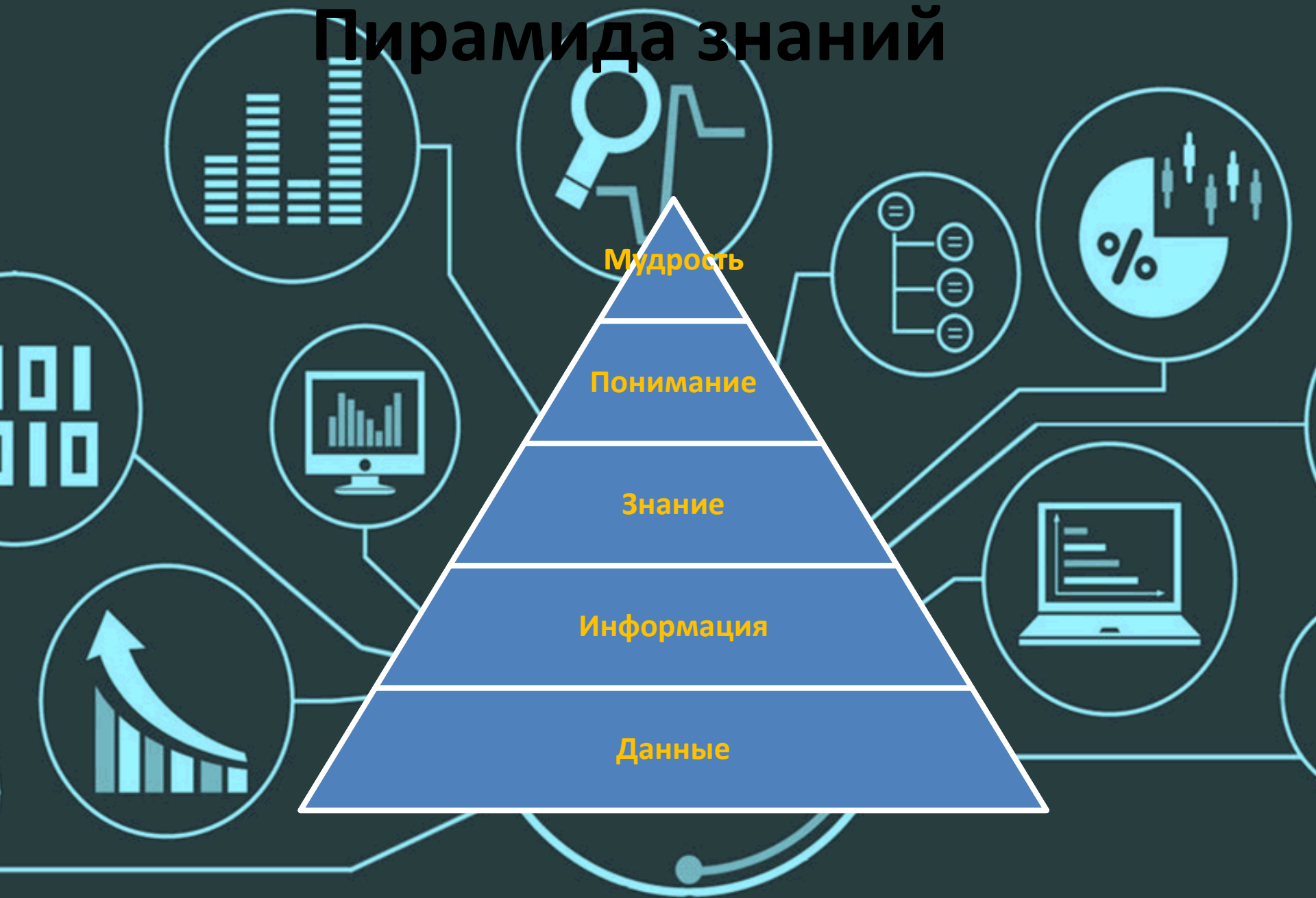# Методы работы с большими данными

Киреев Василий Сергеевич,

к.т.н., доцент

**Москва, 2020**

# Пирамида знаний

# DATA-driven

- **Data Driven (дословно — «управляемый данными»)** — это подход к управлению, основанный на данных. Его главный постулат: решения нужно принимать, опираясь на анализ цифр, а не интуицию и личный опыт.
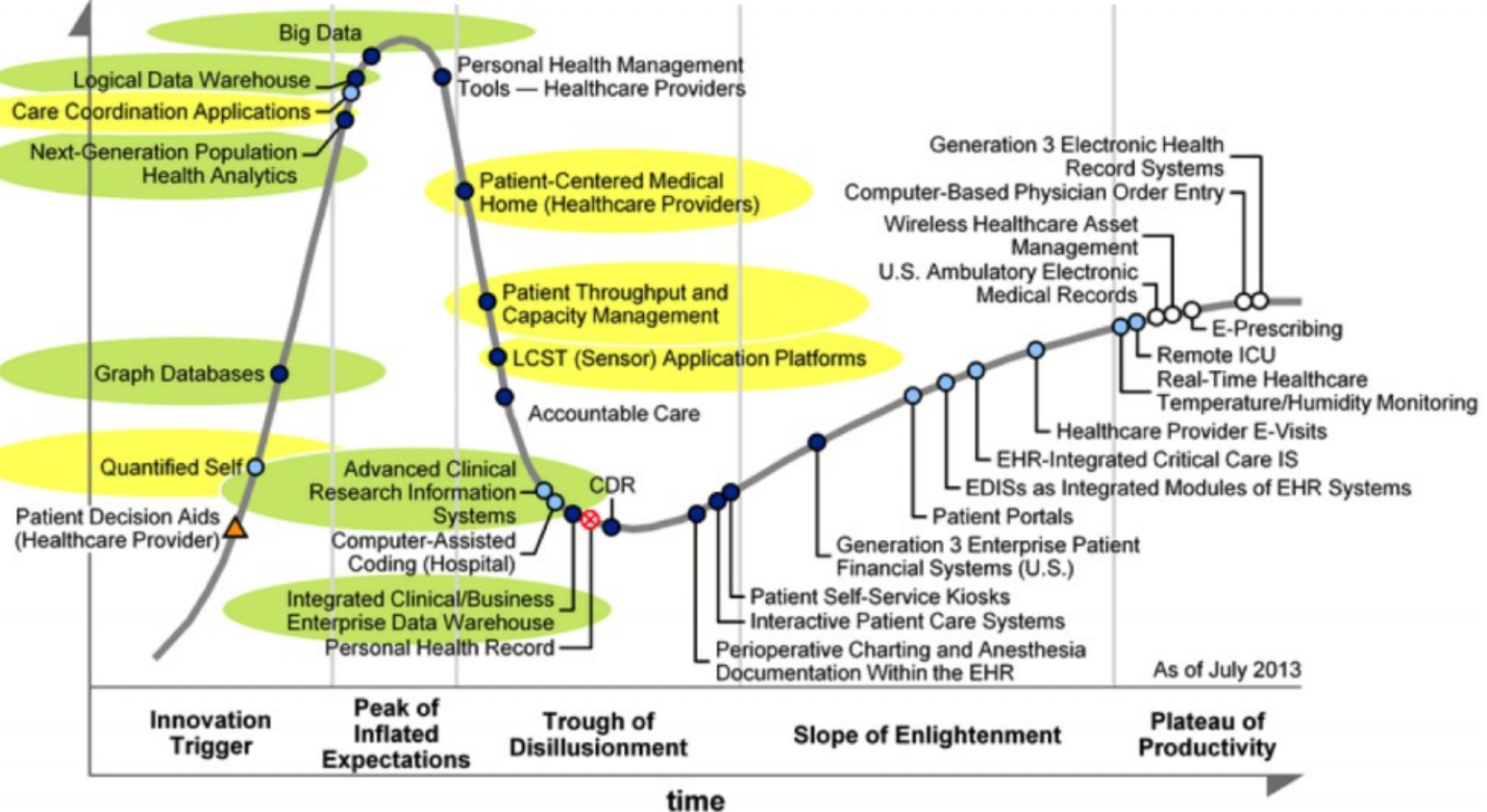
# Большие Данные

- Big Data is a term used to describe the large amount of data in the networked, digitized, sensor-laden, information-driven world *(NIST)*

- Термин «Большие данные» относится к наборам данных, размер которых превосходит возможности типичных баз данных (БД) по занесению, хранению, управлению и анализу информации

# Большие Данные

# Что такое большие данные?

# Большие данные в отраслях

**3,8** ПтБ данных приходится в среднем на каждую компанию

**Банки**

СОСТОЯНИЕ

**25%**
Big Data владеет финансовая индустрия

ЗАДАЧИ

сбор данных о:
- использовании кредитных карточек
- залогах
- кредитах
- профилях клиентов
- сбережениях клиентов

ЭФФЕКТИВНОСТЬ

**76%**
банков заявляют

Big Data позволяет привлекать новых клиентов, лучше взаимодействовать с ними и поддерживать их лояльность

---

**В 2,24** раза меньше было таких компаний 2 года назад

**Телеком**

СОСТОЯНИЕ

**45%**
компаний ведут активные проекты по Big Data

ЗАДАЧИ

- аналитика
- «умные» маркетинговые кампании
- выявление мошенничества
- улучшение качества связи

ЭФФЕКТИВНОСТЬ

**НА 92%** технологии Big Data могут уменьшить время обработки запроса пользователя

**85%**
компаний получают конкурентное преимущество

благодаря аналитике на базе технологий Big Data

---

**2,5** ПтБ
данных в час собирает сеть супермаркетов Walmart

**Ритейл**

СОСТОЯНИЕ

ЗАДАЧИ

- прогнозирование трендов покупательского спроса
- подготовка к резкому росту спроса на отдельные товары
- оптимизация акций и цены
- целевые маркетинговые кампании

ЭФФЕКТИВНОСТЬ

**НА 60%**

может увеличить операционную рентабельность максимально эффективное использование Big Data

# Уровни архитектуры по обращению с Большими данными

# Уровень сбора данных

Этот уровень отвечает за отделение шума от соответствующей информации, а также регулирование объема, скорости и разнообразия данных. Он должен иметь возможность проверять, очищать, преобразовывать, уменьшать и интегрировать данные в стек технологий больших данных для дальнейшей обработки. Это новое программное обеспечение, которое должно быть масштабируемым, устойчивым, отзывчивым и регулирующим в архитектуре больших данных.

# Уровень инфраструктуры

На данном уровне располагается физическая инфраструктура, необходимая для функционирования и масштабируемости архитектуры больших данных. Фактически наличие надежной и недорогой физической инфраструктуры привело к появлению таких важных тенденций, как big data. Для поддержки непредвиденного или непредсказуемого объема, скорости или разнообразия данных физическая инфраструктура для больших данных должна отличаться от инфраструктуры для традиционных данных.

# Уровень хранения

Использование массового распределенного хранилища и обработки является фундаментальным изменением в способе обработки больших данных предприятием. Распределенная система хранения данных обещает отказоустойчивость, а распараллеливание позволяет высокоскоростным алгоритмам распределенной обработки выполнять крупномасштабные данные. Для работы с «большими данными» используется несколько реализаций распределенных файловых систем. Основными из них можно считать реализацию от open-source проекта Hadoop (Hadoop Distributed File System - HDFS) и реализацию от Google (Google File System - GFS).

# Уровень управления и обработки

На уровне управления и обработки находятся инструменты и языки запросов для доступа к базам данных NoSQL с помощью файловой системы хранения HDFS, находящейся поверх уровня физической инфраструктуры Hadoop. С развитием вычислительной техники, теперь можно управлять огромными объемами данных, которые ранее могли бы быть обработаны только суперкомпьютерами за большие деньги. Цены на системы (ЦП, ОЗУ и диск) упали. В результате, новые методы для распределенных вычислений стали основным направлением.

# V-модель Больших Данных

# Volume

- Гигабайты – 1024 мегабайт
- Терабайты – 1024 гигабайт
- Петабайты – 1024 терабайт
- Эксобайты – 1024 петабайт
- Зетабайты – 1024 эксобайт

# Variety

- Тексты
- Изображения
- Видео
- Аудио
- Показания датчиков
- Транзакции
- Гео-данные

BIG DATA

# Velocity

- Пакетная обработка
- Стриминговая обработка
- Обработка в реальном времени

# Veracity

- Отсутствие целостности
- Отсутствие полноты
- Отсутствие определенности

# Value

- Статистическая ценность
- Финансовая ценность
- Ценность для продукта

# Validity

- Прозрачность
- Законность
- Конфиденциальность

# Тренды в Больших Данных



DATA AGE - THE GLOBAL DATASPHERE 2025
TRENDS & DATA-READINESS FROM EDGE TO CORE

## 175 Zettabytes

The global datasphere will grow from 33 zettabytes in 2018 to 175 zettabytes by 2025. IoT devices are expected to create over 90 zettabytes of data in 2025.

## 49%

By 2025, 49% of all data worldwide will reside in public cloud environments as cloud becomes the new core.

## 30%

In 2025 nearly 30% of the world's data will need real-time processing as the role of the edge continues to grow.

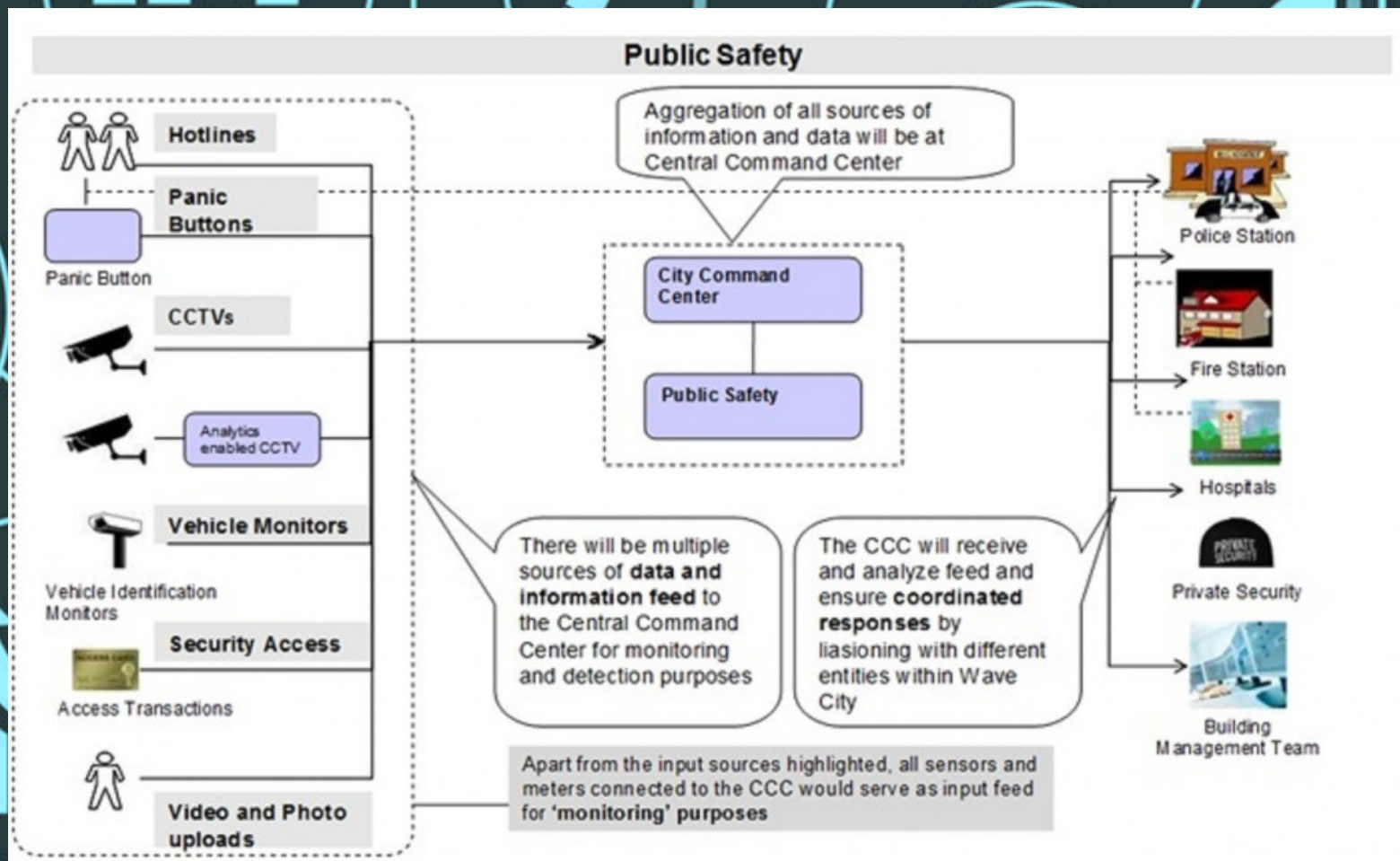IDC & Seagate Data Age 2025 - www.seagate.com/gb/en/our-story/data-age-2025/

# Применение Больших Данных

- Операционная аналитика
- Предиктивная аналитика
- Маркетинг
- Медицина
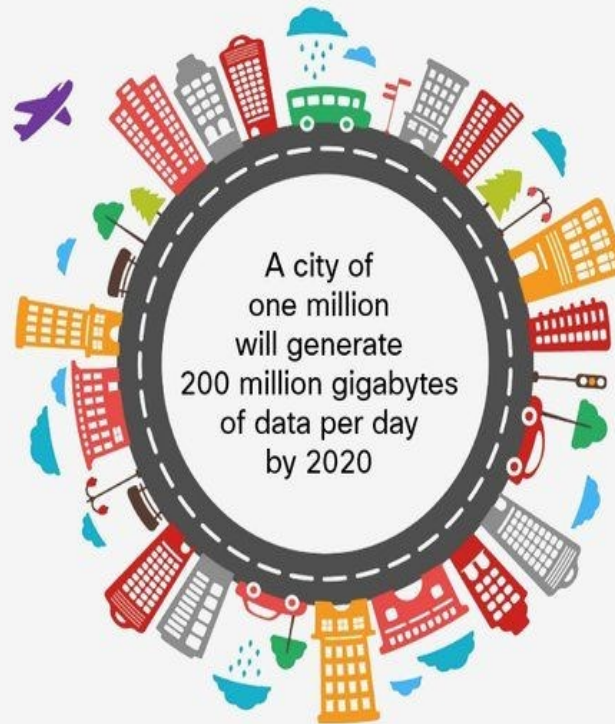- Производство
- Государство
- Городская среда

BIG DATA

# Общественная безопасность



**Public Safety**

Hotlines

Panic Buttons

Panic Button

CCTVs

Analytics enabled CCTV

Vehicle Monitors

Vehicle Identification Monitors

Security Access

Access Transactions

Video and Photo uploads

Aggregation of all sources of information and data will be at Central Command Center

City Command Center

Public Safety

There will be multiple sources of **data and information feed** to the Central Command Center for monitoring and detection purposes

The CCC will receive and analyze feed and ensure **coordinated responses** by liasioning with different entities within Wave City

Apart from the input sources highlighted, all sensors and meters connected to the CCC would serve as input feed for 'monitoring' purposes

Police Station

Fire Station

Hospitals

Private Security

Building Management Team

# Транспорт и городская инфраструктура

# Общественное здравоохранение



**Supply drivers**

**Medical & patient data**
Electronic Health Records (EHR) health sensors, social media, and genomics create rich new data sources for analytics

**Big data analytics**
Cheap computing power and sophisticated analytics drive insights into patient behavior, treatment costs and R&D

**Moblie/mHealth**
Pervasive mobile and smart phone adoption creates new engagement models within daily routines

**Health care professionals digital workflow**
Increasing integration of EHRs and telehealth drives new digitally-enabled coordinated workforce models of care

**Health information technology-enabled opportunities**

**Demand drivers**

Roll out business models tied to patient outcomes that also reduce medical errors and improve quality

Discover and deliver targeted and personalised therapies with real-world evidence of impact

Influence patients behaviors beyond the pill' and sustain engagement outside the traditional care setting

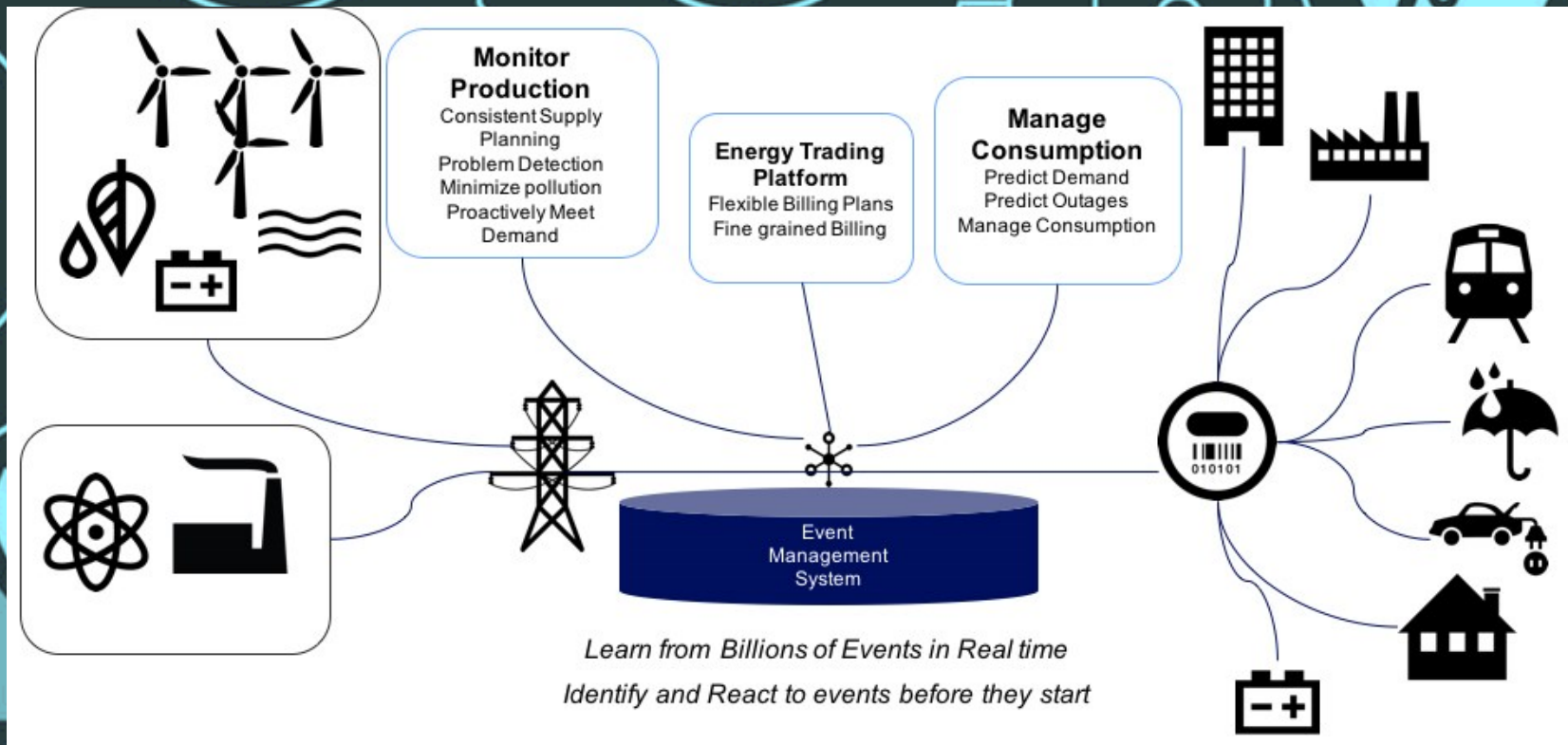Drive population health management, protocol-driven patient risk pool and stratification management

# Общественное здравоохранение

# Энергетика, окружающая среда и ЖКХ

# Цифровые государственные услуги в мире

## Governments Can Greatly Improve Usage and Satisfaction (I)

Very few services score well along both dimensions

**COLOR = % WHO ARE SATISFIED IN 2016**
- ● High satisfaction (70% or more)
- ● Medium satisfaction (50–69%)
- ● Low satisfaction (<50%)

**SIZE = % WHO USED A SERVICE IN THE PAST TWO YEARS**

| | | Denmark | Estonia | France | Germany | Netherlands | Norway | Russia | Sweden | UK | US |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **AVERAGE USAGE OF ALL SERVICES** (% of users) | | 21 | 27 | 21 | 20 | 19 | 21 | 21 | 22 | 19 | 21 |
| **Education** | Interactions with public institutions | 16 | 19 | 14 | 19 | 9 | 18 | 33 | 17 | 13 | 17 |
| **Health** | Health care records | 62 | 37 | 13 | 11 | 45 | 44 | 57 | 34 | 15 | 31 |
| **Immigration** | Passport services | 20 | 28 | 15 | 16 | 26 | 18 | 26 | 21 | 33 | 17 |
| | Visa, residency, and work permits | 6 | 5 | 6 | 9 | 6 | 15 | 5 | 9 | 12 | 8 |
| **Registries** | Address updates | 23 | 19 | 31 | 27 | 21 | 33 | 7 | 30 | 20 | 23 |
| | Company information updates | 26 | 49 | 19 | 15 | 11 | 18 | 26 | 15 | 13 | 17 |
| **Social services** | Subsidy and benefit applications | 23 | 19 | 17 | 18 | 23 | 26 | 13 | 20 | 19 | 28 |
| | Payments to pensions | 18 | 14 | 10 | 16 | 9 | 8 | 14 | 20 | 7 | 9 |
| | Employment services and job searches | 30 | 14 | 32 | 30 | 31 | 36 | 31 | 35 | 32 | 37 |
| | Public-housing services | 23 | 11 | 12 | 12 | 17 | 8 | 8 | 15 | 11 | 11 |
| **Taxes and customs** | Electronic gates at border control checkpoints | 4 | 7 | 6 | 8 | 5 | 10 | 0 | 4 | 16 | 5 |
| | Tax returns | 68 | 92 | 65 | 45 | 69 | 71 | 27 | 75 | 27 | 55 |
| | Tax, rate, and fine payments | 49 | 71 | 57 | 49 | 47 | 26 | 61 | 53 | 33 | 33 |
| **Transportation** | Real-time information | 54 | 72 | 69 | 58 | 36 | 60 | 60 | 65 | 50 | 51 |

**Source:** BCG 2016 Digital Government Satisfaction Survey.

**Note:** Of the 25 services covered in our survey, 15 are shown here. Usage means that a user completed at least some part of a transaction online. Each country may not allow end-to-end transactions online for a service. Survey question (usage): "Have you used the internet for the following interactions with government at least once in the past two years?" Respondents who answered "Yes" have been included. Survey question (satisfaction): "How satisfied are you with the use of the internet in delivering each kind of government service?" Response options range from 1 to 7, where 1 = "Extremely dissatisfied" and 7 = "Extremely satisfied." Respondents who selected 6 or 7 have been included as satisfied.
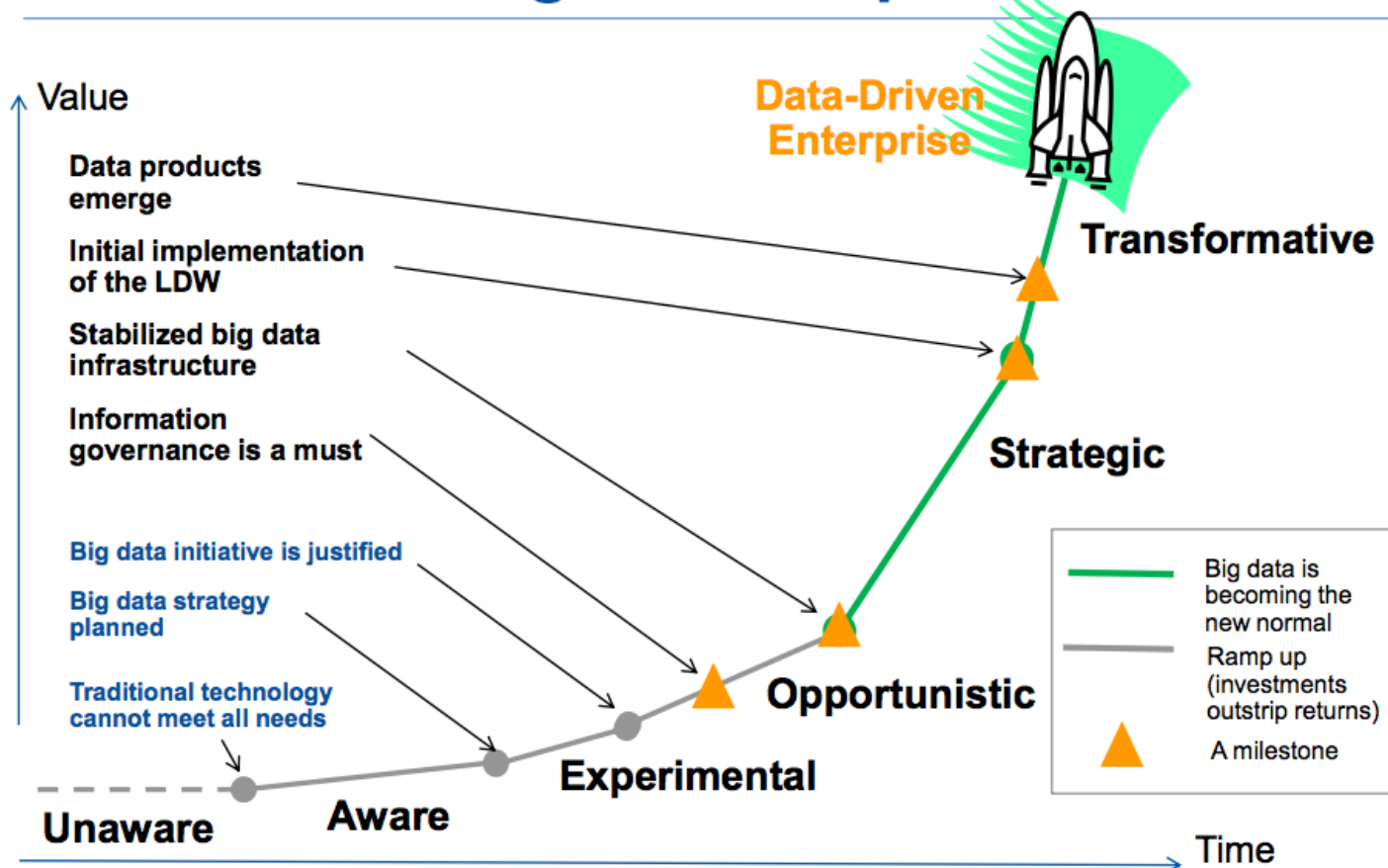
# Большие данные в гос. управлении

- Открытые данные
- Оборона
- Защита прав потребителей
- Общественная безопасность
- Транспорт и городская инфраструктура
- Общественное здравоохранение
- Энергетика, окружающая среда и ЖКХ
- Образование

# Внедрение Больших Данных



The Road Map: Typical Stages and Milestones of Big Data Adoption

# Угрозы и риски использования Больших данных

риск конфиденциальности

риск потери данных

риск переполнения хранилища

риск снижения эффективности больших данных

риск формирования неэффективного набора данных

риск мошенничества

риск неготовности к переменам

риск внешнего консультанта

риск экономической нецелесообразности

риск ошибок бизнес-модели

риск ошибок больших данных

# Снижение риска ошибок Больших данных

- проводить периодические ревизии данных

- контролировать ключевые параметры данных

- вести журнал выявленных ошибок и их устранения

- разрабатывать инструменты и алгоритмы устранения или нивелирования ошибок и некорректных состояний данных

- оценивать результативность инструментов

- проводить независимую оценку и экспертизу

- применять специальные средства тестирования данных и инструментов, которые разрабатываются самостоятельно

- использовать инструменты последовательно, подконтрольно и пошагово с постоянным контролем обрабатываемых данных в целом или по выборкам

# Технологии



Big Data Landscape 2016 (Version 3.0)

Last Updated 3/23/2016 © Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark Capital (@firstmarkcap) FIRSTMARK

# Data Lake

# Процесс работы с данными



**Integrating big data into the traditional IT architecture**

**Multiple data sources**
- Mainframe
- Relational databases
- CRM
- Supply chain management
- ERP
- Documents
- Log files
- External data feeds
- Unstructured data
- Internet and social media
- Audio and video
- Images

**ETL** →

Traditional data warehousing components
- Data extraction
- Report server
- Data warehouse
- Multi-dimensional cubes
- Data marts

**Data mining**

**Business intelligence**

**Offload data and processing** ↓    ↑ **Move high-value results**

**ELT** →

New big data components
- Big data catalogues
- Processing
- Advanced analytics

Data science teams

**Visualization and enterprise reporting**
- Dashboards
- Standard reports
- Ad hoc reports
- Real-time analytics
- Infographics
- Export (pdf, xls, doc)
- Collaboration

**Analytics and segmented users**
- Business analysts
- Business stakeholders
- Executives
- Decision makers
- Database administrators
- Data warehouse users
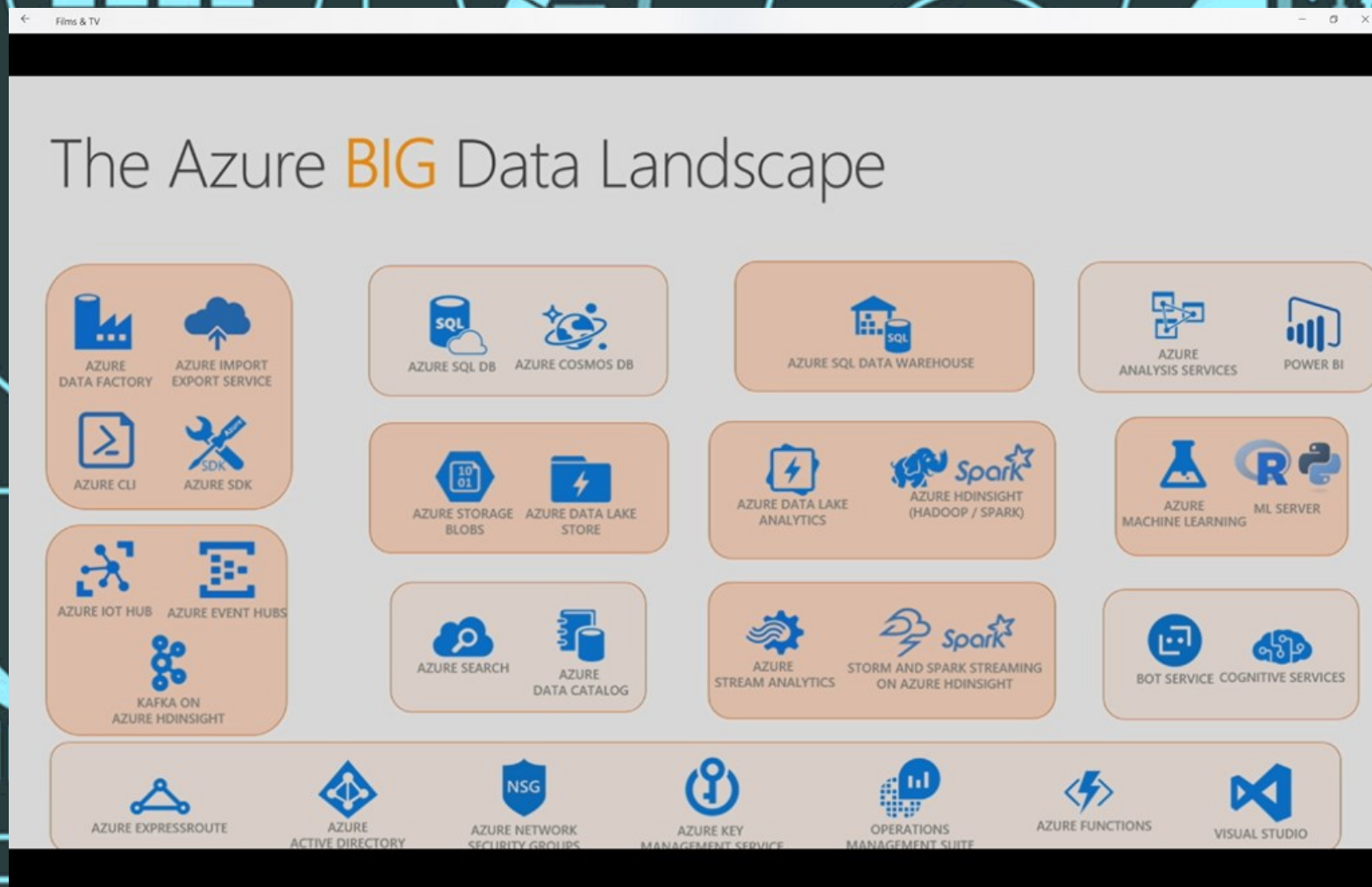- Administrators

# Процесс Apache

# Процесс Google



BigQuery: How Google handles BigData

**A fully-managed data analytics service in the cloud.**
Blazingly fast. Unlimited storage. Interactive analysis on multi-terabyte datasets.

Corporate data
3rd party data

Analytics
Adwords
DoubleClick
AdSense

BigQuery
SQL   API

tableau
QlikView   Other B.I. Tools

Google Spreadsheets

Co-Workers

**Host your big data in the cloud**

**Analyze interactively Mash it up**

**Securely share & distribute the results**

# Процесс Google



BigQuery: How Google handles BigData

# Процесс Microsoft

# Процесс Amazon

# Процесс IBM