

# Особенности Big Data

Максим Губин

Томск



# Два типа больших данных

Это Big Data:	Это тоже Big Data:
Поток данных телескопа VLA	Сделки биржи NYSE
Поток данных LHC	Действия игроков Eve Online
GPS-треки с общественного транспорта	Facebook
Покупки в супермаркете	YouTube
Wayback Machine	E-Bay
	Транзакции Visa и MasterCard (и Мир)

# «Научные» Big Data

- ❖ Обычно выход датчиков какого-то рода;
- ❖ Может быть ограничение на количество данных, получаемых в единицу времени, но обычно это довольно большое значение;
- ❖ Обогащение данных может привести к тому, что они вырастут до произвольного размера, одновременно увеличивая их полезность.

**При работе с «научным» видом больших данных необходимо помнить одну важную концепцию:**

# «Научные» Big Data

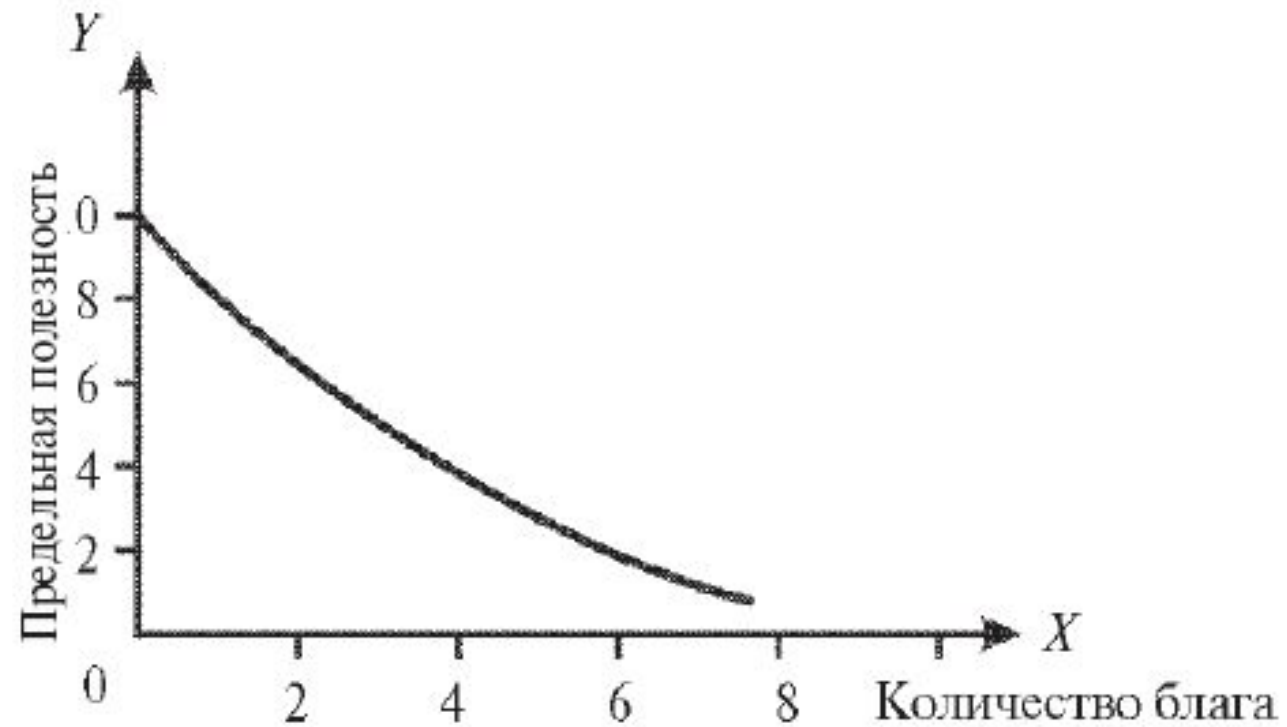


Рис. 2. Кривая предельной полезности

# «Научные» Big Data

## Вопросы, на которые стоит ответить:

- ❖ Насколько растёт точность наших моделей с каждой следующей записью?
- ❖ Сколько стоит получить следующую запись?
- ❖ Являются ли некоторые данные более ценными, чем другие?
- ❖ Можем ли мы пожертвовать некоторыми записями, при этом получая пользу?
- ❖ Когда мы начнем терять полезность из-за накладных расходов?
- ❖ Выйдем ли мы в плюс, если продолжим наращивать объёмы данных?

# «Научные» Big Data

## Важные особенности:

- ❖ Цена данных обычно низкая;
- ❖ Стабильное хранение данных обычно не требуется;
- ❖ Сами данные не ценны и не полезны, ценны и полезны результаты их исследования;
- ❖ Согласованность данных важна, но не критична;
- ❖ Потеря даже 100% данных во многих случаях является лишь незначительной проблемой.\*

# «Бизнес» Big Data

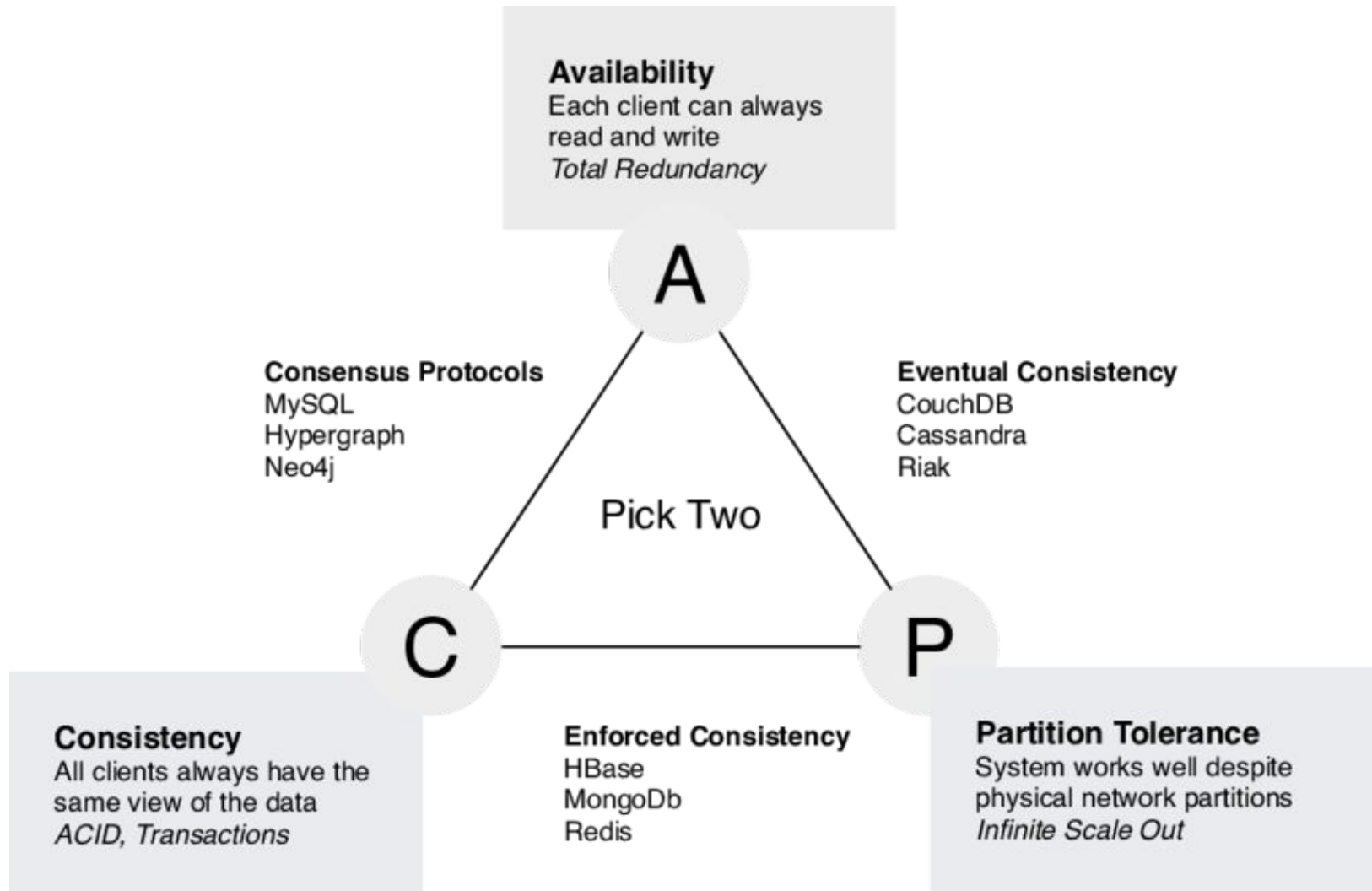
## Важные особенности:

- ❖ Бизнес-данные, выросшие настолько, что с ними уже нельзя работать традиционными подходами.
- ❖ Такие данные обычно важные, дорогие, требуют длительного хранения, и потеря даже малой их части может быть катастрофической.

Здесь компромиссы гораздо менее выражены, потому что такие данные очень плохо переносят деградацию объёма.



# Теорема CAP





# Теорема CAP

## Consistency:

Каждое чтение возвращает самые свежие записанные данные либо ошибку.

# Теорема CAP

## Availability:

Каждый запрос вернет ответ без гарантий, что в ответе содержатся самые свежие данные.

# Теорема CAP

## Partition Tolerance:

Система продолжает работать, несмотря на произвольный уровень потери связности её узлов.

# Теорема CAP, 3 варианта

## Consistency:

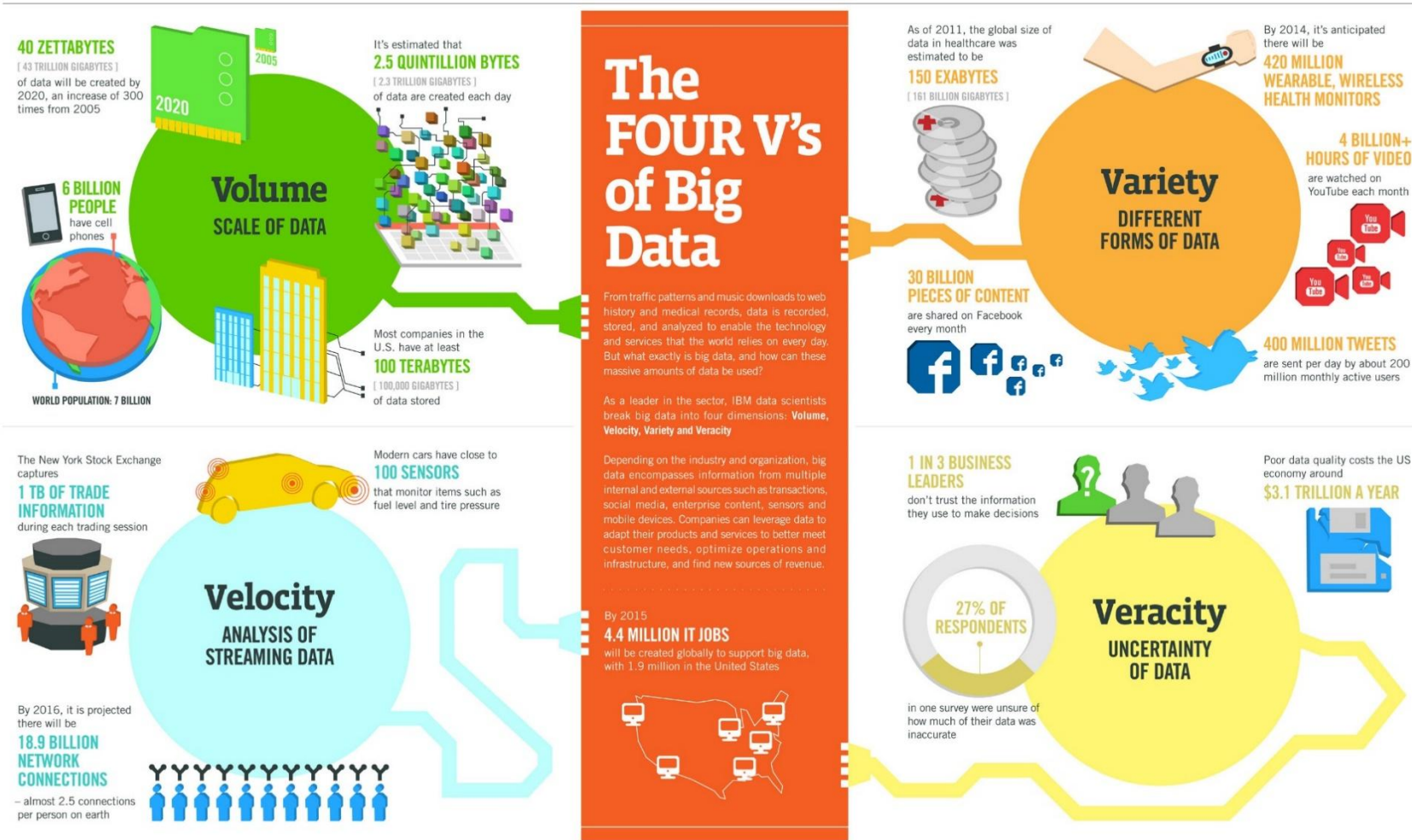
Система вернет **ошибку или таймаут**, если не может гарантировать актуальность данных из-за проблем с сетью.

## Availability:

Система всегда ответит на запрос **самой новой доступной версией** данных, даже если она не может гарантировать актуальности информации из-за проблем с сетью.

Третьего варианта нет. Если БД не фрагментирована, у вас есть все 3 полезных свойства.

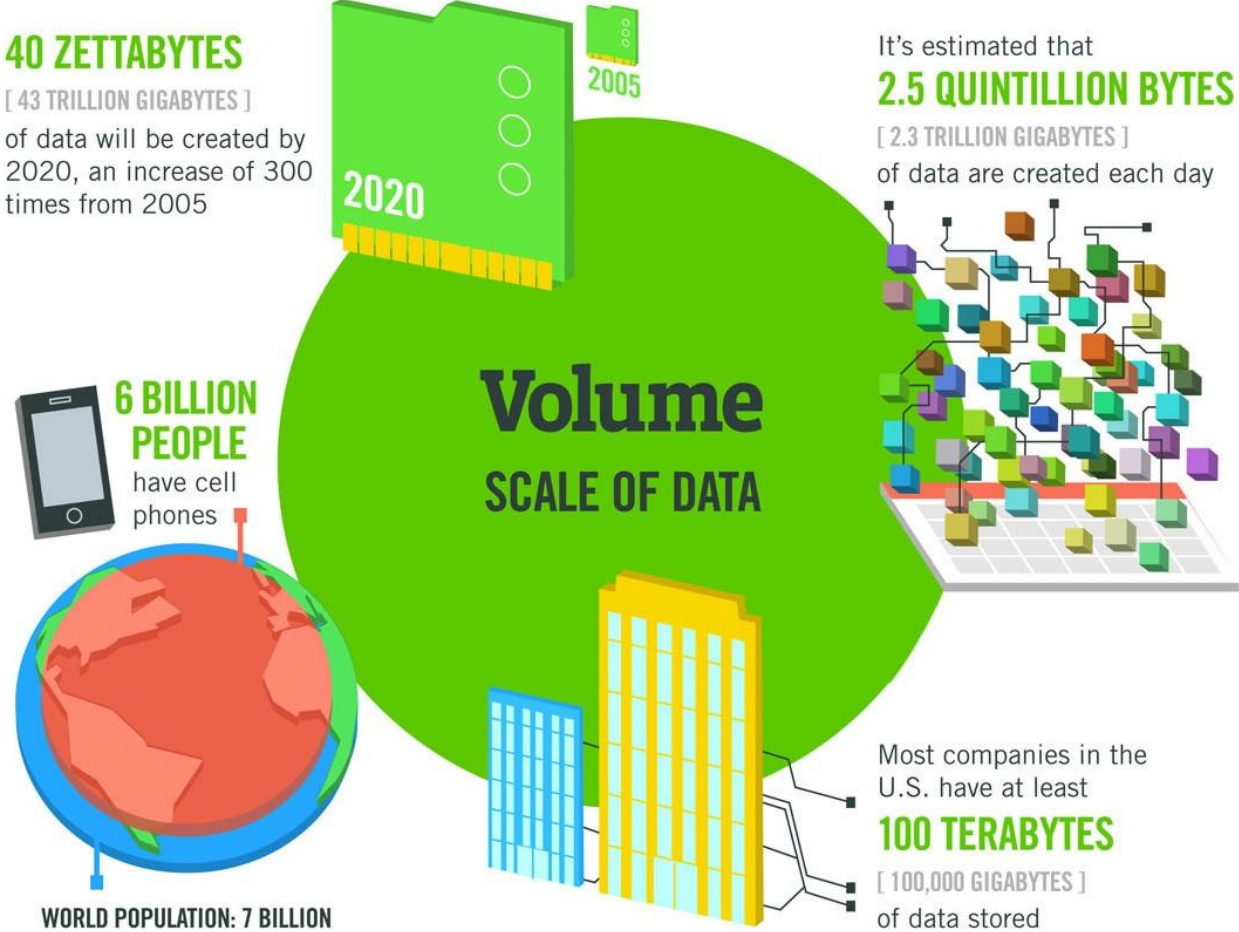
# Volume, Variety, Veracity, Velocity



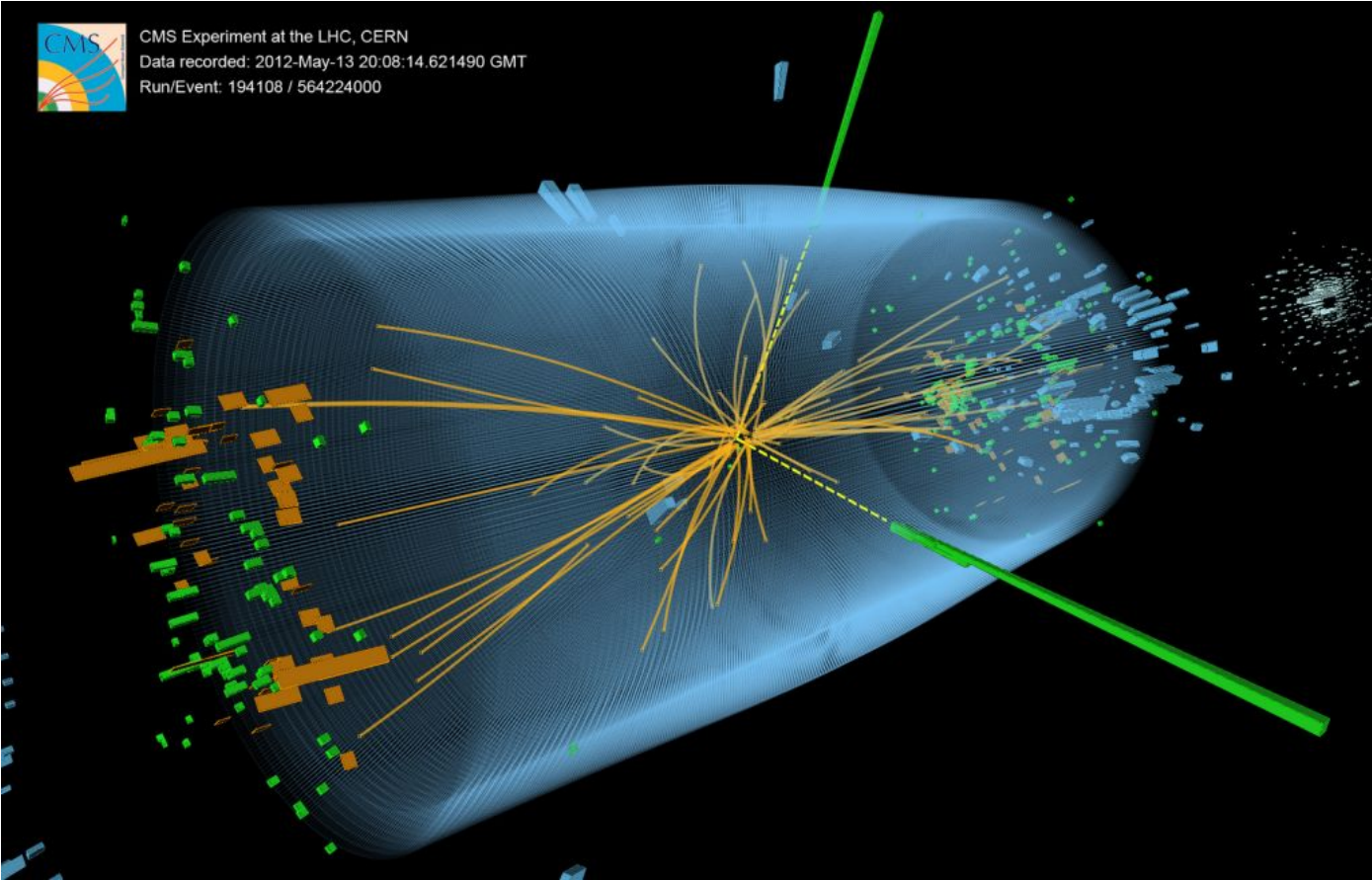
Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTec, QAS



# Volume



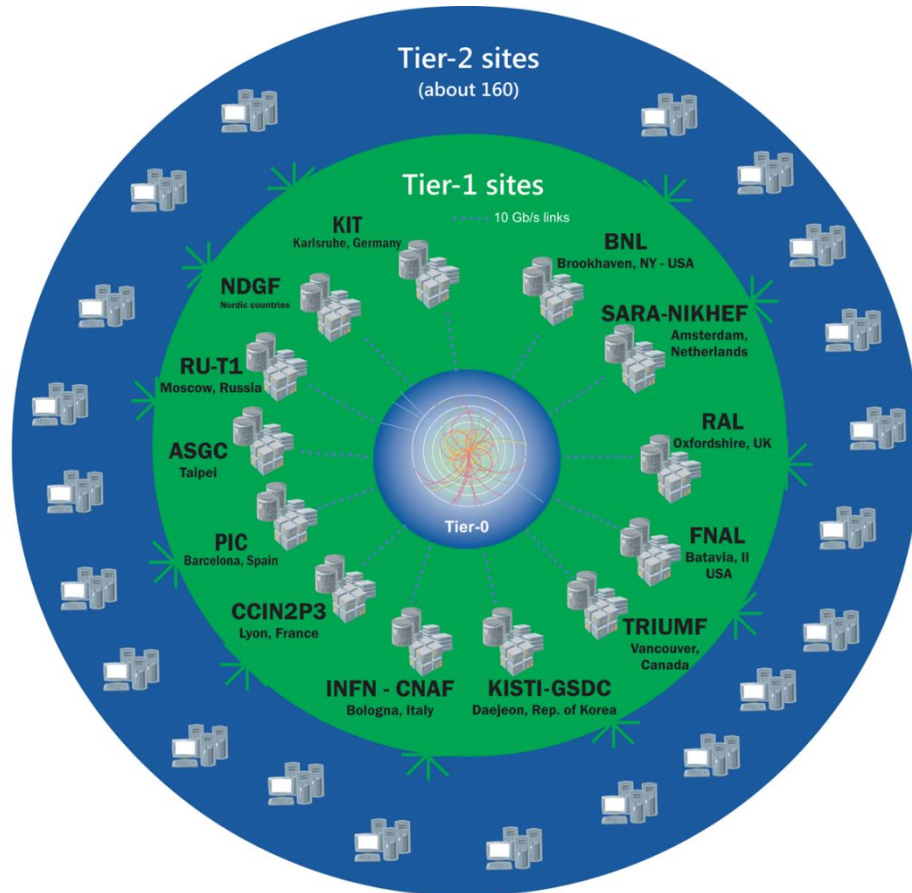
# Volume



CERN



# Volume



## CERN – Rucio

Более 350 PB данных, миллиарды файлов, в более чем 120 ЦОД по всему миру.

Три копии файла на разных континентах и одна на плёнке? Ok.

Стирать непопулярные файлы автоматически? Ok.

# Volume



## CERN

Долговременное хранение на плёнке,  
библиотечные стримерные роботы.

# Volume



Wayback Machine

<http://archive.org>

>15 петабайт данных,  
объём растёт на 20 Тб в неделю.

# Volume

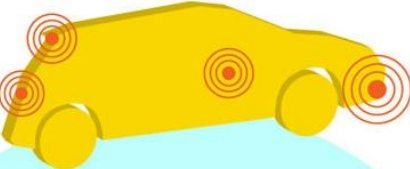
## Подходы к проблеме:

- ❖ **Управление объёмом данных;**
- ❖ **Вложения в хранилища;**
- ❖ **Разработка специализированного ПО и АО;**
- ❖ **Использование специализированного ПО от других разработчиков.**

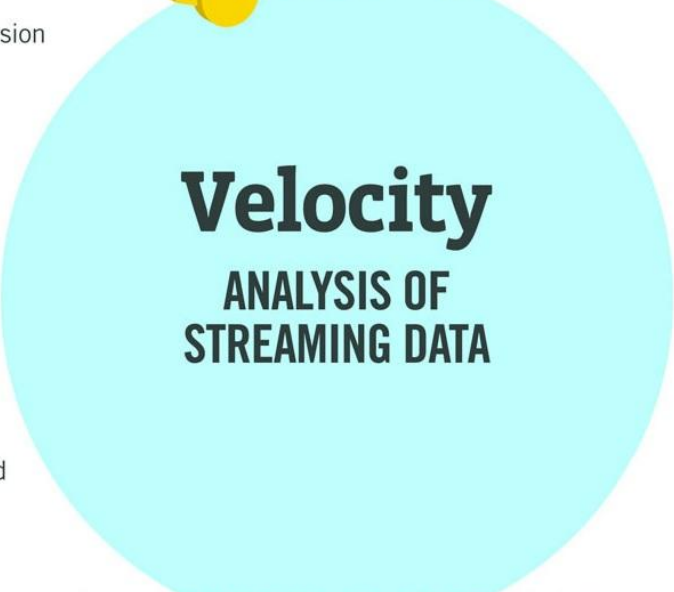


# Velocity

The New York Stock Exchange captures  
**1 TB OF TRADE INFORMATION**  
during each trading session



Modern cars have close to  
**100 SENSORS**  
that monitor items such as  
fuel level and tire pressure



By 2016, it is projected  
there will be  
**18.9 BILLION NETWORK CONNECTIONS**  
– almost 2.5 connections  
per person on earth



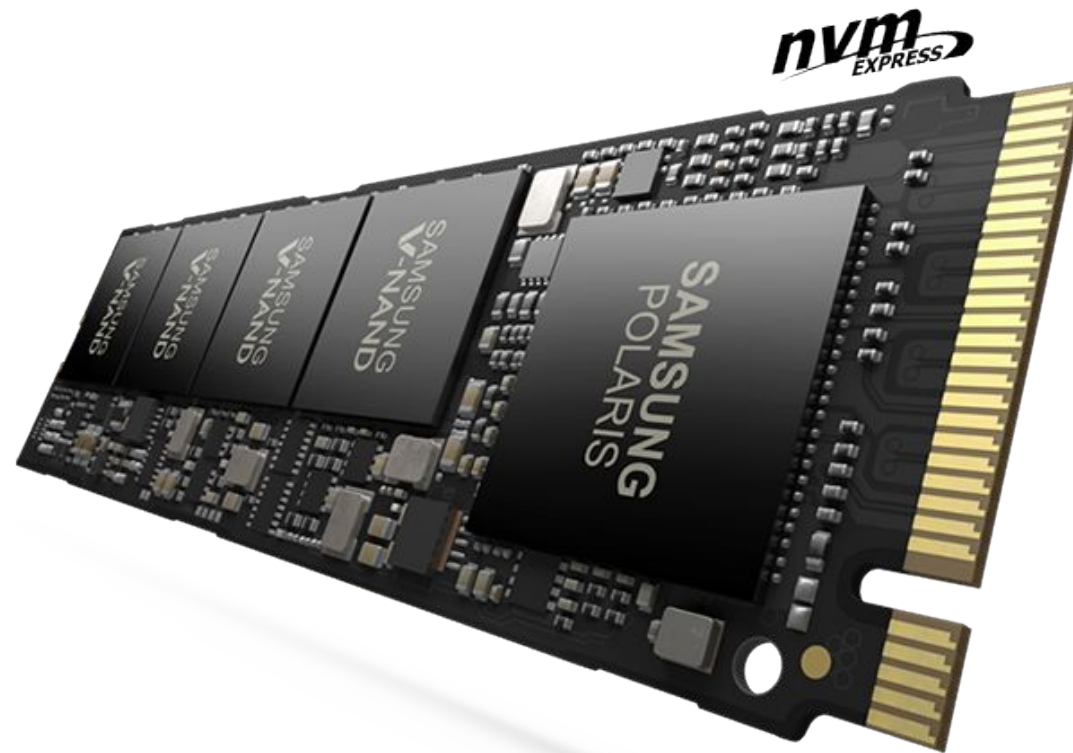
# Velocity

Возможное решение проблемы:



# Velocity

Причина рождения «больших данных»:





# Velocity

## CERN:

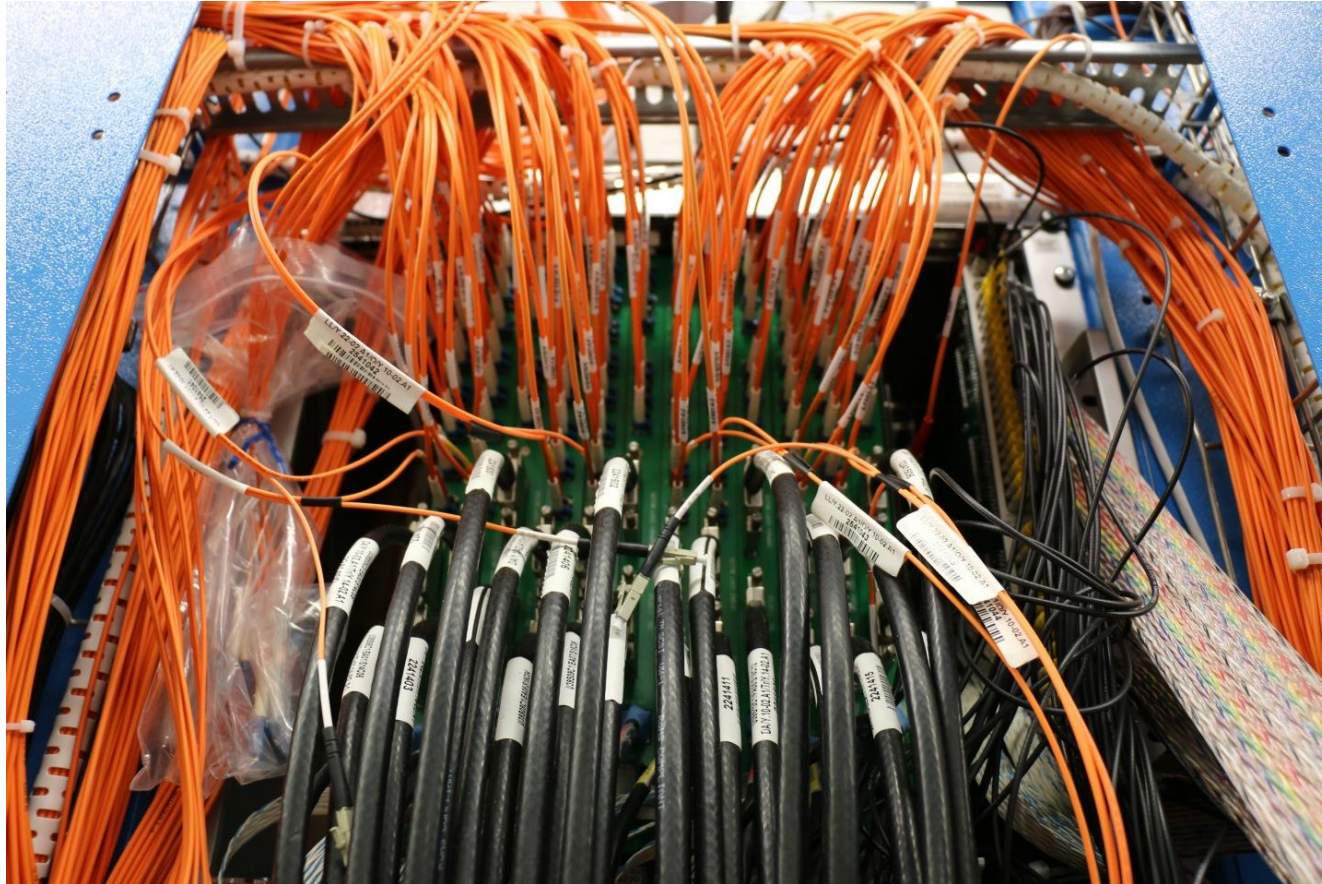
ATLAS выдаёт ~100 TB/с со своих датчиков.

Лишь ~1 GB/с сохраняется для дальнейшей обработки.

## Trigger:

The data reduction is carried out in two stages: first, **custom electronics** performs an initial level of data rejection for each bunch crossing based on partial and localized information. Only data corresponding to collisions passing this stage of selection will be actually read-out from the on-detector electronics. Then, a large computer farm (**~17 k cores**) analyses these data in real-time and decides which ones are worth being stored for Physics analysis.

# Velocity



CERN Trigger hardware

# Velocity



**NYSE:**

2Gbps link to allow for latency reduction,  
to allow for even faster high-speed  
trading.

**Передача данных по лазеру**

# Velocity

Пути решения проблемы:

Kafka, Flume и Logstash дают  
возможность потокового сбора данных  
и совместимы с множеством разных  
источников и приёмников данных.





# Variety

As of 2011, the global size of data in healthcare was estimated to be

**150 EXABYTES**

[ 161 BILLION GIGABYTES ]



By 2014, it's anticipated there will be

**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**



**4 BILLION+ HOURS OF VIDEO**

are watched on YouTube each month



**30 BILLION PIECES OF CONTENT**

are shared on Facebook every month



**Variety**  
DIFFERENT FORMS OF DATA



**400 MILLION TWEETS**

are sent per day by about 200 million monthly active users

# Variety

---

Около 1500 единиц общественного транспорта, оборудованных GPS

GSM

Данные о местоположении абонентов от сотовых операторов

IR;  
CCTV

Карты, данные аэрофотосъёмки и спутниковой съёмки

Видеорегистраторы

Поток видео с веб камер

Дорожные камеры

GPS

Спутниковое видео

---

# Variety

## Подходы к проблеме:

- ❖ Обработка естественного языка;
- ❖ Текст-в-речь;
- ❖ Классификация изображений;
- ❖ Machine Learning-классификация всего;
- ❖ Семантические технологии, Web 2.0;
- ❖ NoSQL.

Универсального решения не существует, каждая предметная область требует своего подхода.



# Veracity

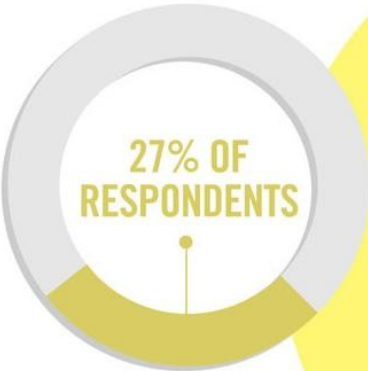
## 1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions



Poor data quality costs the US economy around

**\$3.1 TRILLION A YEAR**



in one survey were unsure of how much of their data was inaccurate

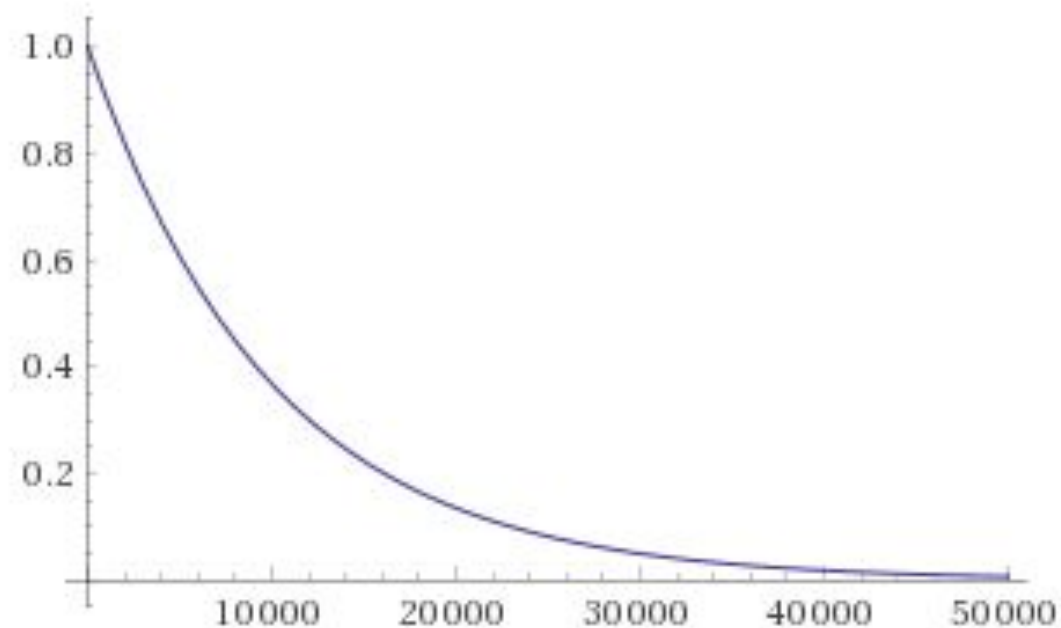
**Veracity**  
**UNCERTAINTY OF DATA**

# Veracity

## Следствие из теории вероятностей

Когда объём выборки стремится к бесконечности, вероятность ошибки в данных возрастает до определённости.

$$P(A \text{ and } B) = P(A) * P(B)$$



Computed by Wolfram|Alpha

$$0.9999^{**}x$$

# Veracity: AAA

Anyone can say anything about anything at any moment.

## Разрешение неопределённости:

- ❖ Проверка, либо
- ❖ Реификация.

.

# Спасибо за внимание!

mgubin@tpu.ru

econophysics 