

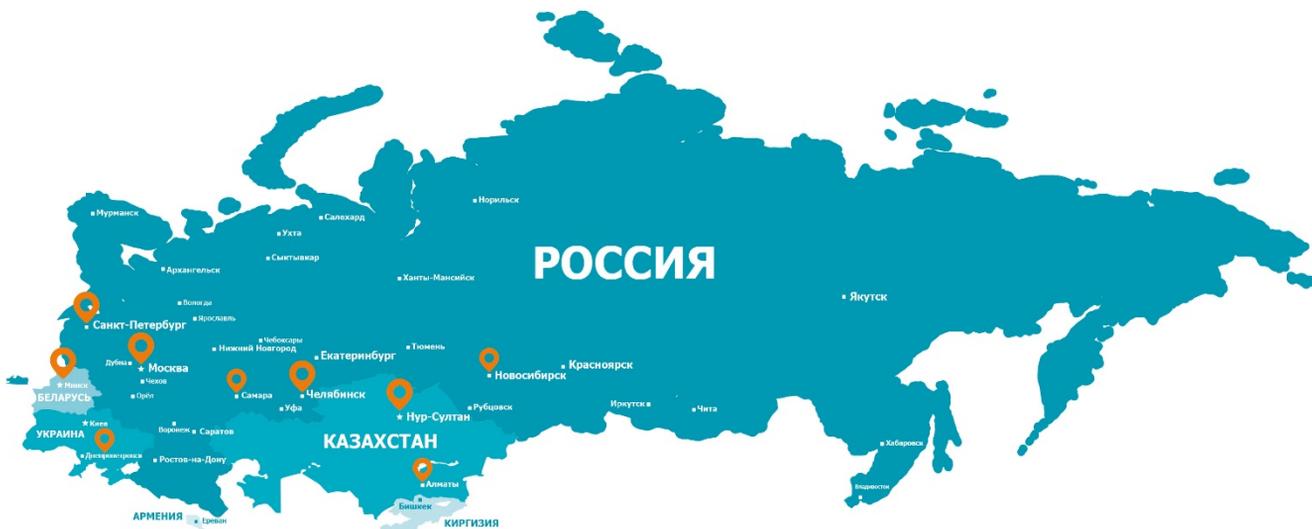
INVITRO

Данные тоже должны быть чистыми

Анастасия Рязанова,
ведущий аналитик-программист

Москва, 2019

Частная медицинская компания в России, специализирующаяся на лабораторной диагностике и оказании медицинских услуг. Представлена в 6 странах



ИНВИТРО – это

- 12,5+ млн пациентов в год
- 1300+ офисов
- 500+ городов присутствия
- 800+ серверов
задействованы в работе
- 9000+ сотрудников в
компании



Проблема

- Огромный объем поступающей информации
 - Общий размер клиентской базы 32 млн. записей (каждый день изменяются 64 тыс.)
- Разный уровень грамотности у администраторов медицинских офисов
- Некорректные персональные данные = недостоверные результаты
- Некорректные контактные данные = отсутствие информации о последних «вкусных» акциях и нововведениях

Проблема. Много точек входа при работе с клиентскими данными

Наполнение клиентской базы:

- администраторы в медицинских офисах
- операторы контакт-центра
- интеграционные автоматизированные сервисы клиентов
- сами пациенты (предзаказы в личном кабинете на сайте)

Обработка клиентской базы:

- операторы контакт-центра (при звонке на горячую линию)
- операторы баз данных (регулярно проводят анализ данных)

Проблема. Примеры «грязных» данных

Ошибки в ФИО

- **Мариа** (вместо Мария)
- **Светлана** (первый символ был записан латиницей)
- **Саша** (вместо Александр)

Некорректно указан пол пациента или не указан совсем

Неполный номер телефона

- 454545 (без указания кода страны и города)
- +7 (только код страны)

Некорректный формат email

- test.mail.ru (test@mail.ru)
- test@gmail.ru (test@gmail.com)

Проблема. Дубликаты

Источники дублей:

- Особенности корпоративных правил и защиты персональной информации
 - администраторы медицинских офисов видят только ограниченную зону клиентов по региону размещения, а не всю базу
- Ошибки в персональных данных существующего в базе пациента
- Невнимательность администраторов при поиске пациента в базе

Проблема. Дубликаты. Пример

Объединение контактов

Выбор главной записи							
Главная запись:	<input checked="" type="radio"/> 201076МЕЕШЦ 55ae10ba-406d-47f6-ba52-56e330b17d49	<input type="radio"/> 180181ЖЕОПС b441a7f1-f89b-11e7-cc39-54de8422b815	<input type="radio"/> 230391ЖЕ98Ъ 82f1407d-c1f4-11e6-bcf4-54a0507a49b2	<input type="radio"/> 030480ЖОЧР5 1e161d56-1dda-11e7-9a0c-54a0507a49b2	<input type="radio"/> 280181ЖЕМНШ 3fa12757-7a1c-11e8-8b9b-54a0507a1a0e	<input type="radio"/> 220588ЖЕЗЖЗ 22467bec-9477-11e6-48b0-54a0507a1a0e	<input type="radio"/> 100284ЖА0ЪХ a7177fec-dd93-11e6-2707-54a0507a1a0e
Основная информация	<input checked="" type="radio"/> Выбрать все поля	<input type="radio"/> Выбрать все поля					
НСС:	201076МЕЕШЦ	180181ЖЕОПС	230391ЖЕ98Ъ	030480ЖОЧР5	280181ЖЕМНШ	220588ЖЕЗЖЗ	100284ЖА0ЪХ
Фамилия:	<input checked="" type="radio"/> [redacted]	<input type="radio"/> [redacted]					
Имя:	<input checked="" type="radio"/> [redacted]	<input type="radio"/> [redacted]					
Дата рождения:	<input checked="" type="radio"/> 20.10.1976	<input type="radio"/> 18.01.1981	<input type="radio"/> 23.03.1991	<input type="radio"/> 03.04.1980	<input type="radio"/> 28.01.1981	<input type="radio"/> 22.05.1988	<input type="radio"/> 10.02.1984
Отчество:	<input checked="" type="radio"/> [redacted]	<input type="radio"/> [redacted]					
Пол:	<input checked="" type="radio"/> Женский	<input type="radio"/> Женский					
Общие сведения	<input checked="" type="radio"/> Выбрать все поля	<input type="radio"/> Выбрать все поля					
Город:	<input checked="" type="radio"/> не заполнено	<input type="radio"/> не заполнено					
Регион:	<input checked="" type="radio"/> Москва	<input type="radio"/> Москва					
Дата регистрации:	<input checked="" type="radio"/> 19.03.2013	<input type="radio"/> 22.01.2018	<input type="radio"/> 14.12.2016	<input type="radio"/> 10.04.2017	<input type="radio"/> 27.06.2018	<input type="radio"/> 17.10.2016	<input type="radio"/> 18.01.2017
Электронная почта:	<input checked="" type="radio"/> [redacted]@mail.ru	<input type="radio"/> [redacted]@mail.ru					
Мобильный телефон:	<input checked="" type="radio"/> +7 [redacted] 7835426	<input type="radio"/> +7 [redacted] 7835426					

Группа из 9 дублей

Проблема. Дубликаты. Пример

Общие сведения	<input checked="" type="radio"/> Выбрать все поля	<input type="radio"/> Выбрать все поля	<input type="radio"/> Выбрать все поля
Город:	<input checked="" type="radio"/> не заполнено	<input type="radio"/> не заполнено	<input type="radio"/> не заполнено
Регион:	<input checked="" type="radio"/> Москва	<input type="radio"/> Москва	<input type="radio"/> Москва
Дата регистрации:	<input checked="" type="radio"/> 19.03.2013	<input type="radio"/> 22.01.2018	<input type="radio"/> 14.12.2016
Электронная почта:	<input checked="" type="radio"/> █████290@mail.ru	<input type="radio"/> █████290@mail.ru	<input type="radio"/> █████290@mail.ru
Мобильный телефон:	<input checked="" type="radio"/> +7 █████ 7835426	<input type="radio"/> 7 █████ 7835426	<input type="radio"/> 7 █████ 7835426

В группе полностью идентичны фамилия, пол, номер мобильного телефона, адрес эл. почты.

Незначительные различия в именах и дате рождения.

Задачи

- Наладить динамический процесс проверки информации, допускающий внесение изменений в функционал сервиса
- Снизить влияние человеческого фактора при вводе и обработке контактных данных
- Автоматизировать обработку, очистку и дедупликацию
- Поддерживать гарантированное качество персональных данных клиентов
- Учитывать специфику бизнес-процессов компании

Решение

В 2013 году была создана система по очистке и дедупликации данных на базе аналитической платформы **Deductor**. В 2019 решение было перезапущено на базе новой версии – **Loginom**.

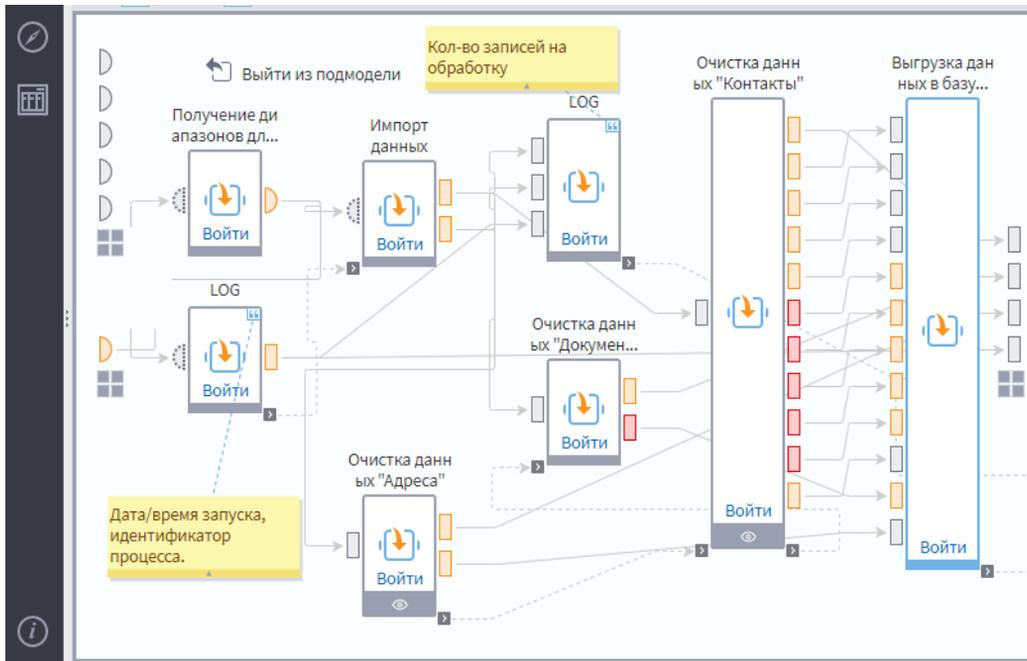
Предпосылки перехода на новую версию

- Заканчивалось технологическое развитие старой системы, а бизнес-процессы требовали изменений
- Система уже не справлялась с объемом измененных данных и просто «захлебывалась» в их количестве

О системе очистки и дедупликации

Loginom

Loginom – очень гибкая система, которая позволила нам максимально «подстроить» ее под наши технологические и бизнес-процессы.



Состав системы

В основе системы использовано решение **Loginom Data Quality**

- **Очистка**
 - Исправление ошибок, проверка по справочникам, стандартизация, обогащение
- **Дедупликация**
 - Проверка дублей по стратегиям
 - Объединение дублей в группы, автоматическое слияние или передача операторам для уточнения информации и принятия решения о слиянии

Входные данные

Система работает с несколькими справочниками:

- Пациенты (ФИО, пол, дата рождения, контактные данные)
- Адреса пациентов
- Документы пациентов

Очистка и обогащение на примере ФИО

- Удаление лишних символов в ФИО (, ; “; –)
- Проверка основного языка написание ФИО, замена символов другого языка аналогами основного (латинская **C** вместо кириллической)
- Сверка имени, отчества с постоянно обновляемыми справочниками имен, отчеств, фамилий, встроенных в Loginom
 - Исправление орфографических ошибок
 - Восстановление до полного имени
- Восстановление пола по ФИО в случае его отсутствия

Дедупликация

- Только уже очищенные данные
- Поиск дублей по полной базе пациентов по заложенным 6 четким и 3 нечетким стратегиям
- Объединение в группы полных и потенциальных дублей
- Передача групп потенциальных дублей операторам
- Выделение «золотой» мастер-записи в группах полных дублей
- Слияние дублей в группе на мастер-запись
- Публикация изменения в шину данных

Дедупликация. Пример стратегии

СТРАНА (непротиворечивость=true, точное совпадение) И

ИМЕНА (непротиворечивость=false, нечеткость 2/20) И

ФАМИЛИИ (непротиворечивость=false, нечеткость 2/20)

И

ИНИЦИАЛЫ ОТЧЕСТВА (непротиворечивость=false, точное совпадение) И

ТИП ДОКУМЕНТА (непротиворечивость=false, точное совпадение) И

СЕРИЯ ДОКУМЕНТА (непротиворечивость=false, точное совпадение) И

НОМЕР ДОКУМЕНТА (непротиворечивость=false, точное совпадение)

Дедупликация. Пример стратегии

Нечеткие совпадения базируются на расчете расстояния редактирования Дамерау-Левенштейна (PP)

Нечеткость **2/20** означает, что допускается PP не больше **2 символов** и процент PP от общей длины не больше **20%**

НАТАЛ**Ь**Я = НАТАЛ**И**Я (разница в 1 символ и 14% от общей длины) Нечеткость 1/14 – будут признаны идентичны.

НАТАЛ**Ь**Я = НАТАЛ**И** (разница в 2 символ и 28% от общей длины) Нечеткость 2/28 – **не будут** признаны идентичны

Процесс обработки данных. Шаг 1

Петров
 Ваня
 Г
 01.01.1986
 Муж
 89111111111
 1 визит: 09.10.2017

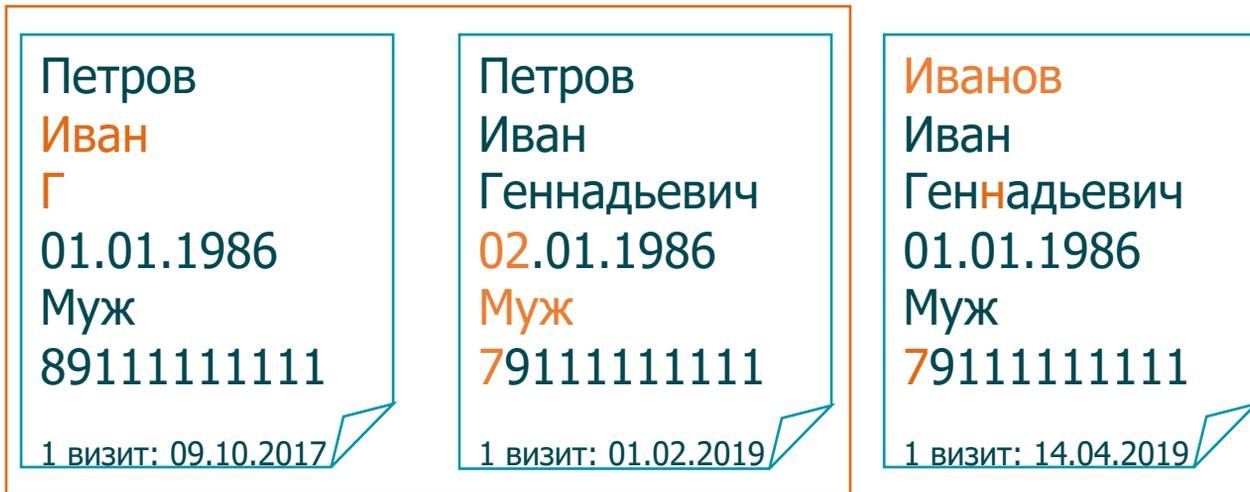
Петров
 Иван
 Геннадьевич
 02.01.1986
 Жен
 79111111111
 1 визит: 01.02.2019

Иванов
 Иван
 Генадьевич
 01.01.1986
 Муж
 79111111111
 1 визит: 14.04.2019



Очистка данных

Процесс обработки данных. Шаг 2

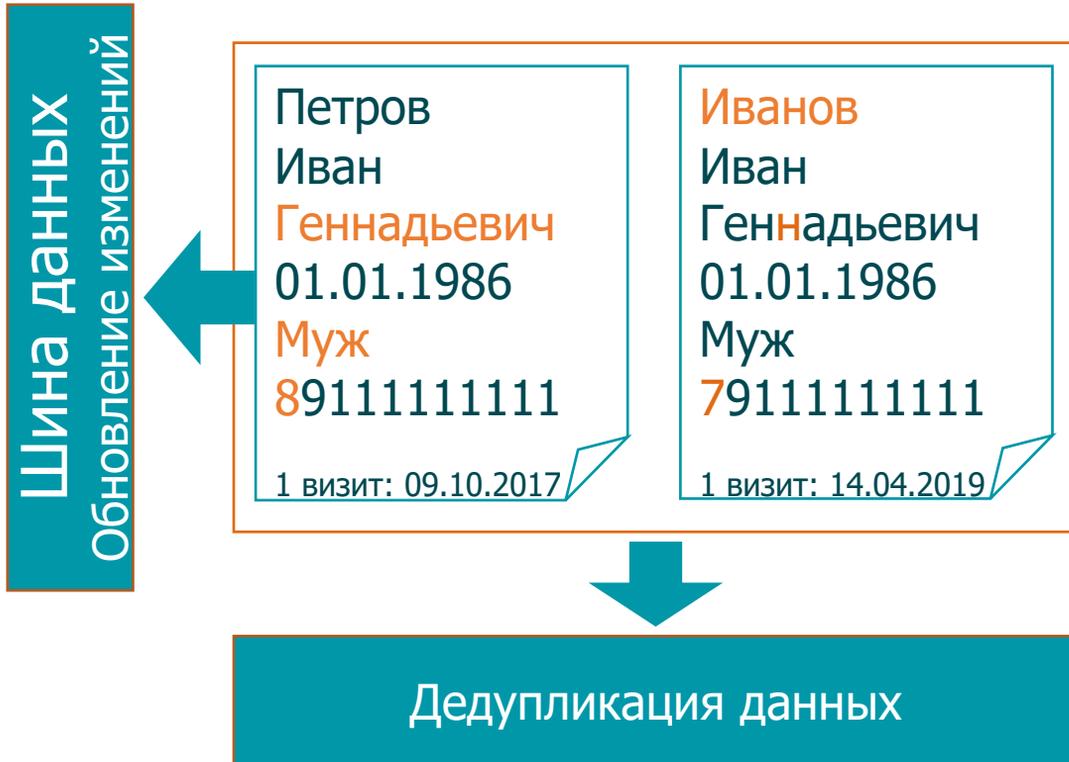


Полные дубли

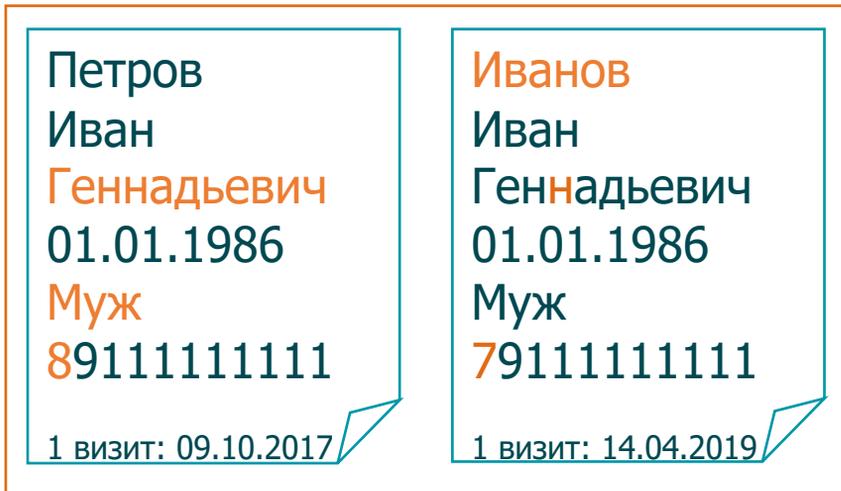


Дедупликация данных

Процесс обработки данных. Шаг 3



Процесс обработки данных. Шаг 4



Потенциальные дубли

Оператор
Ручная правка

Взаимодействие систем

Два направления взаимодействия:

Частичная репликация клиентской базы и базы Loginot

- Для обработки выбираются все измененные за **5 минут** записи
- Записывает потенциальные дубли в отдельную таблицу, доступную только оператору

Шина данных

- После очистки и дедупликации Loginot публикует изменения в шину данных
- Остальные сервисы, работающие с клиентской базой, «разбирают» изменения из шины

Статистика по работе

- Ежедневно
 - Очищаются около 50 тыс. записей
 - Обогащаются около 15 тыс. записей
 - Сливаются около 4 тыс. дублей
 - Ручная обработка оператором 50-100 записей
- За все время работы сервиса
 - Очищены более 5,5 млн записей
 - Проверены на дубли около 3 млн записей

Состав команды ИНВИТРО

В процессе внедрения Logiном участвовали:

- 1 сотрудник группы поддержки лабораторных систем
- 1 сотрудник отдела разработки
- 1 архитектор баз данных

Перспективы в развитии системы

- Подключить к системе другие справочники
- Использовать более широкий спектр возможностей LogiNot в глубокой аналитике данных

Результаты

- **Результативность внедренной системы подтвердила планы**
 - Снижены издержки на поддержку качества клиентских данных
 - Существенно уменьшен процент попадания некорректных и неполных данных в клиентскую базу
 - Успешно запущен процесс «избавления» от дублей

INVITRO

**Спасибо
за внимание!**

www.invitro.ru

8 (800) 200-363-0