



DevCon School

Технологии будущего

Большие данные на платформе Microsoft Azure

Дмитрий Сошников
Technical Evangelist, Microsoft

Что такое Big Data

Откуда берутся большие
данные и что с этим делать

Обзор технологий

Что есть в Azure для работы
с большими данными

Орг.вопросы

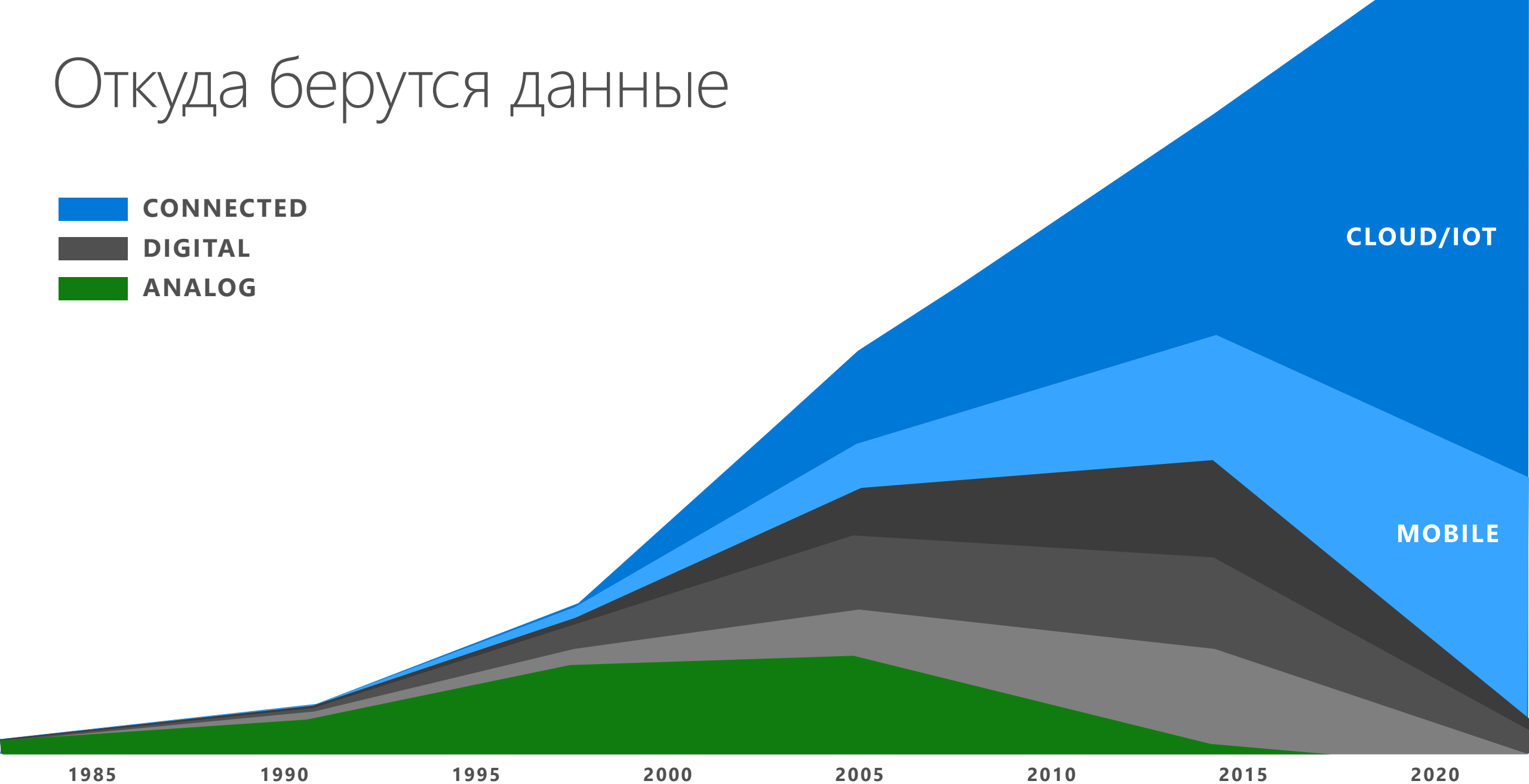
Что мы будем делать на
школе

Что такое большие данные

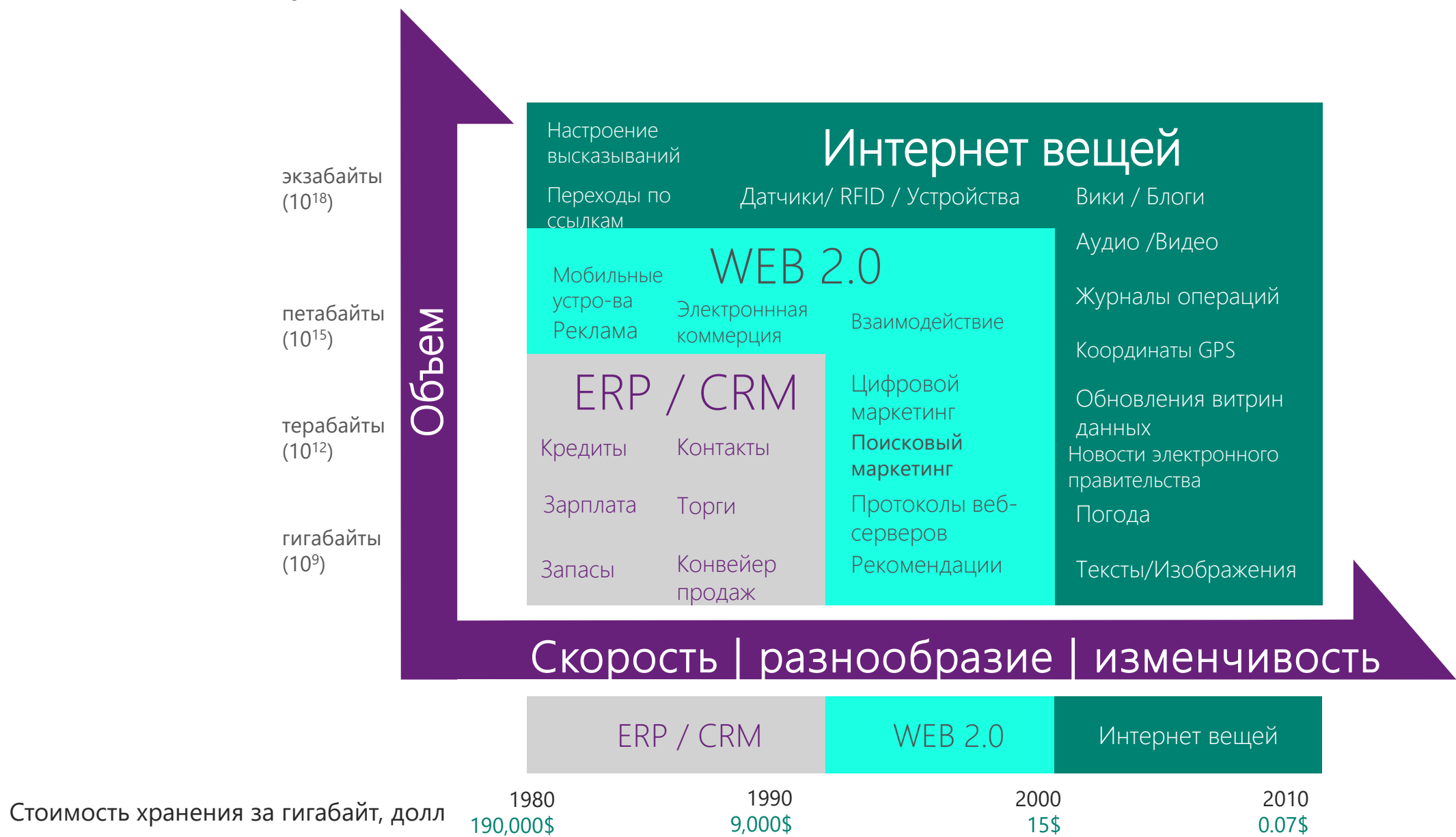
#msdevcon

Откуда берутся данные

CONNECTED
DIGITAL
ANALOG



Классификация больших данных



Зачем собирать большие данные



Получение полной картины бизнеса

- Удалённый мониторинг
- Предсказание показателей
- Управление рисками



Распределение продуктов

- Inventory management
- Supply chain optimization



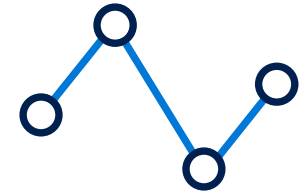
Персонализация для конечных пользователей

- Personalized offers
- Product recommendations



Проактивное устранение проблем

- Predictive maintenance
- Operational efficiency
- Customer service improvement



Новые возможности бизнеса

- Cross-sell and upsell
- Product-as-a-service
- New data-driven services

Современный подход к большим данным

	Традиционный	Big Data
 Структура данных	Реляционные с предопределенной схемой	Слабо-структурированные с гибкой схемой
 Стоимость	Специализированное оборудование	Общего назначения
 Культура	Операционные отчеты Смотрим в прошлое	Пробы и эксперименты Машинное обучение, a/b тестирование и т.д.

Опыт Microsoft

Microsoft столкнулся с большими данными при создании своих продуктов.

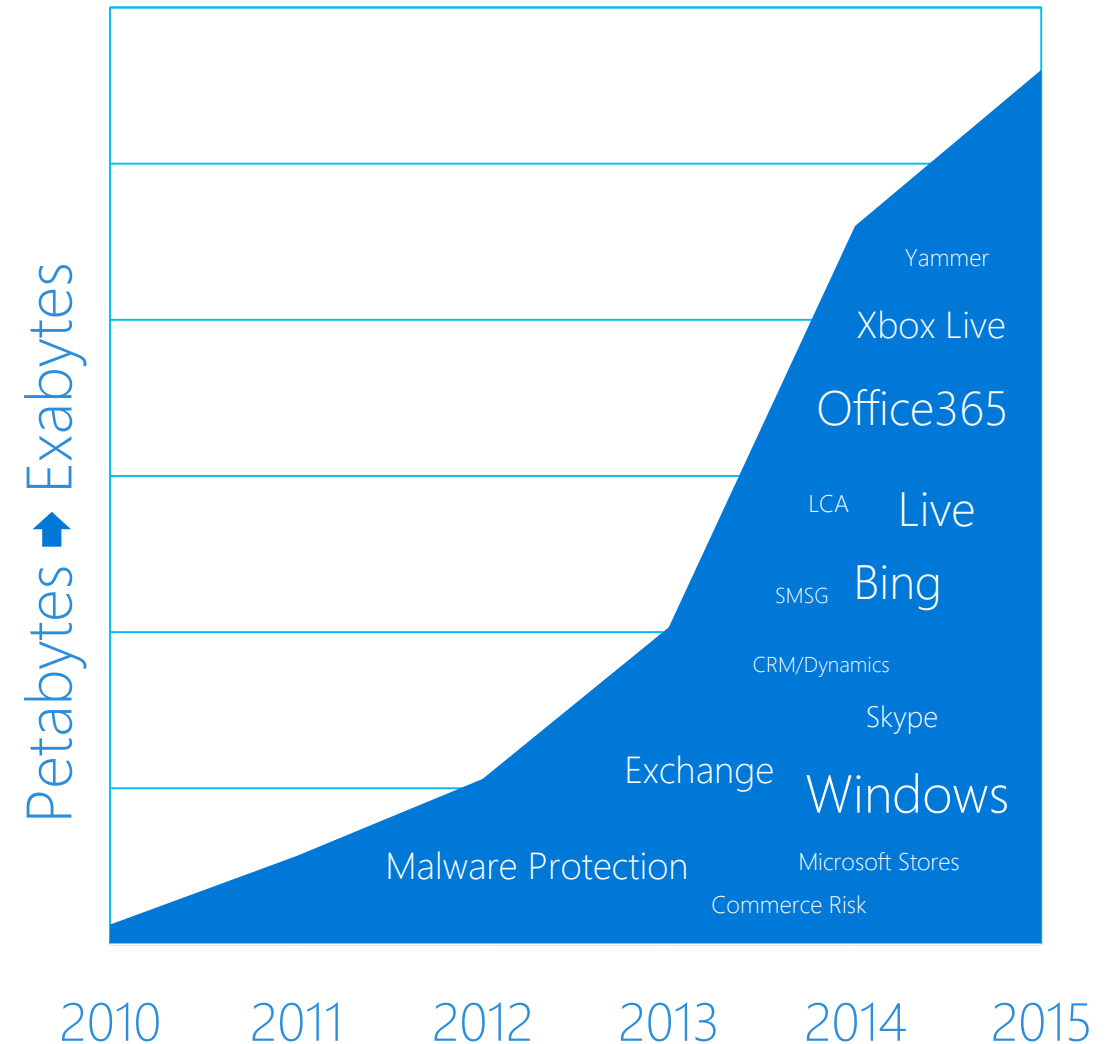
Что мы сделали:

- Общедоступные сервисы для больших данных
- Инструменты, доступные разработчикам
- Доступные инструменты машинного обучения

В результате:

- Разработчикам доступны те же инструменты, которые используются внутри Microsoft

Growth of data @ Microsoft



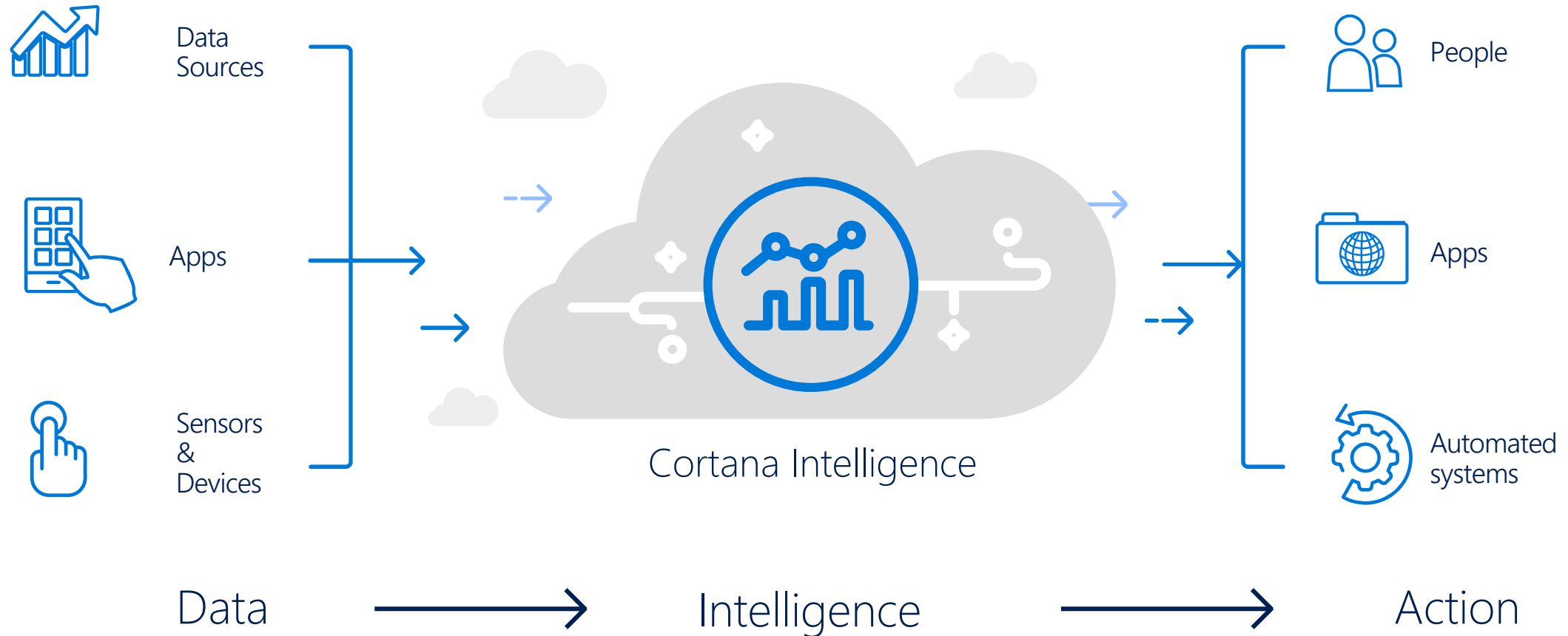
* Microsoft. FY16 Q4 Results, URL: <http://www.microsoft.com/en-us/Investor/earnings/FY-2016-Q4/press-release-webcast>

Cortana Intelligence Suite

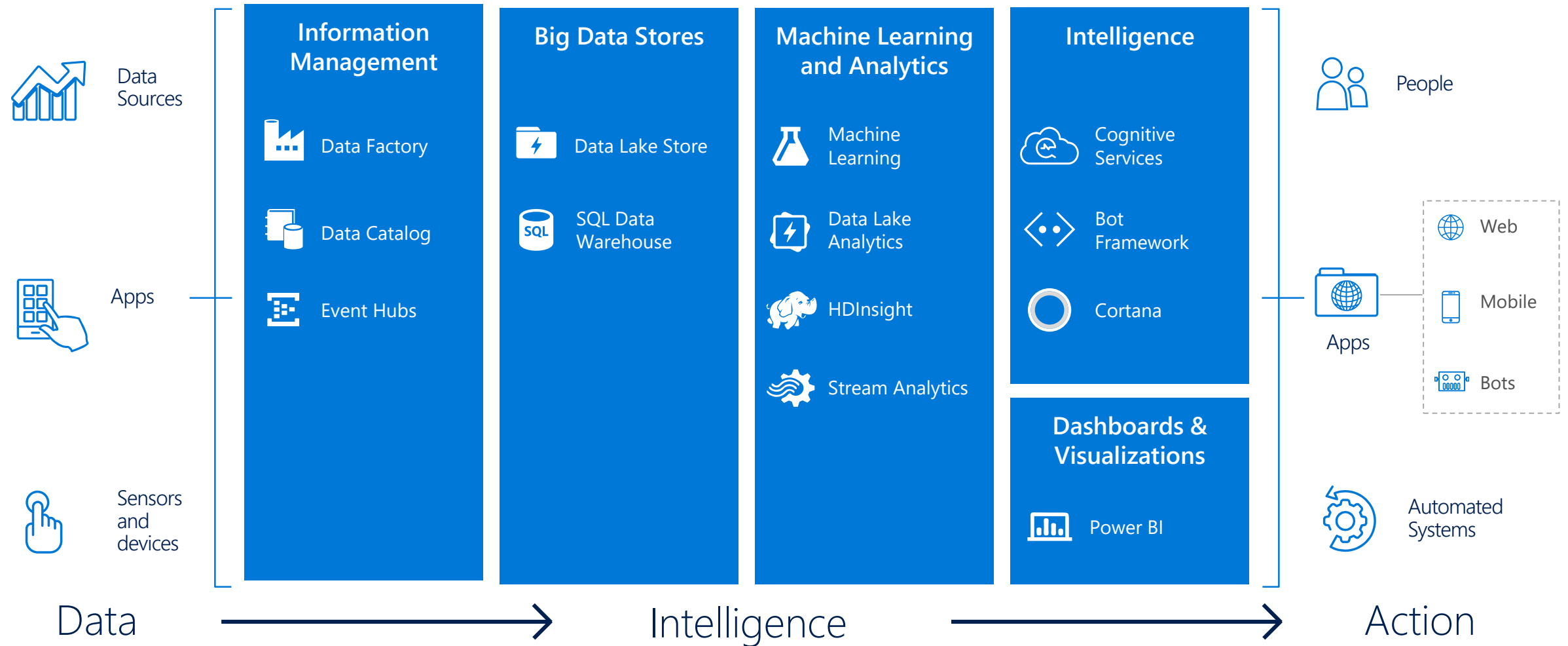


Cortana Intelligence Suite

Превращает данные в умные действия

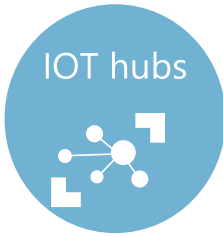


Cortana Intelligence Suite



Получение данных

Потоковые данные



Stream Analytics

Azure Function

Перенос данных



Data Factory

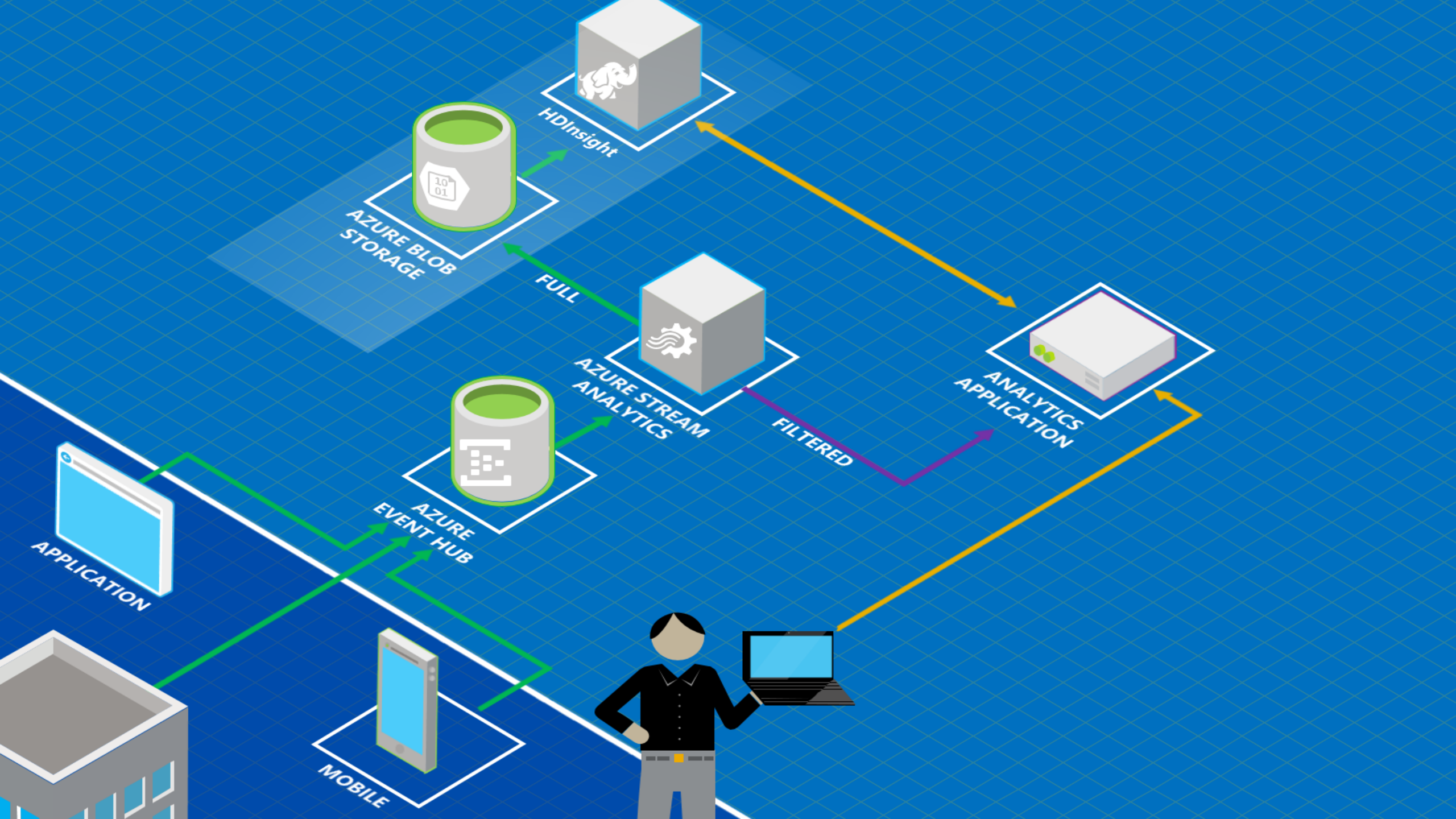
Хранилище

 Демонстрация

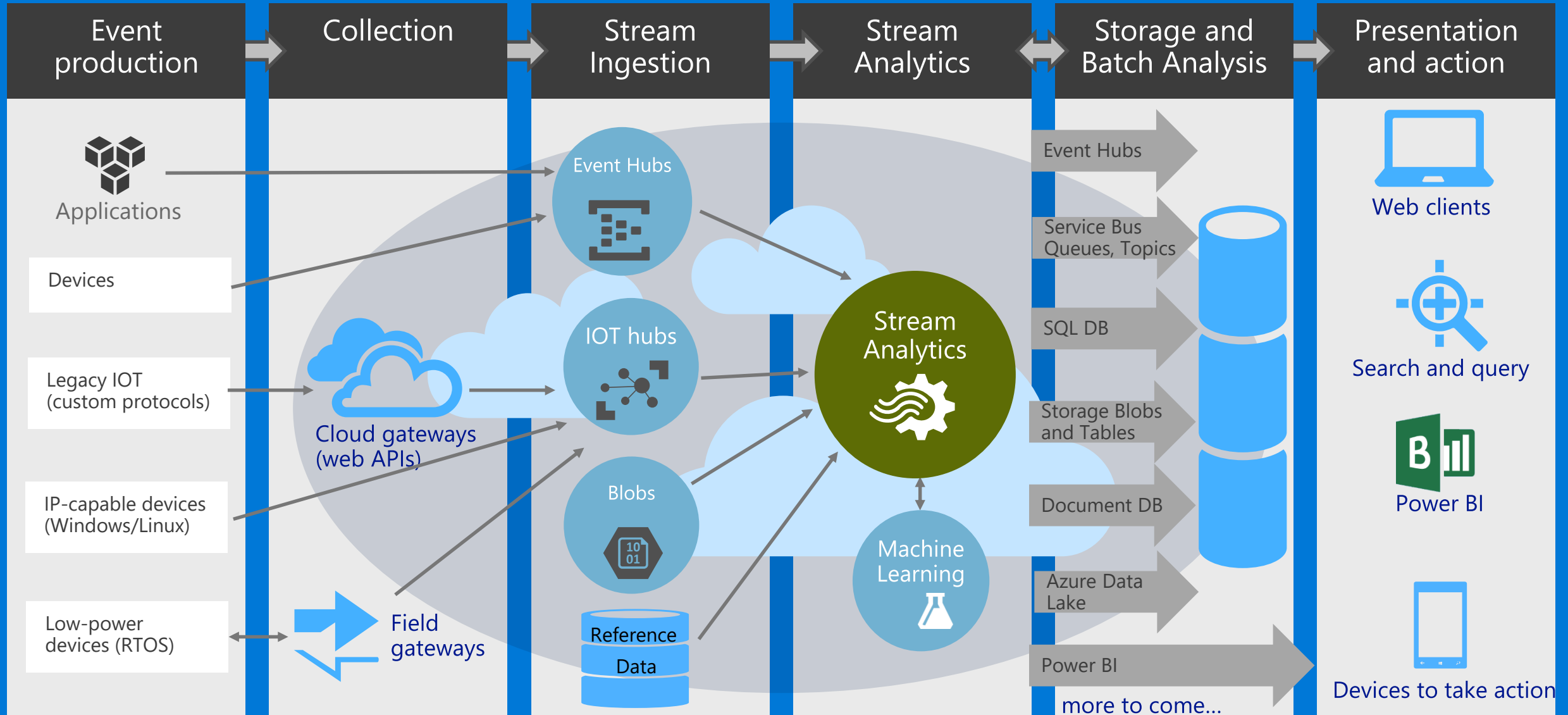
Собираем IoT-данные

<http://github.com/evangelism/FaceRecognitionTracker>

#msdevcon



Типичная схема IoT-решения



Хранение данных

Масштабная обработка



Azure
BLOB



Azure
Table



CosmosDB
(DocumentDB)



DataLake Store
(HDFS)



Data Warehouse



Data
Catalog



Структурированность



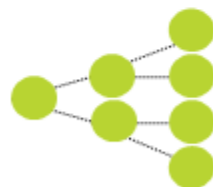
Key-Value



Column-family



Documents



Graph



Global distribution

Elastic scale out

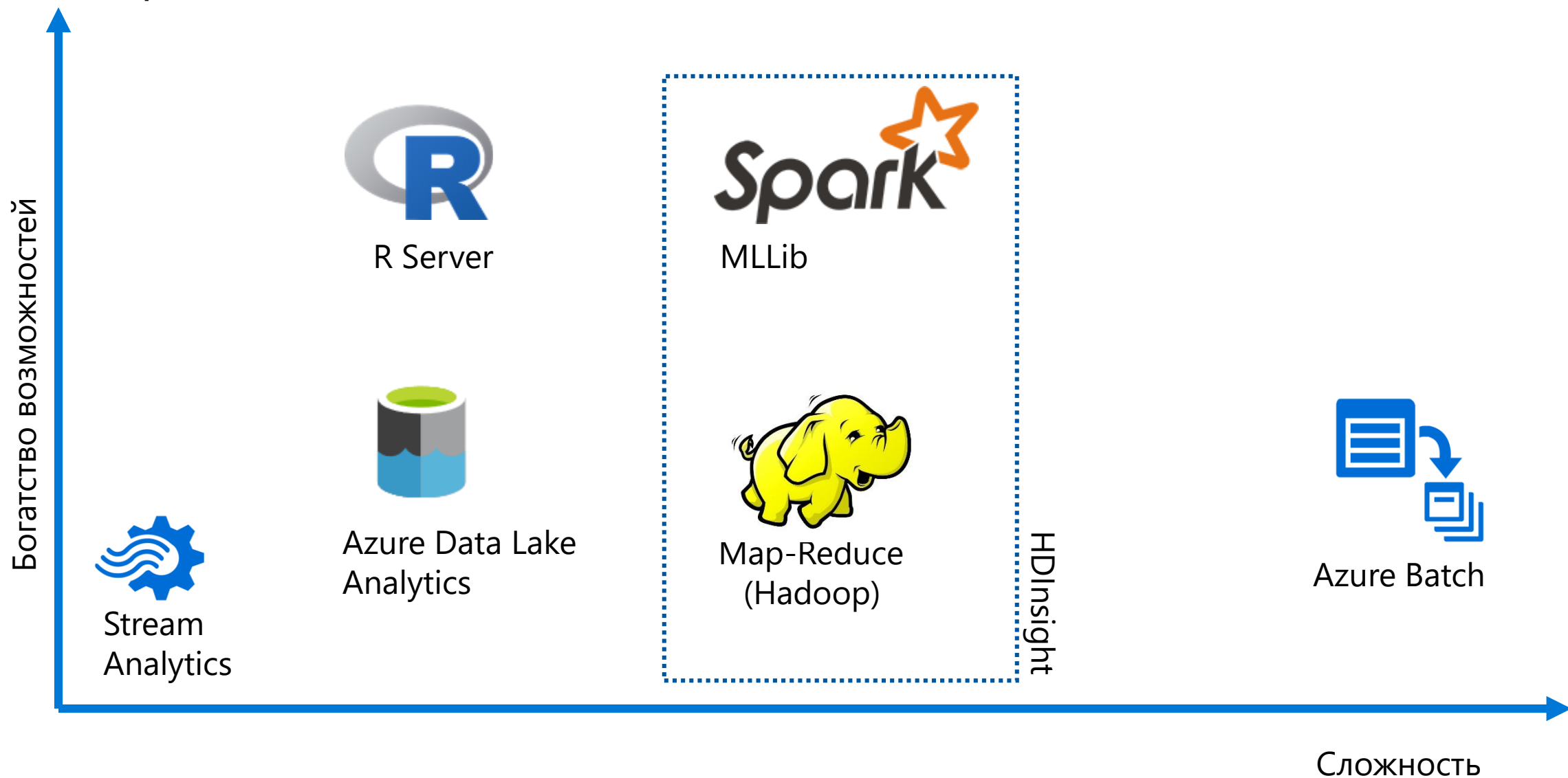
Guaranteed low latency

Five consistency models

Comprehensive SLAs

A globally-distributed, multi-model database service

Обработка данных

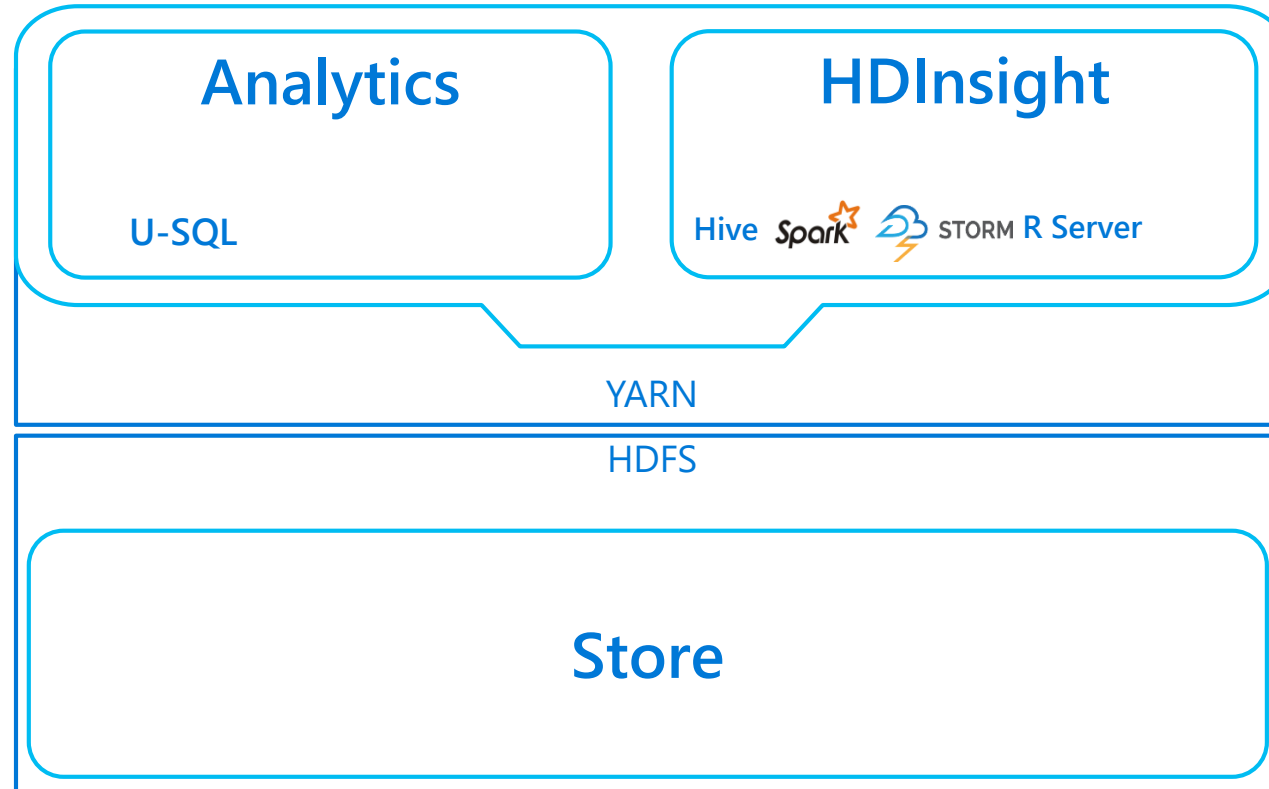


Azure Data Lake

Не думаем об инфраструктуре

Оплата за выполнение запросов

SQL-подобный язык



Контроль над инфраструктурой

Широко используемые открытые технологии

 Демонстрация

Обрабатываем тексты

<http://gutenberg.org> -> Azure Data Lake Analytics

#msdevcon

Azure Data Lake Analytics

```
@t = EXTRACT date string
      , time string
      , author string
      , tweet string
FROM "/Input/MyTwitterHistory.csv"
USING Extractors.Csv();

@res = SELECT author AS author
      , COUNT(*) AS tweetcount
FROM @t
GROUP BY author;

OUTPUT @res TO "/Output/MyTwitterAnalysis.csv"
ORDER BY tweetcount DESC
USING Outputters.Csv();
```

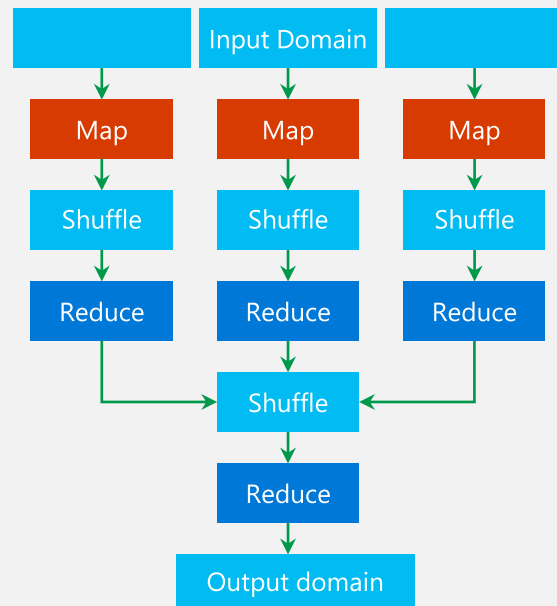
- U-SQL: простой и расширяемый язык запросов
- Декларативная природа SQL с мощностью C#
- Использование библиотек .NET, R и Python
- Параллелизация кода для различных задач (ETL, ML, image tagging, facial detection)

MapReduce



Batch	Script	SQL	Nosql	Stream	Machine Learning	Others
Map reduce	Map	Map Spark SQL	Map Accumulo Phoenix	Kafka Storm Spark	Spark MLlib Mahout	TV engines
YAKIN: data operating system						
HDFS (Hadoop Distributed File System)						

- Обработка происходит там, где хранятся данные
 - Данные хранятся распределенно (HDFS)
 - После первичных вычислений агрегируются только результаты
 - Линейное масштабирование при добавлении узлов
- Фреймворк берет на себя распределение функций по узлам



```
var map = function (key, value, context) {  
  var words = value.split(/[^\a-zA-Z]/);  
  for (var i = 0; i < words.length; i++) {  
    if (words[i] !== "")  
      context.write(words[i].toLowerCase(),  
1);  
  }  
};  
  
var reduce = function (key, values, context) {  
  var sum = 0;  
  while (values.hasNext()) {  
    sum += parseInt(values.next());  
  }  
  context.write(key, sum);  
};
```


Hadoop / HDInsight

Tools

Zeppelin

Ambari User Views

Data access

Batch

Map reduce

Script

Pig

SQL

Hive
Spark SQL

Nosql

Hbase
Accumulo
Phoenix

Stream

Kafka
Storm
Spark

Machine Learning

Sparkl Mlib
Mahout

Others

ISV engines

YARN: data operating system

HDFS (Hadoop Distributed File System)

Data management

 Демонстрация

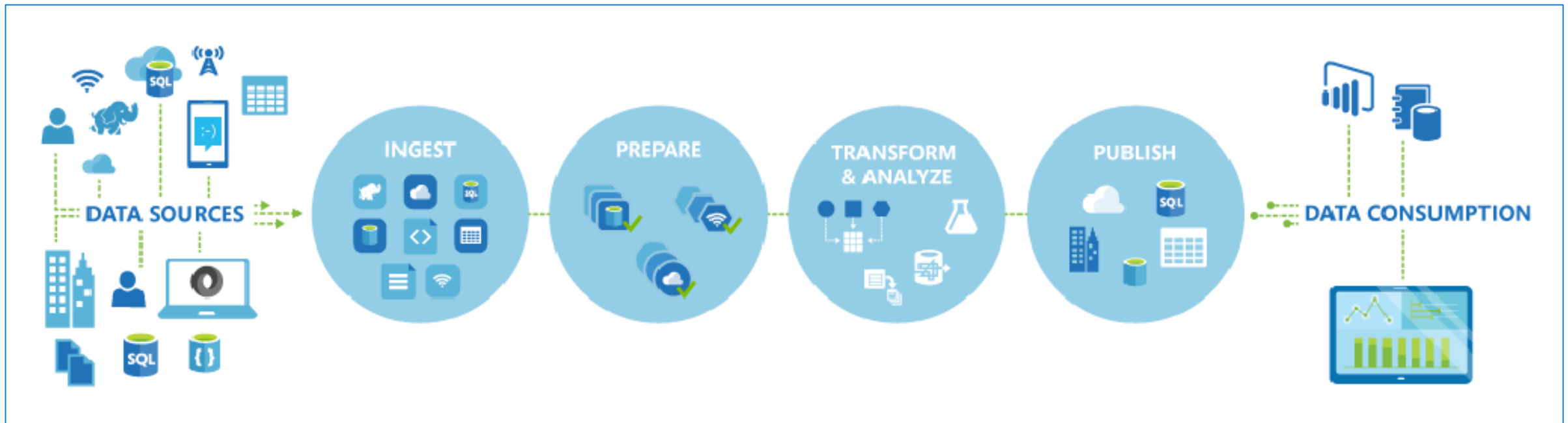
Обрабатываем тексты

<http://gutenberg.org> -> Spark

#msdevcon

Azure Data Factory

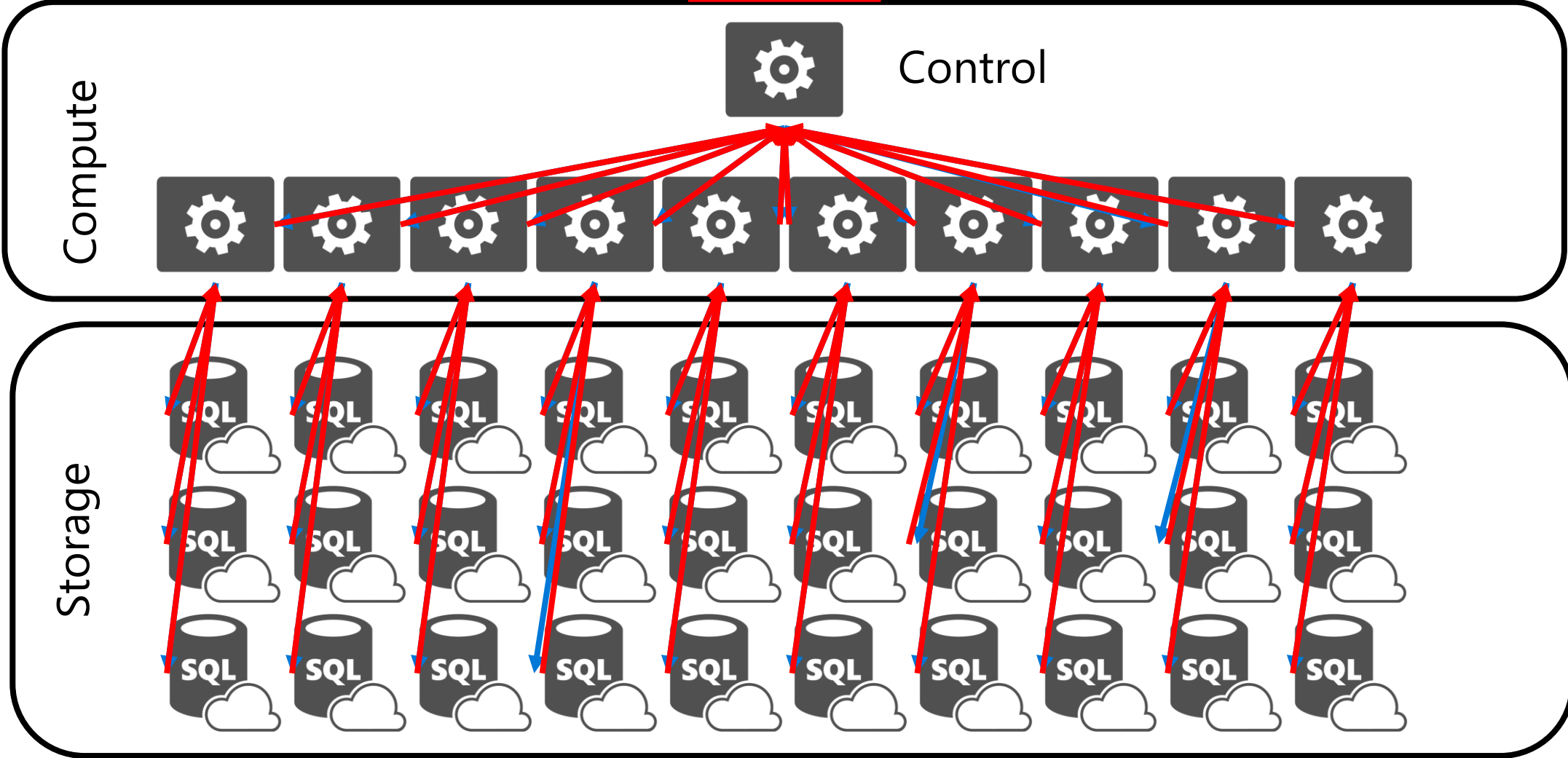
- Оркестрирование всех задач обработки данных
- Pipelines, которые запускаются по расписанию
- Различные источники данных на входе и на выходе
- Внутри pipeline – запросы на U-SQL, Spark/Hadoop и т.д.



SQL Data Warehouse



Result



Простой пример

```
SELECT  COUNT_BIG(*)  
FROM    dbo.[FactInternetSales]  
;
```



```
SELECT  SUM(*)  
FROM    dbo.[FactInternetSales]  
;
```



Control

Compute



```
SELECT  COUNT_BIG(*)  
FROM    dbo.[FactInternetSales]  
;
```



```
SELECT  COUNT_BIG(*)  
FROM    dbo.[FactInternetSales]  
;
```



```
SELECT  COUNT_BIG(*)  
FROM    dbo.[FactInternetSales]  
;
```



```
SELECT  COUNT_BIG(*)  
FROM    dbo.[FactInternetSales]  
;
```



Задачи для SQL Warehouse

Хорошо подходит

- Хранение больших объемов данных
- Отчеты по большим объемам данных
- Агрегирование данных

Плохо подходит

- Подготовка данных (построчная обработка)
- OLTP
- Большое количество индивидуальных запросов

AZURE ML STUDIO

ML Studio

Enter feedback here

User

Menu

Sample Experiment: Recommender System

Locked

Properties

Experiment Properties

START TIME	6/26/2014 1:07:1...
END TIME	6/26/2014 1:07:5...
STATUS CODE	Finished
STATUS DETAILS	None
<input type="checkbox"/> Disable upgrades	

[Prior Run](#)

Quick Help

```
graph TD; A[Restaurant ratings] --> B[Split<br/>Split data to train and test sets]; C[Restaurant customer data] --> D[Project Columns<br/>Select customer features to use in recommender model]; E[Restaurant feature data] --> F[Project Columns<br/>Select restaurant features to use in recommender model]; B --> G[Matchbox Recommender Train<br/>Train recommender model]; B --> H[Recommender Scorer<br/>Recommend restaurants to customers in test set]; D --> G; E --> G; G --> H; F --> H; H --> I[Recommender Evaluator<br/>Evaluate restaurant recommendations against test data];
```

NEW

VIEW RUN HISTORY

SAVE

SAVE AS

DISCARD CHANGES

REFRESH

CANCEL

RUN

PUBLISH WEB SERVICE

Эксперты



**Евгений
Григоренко**

Главный эксперт
по Big Data



**Ольга
Тележная**

Python Guru



**Дмитрий
Сошников**

Обработка данных, F#
IoT, ML, CNTK



**Jan
Pospisil**

Приглашенный эксперт

Сегодня: Azure и Big Data (в астрофизике)



Сергей Герасимов

ВМК МГУ



Александр Мещеряков

Институт космических исследований РАН

В четверг



Дмитрий Зобнин

Azure Batch: Big Compute без
головной боли (или почти без)



Jan Pospisil

Потоковые технологии обработки
больших данных на реальных примерах

В субботу



Андрей Устюжанин

**Большие данные в большом
адронном коллайдере**

Руководитель лаборатории больших
данных ФКН НИУ ВШЭ
Руководитель совместных проектов
Яндекс-CERN

