

Большие данные и машинное обучение в HeadHunter: умный поиск, рекомендации, прогнозирование

Александр Сидоров

The logo consists of the lowercase letters 'hh' in white, centered within a solid red circle.

hh

У каждого свои большие данные

- Дисковые полки в 2002
- Медицинские экспертные системы
- Антивирусы
- Веб-антивирусы
- Веб-поиск и признаки по поведению пользователей
- Classified

Big data - не только про объём данных

Достаточный объём
и качество
данных

Подходящие
инструменты

Задачи, выполнимые
и экономически
целесообразные

Команда, способная
решить эти задачи

Автоматизированная модерация резюме

Задача

- Проверять качество заполнения резюме на самом верху воронки подбора

Проблема

- Значительная часть соискателей плохо и неполно заполняют резюме
- Модерировать вручную – дорого

Решение

- Продолжать модерировать вручную
- Машинное обучение

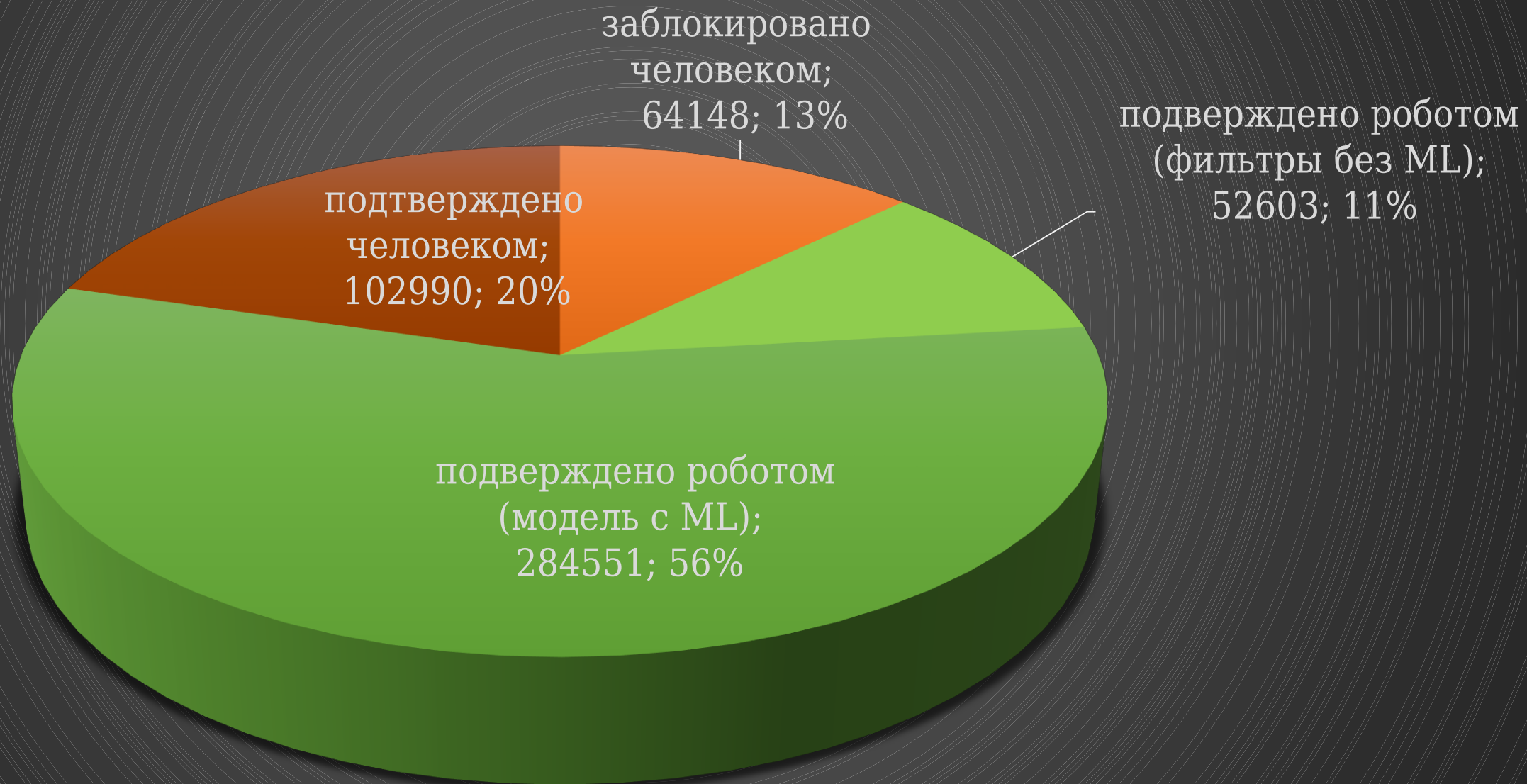
Схема работы автомодерации

2016: 20к новых резюме в день, 20 модераторов

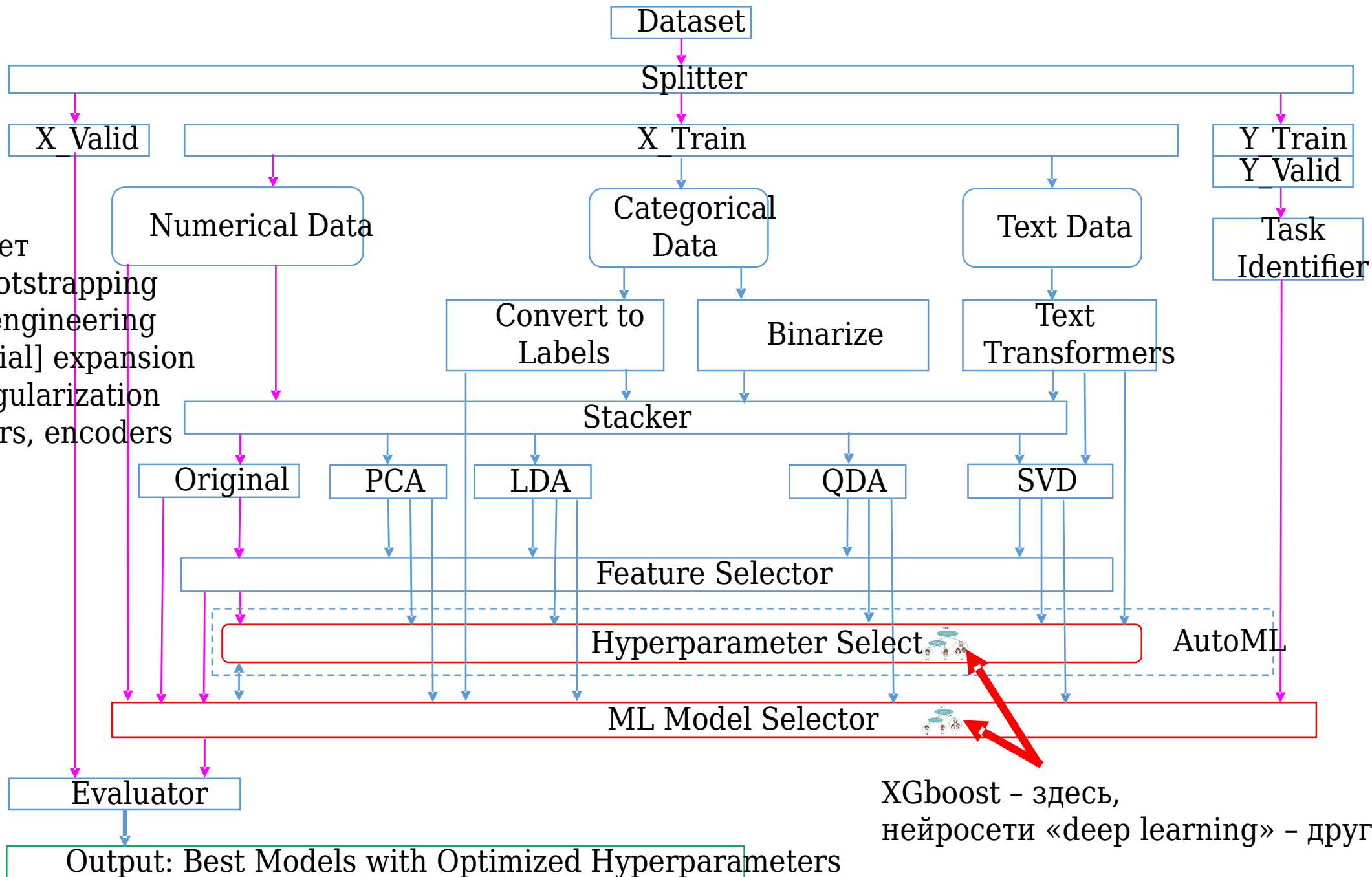
1. Результаты ручной модерации - исторические данные, 200к рез
2. Научили модель автомодерации, 1к признаков, auc_roc 0.94
3. Новые резюме проверяются автомодератором
4. Качественные резюме автоматически подтверждаются
5. Спорные случаи отправляются людям

2017: >30к новых резюме в день, ↑, 12 модераторов и качество мо
const

Экономия времени модераторов-людей



Не хватает
Kfold, bootstrapping
Feature engineering
[Polynomial] expansion
[l1/l2] regularization
Vectorizers, encoders
...



XGboost - здесь,
нейросети «deep learning» - другой вариан

Кейс 2: Прогнозирование времени работы KPI, риска ухода для конкретной компании

Задача

- Снизить текучесть персонала, сэкономить ФОТ
- Увеличить производительность и выручку

Решение

- Больше работы HR'ов и руководителей
- Машинное обучение

Функционал

- Дописываем к резюме кандидатов прогнозы по времени работы и KPI

Схема прогнозирования времени работы и I



26 000 000 резюме

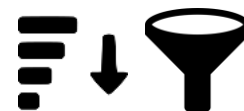
Отбор по проф. области
и дополнительным критериям



3-5 признаков

Модель, разделяющая по
времени и KPI

Ранжирование по
вероятности,
времени, KPI



≈ 1000 признаков

10-200k резюме в обучающих выборках

Результаты предпилотных проектов

Компания	Банк 1	Банк 2	Банк 3	FMCG 1
Процент новых сотрудников, которые...				
...работают более 2 года сейчас	36,5% ↓	34,0% ↓	40,3% ↓	30,5% ↓
...будут работать более 2 лет с нашим продуктом	48,6%	68,2%	75,9%	70,6%
Уменьшение текучести	1.33 раза	2,05 раза	1,88 раз	2,31 раза
Экономия ФОТ на текучести	2,5%	5,1%	4,68%	5,7%

- Сейчас технология рекрутмента у клиентов такая, что они нанимают кандидатов, из которых только **1/3** проработают больше 2 лет
- Мы отбираем и предоставляем резюме кандидатов, из которых до **2/3** проработают больше 2 лет

Кейс 3: Поиск, рекомендации вакансий, реклама вакансий

Задача

- Подобрать по резюме вакансии, интересные соискателю


Решение

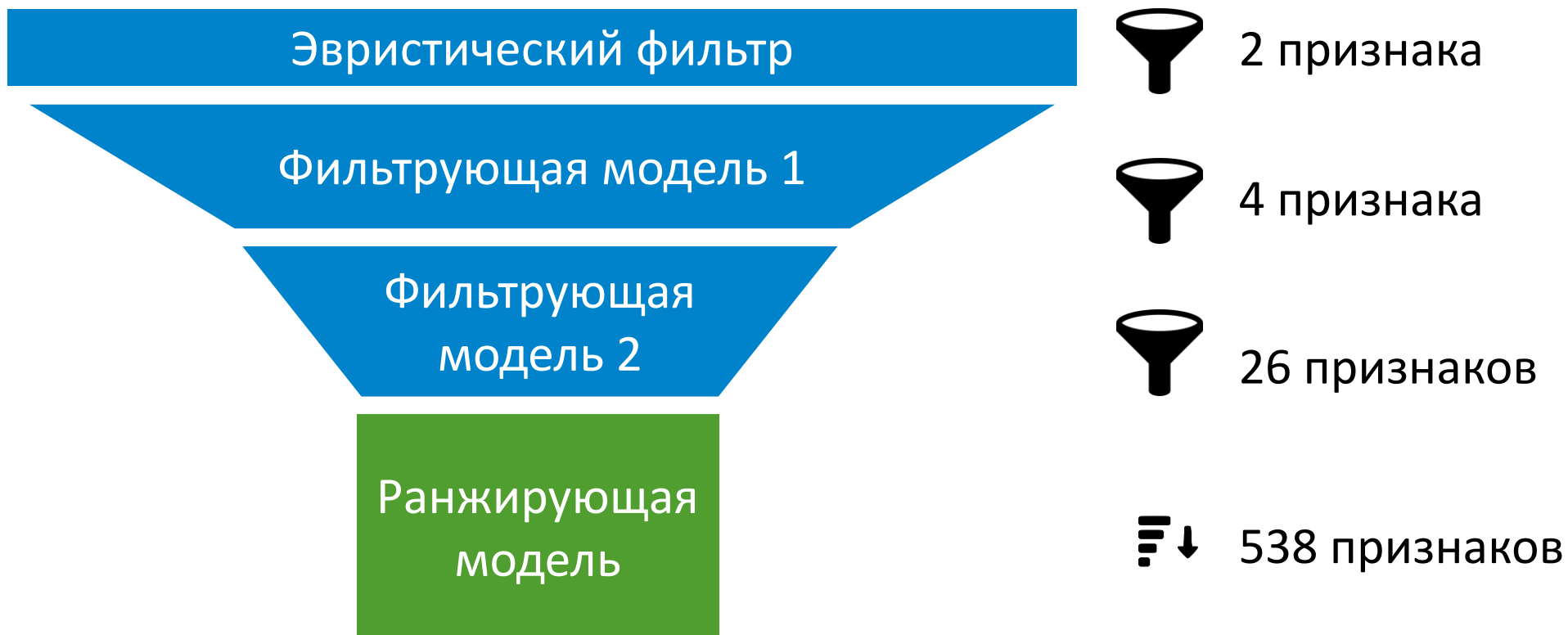
- Фильтр по параметрам вакансии, ручной таргетинг
- Машинное обучение

Функционал

- Главная страница, рекомендуемые вакансии в списке резюме, рассылки с подходящими вакансиями

Схема работы рекомендательной системы

 490 000 вакансий



44-500m пар в обучающих множествах, переподбор за ночь

Схема работы поиска

 490 000 вакансий

Индекс Lucene

Разделяющая модель +
простое ранжирование

Сложное
ранжирование



17 признаков



≈200 признаков

44-500m пар в обучающих множествах, переподбор за ночь

С ML и Big Data связана только часть проектов

Постановка задачи,
выбор приоритетных
процессы, команда

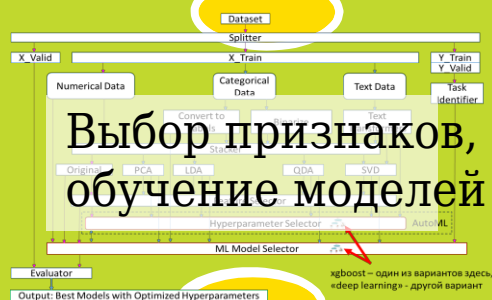
Сбор
исторических
данных

Мониторинг
ресурсоёмкости,
контролируемое
поведение при
превышении нагрузки

Оптимизация
использования
ресурсов

Прототипирование

Выбор признаков,
обучение моделей



Автоматизированное
тестирование

Работа с
экспериментальным
кодом

Регулярный пересчёт
статических признаков,
повторное обучение
моделей

Расчёт признаков
и работа моделей
во время обработки
поисковых запросов

A/B-тесты,
соблюдение баланса
«польза/ресурсы/
трудозатраты
на оптимизацию»

Синхронное, без остановок
обновление
на серверах
в production





Frontend (обычный и мобильный сайт, также API, через который работают мобильные приложения)
20 instances

Интеграционный слой: 15 instances

Backend
остальных
сервисов:
30 instances

Backend поиска и рекомендаций по вакансиям: 30+ instances

Meta (балансировщик, индексатор) 3 instances

Master (обновляет сегменты индексов) 2 instances

Baserearch (Lucene, расчёт pairwise features, модели) 20+ instances

Kardinal
(извлечение
статических
признаков)
7 instances

Data storage: ≈50 instances

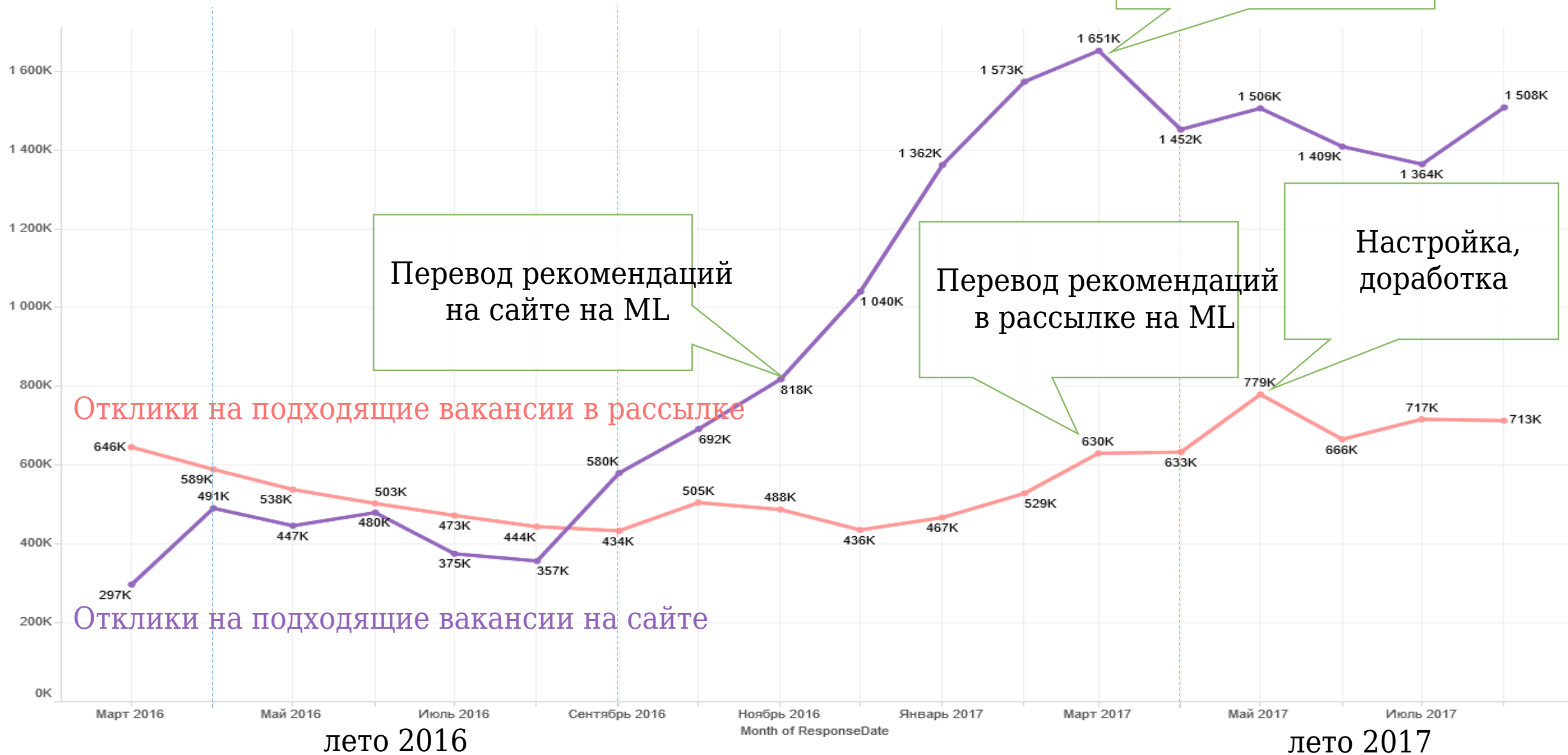
PostgreSQL 7 servers

Hadoop cluster 30+ servers

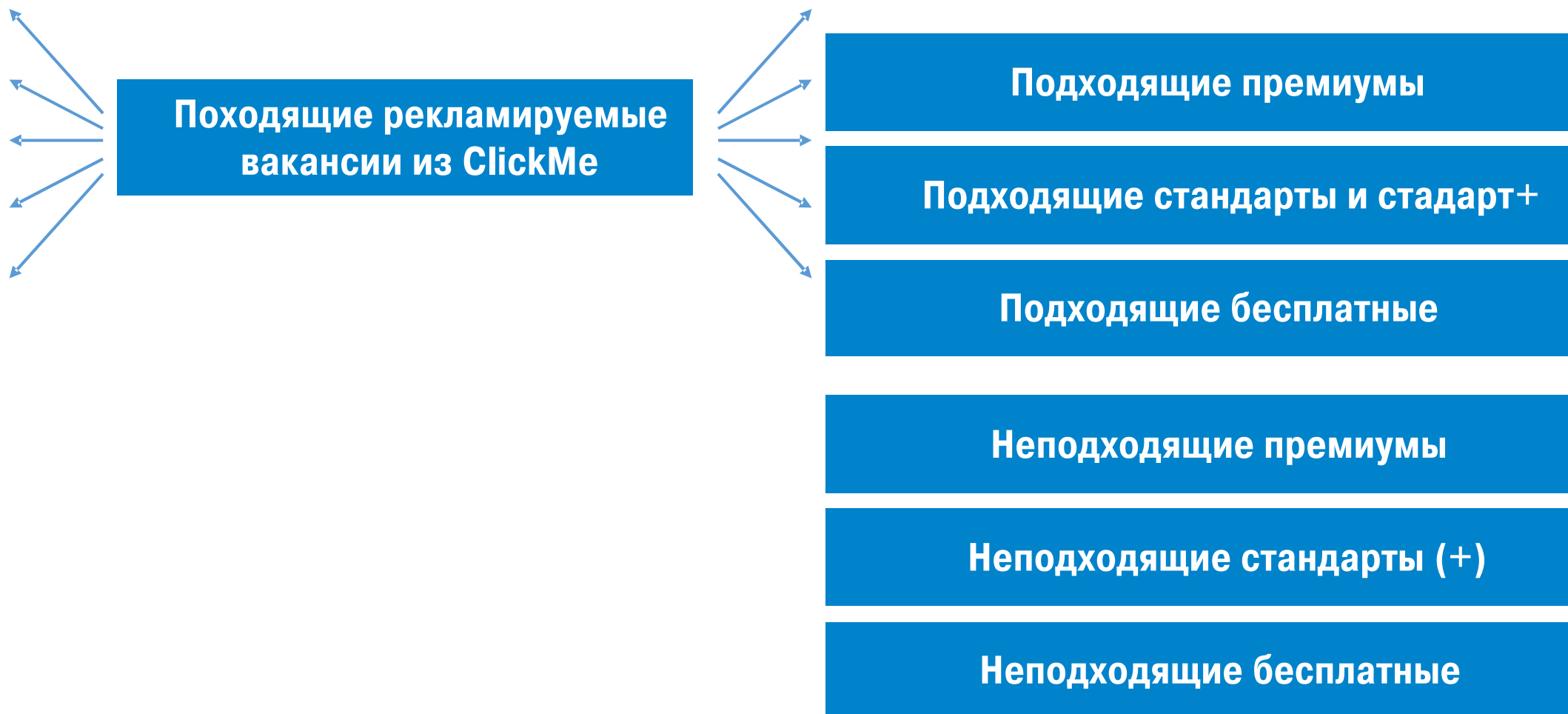
Cassandra 7 servers

Instance это docker container, vm или сервер

Отклики на вакансии



Умный поиск 1.0 + умная реклама вакансий в С



Прогнозируемая вероятность отклика выше заданного порога

Воронка найма при обычном поиске, **было**



Воронка найма при умном поиске, становитс



Воронка при умном поиске и **clickme**, скоро



Big data - не только про объём данных

Достаточный объём
и качество
данных

Подходящие
инструменты

Задачи, выполнимые
и экономически
целесообразные

Команда, способная
решить эти задачи

Большое спасибо!

Вопросы можно сейчас или потом, по почте