

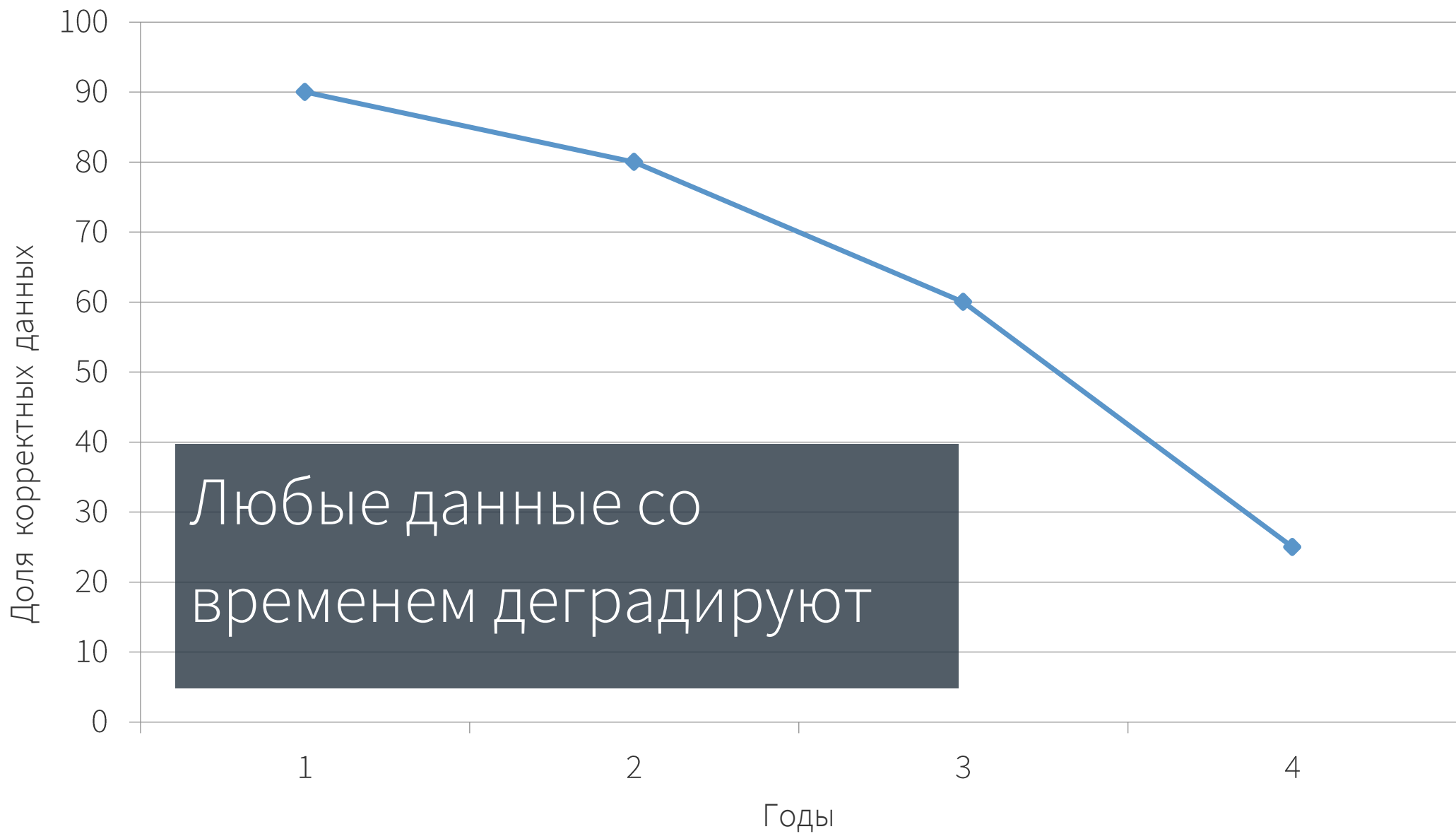


# Loginom Data Quality: очистка данных



Люди ошибаются, поэтому проблемы с качеством данных будут всегда:

- Пропуски
- Опечатки
- Дубли и противоречия
- Фиктивные сведения
- Устаревшие данные



# Последствия:

- Безадресный маркетинг
- Раздражающие коммуникации
- Потерянные контакты
- Отток клиентов
- Искаженная отчетность

Поле	Значение	Ошибка
Имя	Сергей	Первая буква - латинская
Фамилия	Петрович	Неверное поле
Страна	Лимония	Нет такой страны
Город	Мсква	Опечатка
Телефон	0000000	Фиктивный номер
E-mail	<a href="mailto:1@siteforspam.com">1@siteforspam.com</a>	Нет такого домена
Адрес	Новые Васюки 10	Нет адреса
Паспорт	АБВГДЕЖЗ	Не стандарт



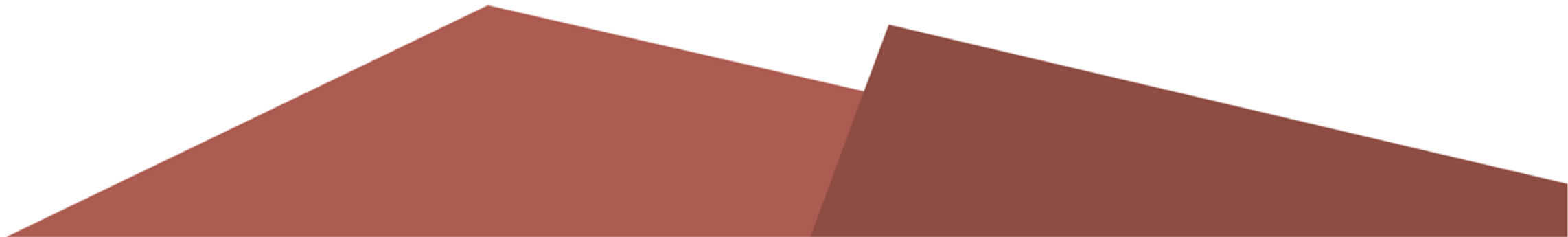
# Решение Loginom Data Quality

Автоматически  
исправляет ошибки,  
приводит к  
стандартному виду,  
восстанавливает и  
обогащает  
клиентские данные

Задача	Описание
Очистка	Исправление ошибок, опечаток, фиктивных данных
Обогащение	Дополнение полезными данными, заполнение пропусков, актуализация
Стандартизация	Унификация представления данных
Дедупликация	Обнаружение дублей и «похожих» записей, объединение их в группы
Формирование «золотой записи»	Создание профиля клиента, включающей актуальные данные всех систем *

\* В рамках проекта

# Очистка, обогащение и стандартизация

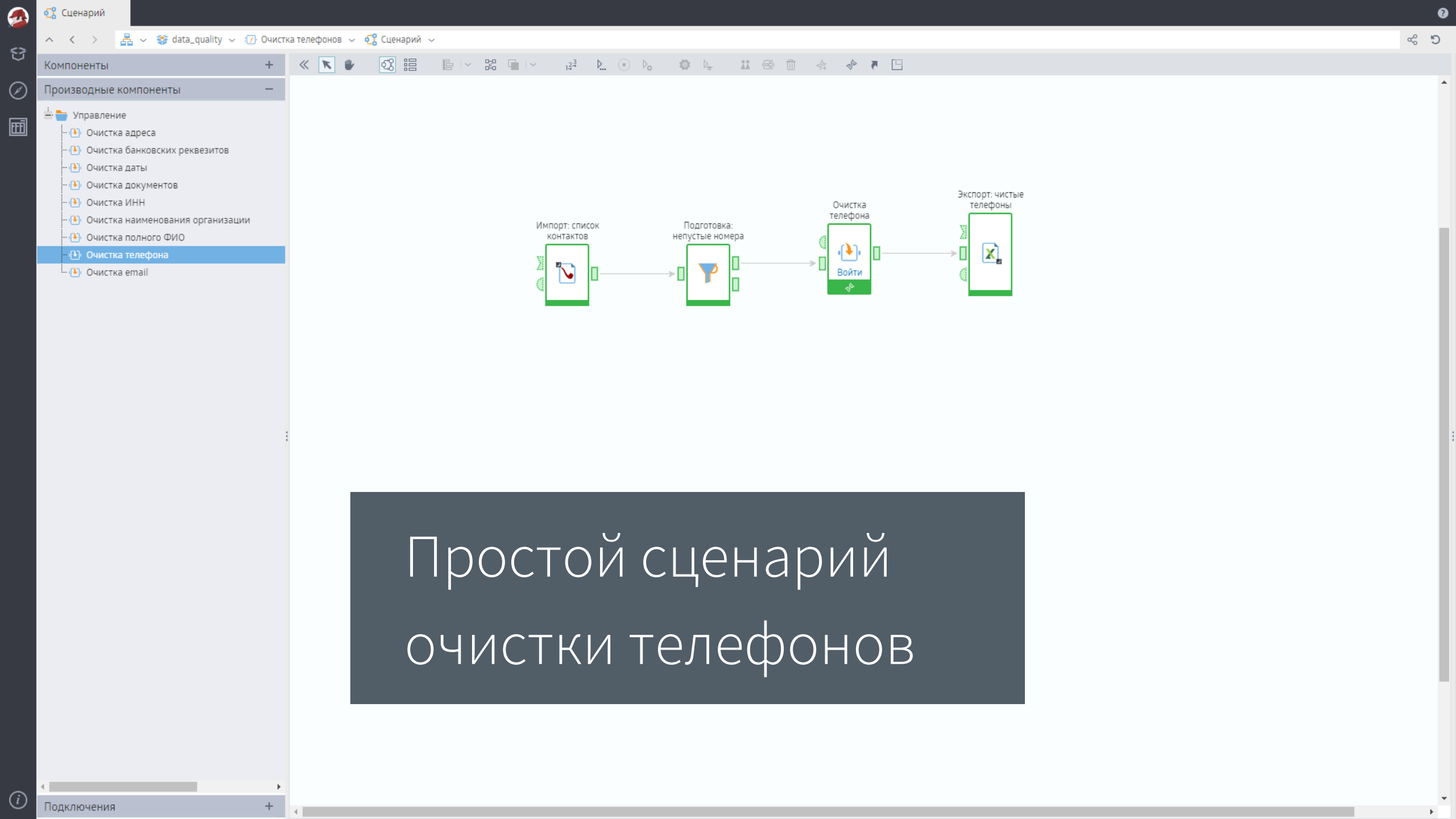




Очищаемые домены	Физ. лица	Юр. лица
Фамилия, имя, отчество	●	
Название организации		●
Почтовый адрес	●	●
Телефоны	●	●
Электронная почта	●	●
Удостоверения личности	●	
Реквизиты организаций		●
Даты	●	●
Банковские реквизиты		●

# Пример схемы очистки ФИО

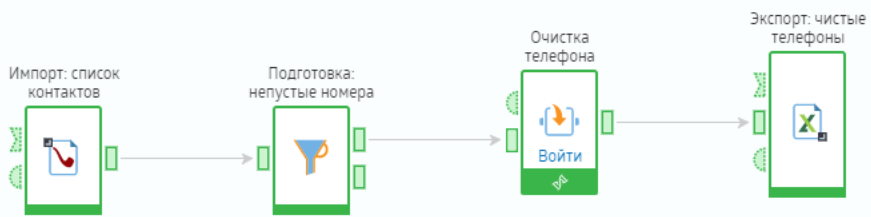




Компоненты +

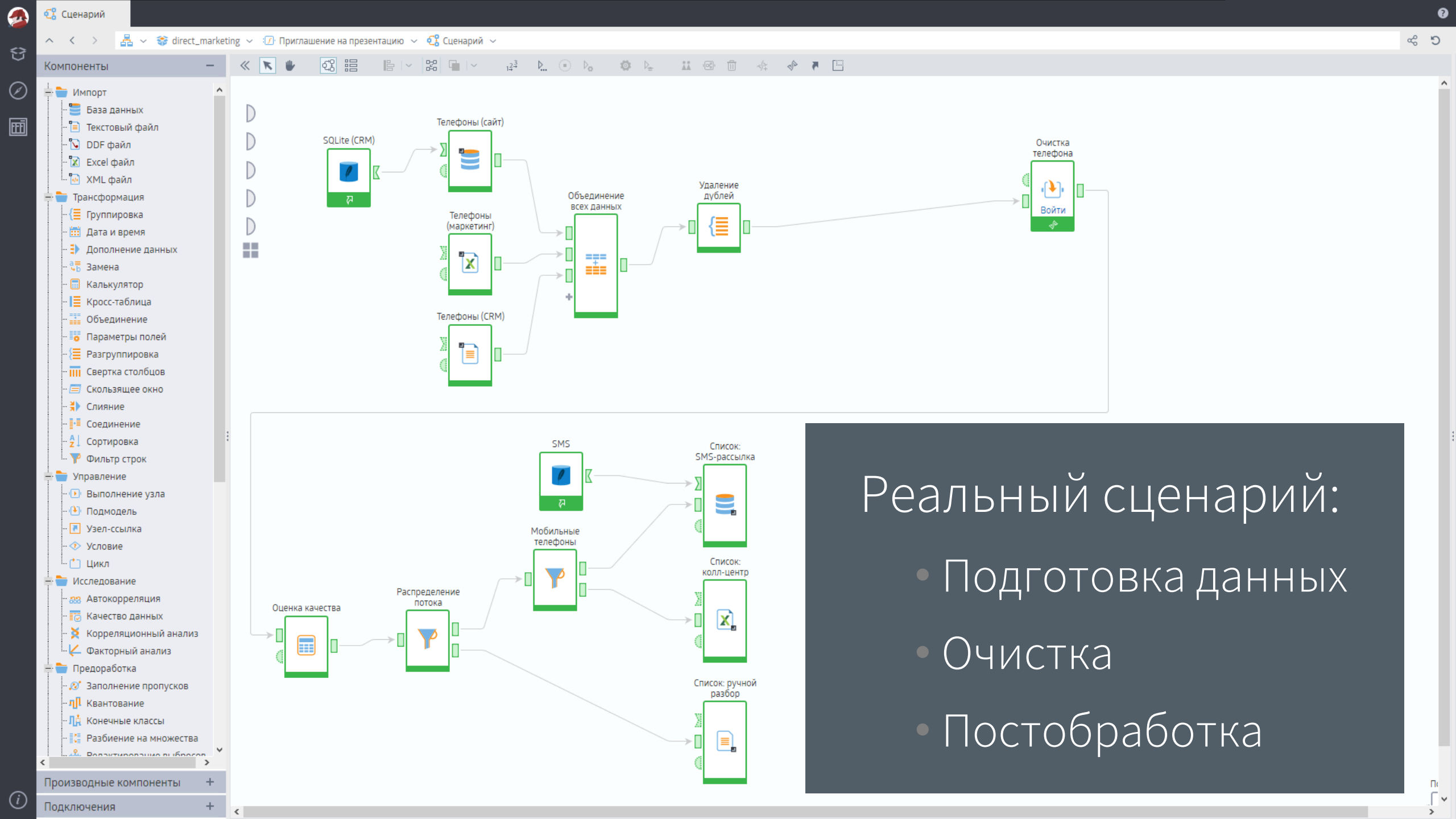
Производные компоненты -

- Управление
  - Очистка адреса
  - Очистка банковских реквизитов
  - Очистка даты
  - Очистка документов
  - Очистка ИНН
  - Очистка наименования организации
  - Очистка полного ФИО
  - Очистка телефона**
  - Очистка email



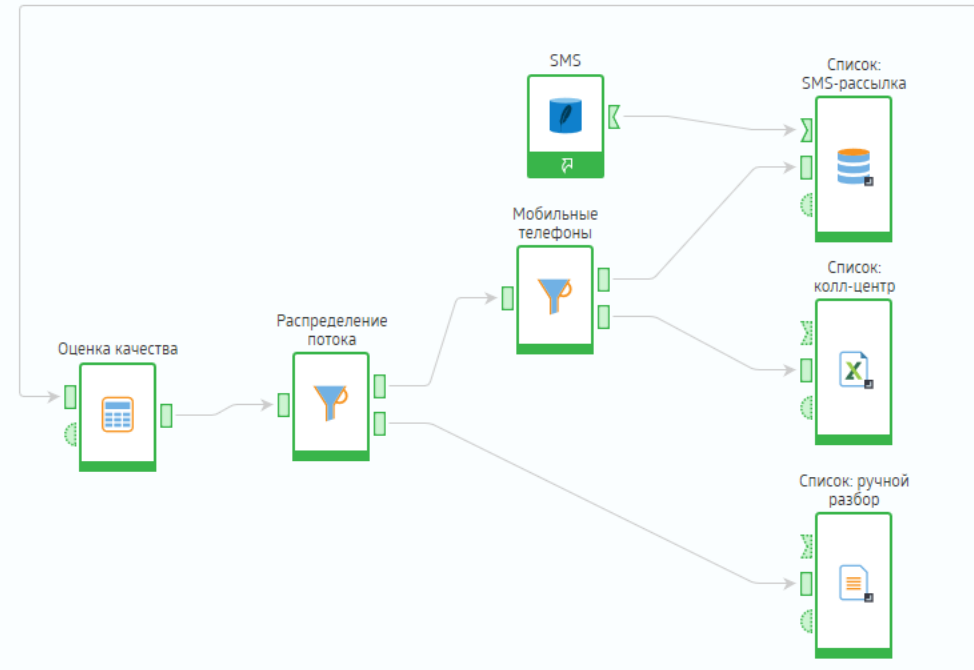
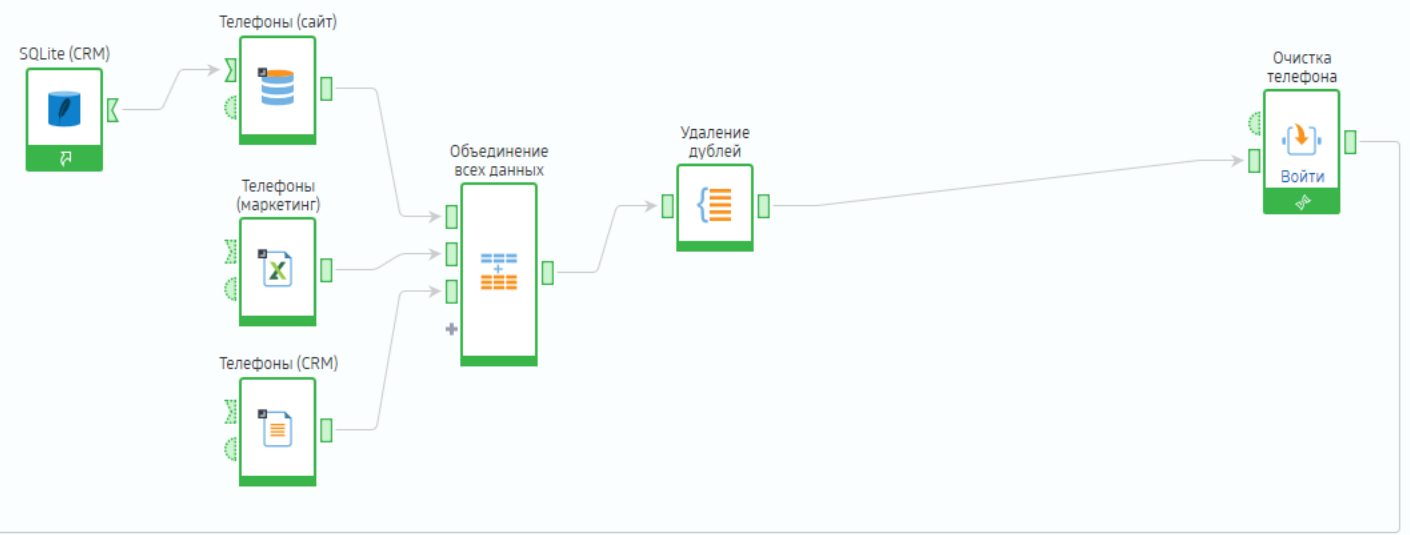
Простой сценарий  
очистки телефонов

Подключения +



Компоненты

- Импорт
  - База данных
  - Текстовый файл
  - DDF файл
  - Excel 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
  - Качество данных
  - Корреляционный анализ
  - Факторный анализ
- Предобработка
  - Заполнение пропусков
  - Квантование
  - Конечные классы
  - Разбиение на множества
  - Восстановление выбросов



Реальный сценарий:

- Подготовка данных
- Очистка
- Постобработка

# Телефон – пример очистки

ab ID	8
ab Исходный номер телефона	8 908 92 24933 справочная
ab Очищенный номер телеф...	+7 (908) 9224933
ab Добавочный номер теле...	
ab Тип номера телефона	Мобильный
ab Остаток строки	справочная
ab Код страны	7
ab Код ABC/DEF	908
ab Номер телефона	9224933
ab Страна	Россия
ab Регион	Свердловская область
ab Город	
ab Оператор	ООО "ЕКАТЕРИНБУРГ-2000"
ab Часовая зона 1	МСК+2
ab Часовая зона 2	московское время плюс 2 часа
ab Часовая зона 3	UTC+5
ab Лог	В номере 11 цифр и первая 8, она заменяется на 7
7 Дата очистки	20.10.2017 18:00:20

Исходные данные

Результат обработки

# Почтовый адрес – пример очистки

ab ID	252162
ab Исходный адрес	ЭНГЕЛЬС-23,УЛ.,АТКАРСКАЯ,24
ab Найденный адрес	413110, Саратовская область город Энгельс улица Аткарская д. 24
90 Доля найденных букв запроса	0,783
12 Масштаб адреса	4
ab Код КЛАДР	640000130000006
ab Ключ ФИАС	cd6746ba-1dab-4dc7-a86e-c84b71fdaf5b
ab Почтовый индекс	413110
0/1 Признак актуальности	Ложь
ab Полный адрес актуального АЭ	Саратовская Область Энгельсский Район Рабочий поселок Приволжский Улица Аткарская
ab Код КЛАДР актуального АЭ	640390000770037
ab Ключ ФИАС актуального АЭ	6272133e-3bb0-4dff-acee-8e0293387f20
ab Почтовый индекс актуально...	413110
ab Тип субъекта РФ	Область
ab Наименование субъекта РФ	Саратовская
ab Тип "города"	Город
ab Наименование "города"	Энгельс
ab Тип "улицы"	Улица
ab Наименование "улицы"	Аткарская

Исходные  
данные

Результат  
обработки

# Почтовый адрес – результат очистки

3 группы адресов по качеству:

- **Белая**, можно использовать – адрес подтвержден и разобран до дома
- **Серая**, на ручную проверку – не найден дом, есть коллизии, нет уверенности
- **Черная**, не пригоден для работы – не почтовый адрес, невозможно идентифицировать адрес

Сырой адрес	Результат очистки	Зона
СОЧИ, ул. ШОССЕЙНАЯ д. 5В кв. 25, 354037	354037, Краснодарский край город Сочи улица Шоссейная д. 5В, кв. 25	Белая
614000;Пермский край;Город ПЕРМЬ;Улица СЕМФИРОПОЛЬСКАЯ;д. 86	614089, Пермский край город Пермь улица Симферопольская д. 86	
644117 обл Омская Г. ОМСКУЛ. ЗМОЛОДЕЖНАЯ дом 62 корп. СЕК. 5 кв. 64	644117, Омская область город Омск улица Молодежная 3-я д. 62, кв. 64	
414009 Астрахань Атарбекова 17 12	414009, Астраханская область город Астрахань улица Атарбекова д. 17, кв. 12	Серая
	414009, Астраханская область город Астрахань площадь Атарбекова д. 17, кв. 12	
	414009, Астраханская область город Астрахань переулок Атарбекова д. 17, кв. 12	
край Краснодарский Г. КРАСНОДАР ПР-КТ. И-М. ПИСАТЕЛЯЗНАМЕНСКОГО дом 14 кв. 157	350065, Краснодарский край город Краснодар проспект им писателя Знаменского д. 14, кв. 157	
Самара в/ч 12454	443000, Самарская область город Самара в/ч 12454	
Калужская область Д ПЕРЕДОЛЬ улица НЕТ стр. д. 10 кв.	249160, Калужская область Жуковский район деревня Передоль д. 10	
2 а ул. Усть-Курдюмская	410018, Саратовская область город Саратов улица Усть-Курдюмская д. 2А	Черная
10А ул.м-он Солопова		
10-Я КРАСНОАРМЕЙСКАЯ,22А		
Всеволожск нет Промзона Кирпичный завод		



Домен	Записей/сек.
Фамилия, имя, отчество	540
Название организации	262
Почтовый адрес	5
Телефоны	714
Электронная почта	9 091
Удостоверения личности	1 785
ИНН	2 272
Даты	7 142
Банковские реквизиты	562

Набор:

100 000 записей

Режим:

Пакетная обработка

Сервер:

Windows Server 2012 R2

x64/ Intel Xeon E5, 4 ядра

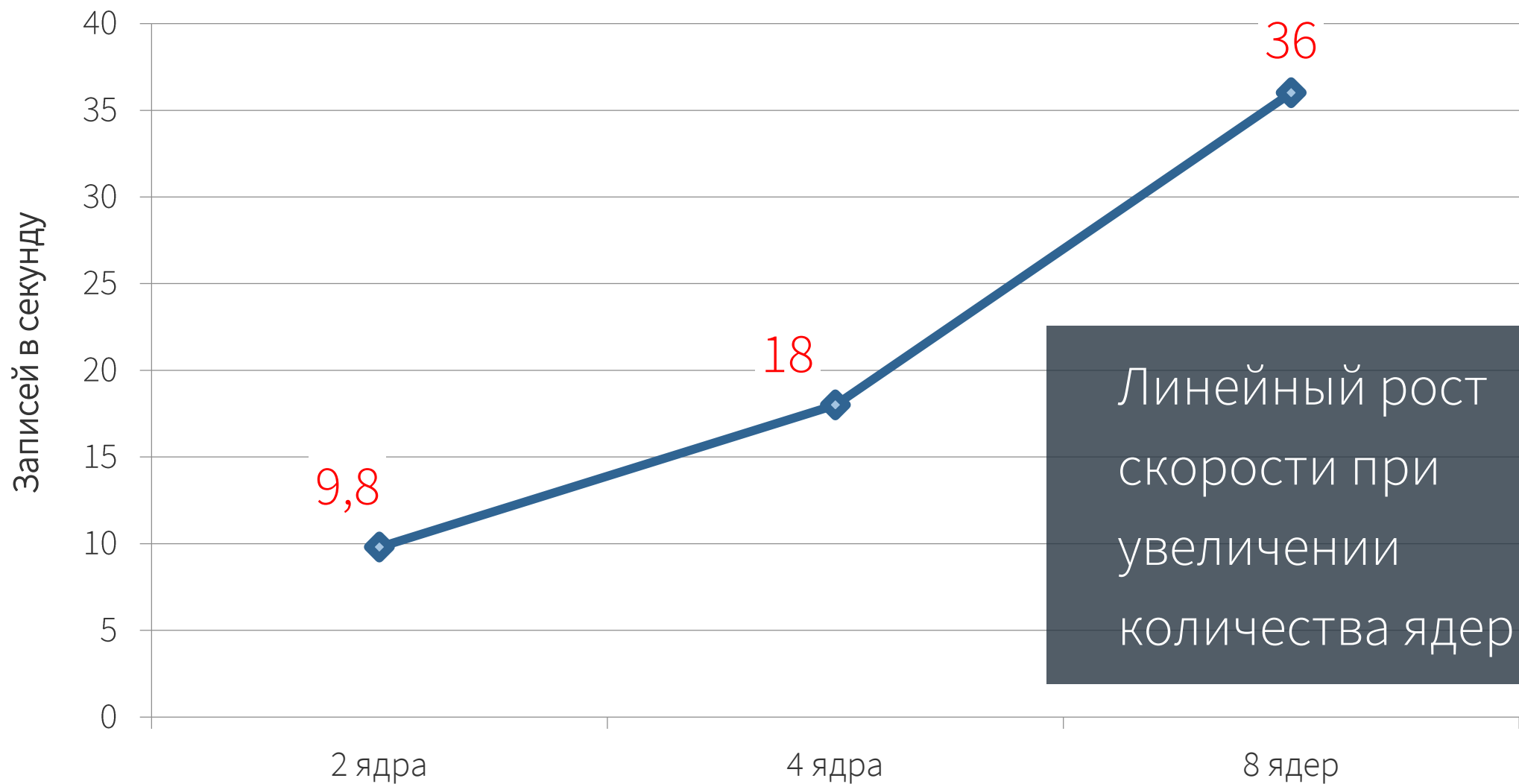
2.60 GHz/16 Gb/1 Tb

Поставляемые справочники	Записей
Почтовые адреса: до улиц	1 909 597
Почтовые адреса: дома	32 592 474
Почтовые адреса: квартиры	54 280 979
Фамилии	752 367
Имена	10 088
Отчества	28 305
Операторы связи/регионы	396 037
Недействительные паспорта	128 631 627
E-mail домены	32 440

8

поставляемых  
справочников

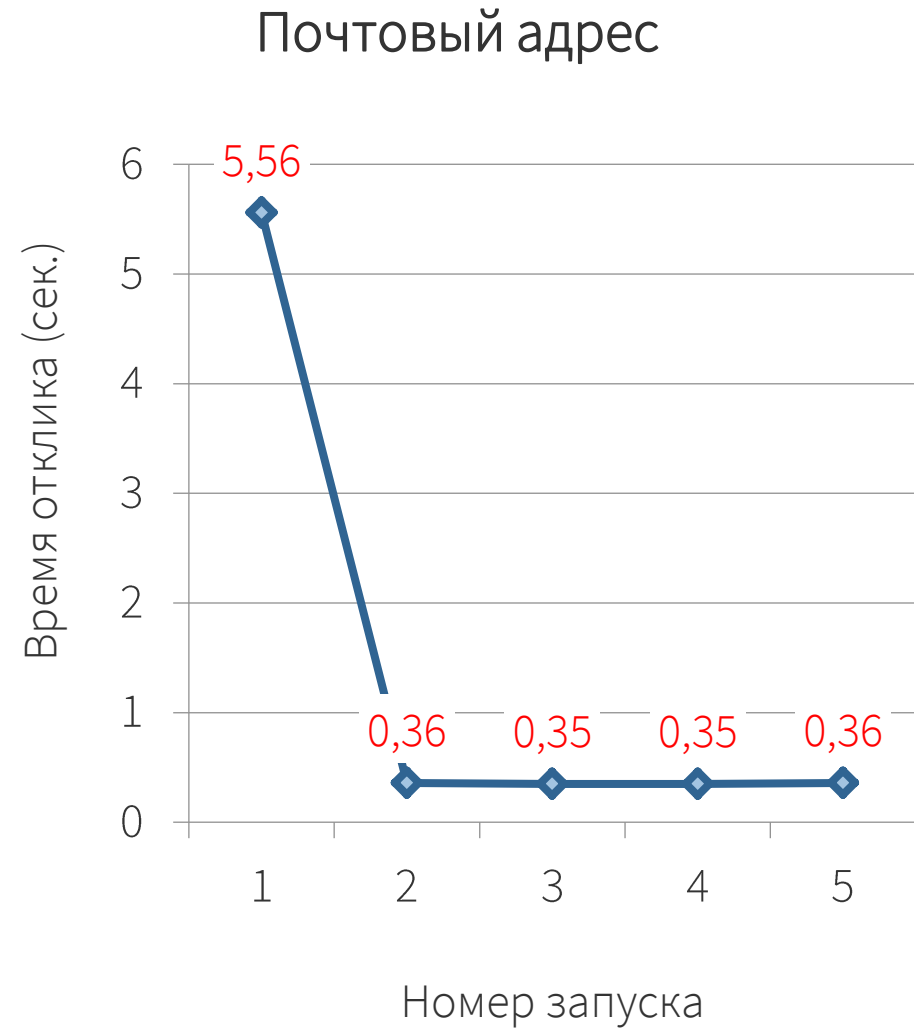
# Пример: Очистка почтовых адресов



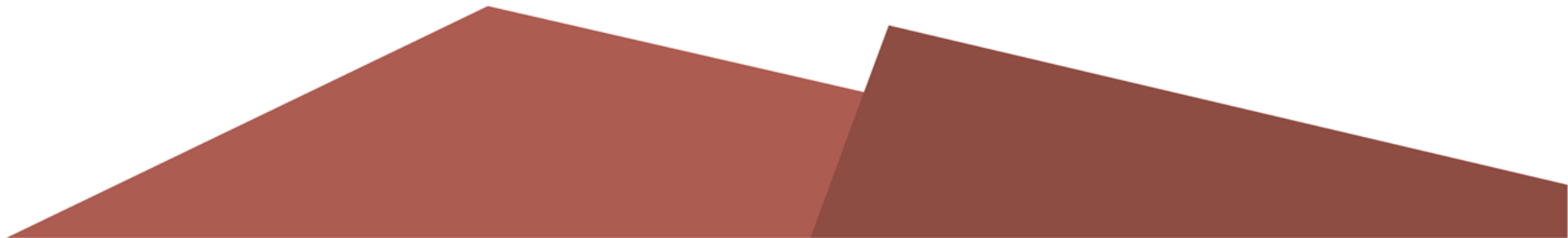
Линейный рост  
скорости при  
увеличении  
количества ядер


# Веб-сервис:

1. Работает медленнее из-за накладных расходов
2. Сильнее зависимость от сетевой инфраструктуры
3. Нужно время на первое обращение – «разогрев»



# Дедупликация и создание «ЗОЛОТОЙ» записи





Конечная цель  
очистки клиентских  
данных – создание  
«золотой записи»:  
единого, точного,  
актуального и  
полного варианта  
записи клиента.

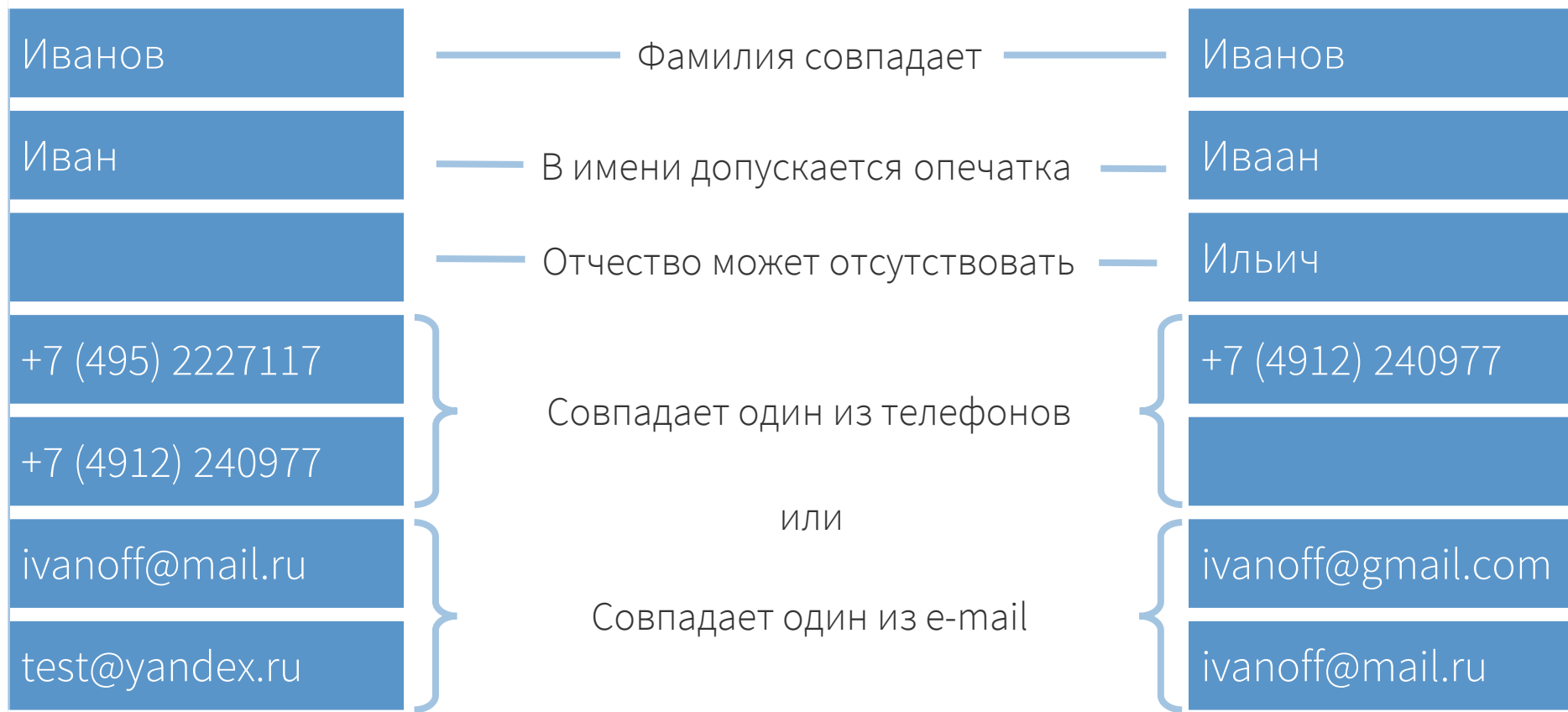
Для этого требуется:

1. Найти дубли и  
похожие записи
2. Объединить  
данные из всех  
источников

# Возможности дедупликации:

1. Формирование групп дублей:
  - Задание сочетаний атрибутов, определяющих дубль
  - Нечеткий поиск с заданной точностью
  - Учет наличия пропусков
2. Создание эталонных записей
3. Формирование «серой зоны» групп дублей для ручной обработки

# Пример стратегии дедупликации





# Преднастроенные стратегии дедубликации

Физические лица:

- 5 автоматических (жестких) стратегий
- 3 ручных (мягких)

Юридические лица:

- 5 автоматических (жестких) стратегий
- 3 ручных (мягких)

StrategyList	GroupID	ID	Phone1	Email1	FirstName	MiddleName	LastName	DocSeries	DocNumber	Address
Стратегия 2A	a-22-41-09102019-114852	22	+7 (927) 9807728	ferrumlena@mail.ru	Екатерина	Владимировна	Котова	6525	392289	
		41		ferrumena@mail.ru	Екатерина	Витальевна	Котова	6525	392289	
	a-78-79-09102019-114852	78	+7 (912) 9708511	shap314@mail.ru	Евгений	Алексеевич	Суворов	7120	464665	305007, Курская область город Курск улица Сумская д. 37, ко...
		79	+7 (951) 3264245		Евгений	Алексеевич	Суворов	7120	464665	363246, республика Северная Осетия - Алания Алагирский р...
Стратегия 4A	a-33-34-09102019-114852	33	+7 (929) 7113339	asta.aniki@gmail.com	Анастасия	Алексеевна	Тимухина	6530	588514	140055, Московская область город Котельники микрорайон ...
		34	+7 (937) 3582277	asta.aniki@gmail.com	Анастасия	Александровна	Тимухина			

Пример стратегии дедупликация по жестким стратегиям для автоматической обработки

Атрибуты	Расстояние Дамерау-Левенштейна	Непротиворечивость
Страна	0	True
Имя	2/25	False
Фамилия	2/25	False
Инициалы	0	False
Тип документа	0	False
Серия документа	0	False
Номер документа	0	False

После дедупликации возможно формирование эталонной «золотой записи»\* с учетом:

- Доверия к источнику данных
- Качества очистки
- Актуальности данных
- Заданной метрики

\* В рамках проекта

## Группа дублей

Фамилия	Имя	Отчество	Дата рождения	Телефон	Телефон	E-mail	E-mail
Иванов	Иван			+7 (495) 2227117	+7 (4912) 240977	<a href="mailto:ivanoff@mail.ru">ivanoff@mail.ru</a>	<a href="mailto:test@yandex.ru">test@yandex.ru</a>
<del>Иванов</del>	<del>Иваан</del>	Ильич		<del>+7 (4912) 240977</del>		<a href="mailto:iii@yandex.ru">iii@yandex.ru</a>	<del>ivanoff@mail.ru</del>
<del>Иванов</del>			15.12.1971	+7 (4912) 240699		<del>ivanoff@mail.ru</del>	



Единичные атрибуты	Значение
Фамилия	Иванов
Имя	Иван
Отчество	Ильич
Дата рождения	15.12.1971

Множественные атрибуты	Значение
Телефон 1	+7 (495) 2227117
Телефон 2	+7 (4912) 240977
Телефон 3	+7 (4912) 240699
E-mail 1	<a href="mailto:ivanoff@mail.ru">ivanoff@mail.ru</a>
E-mail 2	<a href="mailto:test@yandex.ru">test@yandex.ru</a>
E-mail 3	<a href="mailto:iii@yandex.ru">iii@yandex.ru</a>

«Золотая» запись

# Поставка решения



Модуль решения	Компонент Loginom	Веб-сервис
Очистка ФИО	●	●
Очистка названия организации	●	●
Очистка почтового адреса	●	●
Очистка телефонов	●	●
Очистка электронной почты	●	●
Очистка удостоверений личности	●	●
Очистка реквизитов организаций	●	●
Очистка дат	●	●
Очистка банковских реквизитов	●	●
Дедупликация	●	
Формирование «золотой записи»	●	

---

Режим	Назначение
Пакетная обработка	<ol style="list-style-type: none"><li data-bbox="805 418 1913 496">1. Первичная очистка данных</li><li data-bbox="805 518 2175 596">2. Регулярная регламентная очистка</li><li data-bbox="805 618 2117 796">3. Эпизодическая очистка больших наборов данных</li></ol>
Веб-сервис	<ol style="list-style-type: none"><li data-bbox="805 875 1447 953">1. Online очистка</li><li data-bbox="805 975 2252 1053">2. Очистка небольших наборов данных</li><li data-bbox="805 1075 2142 1160">3. Интеграция с другими системами</li></ol>

---

# Техническая поддержка:

1. Ежемесячное обновление и расширение справочников
2. Регулярные улучшения логики очистки данных



# Выгоды решения

1. **Быстрый запуск** – минимальное время от получения данных до первого результата
2. **Гарантии качества** – исправление сотен ошибок, актуальные справочники
3. **Визуальный конструктор** – интеграция и настройка логики без кодирования

[loginom.ru](http://loginom.ru)

