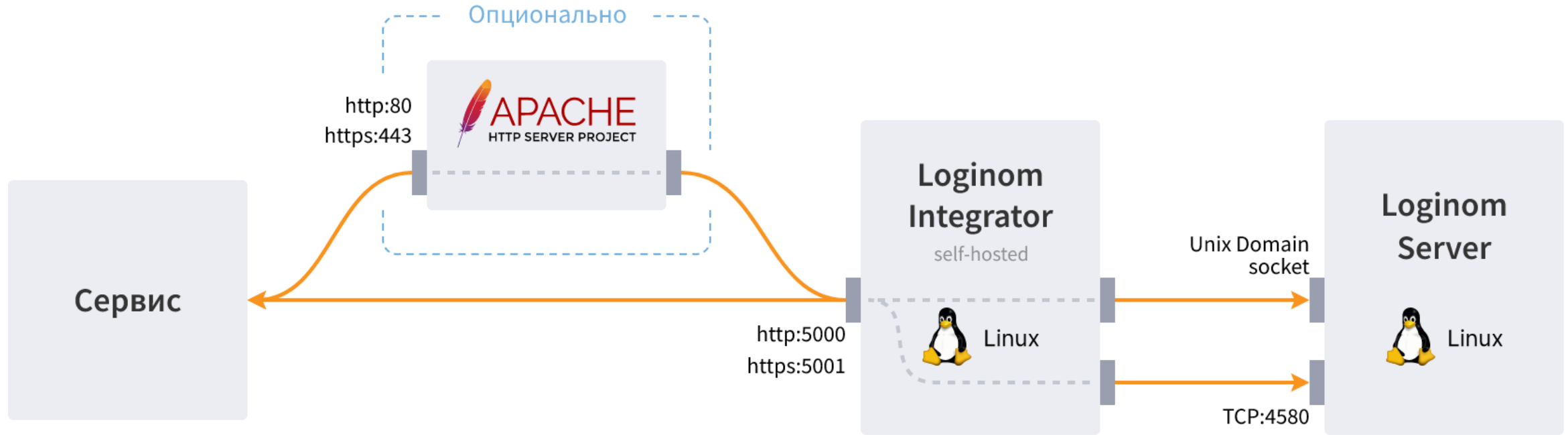
## Оптимизация формата Loginom Data File:

- Замена алгоритма сжатия LZ0 на LZ4 с добавлением расчета контрольной суммы.
- Размер файла увеличился примерно на 10%, но запись быстрее на 13%, а чтение на 39%.
- В процессе экспорта в файле сохраняются уникальные строки, из-за чего скорость записи уменьшилась. При этом при импорте скорость чтения увеличилась, а расход ОЗУ снизился.



# Интеграция компонентов платформы

## Обмен данными между Server-ом и Integrator-ом



Обычно Integrator и Server общаются по TCP, но при работе на одной машине можно использовать Unix domain socket. Отклик быстрее на ~15%.

## Встраивание Python в сценарий

Для параллельного выполнения узла Python анализируемые выборки выгружаются в файлы и запускается несколько интерпретаторов.

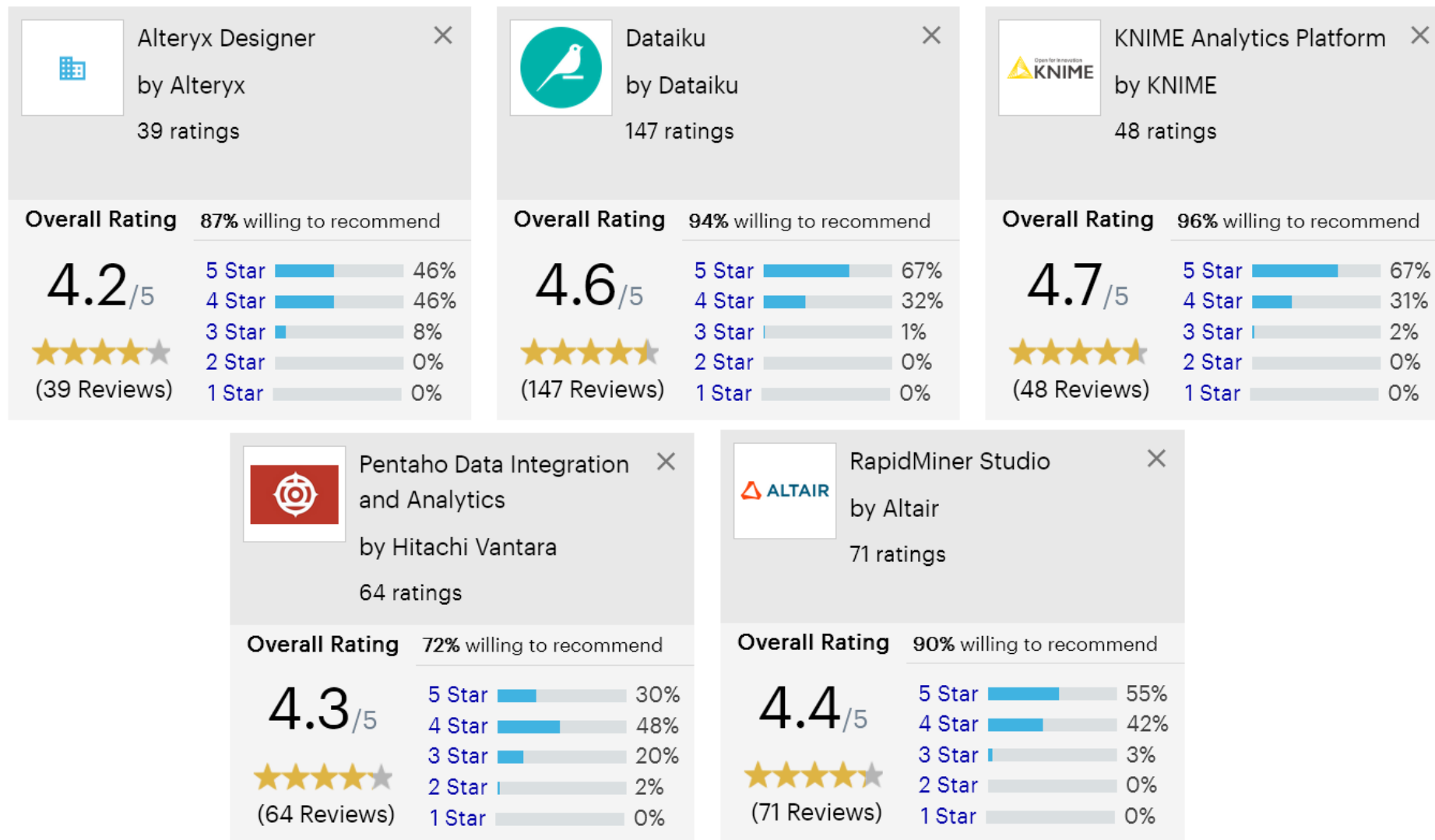
Наборы записываются в файлы бинарного формата с хранением по столбцам и прозрачно для аналитика загружаются в Python. Это повышает скорость обмена данными между Python и другими узлами сценария.

# Сравнение с конкурентами

# Сравнение настольных редакций популярных платформ low-code-аналитики

1. Alteryx
2. Dataiku
3. KNIME
4. Loginom
5. Pentaho
6. RapidMiner

# С Loginom сравниваются лучше мировые продукты по мнению Gartner



## Задача 1. Простой сценарий обработки.

Построение отчетности по продажам, ABC и XYZ-анализ с минимальными настройками. Тест производительности базовых операций работы с данными:

1. Импорт csv-файлов
2. Слияние наборов данных
3. Расчет по формулам
4. Группировки, сортировки
5. Расчет простых статистик
6. Квантование
7. Экспорт csv-файлов

Параметры тестового компьютера: AMD Ryzen 5 5600G 3.9 GHz (6 ядер, 12 потоков), DDR 32 Gb 3200 MHz, SSD ADATA SP90

Workflow - Configuration

Canvas Workflow Runtime Events Meta Info

CANVAS OPTIONS

Layout Direction

Horizontal

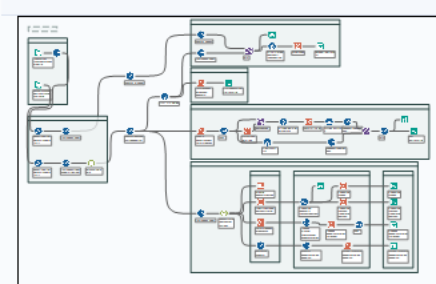
Annotations

Show

Connection Progress

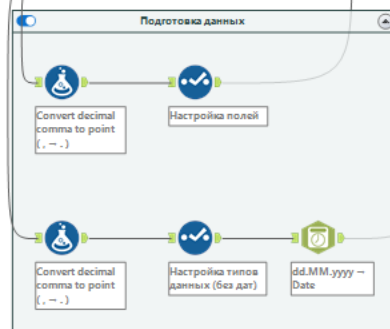
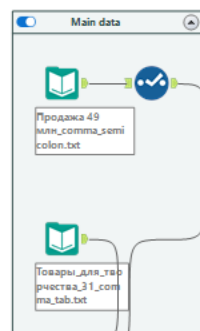
Show Only When Running

Overview



sales.yxmd

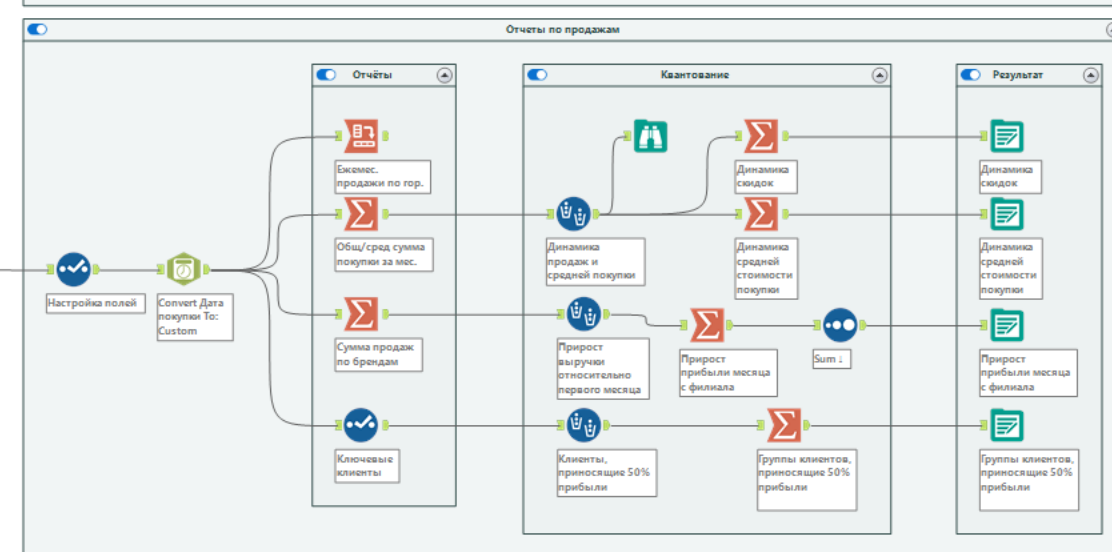
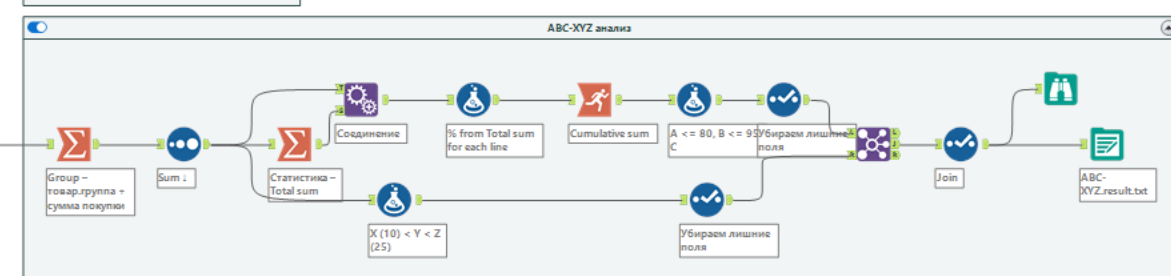
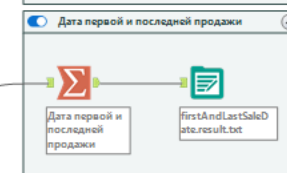
Sample data



Транзит Продаж

First 1.5 mln rows

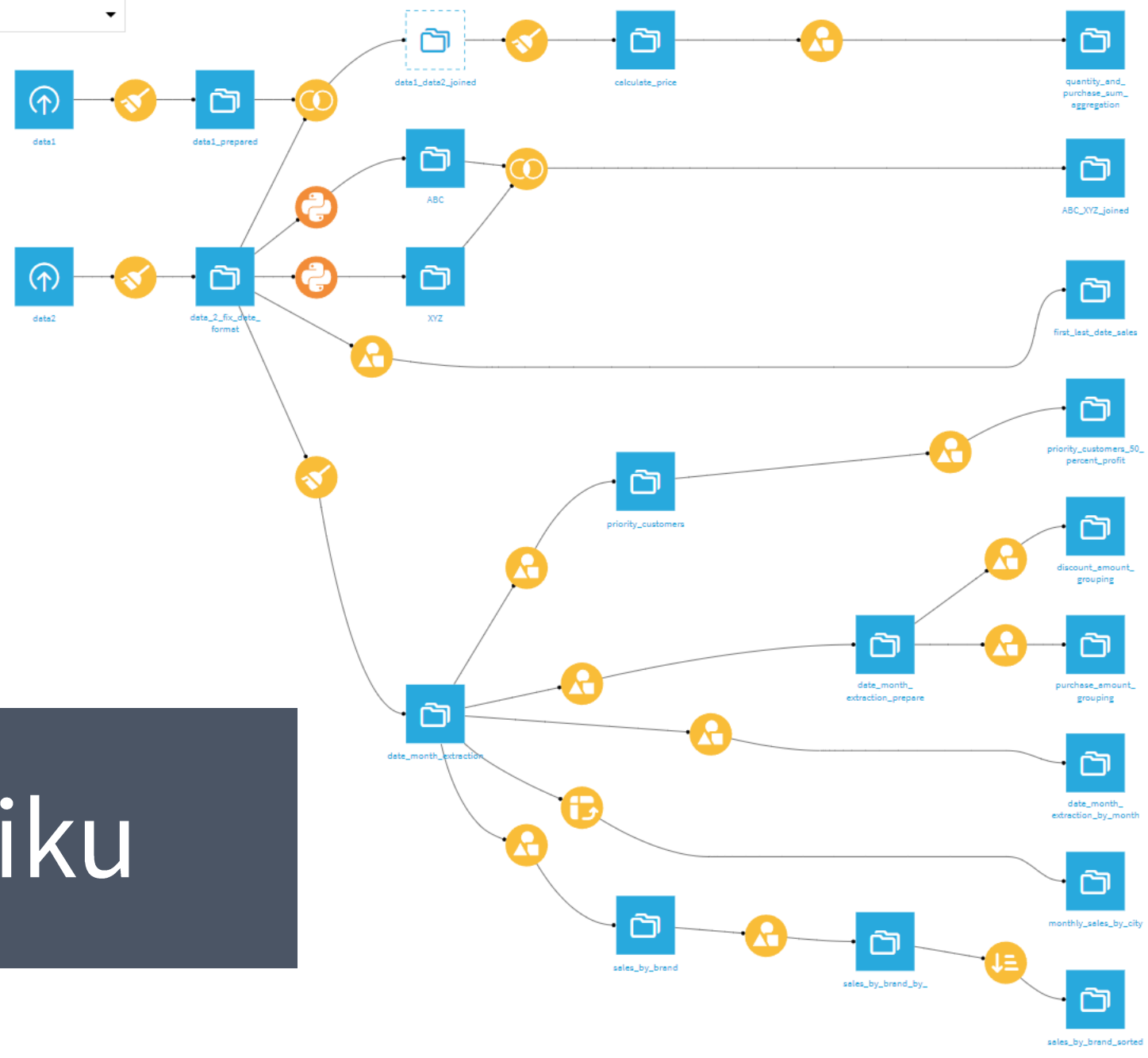
Транзит Товаров для творчества



# Alteryx



22 datasets 20 recipes



# Dataiku

KNIME Analytics Platform

Homecompare

Execute allCancel allReset all

compare > Metanode > Отчеты по продажам > Квантование-1

100%

Auto-Binner

Динамика продаж и средней покупки

Auto-Binner

Динамика продаж и средней покупки

Auto-Binner

Прирост выручки относительно первого месяца

Auto-Binner

Клиенты, приносящие 50% прибыли

Auto-Binner

Динамика скидок

Auto-Binner

Динамика средней стоимости покупки

Auto-Binner

Группы клиентов, приносящие 50% прибыли

groupBy

groupBy

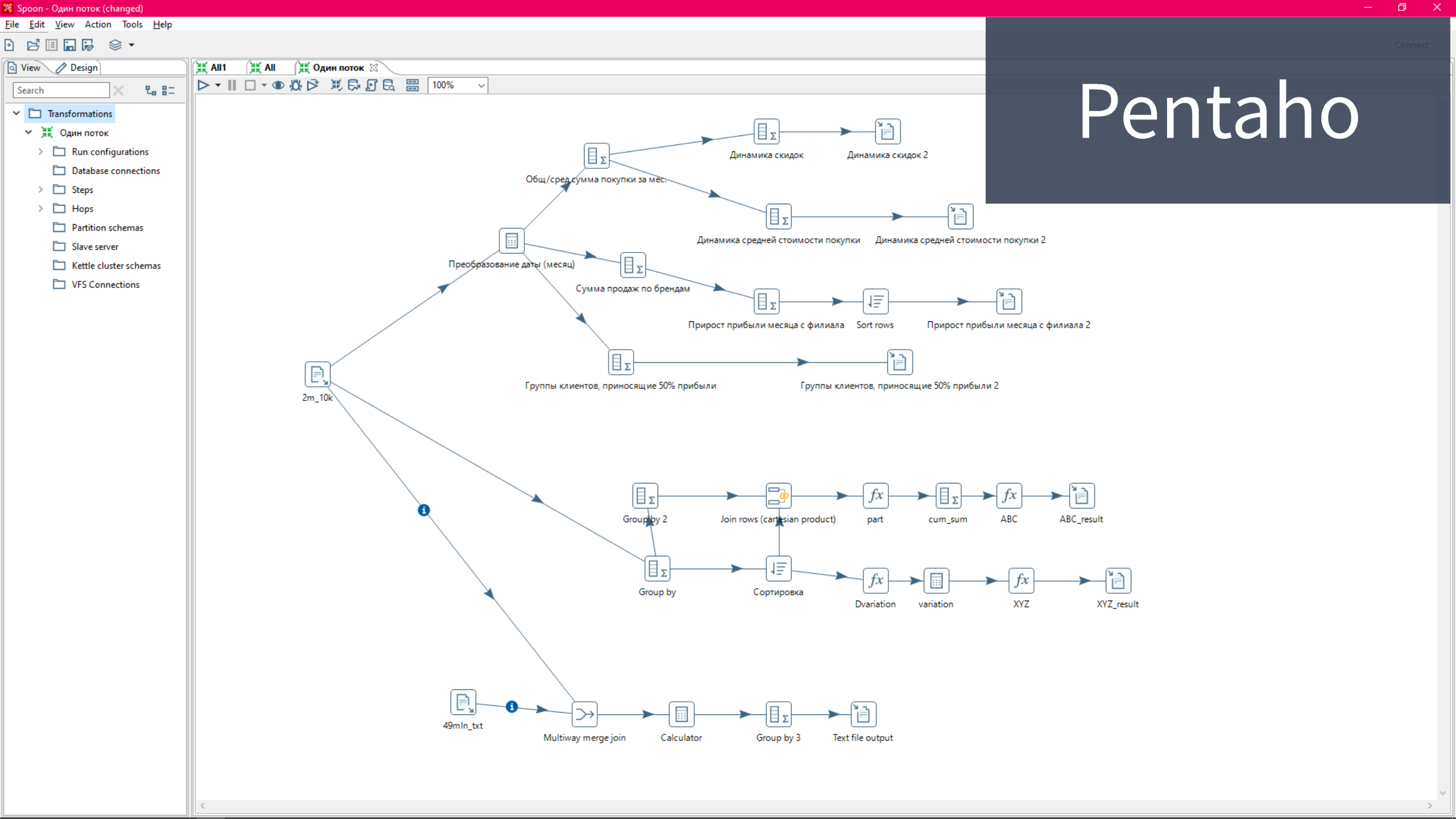
groupBy

Sorter

groupBy

KNIME

To show the node output, please select a configured or executed node.



Pentaho



Views:

Design

Results

Turbo Prep

Auto Model

Interactive  
Analysis

Find data, operators...etc



All Studio

## Repository

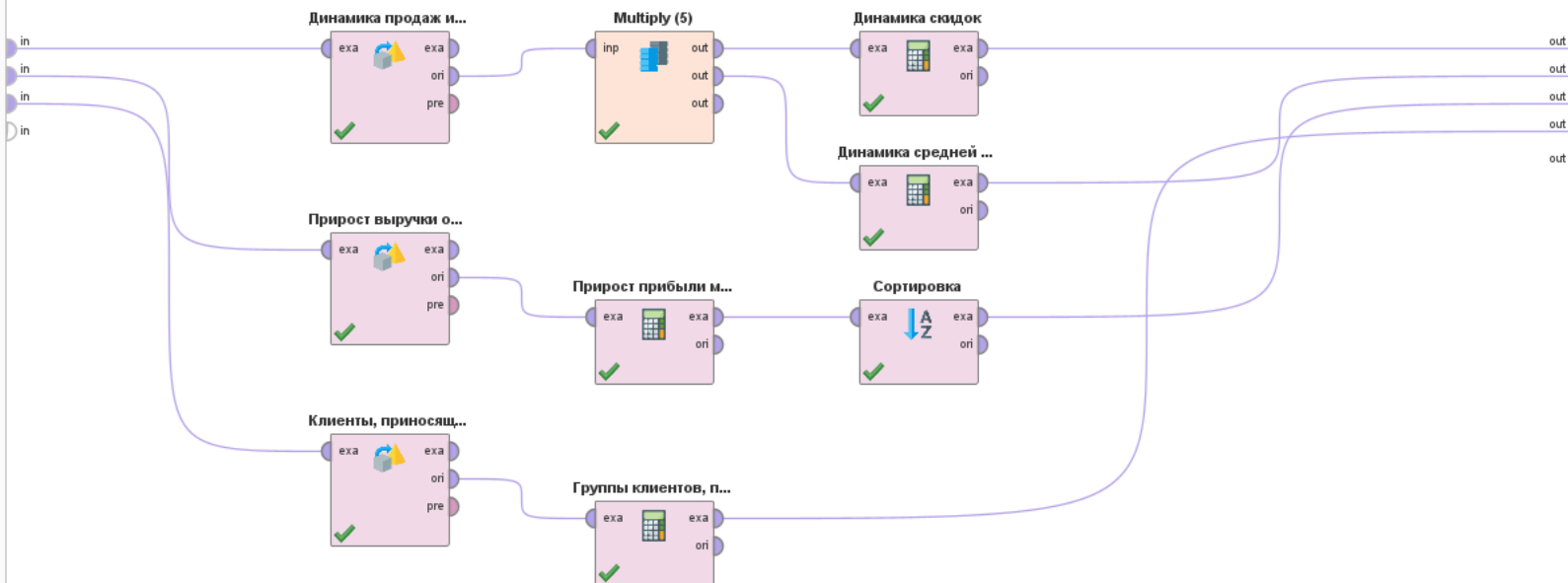
+ Import Data

- Training Resources (connected)
- Samples
- Community Samples (connected)
- Local Repository (Local)
- DB (Legacy)

## Process

Process > ... Отчеты по продажам > Квантование-1

Квантование-1



## Parameters

Квантование-1 (Subprocess)

No parameters to display.

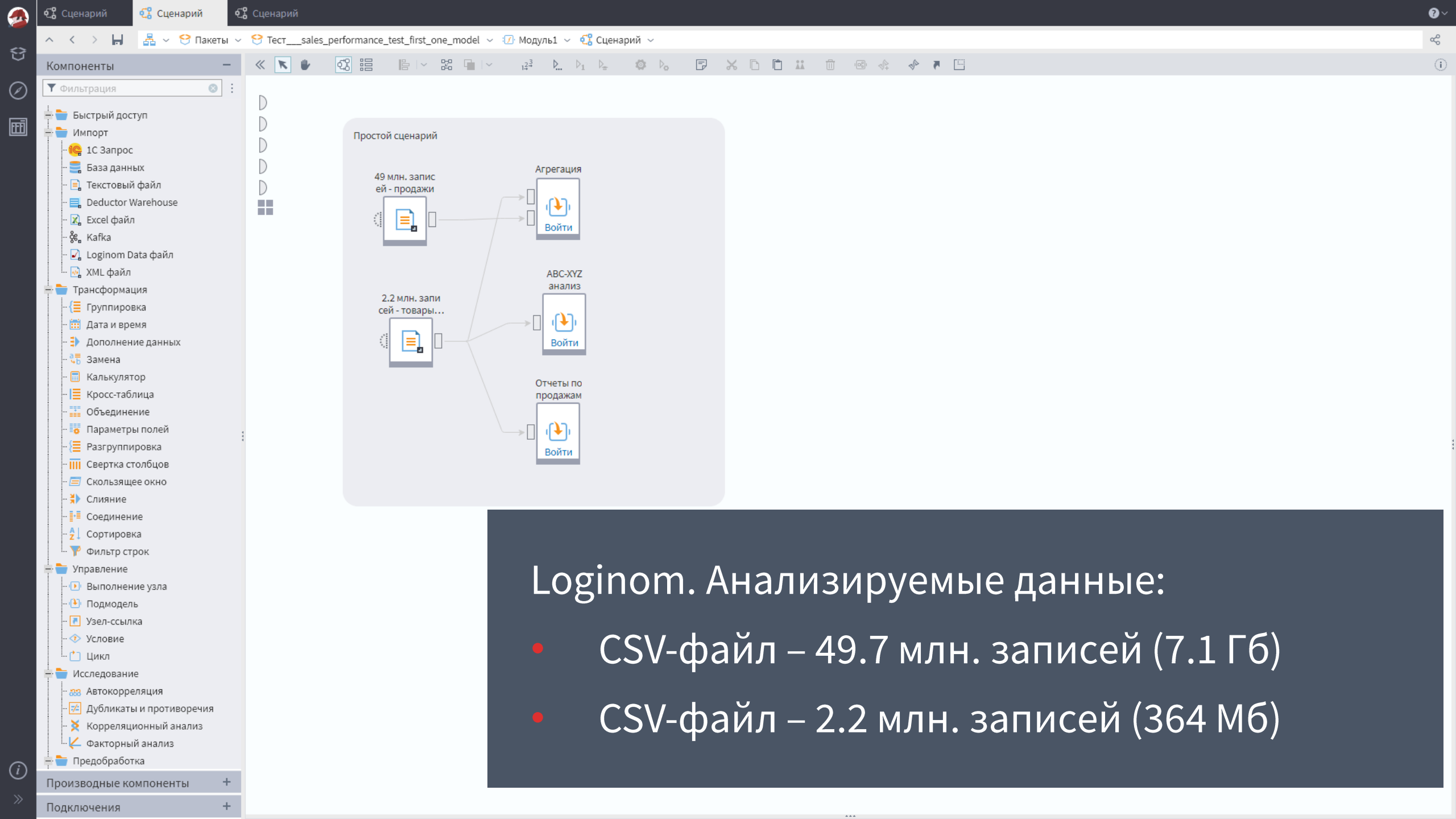
## Operators

Search for Operators







- Data Access (59)
- Blending (81)
- Cleansing (28)
- Modeling (167)
- Scoring (13)
- Validation (30)
- Utility (85)
- Extensions (116)

Get more operators from the Marketplace

# RapidMiner



# Сравнение платформ

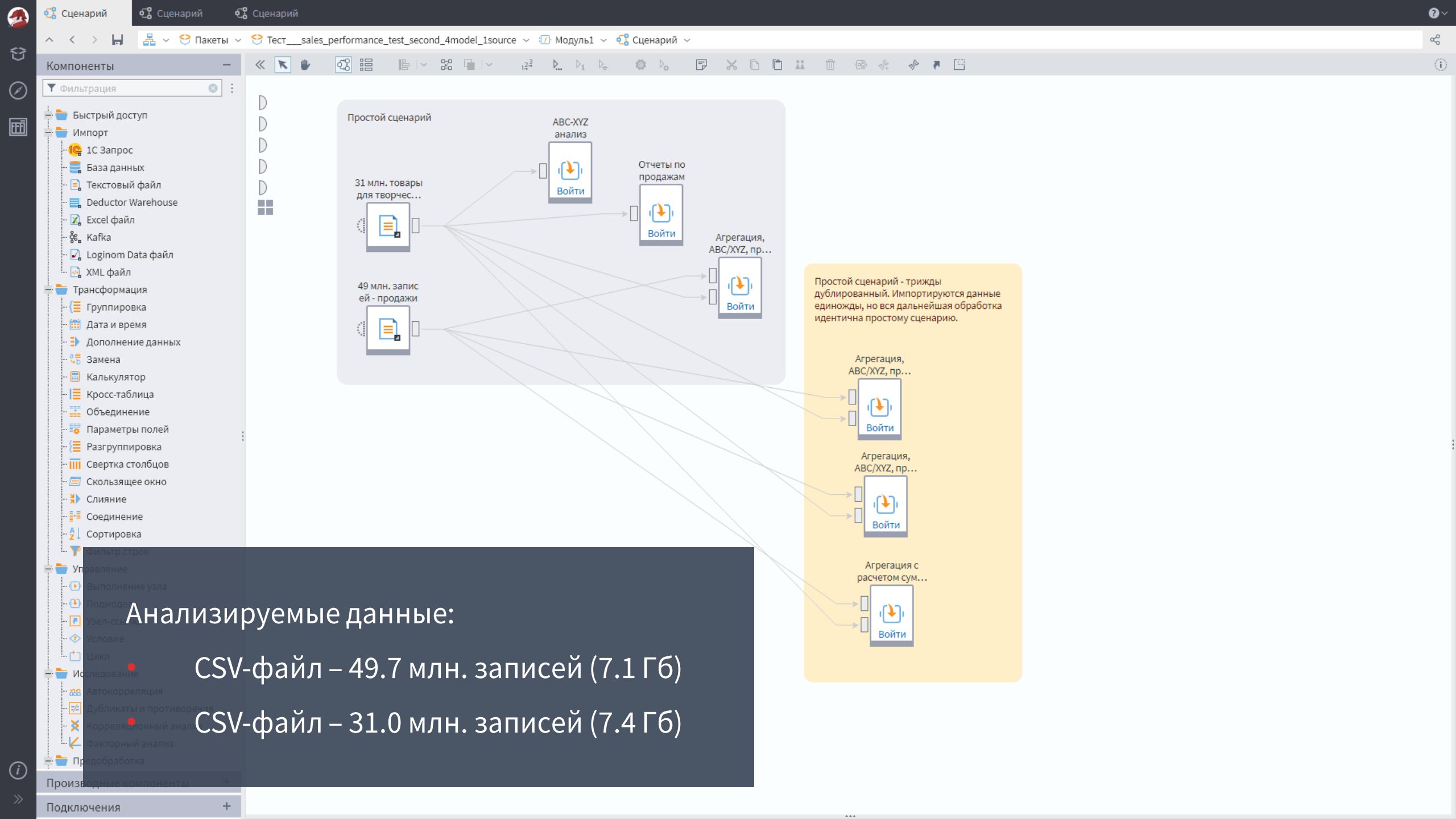
Продукт	Время (сек.)	RAM (Мб)	Загрузка CPU (%)
 Loginom	66	4 839	15
 alteryx	73	10 000	50
 pentaho®	170	1 800	15
 RAPIDMINER	801	8 102	40
 <small>Open for Innovation</small> KNIME	4560	22 500	30
 dataiku	Более 2-х часов или не прогрузилось		

Данные необходимо загрузить в репозиторий, время загрузки суммировалось

## Задача 2. Большие объемы данных, много расчетов







Трижды продублированы ранее разработанные простые сценарии для демонстрации возможностей параллельной обработки.

Эмуляция ситуации при которой большие объемы данных извлекаются из нескольких источников, а затем обсчитываются множеством не очень сложных сценариев.





# Сравнение платформ

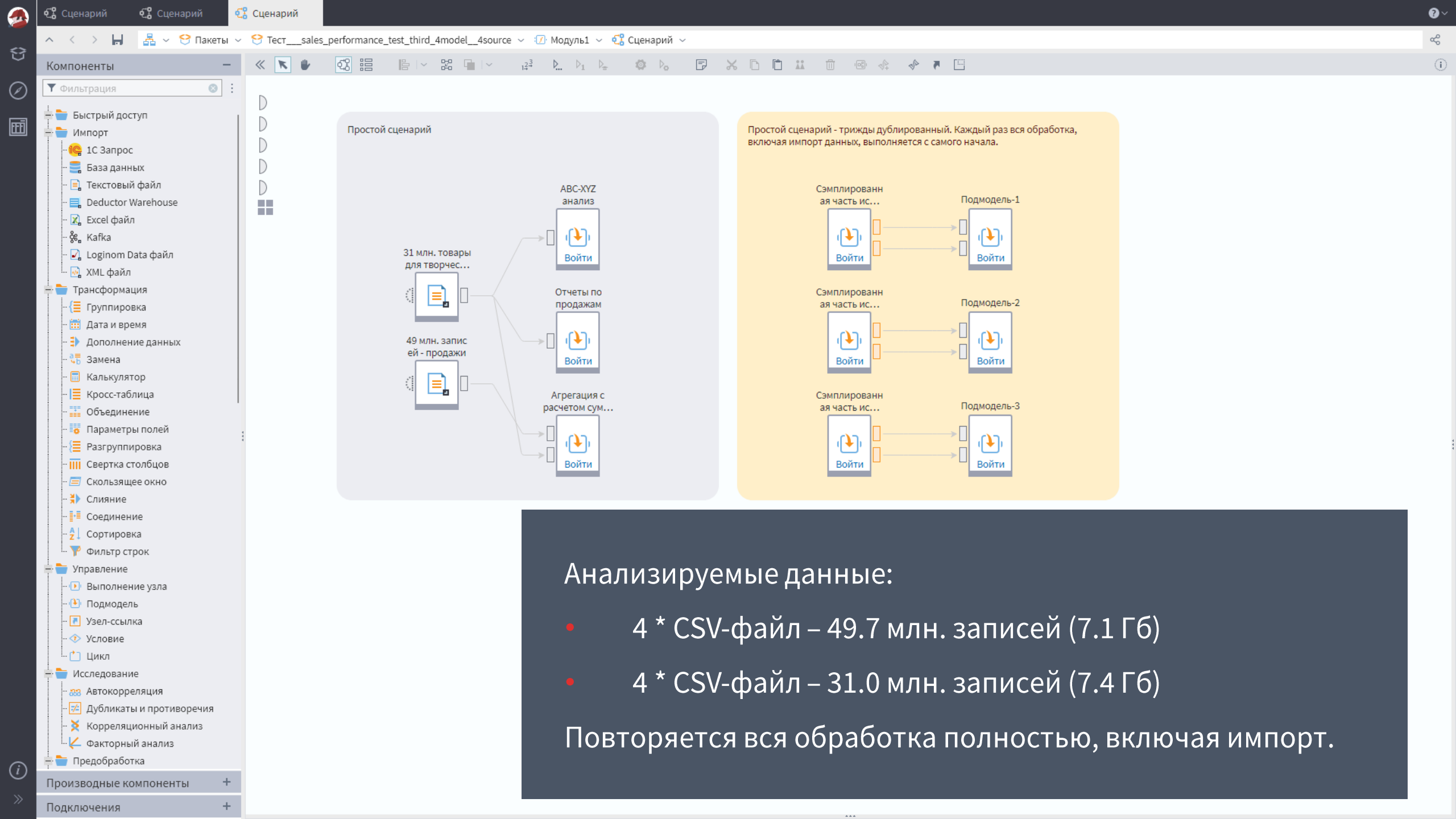
Продукт	Время (сек.)	RAM (Мб)	Загрузка CPU (%)
 Loginom	<div><div></div><div>117</div></div>	<div><div></div><div>13 700</div></div>	<div><div></div><div>85</div></div>
 pentaho®	<div><div></div><div>197</div></div>	<div><div></div><div>1 800</div></div>	<div><div></div><div>35</div></div>
 RAPIDMINER	<div><div></div><div>1 495</div></div>	<div><div></div><div>8 500</div></div>	<div><div></div><div>30</div></div>
 alteryx	Более 2-х часов или не прогрузилось		
 <small>Open for Innovation</small> KNIME	Более 2-х часов или не прогрузилось		
 dataiku	Более 2-х часов или не прогрузилось		

Данные необходимо загрузить в репозиторий, время загрузки суммировалось

### Задача 3. Высоконагруженная обработка

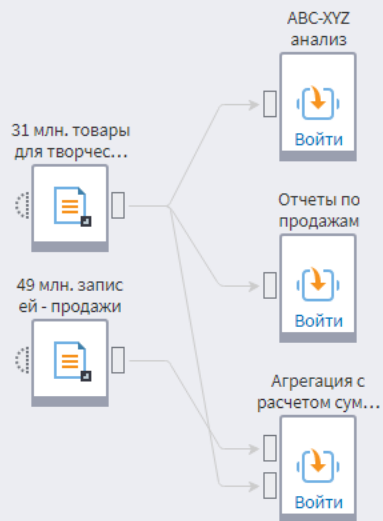
Трижды продублированы все этапы обработки от импорта данных до выгрузки результатов.

Эмуляция ситуации при которой регулярно поступают большие объемы новых данных, которые необходимо прогнать через все этапы обработки – типичный кейс работы высоконагруженных систем.

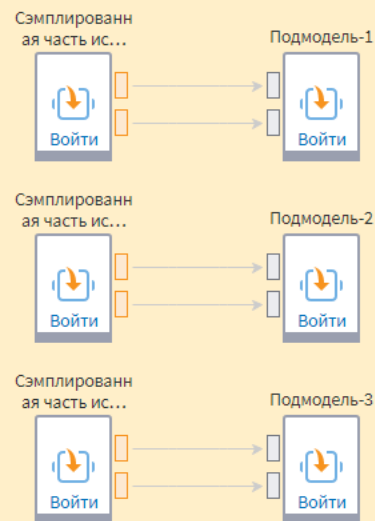


- Быстрый доступ
- Импорт
  - 1С Запрос
  - База данных
  - Текстовый файл
  - Deductor Warehouse
  - Excel файл
  - Kafka
  - Loginom Data файл
  - XML файл
- Трансформация
  - Группировка
  - Дата и время
  - Дополнение данных
  - Замена
  - Калькулятор
  - Кросс-таблица
  - Объединение
  - Параметры полей
  - Разгруппировка
  - Свертка столбцов
  - Скользящее окно
  - Слияние
  - Соединение
  - Сортировка
  - Фильтр строк
- Управление
  - Выполнение узла
  - Подмодель
  - Узел-ссылка
  - Условие
  - Цикл
- Исследование
  - Автокорреляция
  - Дубликаты и противоречия
  - Корреляционный анализ
  - Факторный анализ
- Предобработка

### Простой сценарий



### Простой сценарий - трижды дублированный. Каждый раз вся обработка, включая импорт данных, выполняется с самого начала.



## Анализируемые данные:

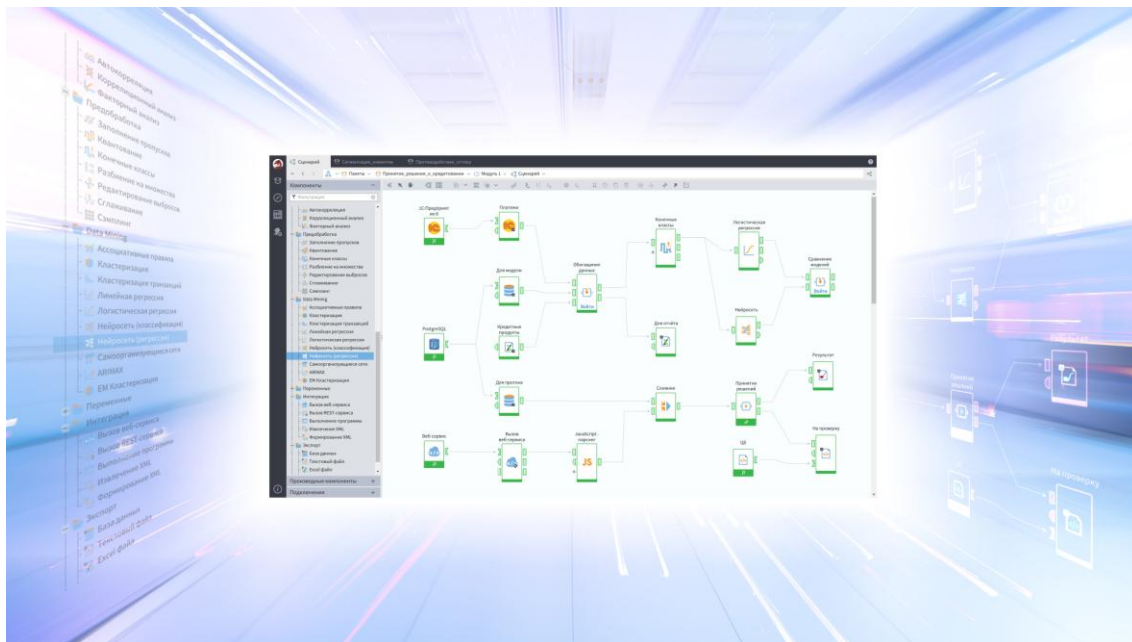
- 4 \* CSV-файл – 49.7 млн. записей (7.1 Гб)
- 4 \* CSV-файл – 31.0 млн. записей (7.4 Гб)

Повторяется вся обработка полностью, включая импорт.

# Сравнение платформ

Продукт	Время (сек.)	RAM (Мб)	Загрузка CPU (%)
 Loginom	226	27 000	85
 RAPIDMINER	5270	21 000	50
 alteryx	Более 2-х часов или не прогрузилось		
 pentaho®	Более 2-х часов или не прогрузилось		
 <small>Open for Innovation</small> KNIME	Более 2-х часов или не прогрузилось		
 dataiku	Более 2-х часов или не прогрузилось		

Данные необходимо загрузить в репозиторий, время загрузки суммировалось



Loginom демонстрирует выдающуюся эффективность и производительность в сравнении с любой аналитической low-code платформой даже при использовании настольной редакции без тонких настроек. Чем больше источников и объемы данных, тем сильнее отрыв.