

Introduction to Big Data

Nelson L. S. da Fonseca
IEEE ComSoc Summer School
Albuquerque, July 17-21, 2017

What are big data?





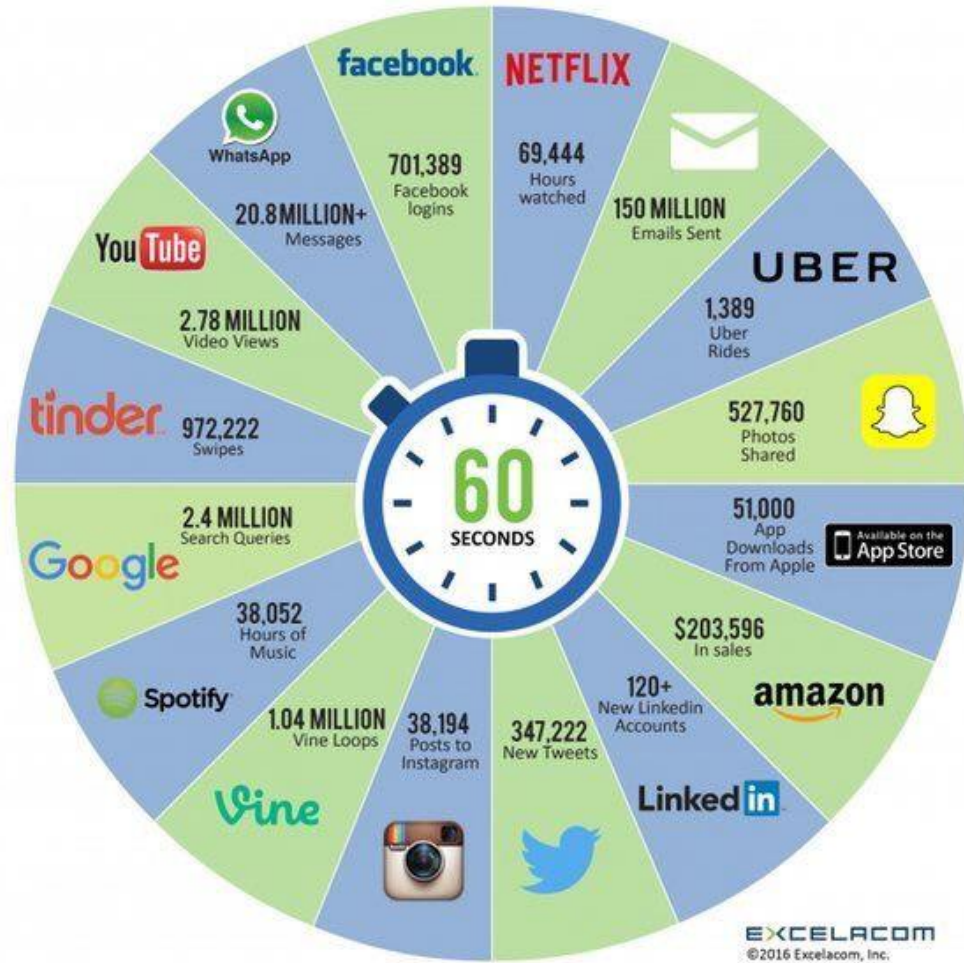
Source: <http://www.spiegel.de/panorama/bild-885931-473266.html>



Source: <http://www.spiegel.de/panorama/bild-889031-473242.html>

In 60 seconds..

2016 What happens in an INTERNET MINUTE?

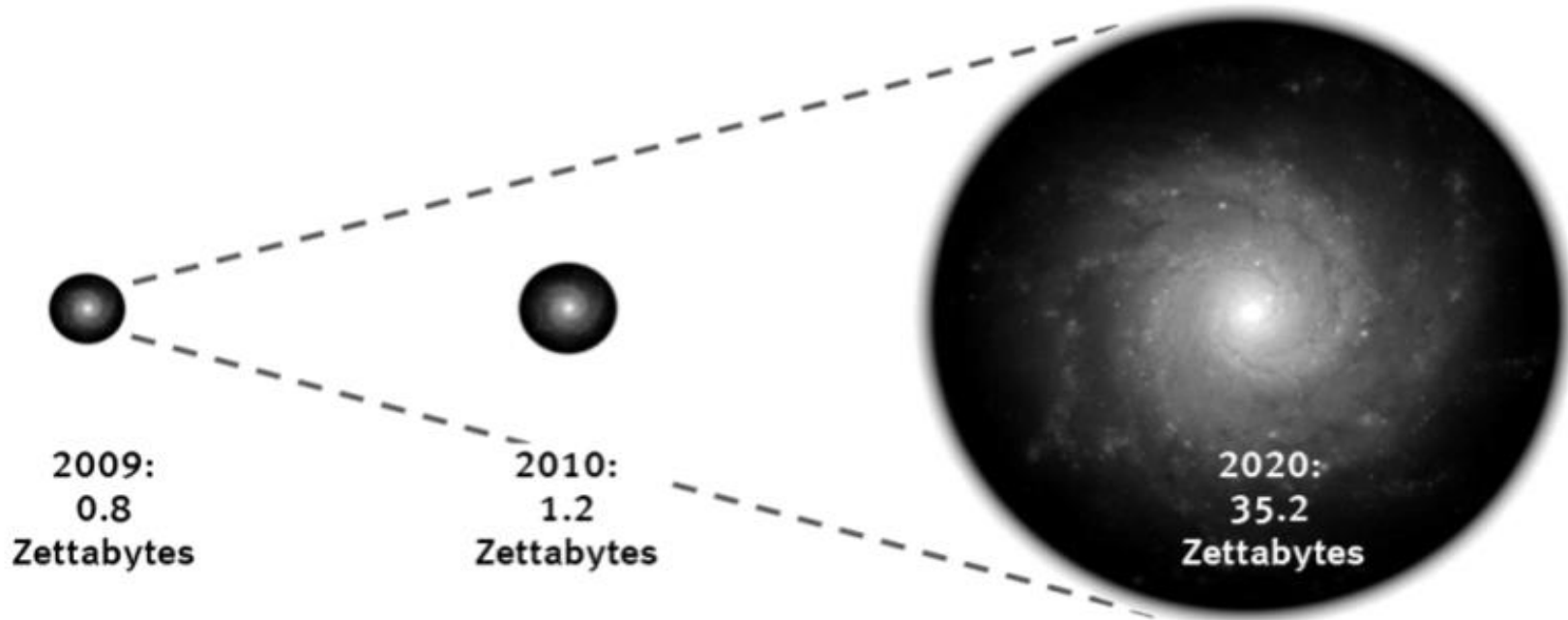


2017 This Is What Happens In An Internet Minute



Big Data Size: The Volume Of Data Continues To Explode

The Digital Universe 2009 - 2020





"If a tree falls in a forest and no one is around to hear it, does it make a sound?" --- George Berkeley

Data ≠ Knowledge

Big Data

Data



Information



Presentation




Knowledge



Processing Big Data

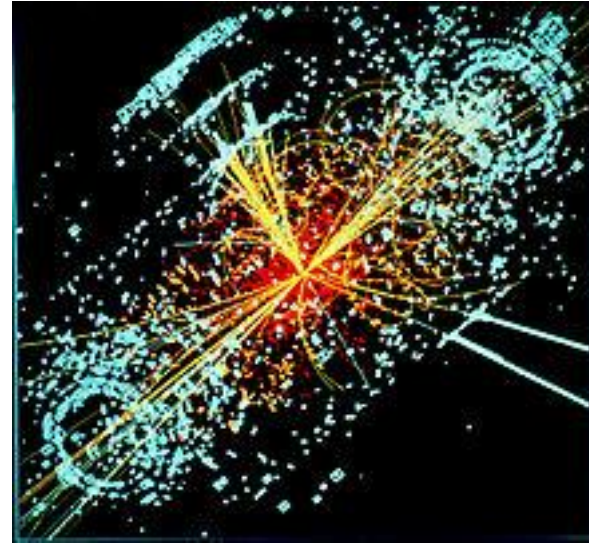
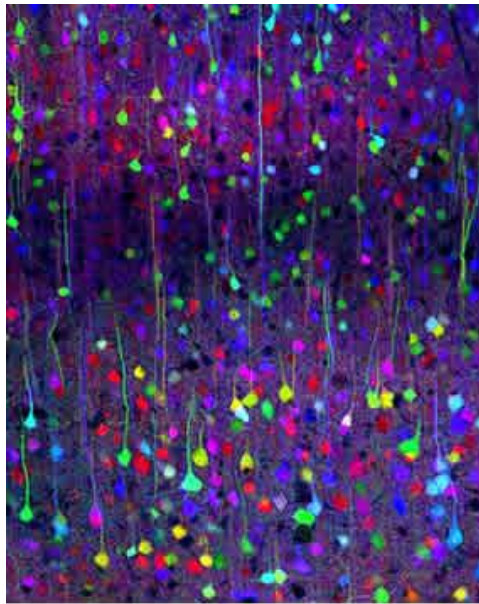
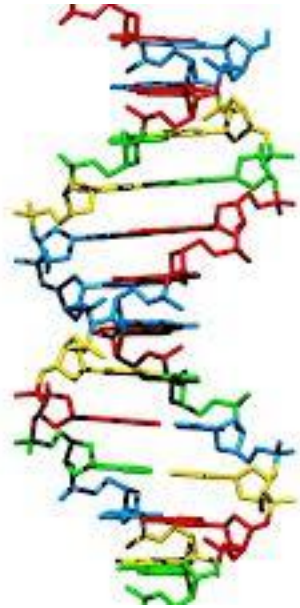




1%

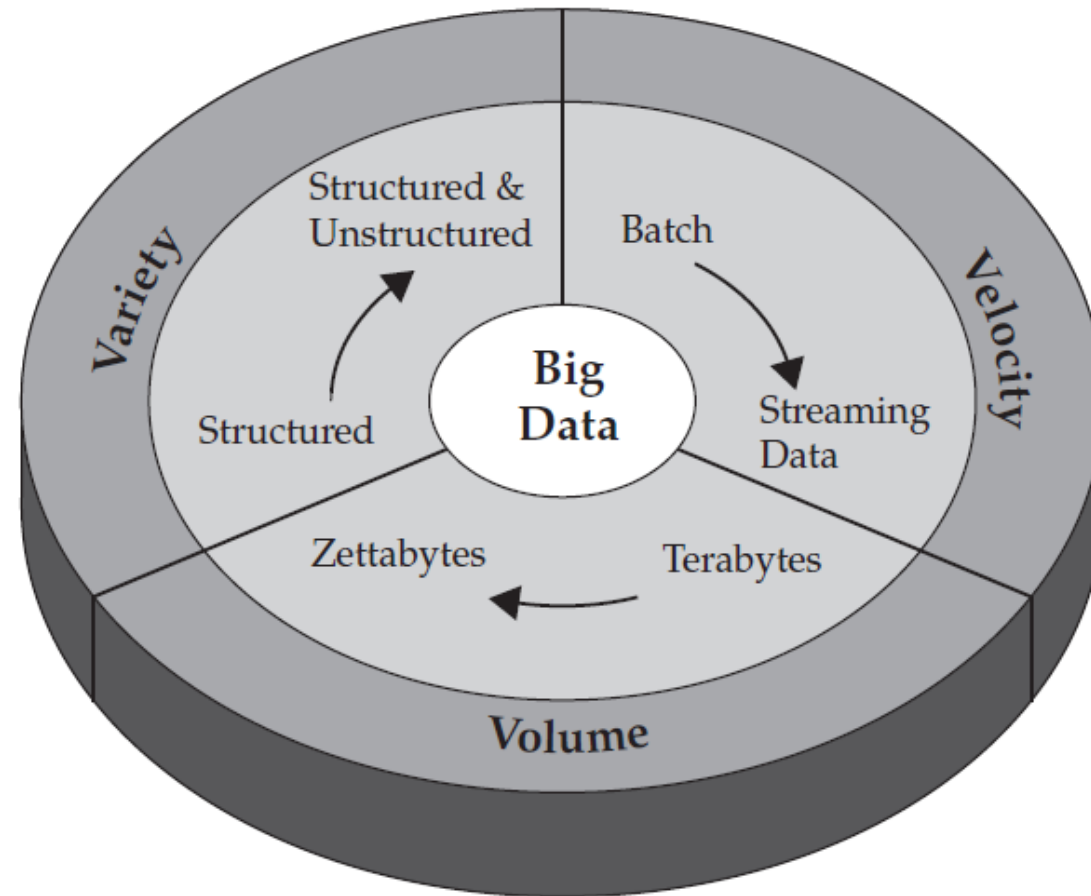
of the world's data is analyzed
today

Big Data: what are they?



✳ It is not just about Volume!

Big Data



Volume

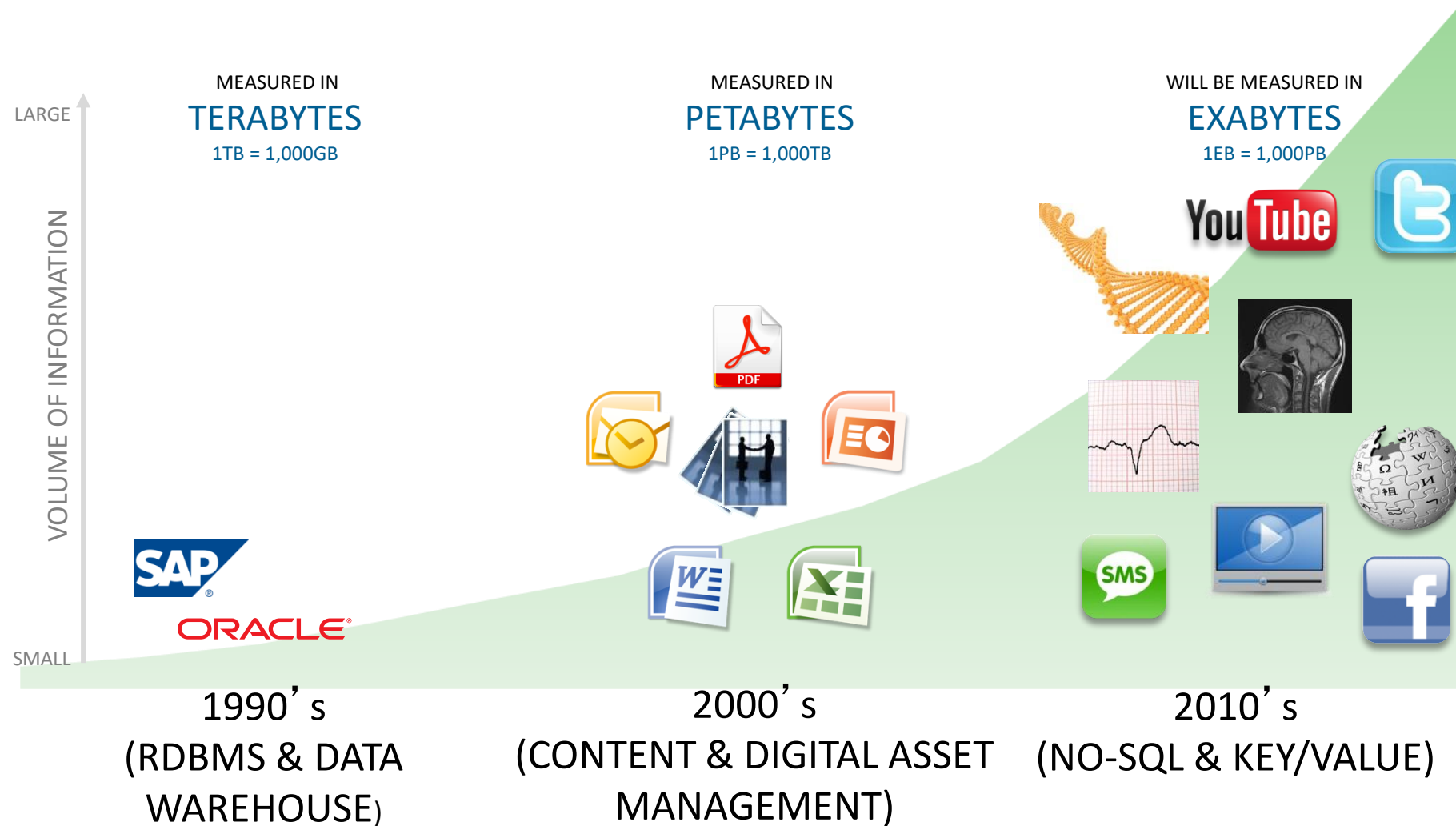
- Data set to be processed at a time is too large
- Data set is not too large but the collection of data set is large
- Volume of data set too large per se, but processing is time consuming perhaps due to too many IO operations



Present of Big Data

Too big to handle

New Applications Driving Data Volume



Velocity

- Data arrive is faster than the processing capacity
- Results must be produced with certain delay bound, processing is limited by disk I/O throughput



Future of Big Data

Drinking from a firehose

Variety

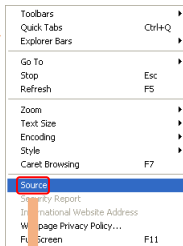
Structured Data

SUMMER FOOD SERVICE PROGRAM 1]				
(Data as of August 01, 2011)				
Fiscal Year	Number of Sites	Peak (July) Participation	Meals Served	Total Federal Expenditures 2]
	-----Thousands-----		--Mil.--	--Million \$--
1969	1.2	99	2.2	0.3
1970	1.9	227	8.2	1.8
1971	3.2	569	29.0	8.2
1972	6.5	1,080	73.5	21.9
1973	11.2	1,437	65.4	26.6
1974	10.6	1,403	63.6	33.6
1975	12.0	1,785	84.3	50.3
1976	16.0	2,453	104.8	73.4
TQ 3]	22.4	3,455	198.0	88.9
1977	23.7	2,791	170.4	114.4
1978	22.4	2,333	120.3	100.3
1979	23.0	2,126	121.8	108.6
1980	21.6	1,922	108.2	110.1

Semi-Structured Data



View → Source

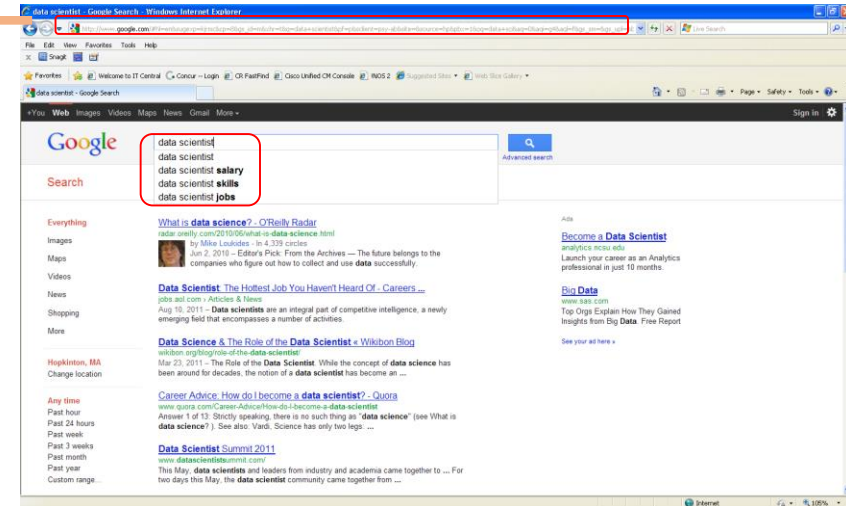


```

1 <!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN" "http://www.w3.org/TR/xhtml1/DTD/xhtml1-trans
2 <html xmlns="http://www.w3.org/1999/xhtml">
3
4 <head>
5
6 <meta http-equiv="Content-Type" content="text/html; charset=UTF-8" />
7 <META name="y_key" content="859b402e1c9a6c">
8 <link rel="canonical" href="http://www.emc.com/index.htm" />
9 <META NAME="verify-v1" CONTENT="yi2c9V0P4eV0jFdiPeVvIFRf32g4qWFE0I2UvTMfSU
10 <title>EMC - Data Recovery, Cloud Computing, and Storage Hardware</title>
11 <META NAME="description" CONTENT="EMC is a leading provider of storage hardware solutions th
12 data recovery and improve cloud computing." />
13 <META NAME="keywords" CONTENT="emc, network storage, data recovery, information manage
14 software, nas storage, information protection, information management" />
15 <!-- Start: stylesheet includes -->
16 <link rel="stylesheet" href="/_admin/css/styles.css" />
17 <link rel="stylesheet" href="/_admin/css/styles_nav.css" />
18 <!--if IE>

```

Quasi-Structured Data

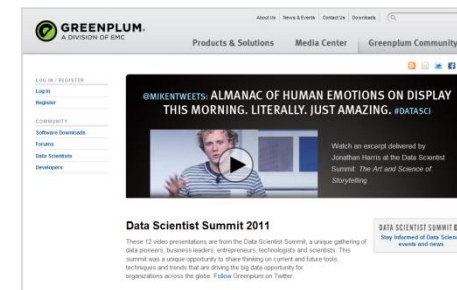


http://www.google.com/#hl=en&sugexp=kjrmc&cp=8&gs_id=2m&xhr=t&q=data+scientist&pq=big+data&pf=p&scient=psyb&source=hp&pbx=1&oq=data+sci&aq=0&aqi=g4&aql=f&gs_sm=&gs_upl=&bav=on.2,or.r_gc_r_pw,cf.osb&fp=d566e0fbd09c8604&biw=1382&bih=651

Unstructured Data

The Red Wheelbarrow, by William Carlos Williams

so much depends
upon
a red wheel
barrow
glazed with rain
water
beside the white
chickens.

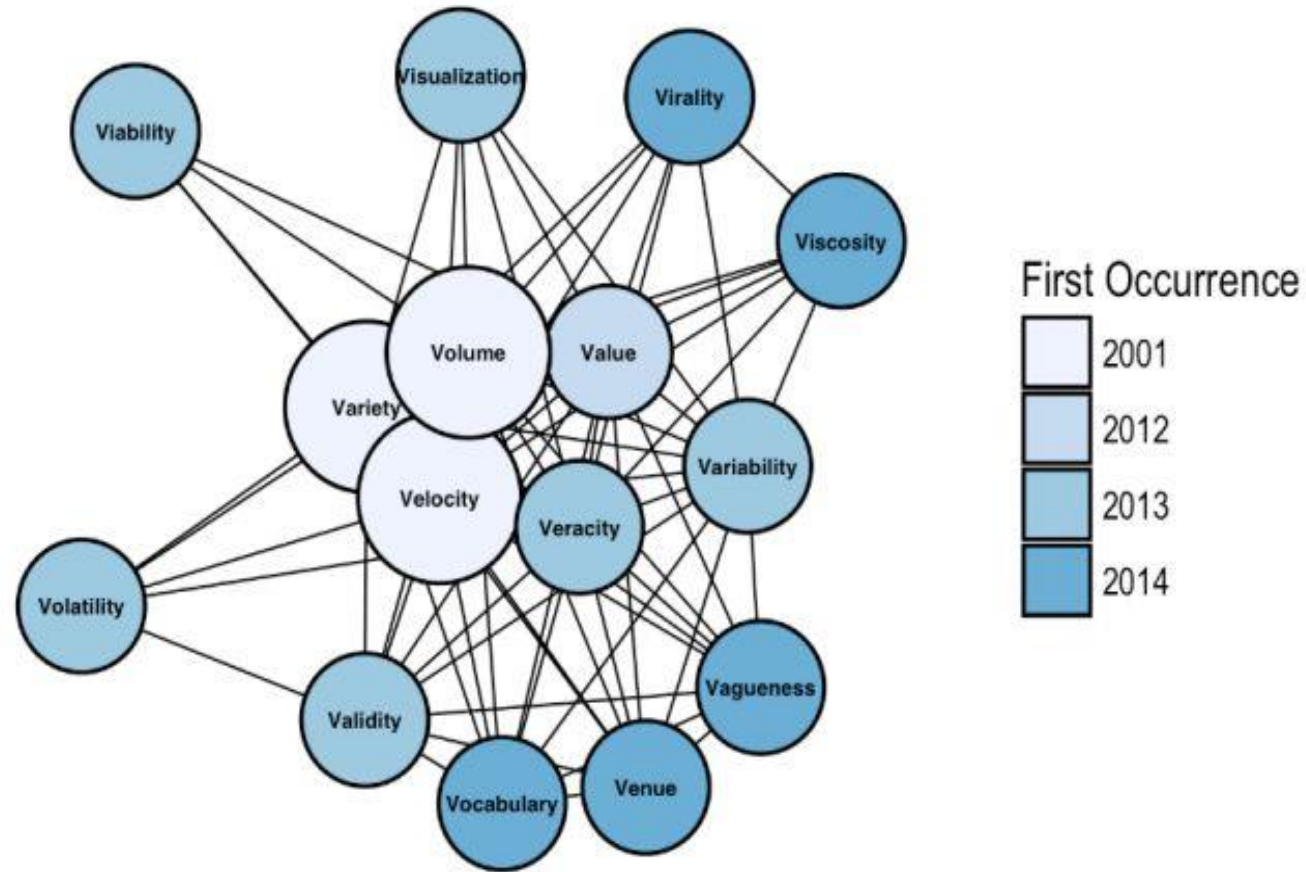


Other V's



- Valence: Non-trivial inter-relatedness of data
- Veracity: The degree of certainty in data
- Variability: variable interpretations

More and more Vs



An example

You are What you Eat and Drink

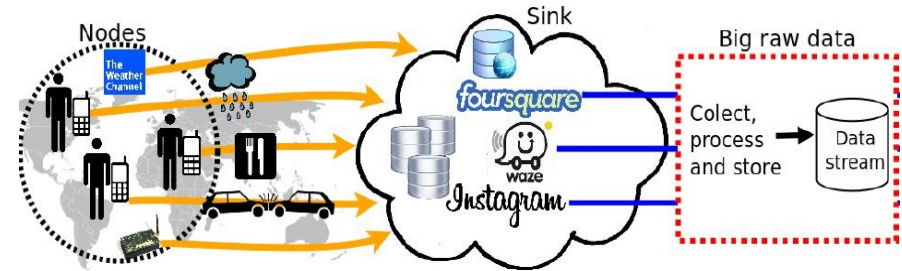


- Food and drink became also a strong cultural aspect, being able to describe strong differences
- Foursquare, created in 2009, registered 5 million users in December 2010 and 45 million users in January 2014



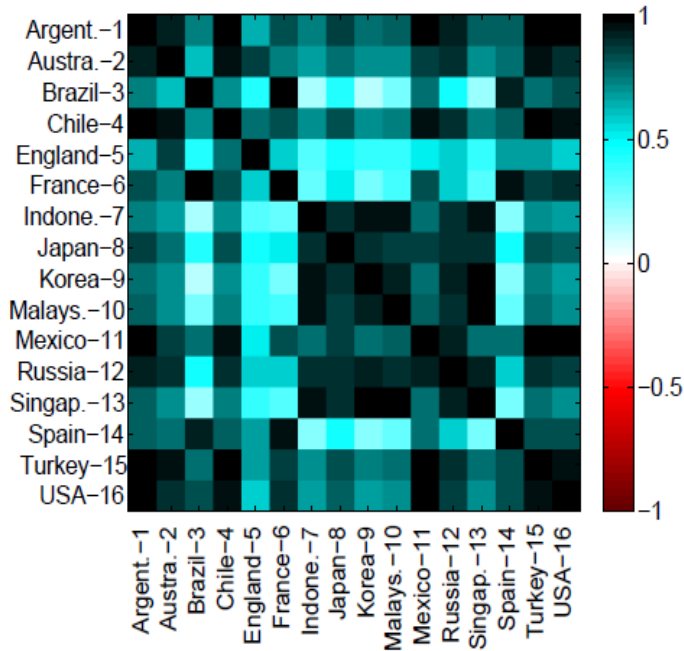
- Possibility to sense human activities related to food and drink practices in large geographical areas
- Delineate and describe regions that have common cultural elements, defining signatures that represent cultural differences between distinct areas around the planet

You are What you Eat and Drink

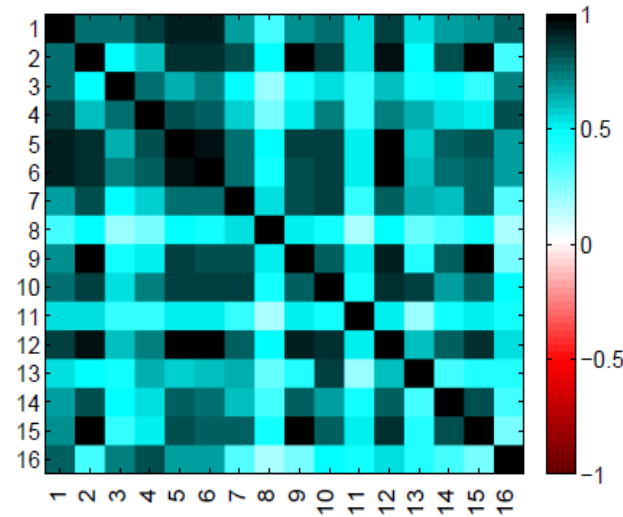


- 4.7 million tweets containing check-ins were gathered, each one providing a URL to the Foursquare website (one-week dataset same order of magnitude of the number of interviews performed in World Values Survey in almost three decades)
- Locationbased social networks (LBSNs)
- Location identified by free reverse geocoding API offered by Yahoo (<http://developer.yahoo.com>)

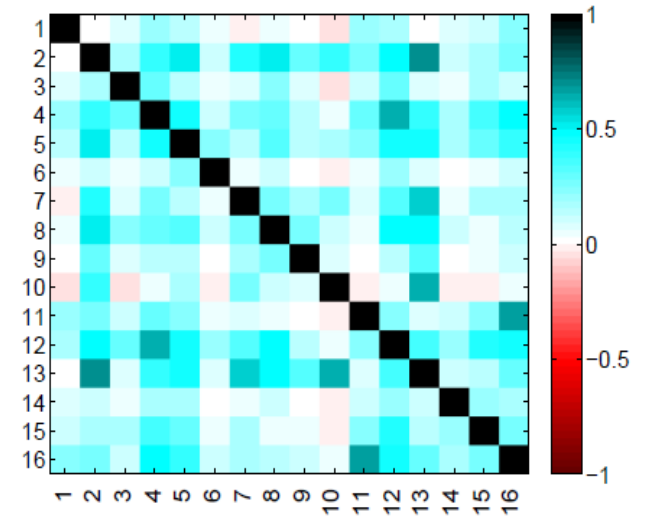
You are What you Eat and Drink



(a) Drink

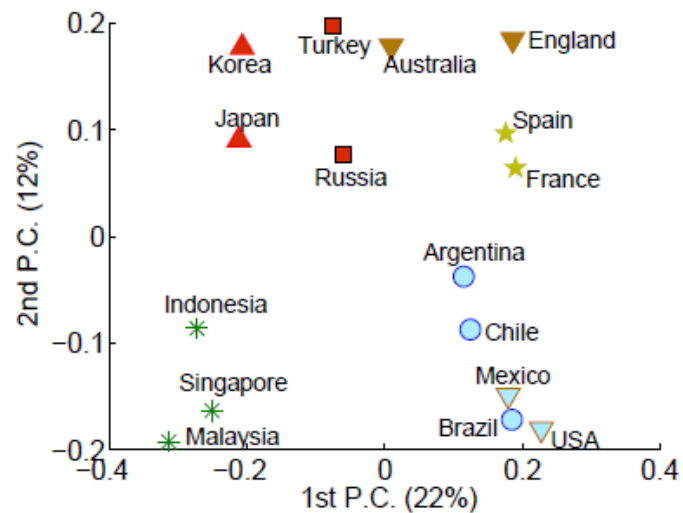


(b) Fast Food

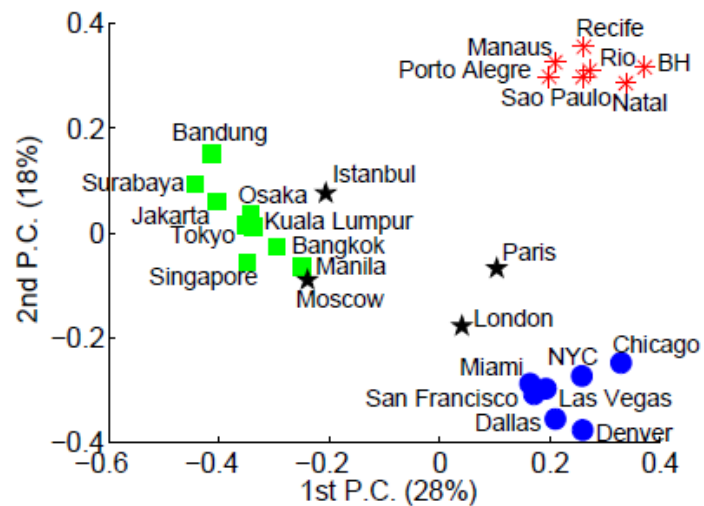


(c) Slow Food

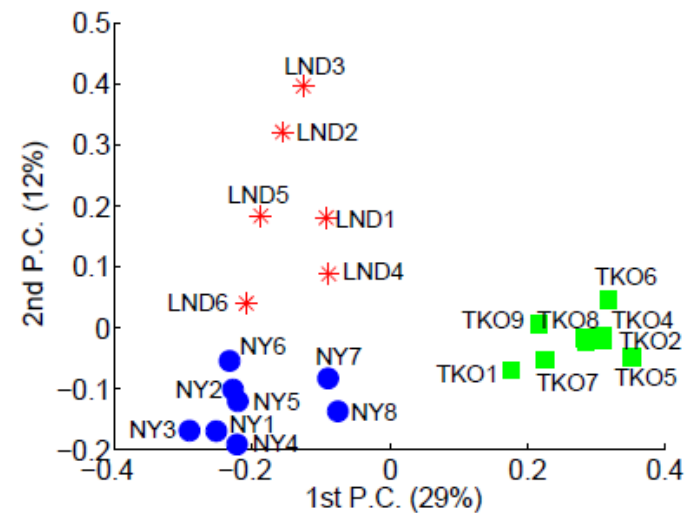
Figure 3: Correlation of preferences between countries.



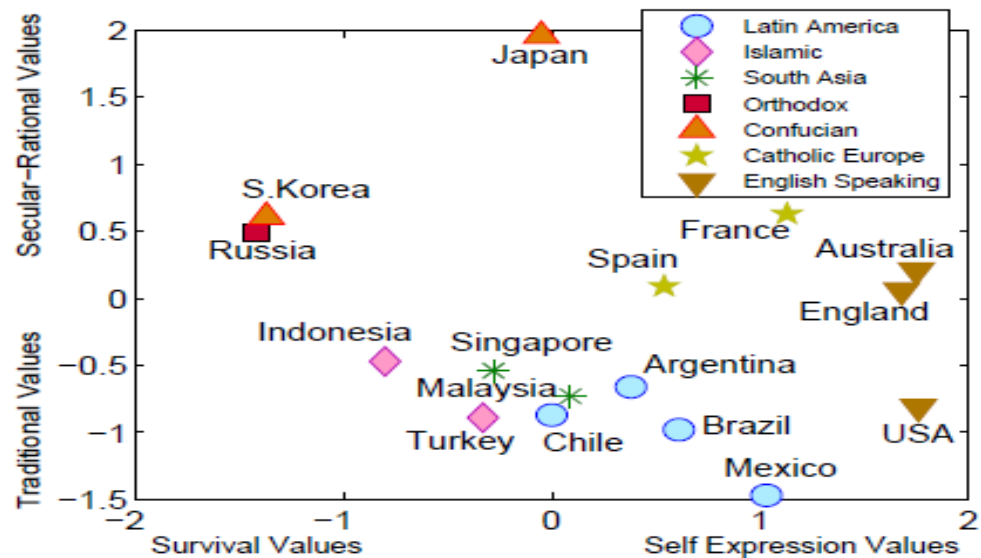
(a) Countries



(b) Cities



(c) Regions

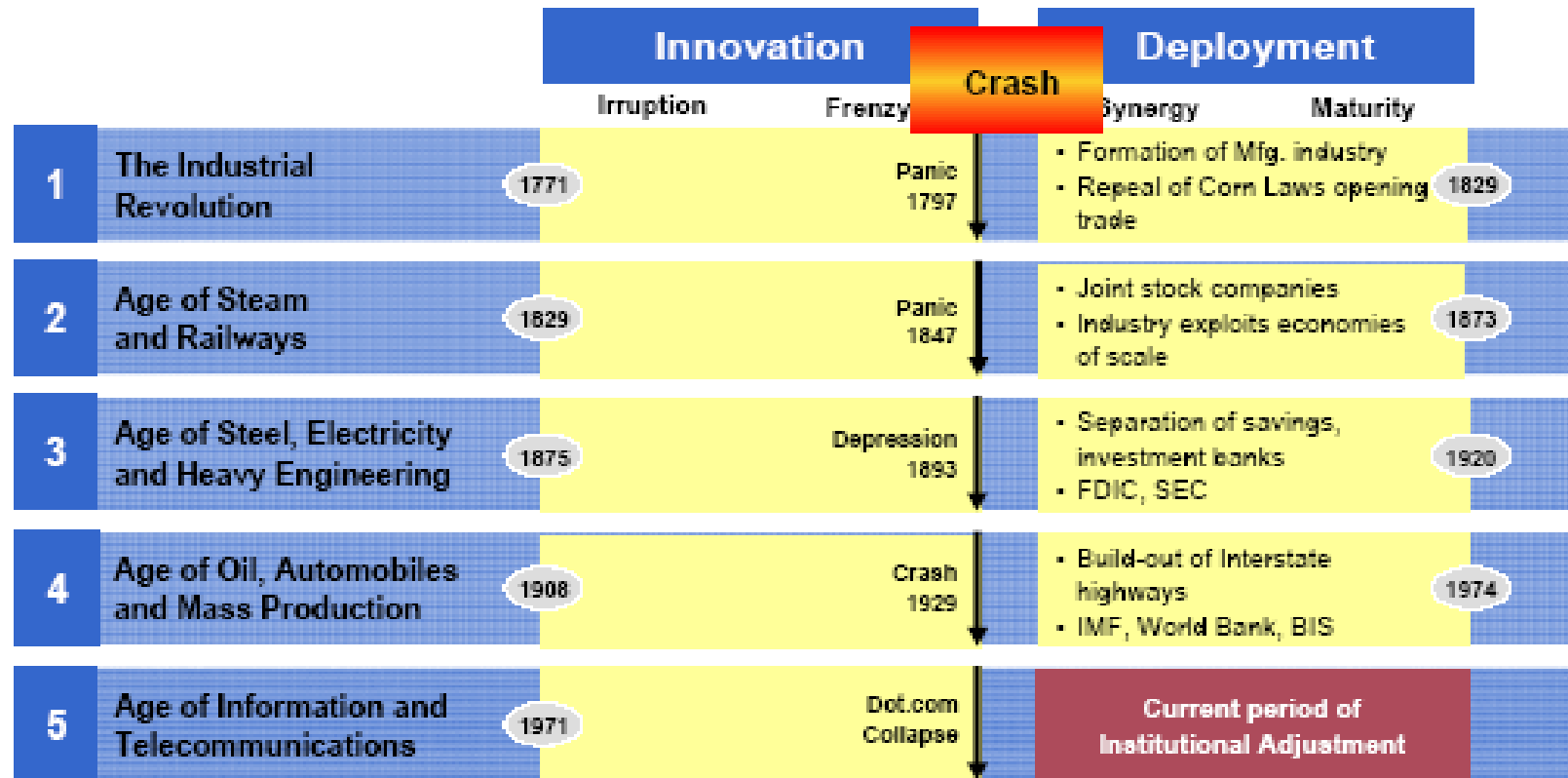


Why Big Data?



Revolutions needs Innovation...

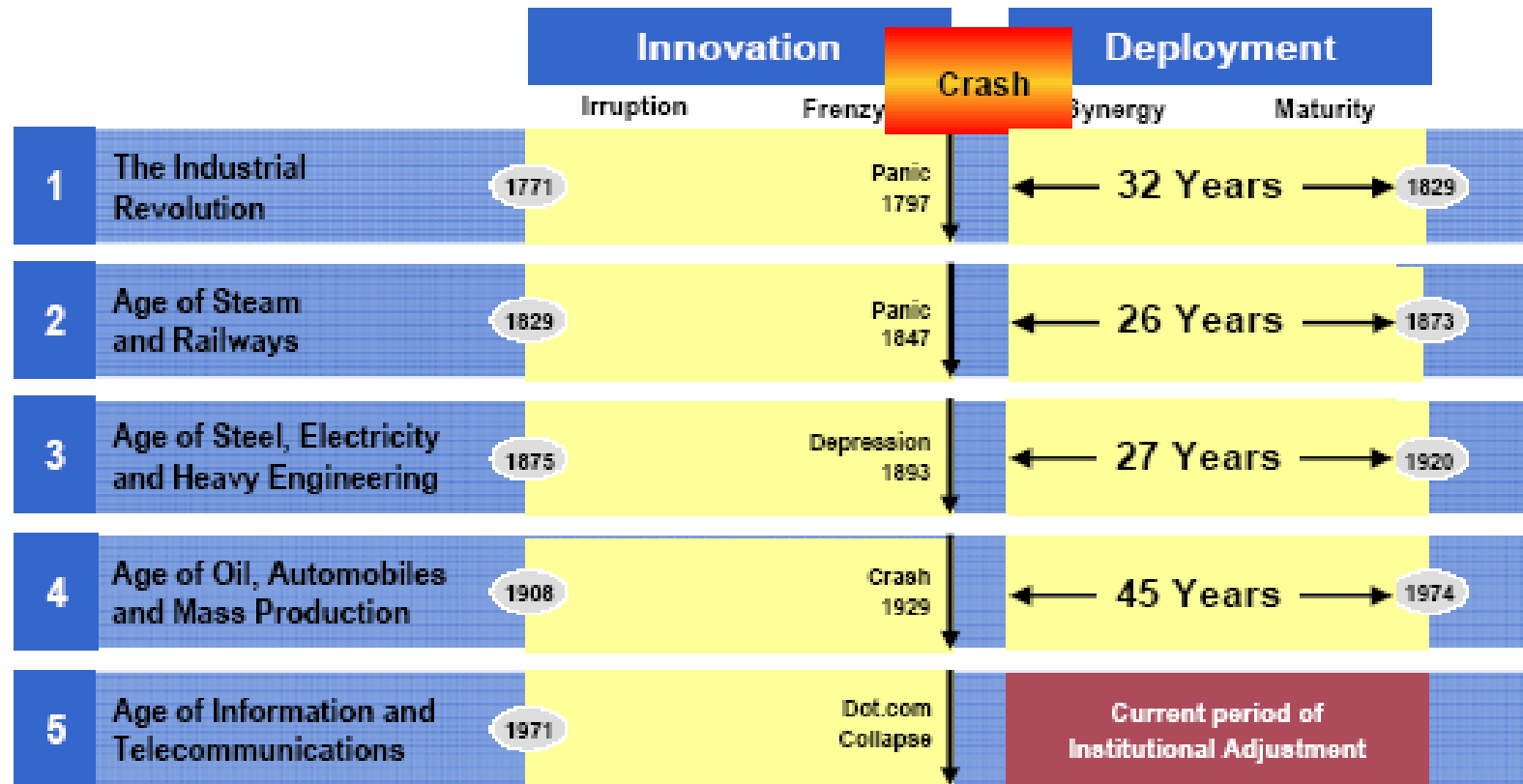
Five historical cycles ...



Source: "Technological Revolutions and Financial Capital, Carlota Perez, 2002

Takes time to deploy...

The deployment phase lasts 26 to 45 years ...

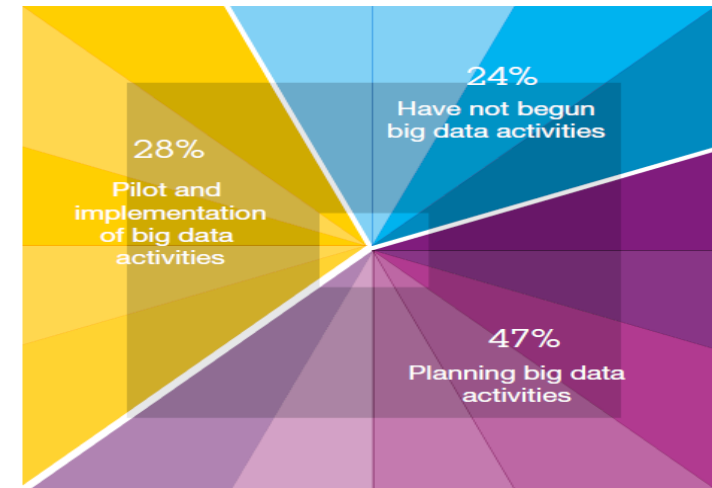
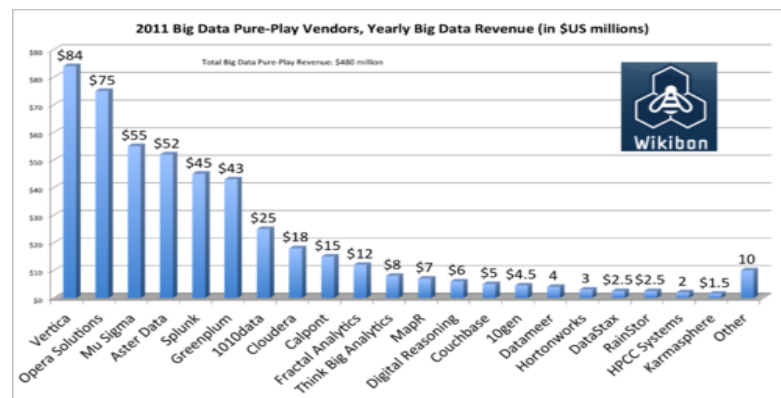
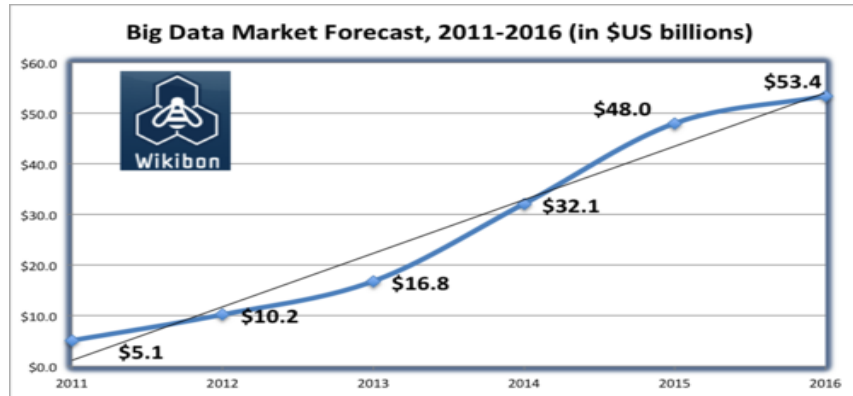


Source: "Technological Revolutions and Financial Capital, Carlota Perez, 2002

Big Data and Enterprise



Big Data and Enterprise



http://wikibon.org/wiki/v/Big_Data_Market_Size_and_Vendor_Revenues

Analytics: The real-world use of big data: How innovative enterprises extract value from uncertain data, Executive REport, IBM Institute for Business Value

THE BIG DATA LANDSCAPE

JANUARY 2014

Apps

Vertical

KNEWTON ellucian practice fusion
SurveyMonkey PROACTIVE PUBLISHING

Operational Intelligence

splunk New Relic PopDynamics
sumologic VITRIA

Data As A Service

factust FICO Kaggle INRIX
GNIP Opigive Placed LOGATE
DATASET TOPSY LexisNexis

Consumer

Google amazon
Walmart Labs ebay BETAFUN in

Ad/Media

METAMARKETS collectiveIQ
rocketHub FLURBY
collective DataXu
Media Science TURN
bloomreach

Business Intelligence

ORACLE | Hyperion
SAP Business Objects | IBM
Microsoft | Business Intelligence
IBM | SAP | BI
pentaho RECOMMIND
Autonomy | blime | GoodData
Chart.io | WATTIVO | Recorded Future

Analytics and Visualization

tableau QlikView
Palantir QUBA TRIFACTA
TERADATA ASTER Knoema
sas TIBCO centrifuge
panopticon UFDRA
Datameer mindata GIND
platforma ClearStory
alteryx visual.ly AVATA
metaLayer Adigo Alpine STScore

Infrastructure

Analytics

cloudera Hortonworks
MAER Gobon HARAPT
Pivotal INFOBRIGHT
NETEZZA VERTICA
ORACLE SQL kognitio

Operational

COUCHBASE mongoDB
AEROSPIKE splice
DATASTRX VoltDB
TERRACOTTA INFORMATICA
MarkLogic

As A Service

Qubole amazon web services
Windows Azure MORTAR
CSC Google BigQuery

Structured DB

ORACLE MySQL
SQL Server PostgreSQL
IBM DB2 SYBASE
memsql TERADATA



Technologies



HBASE

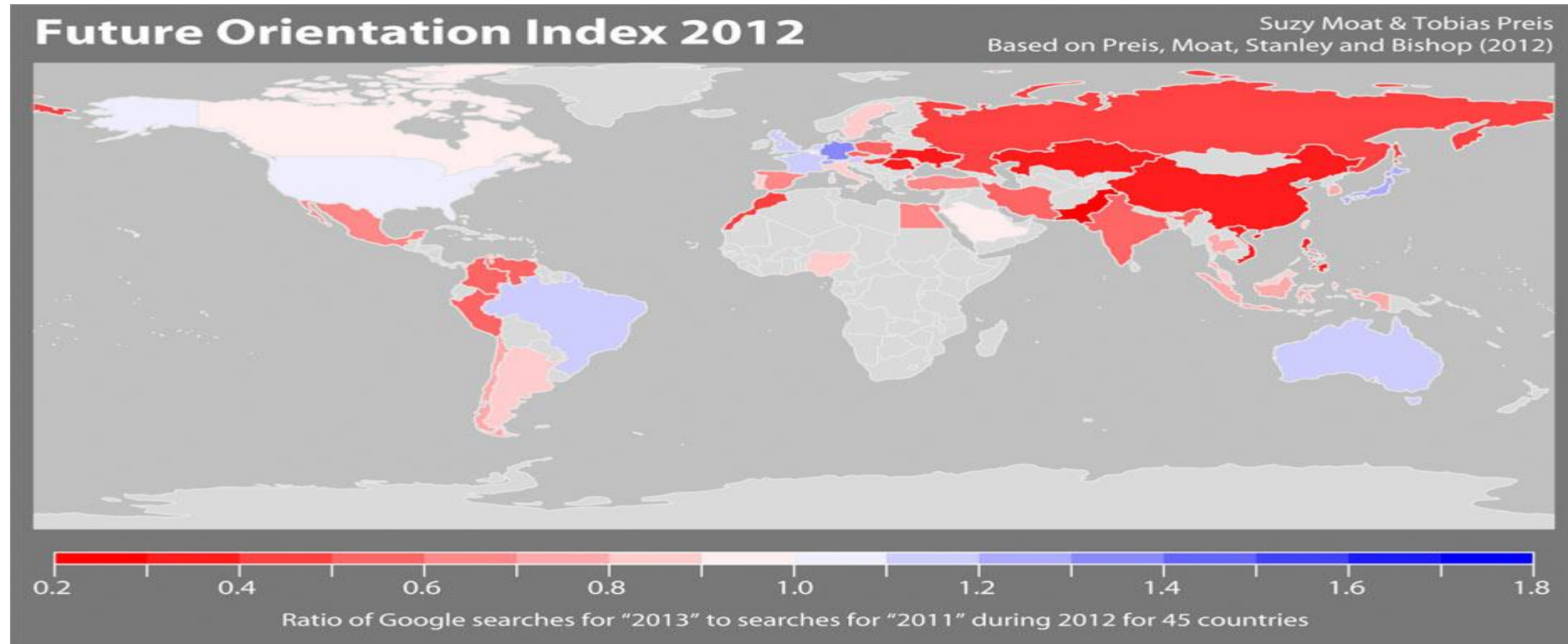


Big Data and the Government



“Can you explain all the emails you’ve received from Russia and Iran?”

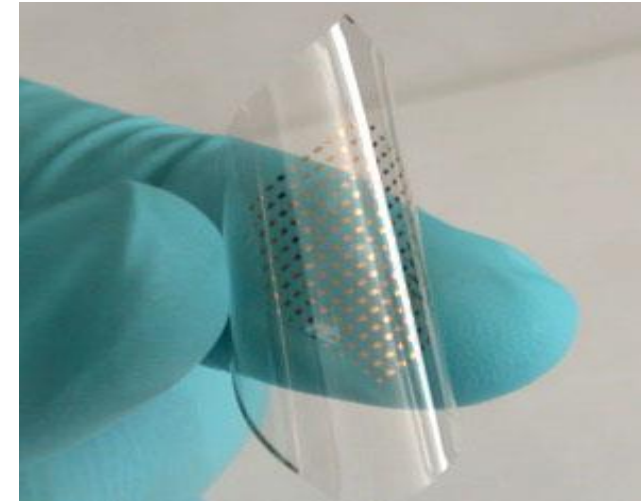
Big Data and Economy



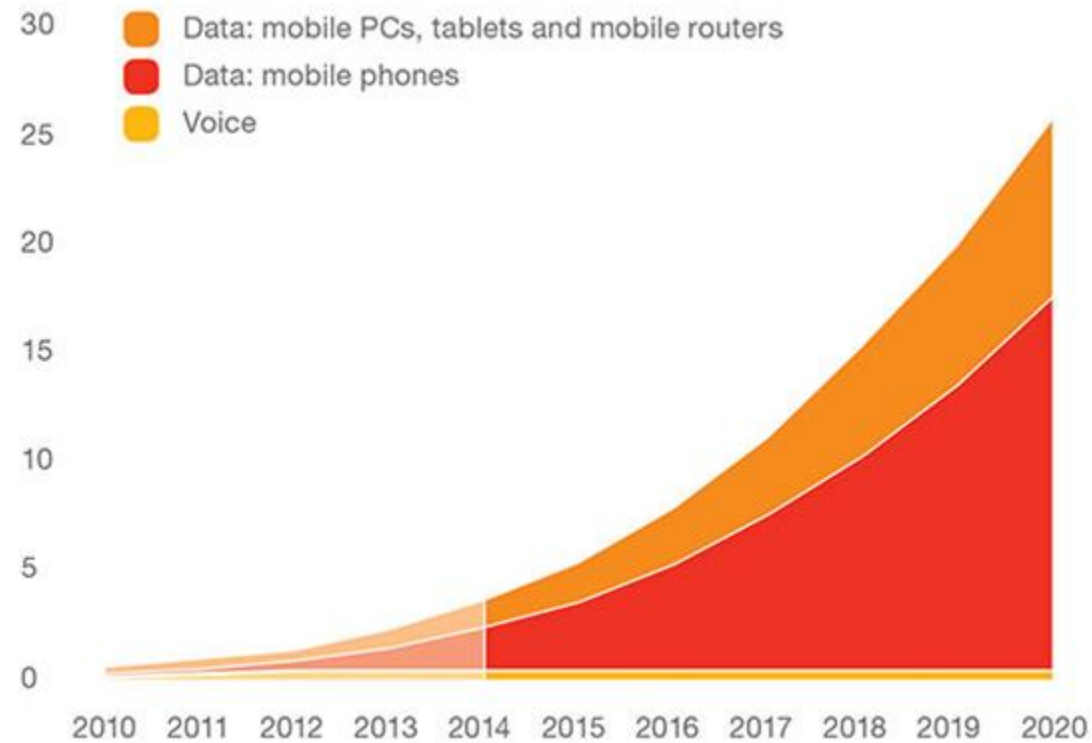
"We see two leading explanations for this relationship between search activity and GDP. Firstly, these findings may reflect international differences in attention to the future and the past, where a focus on the future supports economic success. Secondly, these findings may reflect international differences in the type of information sought online, perhaps due to economic influences on available Internet infrastructure."

What are the sources of data ?

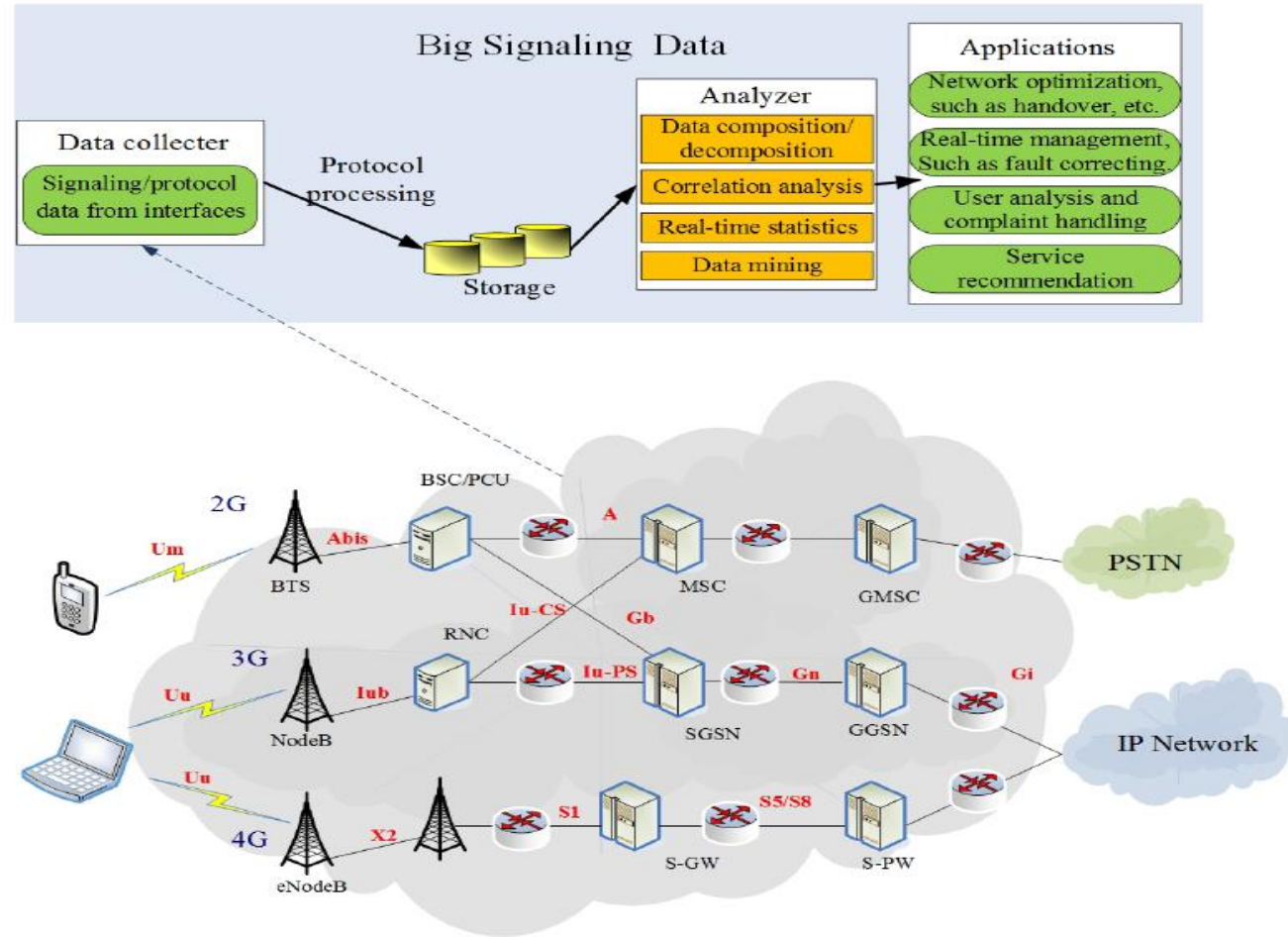




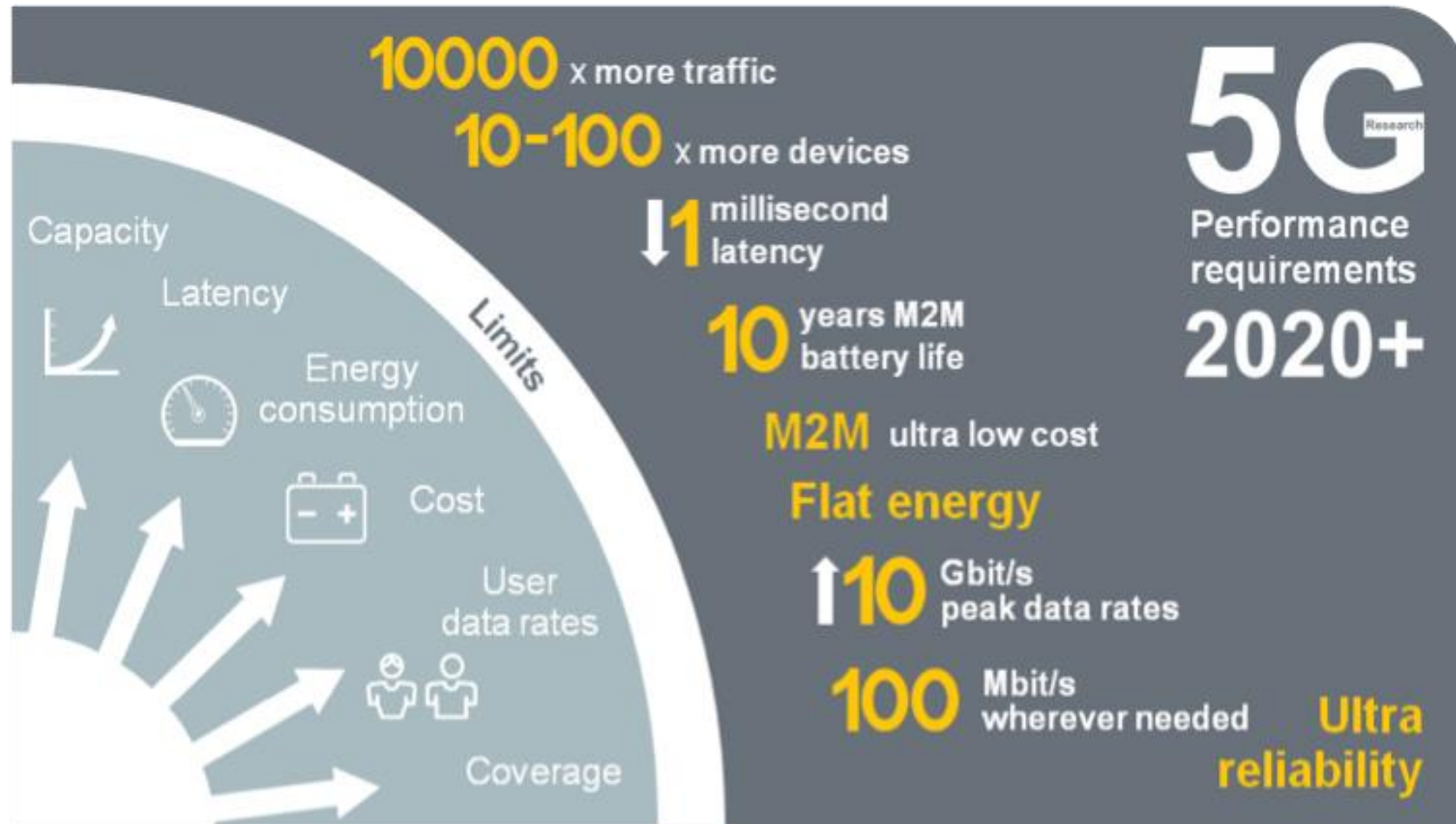
Mobile Traffic Growth (in ExaBytes)



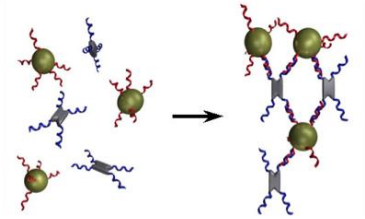
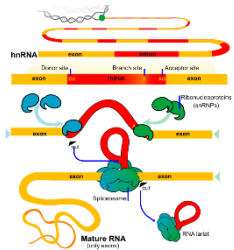
Mobile Traffic



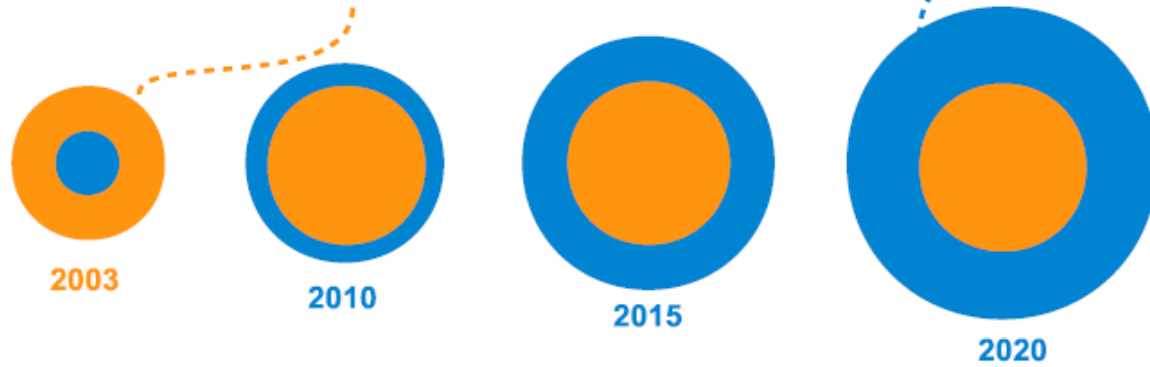
5G Requirements



Sensors



During 2008, the number of things connected to the internet exceed the number of **people** on earth



Internet of Things (IoT)

Slides provided by Dzmitry Kliazovich, University of Luxembourg

What is IoT ?

"The Internet of Things is the network of physical objects that contain embedded technology to communicate and sense or interact with their internal states or the external environment." **Gartner**



Sensors



Linker Intel Group



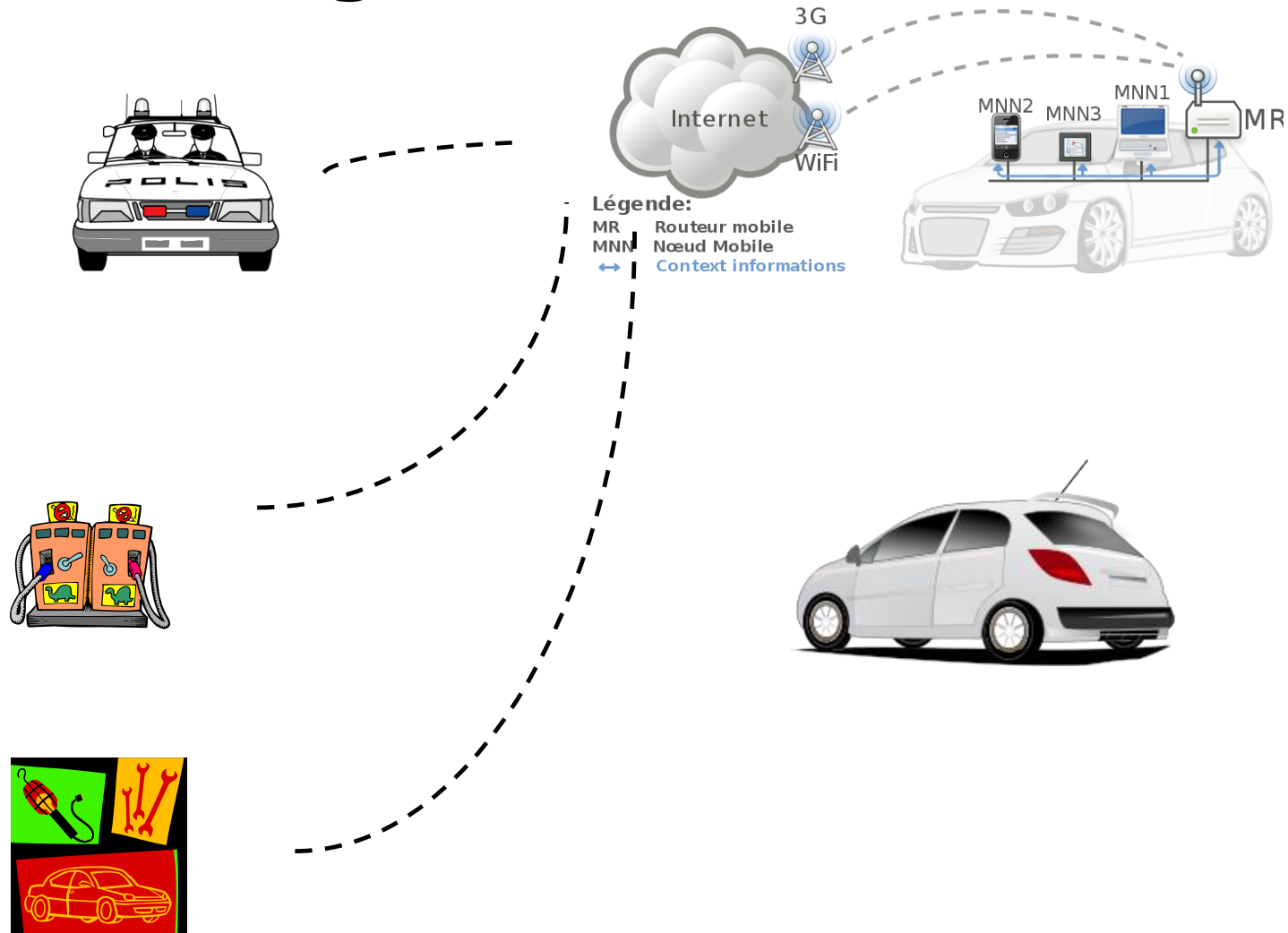
Image Sensor Device



Smart Devices



Things to Things



Human to Things

ECG sensor



Motion sensor

Motion sensor



IoT Applications



Wireless Farming
Sensors: Phytech



Factory Wearable:
ProGlove



Connected Commuter
Bike: Valour

1,249.00



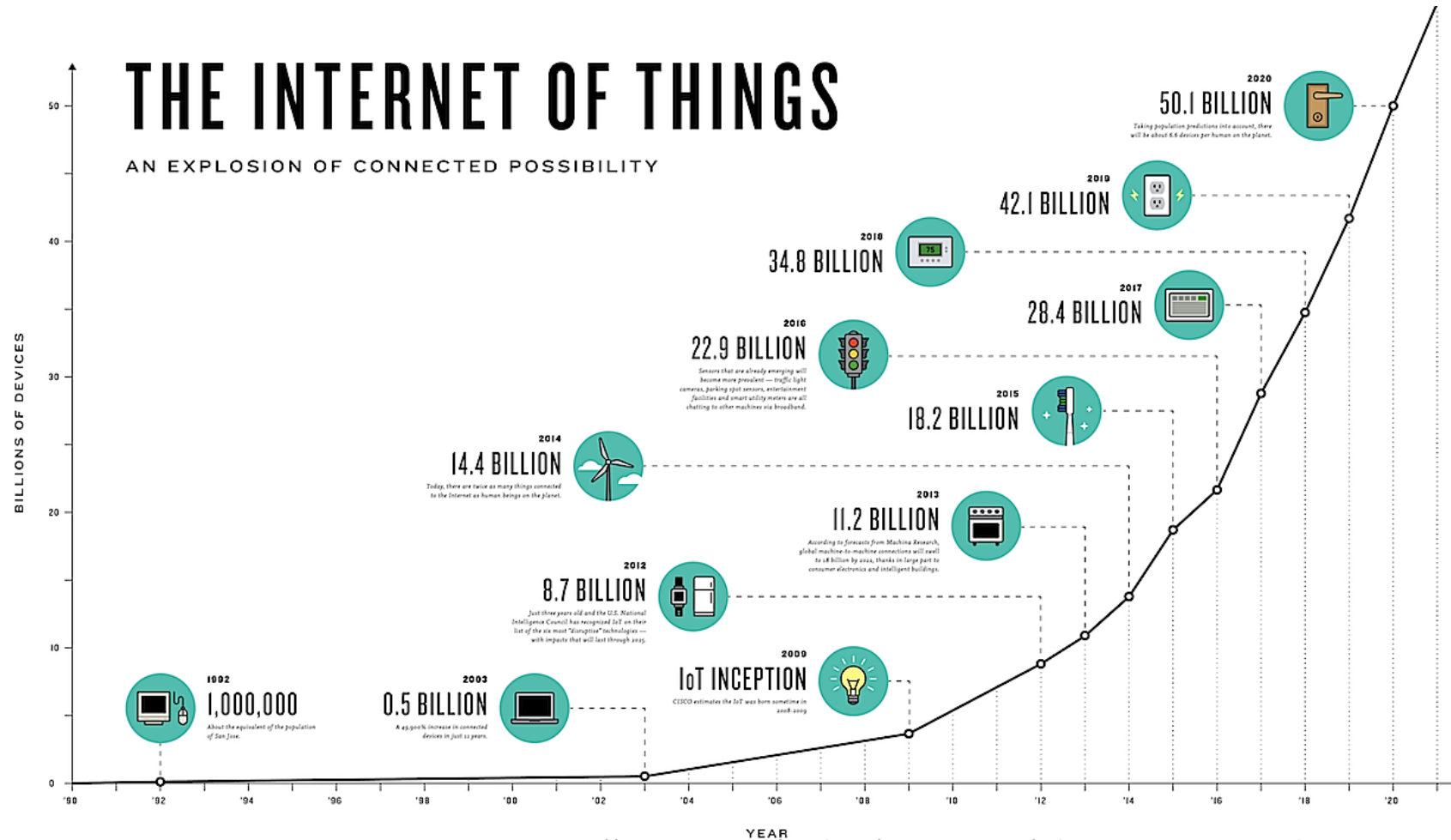
Global Location Device:
iTraq

49.00



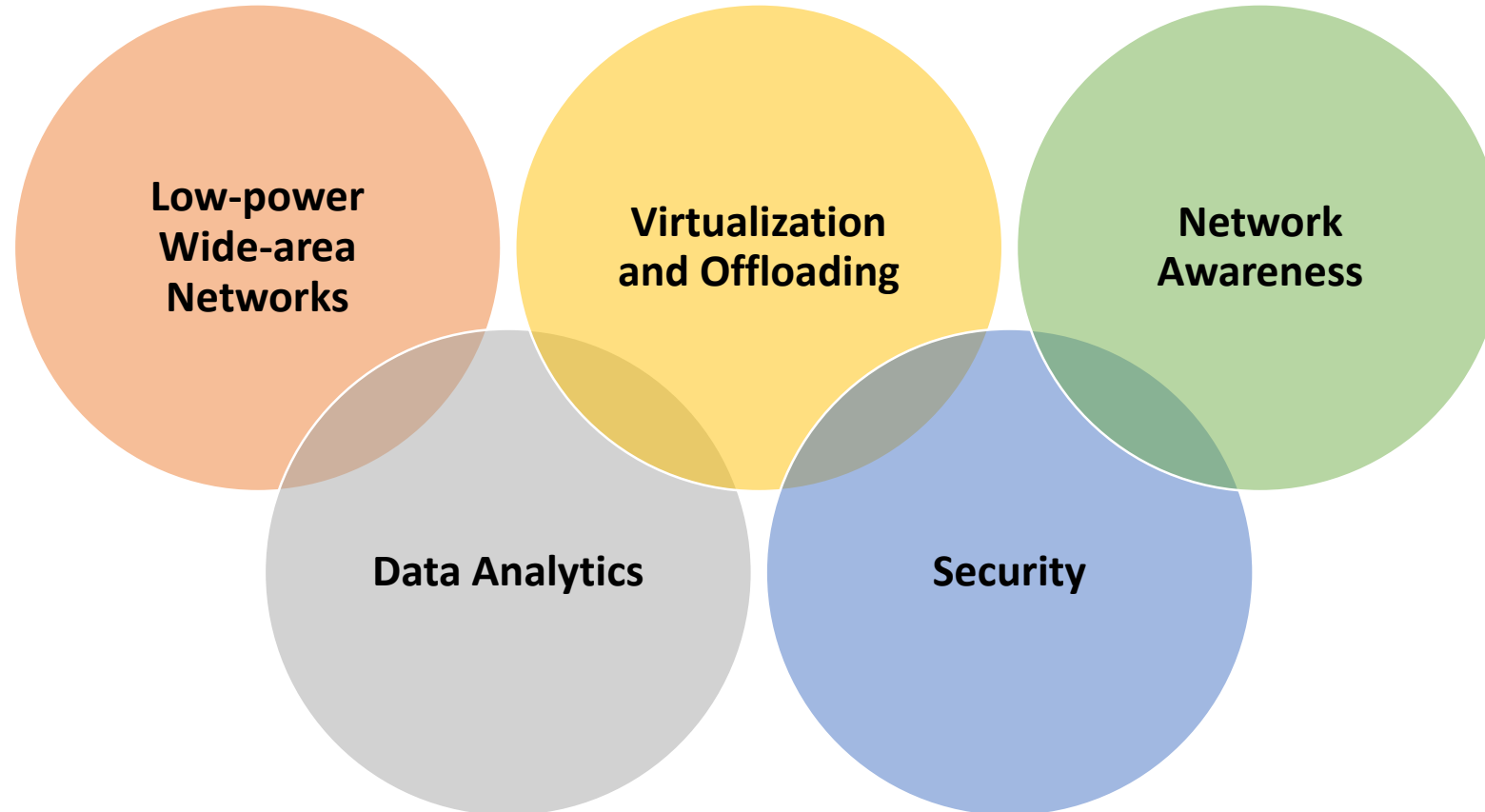
Reemo addresses
mobility challenges with
gesture-based smart
home control

Connected Devices

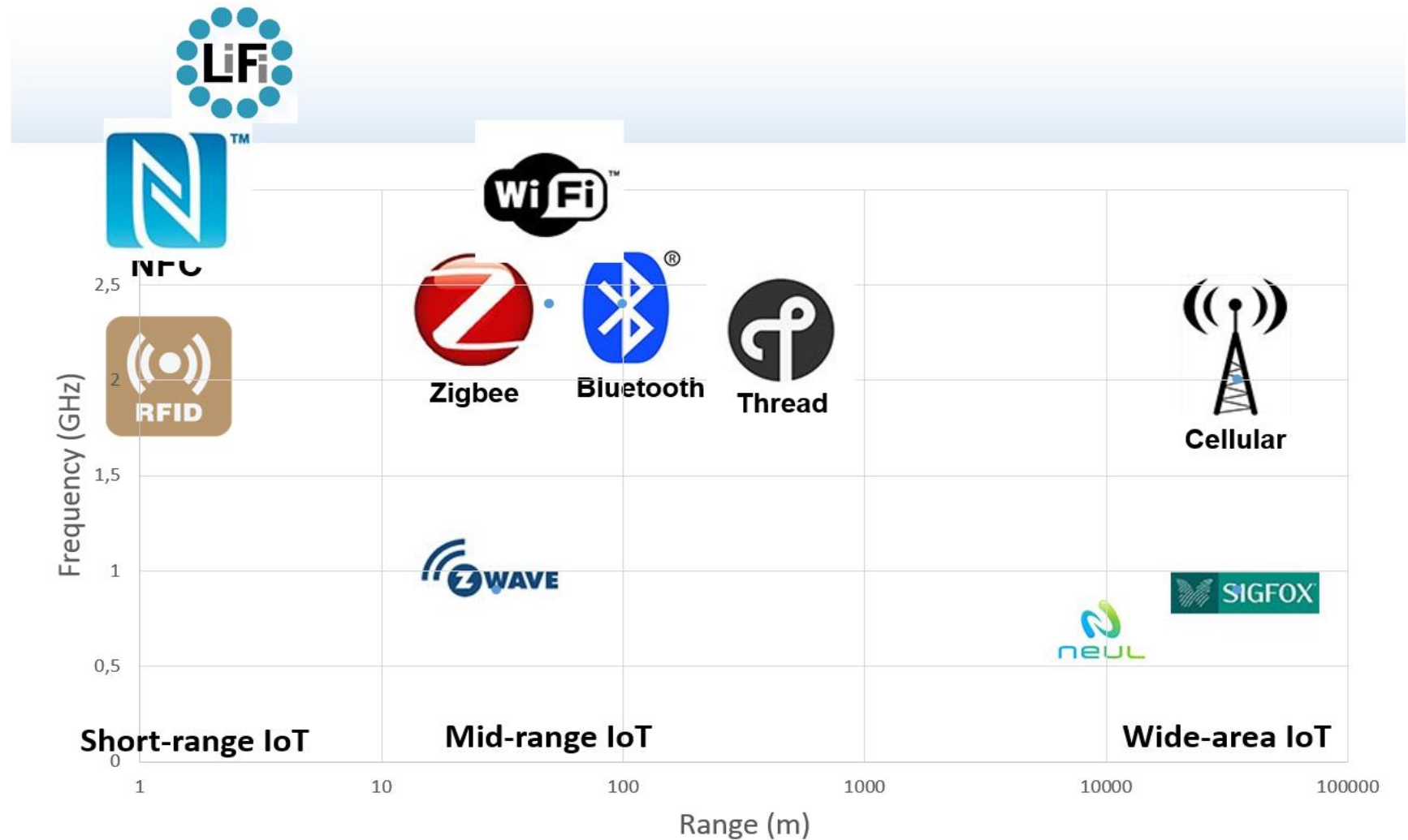


<https://www.ncta.com/platform/industry-news/infographic-the-growth-of-the-internet-of-things/>

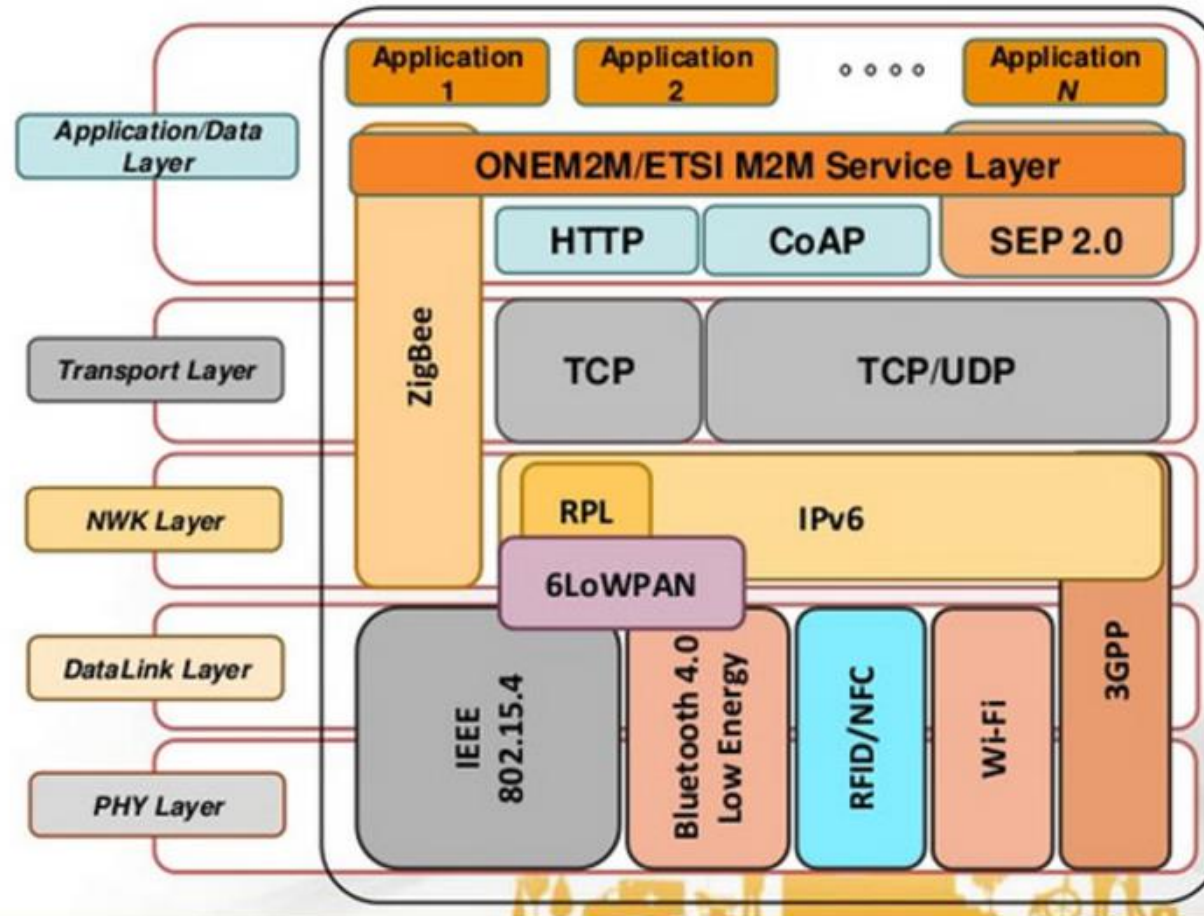
Key Enabling Technologies for IoT



Communication Technologies



Communication Stack



IoT Key Enablers

• IoT Devices Data Collection

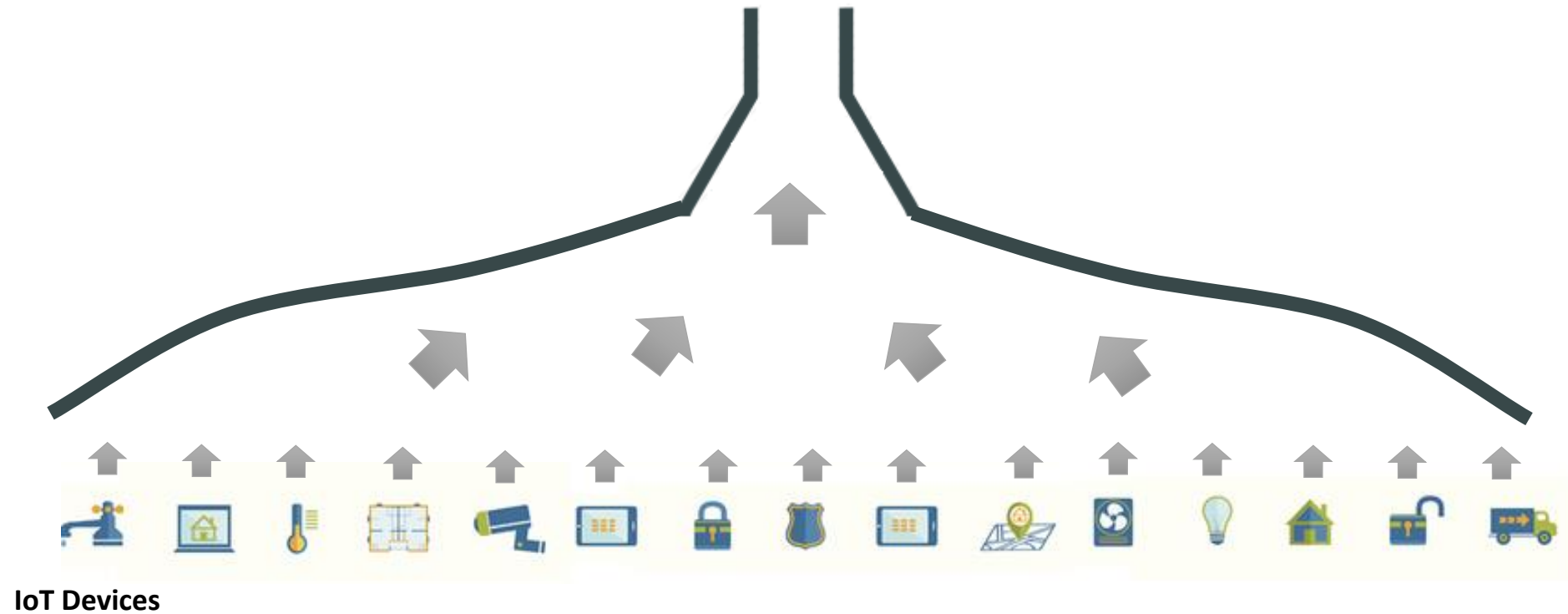
- Daily (Sensors, Meters)
- Hourly (Home appliances)
- Every minute (control)
- Real time (transactions, teleoperation)



IoT Devices

IoT Key Enablers

- Data Fusion and Analytics



IoT Platform

- Openness
- Flexible lightweight virtualization techniques
 - to provide access to shared data, computing, storage and networking resources of IoT infrastructure
 - to augment functionality of the resource-constrained IoT devices and enable new applications
- Effective data collection and management
- Dynamic partitioning of IoT applications in real-time
 - to save energy
 - to augment IoT applications
- Awareness of network dynamics and events
- Security and trust

Int

Applications (Verticals)

Personal Devices Wearable Computing: pebble, cookoo, rucon, iM, IDEYE, GLASS, RINGLY, striiv, APX, MOTÁ, M, wearable intelligence, VUZEX Fitness: GARMIN, Nike, senSonic, amigo, iFIT, JAWBONE, HISFIT, VALEN, BASIS, fitbit, tomtom, ARA, LifeBEAM, wahoo Health: LUMO, manowear, vcsay, kinsq, HAPIfork, soundshare, SPIRE, Withings, QUANTUS, Lively, SCANDIA, RENK, ΔMCTO, Kuanatomy, hello, AliveCor, iHealth, HURON, remee, acion, corventis, TELGARE Family: FILIP, Sproouting, ovuline, AmberAlert, greatcall, Oovoo, Night Lamp, monobaby, OWLET, Secur, BELLABEAT, mimo, Glow, pocketfinder	Lifestyle Sports: SOLO, Brain Sentry, BIKESPIKE, InfoMotion, swingbyte, ZEPP, HAMMERHEAD Cooking: Smart Diet Scale, ANOVA, drop, blossom, iDevices, THE ORANGE CHEF CO., pantry, nomiku Pets: Whistle, PetPace, pintofeod, PetHub, tagg, BISTRO, hoytag, PetBark, PetChatz, Petcube, tractive, Petnet, gibi, Petziko Toys: KAROTZ, UBOOLY, MAKIES, atoms, UBOOLY, seeo Music/Art/Video: ROLI, CATCH, GOPIPO, narrative Garden: plantlink, BITPONICS, radio, EDYN, Greenbox, Koubachi	Connected Home Automation: Quirky, Radiator Labs, netatmo, LEVITON, SmartThings, Ubi, nest, LIFX, gecko, CRESTRON, smarhome, LUTRON, ecobee, Advanced.americs, vivint, SAVANT, INSTEON, CHAMBERLAIN, PROPER, LIGHT, somfy Monitoring: lapka, sense, birdi, BlueMaestro, SUPERMECHANICAL, leeo, knot, CUBESENSORS, tado°, ambient Security: HomeMonitor, canary, ring, dropcam, butterflye, Locktron, Fugust, SCHLAGE, RAVEN, Kwikset, globasense, genie, UNKEY, GOJI, scout, nSmartAlarm, keyfree, SmartCode Tracker: Chipolo, linquet, locca!, TrackR Hub: Homey, revolv, NINJABLOCKS, Control, LOWE'S, zonoff, STAPLES, NEXIA, muzzley, wink	Industries Retail: leapt, PROMOTY, SVILELO, malahan, GIMBAL, boni, PassKA, PERCH, P, stream, stream Payment: Square, shopify, PayPal, ACS, VerFone, LevelUp, belly, payleven, con Healthcare: VISI, Senseonics, STANLEY, AERGENX, VITALITY, MedMinder, GoleTracking, MedSignale, AdhereTech, VERSUS, CENTRAK, intelligent, Sotera, RayGen Automotive: Zobe, fitse, INRIX, nadyDELPHI, Airbiquity, dash, waze, OpenXC Infra-structure: wavelink, iBeam, kisi, Johnson Controls, Trimble, Robin, Schneider Electric, adapt-N, Ag Leader Agri-culture: adapt-N, Ag Leader	Industrial Internet Robotics: Double Robotics, ALDEBARAN, ROBOTIX, EMPIRE, iRobot, KUKA, ABB, LIQUID ROBOTICS, JIBO Drones/Aerospace: SDR, KMO, CYPRIS, Airware, SKYCATCH, spire, Parrot, Skybotix Green-tech: BigBelly, enlightened, Smart, Genevo, compology, AMPY 3D Scary/Print: MakerBot, Stratasys, HREL 3D, formlabs, M3D, matterport, FUEL3D, AIO, NEXTENGINE, RepRap, occipital, DAVID, FSLD, Solidoodle Smart Grid: GRIDNET, e-on, SilverSpring, SMART GRID CLUSTER, Itron, Trilliant Asset Tracking: asip, vloc, NESHSYSTEMS, PHOTON, Colson, ekahau, AccuLync, Pictomatic, Inspire, CUBIC, mactrac, Agency
--	---	--	---	---

Platforms & Enablement (Horizontals)

Connectivity/Dev Platforms spark, kynetx, pinoccio, ioBridge, Ayla Networks, EUROTECH, resin.io, Symplic, TESSEL, bluecity	Software/Data Platforms EXOSITE, iconcontrol, thingsquare, carriers, Keen IO, SeeControl, Uthings, ConnectHO, NewAer, BERG, Axeda, Yaler.net, RacoWireless, SpaceGen, IFTTT, greenWAVE, ARABERAT, swot.io, ZALAR, CyberLighting, altiux, Yo, ThingWorx, DN2K, Xocio, power, IOTER, thingful, CANDI, bugswarm, TempoIQ, evercam.io, covisint, Jasper, Grovemade, ETHERIOS, PubNub, NURA, SensorCloud, Jiveby	Open Source webinos, AllJoyn, openHAB, nimble.com, OPEN INTERCONNECT, ThingSpeak, GRID2HOME	Sensor Networks SAFECAST, placemeter, Motionloft	Personal Interfaces NeuroSky, Raycon, wit.ai, LEAP, gestigon, speech, TRAILBLAZER, LINTUONE, apLai, EMO TIV, MokuLab, Reemo, Oculus	Security inside, SafeNet, utimaco, escrypt, gemalto, BASTILLE NETWORKS, MOCANA	Corporates amazon, hp, LG, intel, htc, PHILIPS, IBM, SAMSUNG, Google, WIND RIVER, MOTOROLA, belkin, DELL, BOSCH, NATIONAL INSTRUMENTS, ARM, LogMeIn, Microsoft, Honeywell, SONY, Atmel, SIEMENS, QUALCOMM, CISCO, TOSHIBA, SHARP
--	---	---	--	---	--	--

Building Blocks

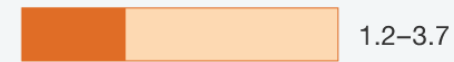
Protocols Bluetooth, Weaved, MQTT, NFC, RuBee, Wi-Fi, ZigBee, OMA, WAVE, enModus, HART, M2M, M-Bus, 2G, 3G, 4G, LTE, CoAP, 6LoWPAN, LWM2M, BITXmI	M2M Networks Helium, SIGFOX, KORE, aeris, HACHEEN, M2M	Portable WIFI Open Garden, GOODSPEED, BRCK, narma	Telecom at&t, boostmobile, Verizon, T-Mobile, Sprint, US Cellular, Vodafone, airtel	M2M arkessa, enOcean, ecanais, seed, WICED, QuSensors, Wireless, seed, WICED, arkessa, enOcean, ecanais, seed, WICED, Wicacomm, Laird, Wicacomm, Laird, Wicacomm, Laird, Wicacomm, Laird			
Cloud Google Cloud Platform, amazon, redhat, ORACLE, Microsoft Azure	Mobile iOS, Windows Phone, BlackBerry	Processors/Sensors WEIO, mCube, BeagleBoard.org, ESP8266, Raspberry Pi, Arduino, MEMS, Freescale	Parts/Kits Makey Makey, SAM, reedymate, littleBits, Wunderlab	Services TINKR, dragon, makeyzi, sculpteo, radstruc, CIRCUIT	Incubators Highway 1, LEANOS Labs, WEARABLE WORLD, R.GA Accelerator, TechShip	Funding KICKSTARTER, indiegogo, MedStart	Distribution angelcam

**Nine settings
where value may accrue**

Size in 2025, \$ trillion¹

■ Low estimate ■ High estimate

Factories—eg, operations management, predictive maintenance



Cities—eg, public safety and health, traffic control, resource management



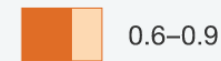
Human—eg, monitoring and managing illness, improving wellness



Retail—eg, self-checkout, layout optimization, smart customer-relationship management



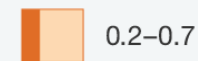
Outside—eg, logistics routing, autonomous (self-driving) vehicles, navigation



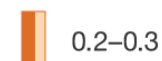
Work sites—eg, operations management, equipment maintenance, health and safety



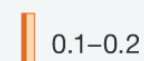
Vehicles—eg, condition-based maintenance, reduced insurance



Homes—eg, energy management, safety and security, chore automation



Offices—eg, organizational redesign and worker monitoring, augmented reality for training



Total \$4 trillion–\$11 trillion

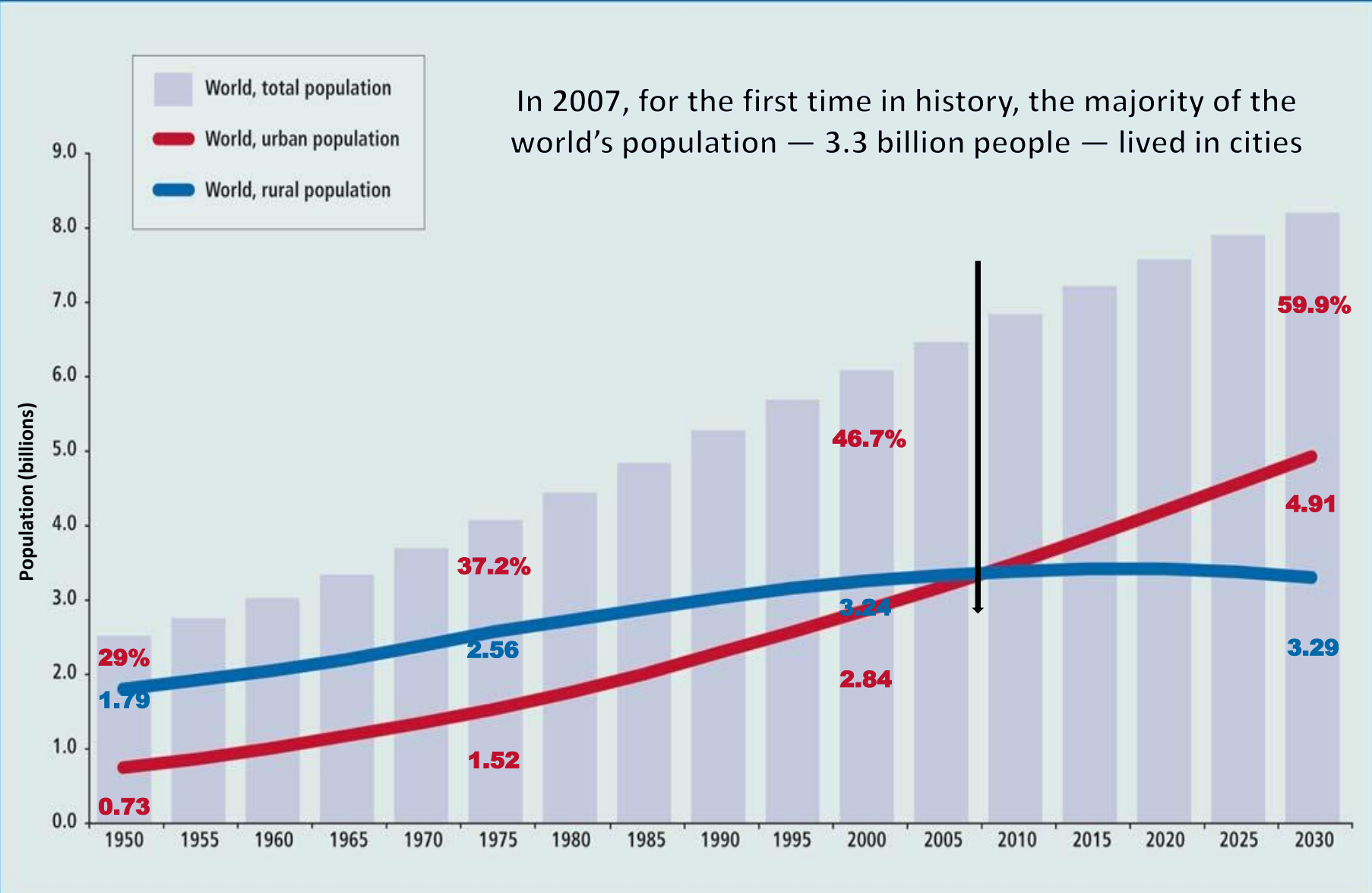
¹Adjusted to 2015 dollars; for sized applications only; includes consumer surplus. Numbers do not sum to total, because of rounding.

Smart Cities

Cities



The urban and rural population of the world, 1950-2030



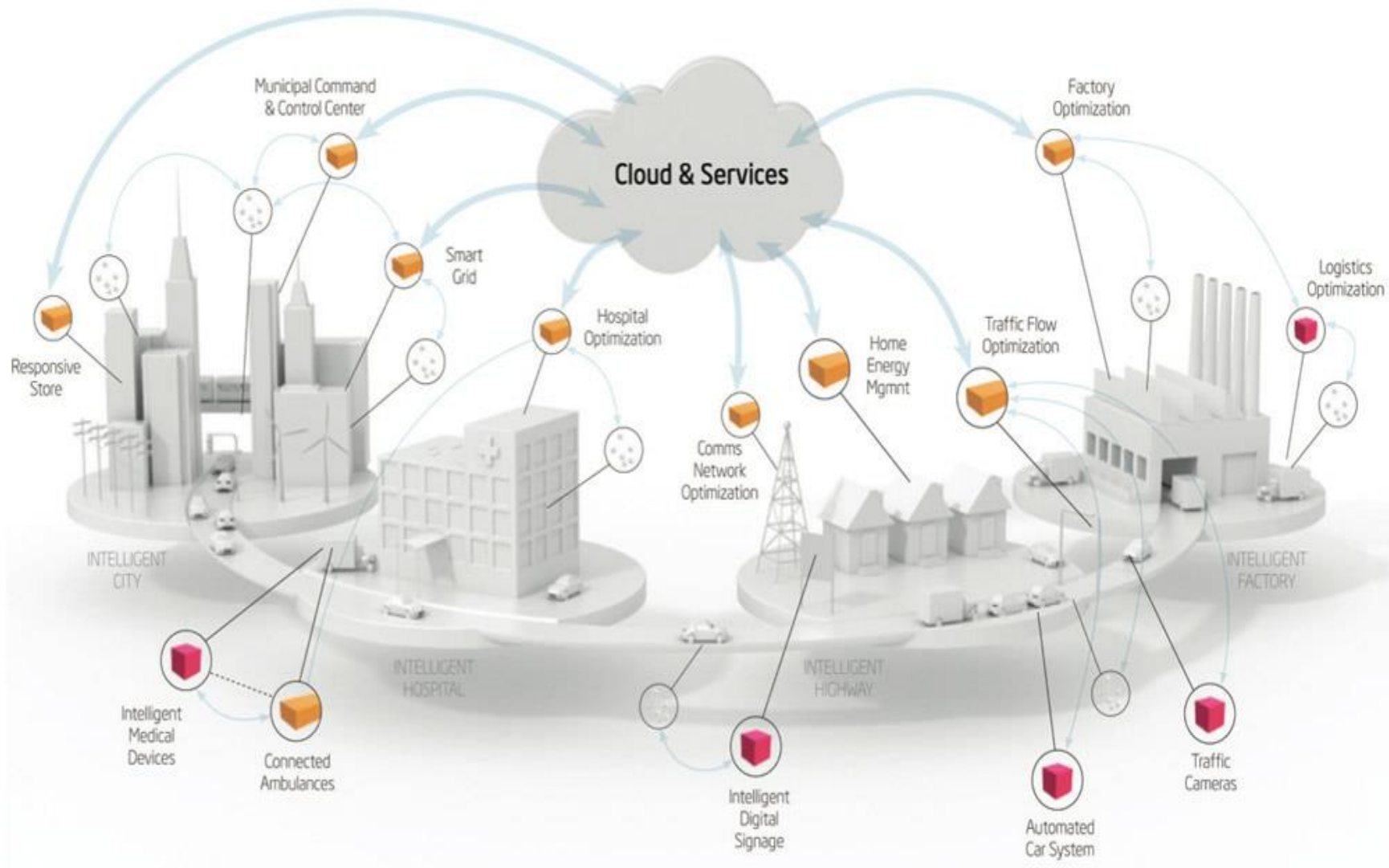
World Urbanization Prospects

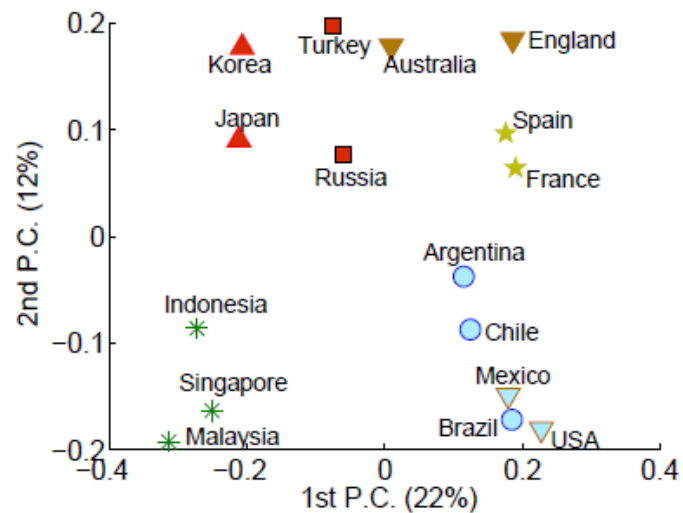




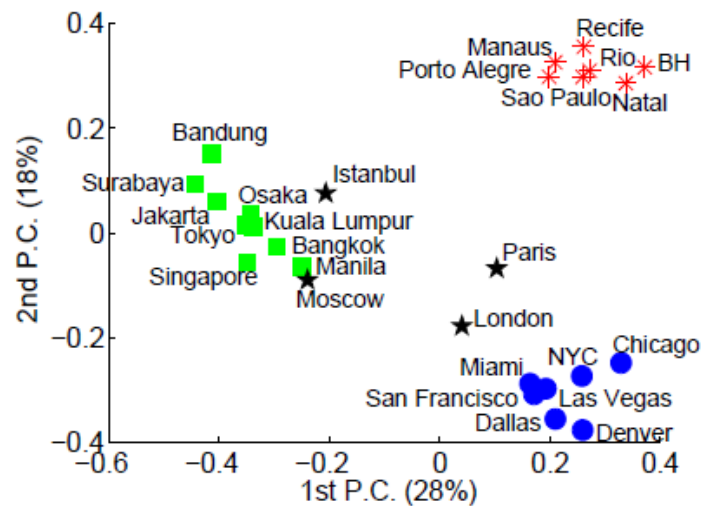
Smart City applications

- **Smart parking:** Monitoring of parking spaces availability in the city.
- **Structural Health:** Monitoring of vibrations and material conditions in buildings, bridges and historical monuments.
- **Traffic Congestion:** Monitoring of vehicles and pedestrian levels to optimize driving and walking routes.
- **Smart lighting:** Intelligent and weather adaptive lighting in street lights.
- **Waste management:** Detection of rubbish levels in containers to optimize the trash collection routes.
- **Smart roads:** Intelligent Highways with warning messages and diversions according to climate conditions and unexpected events like accidents or traffic jams.

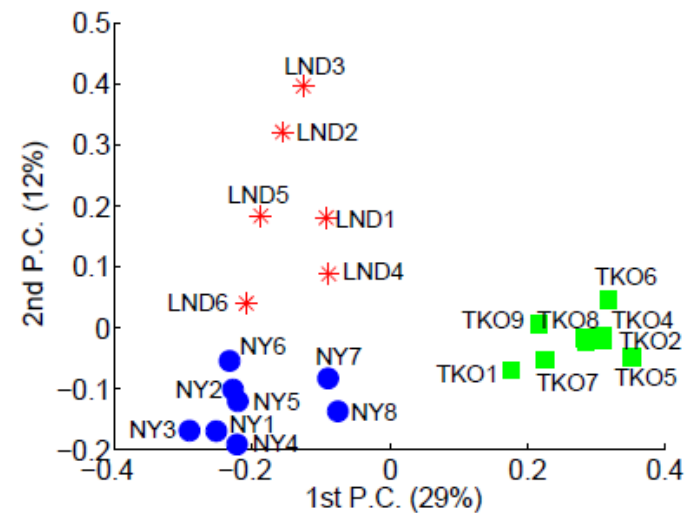




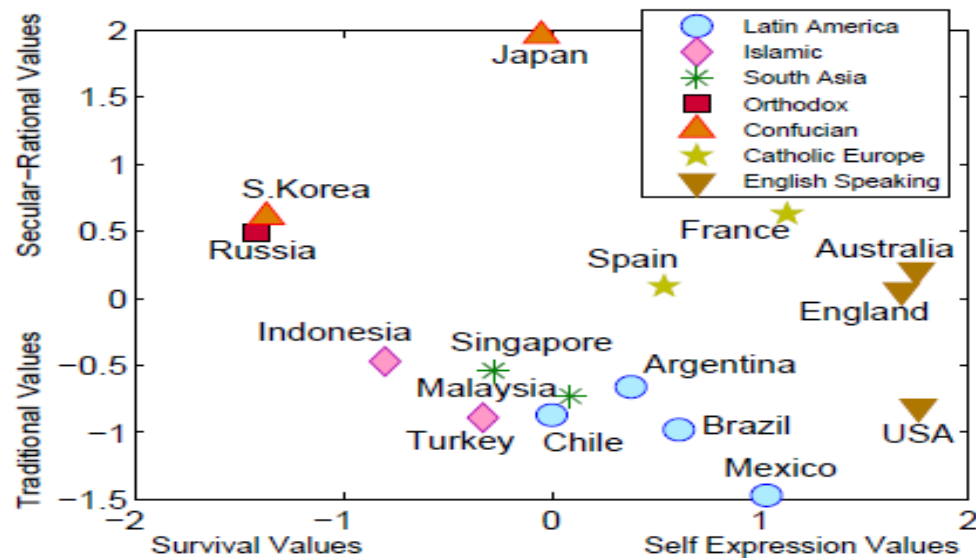
(a) Countries



(b) Cities



(c) Regions

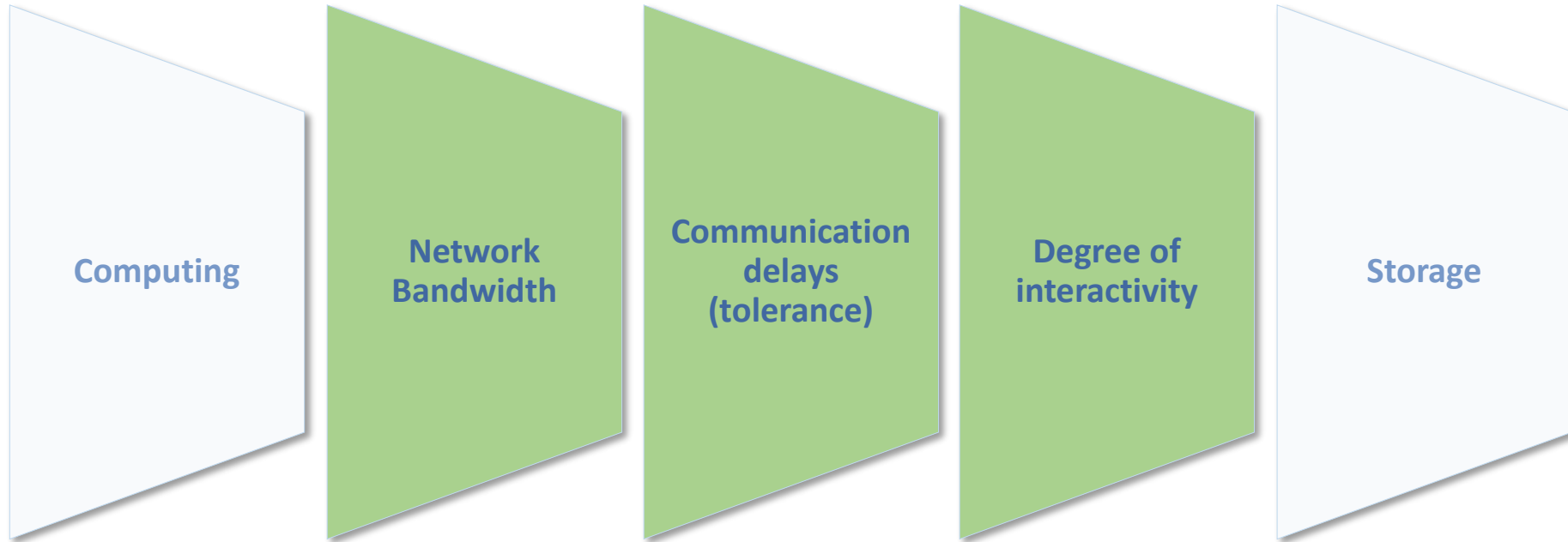


- Latin America
- ◇ Islamic
- * South Asia
- Orthodox
- ▲ Confucian
- ★ Catholic Europe
- ▼ English Speaking

Networking



Networking



Network Virtualization

Networking for Big Data
and
Big Data for Networking

Networking for Big Data

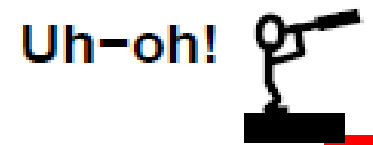
What is the role of networking in Big Data?



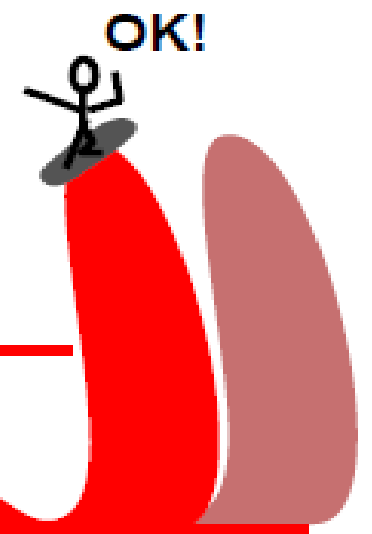
Big Data ... and the Next Wave of **InfraStress**

***John R. Mashey
Chief Scientist, SGI***

***Technology Waves:
NOT technology for technology's sake
IT'S WHAT YOU DO WITH IT
But if you don't understand the trends
IT'S WHAT IT WILL DO TO YOU***



Uh-oh!



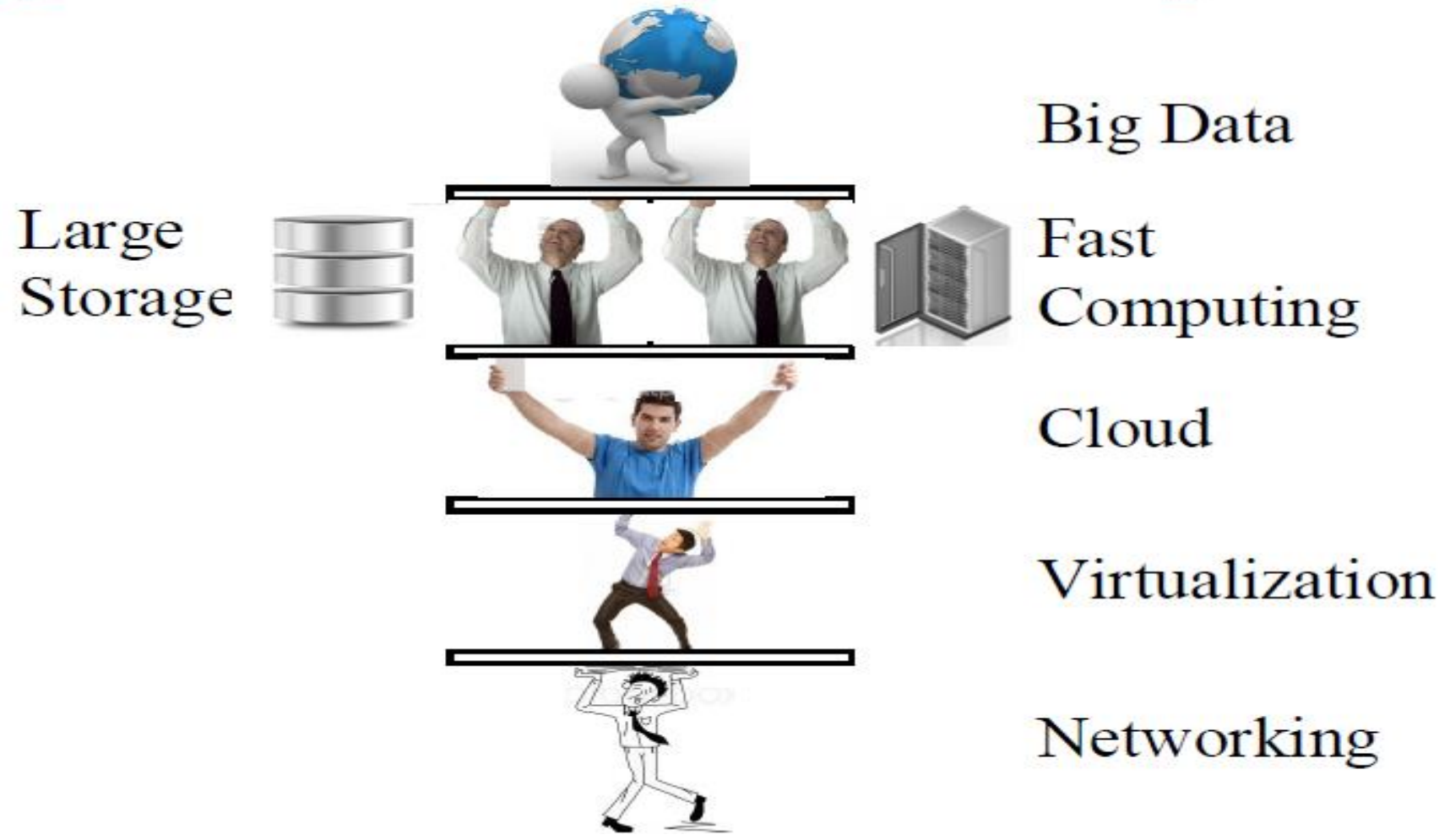
OK!

“Infrastrress”



Alibaba Mall processes in a single day (Nov 11th, 2013) 105.8 million online transactions from 213 million users and 4.1 billion transactions

Big Data Enabled by Networking



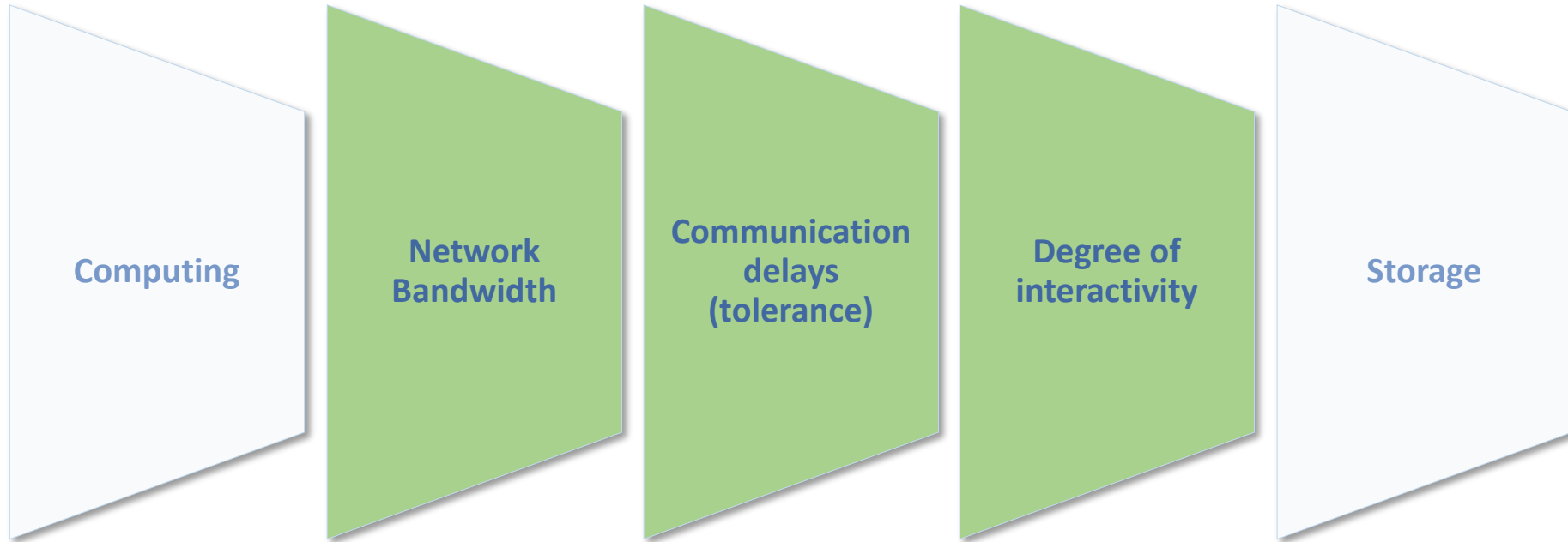
Networking Requirements for Big Data

- Elastic bandwidth: to match the variability of volume
- High Speed data transfer
- Security: Access control privacy, threat detection, all in real-time in a highly scalable manner
- Network partitioning to handle big data
- Network congestion control for big data applications
- Network service consistency

Networking



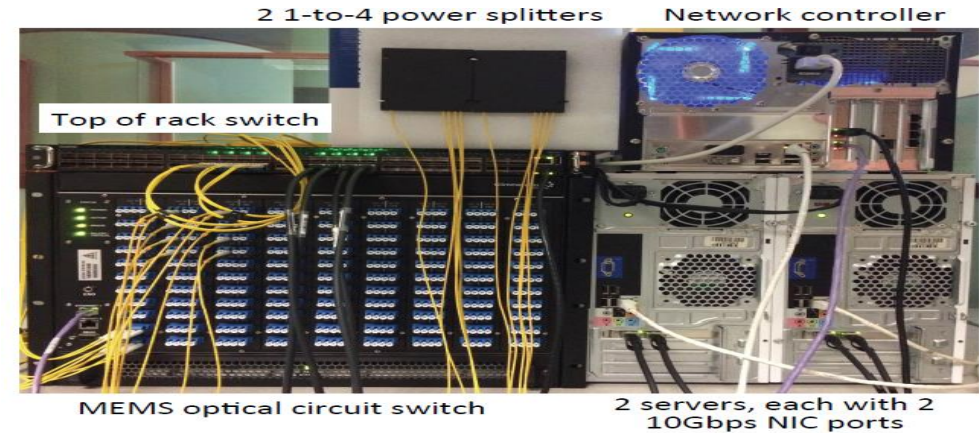
Networking



Optical Multicast

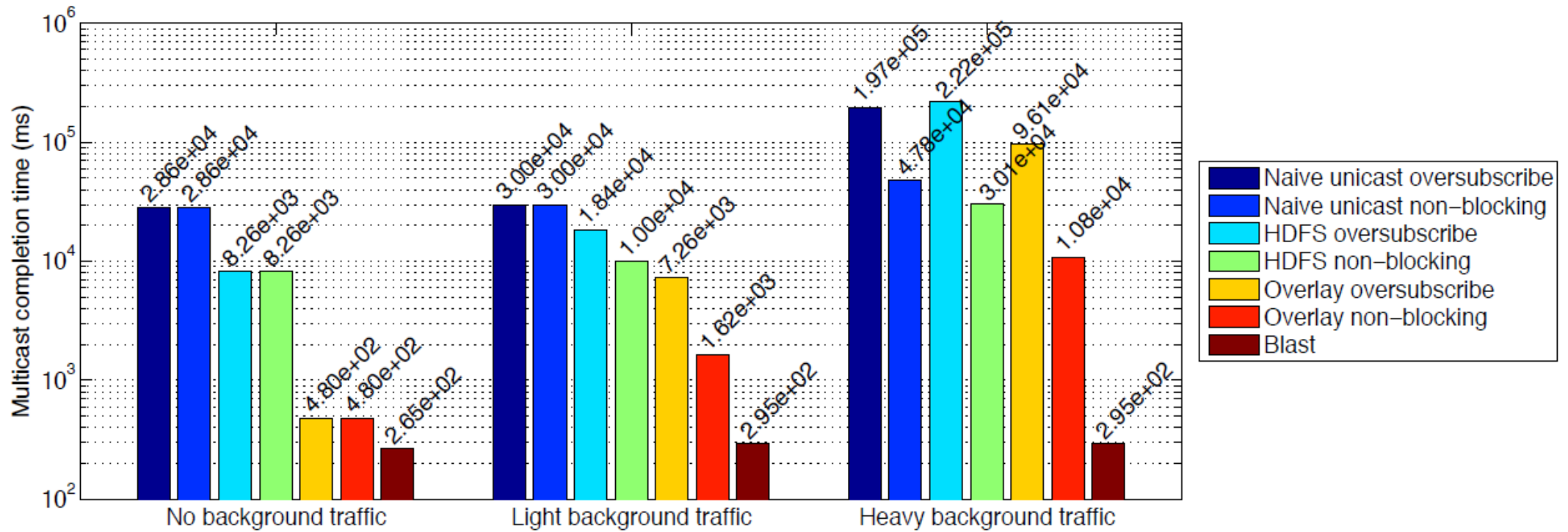
- Data analytics applications routinely need to distribute terabytes of data from a central data source to hundreds of servers for processing
- In Hadoop Distributed File System(HDFS), multicast sender stores the data in HDFS and a multitude of receivers retrieve the data from a few data replicas, creating very high fan-out on the replicas.
- In Spark, a BitTorrent style P2P overlay among the recipient nodes, but BitTorrent suffers from suboptimal multicast trees that render high link stress, performs worse than HDFS

Optical Multicast

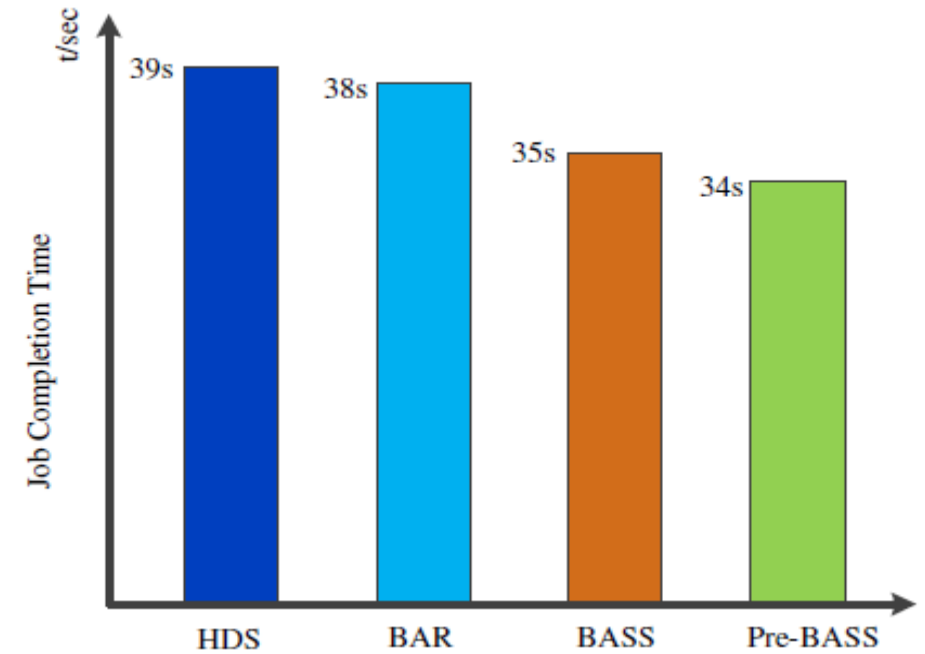
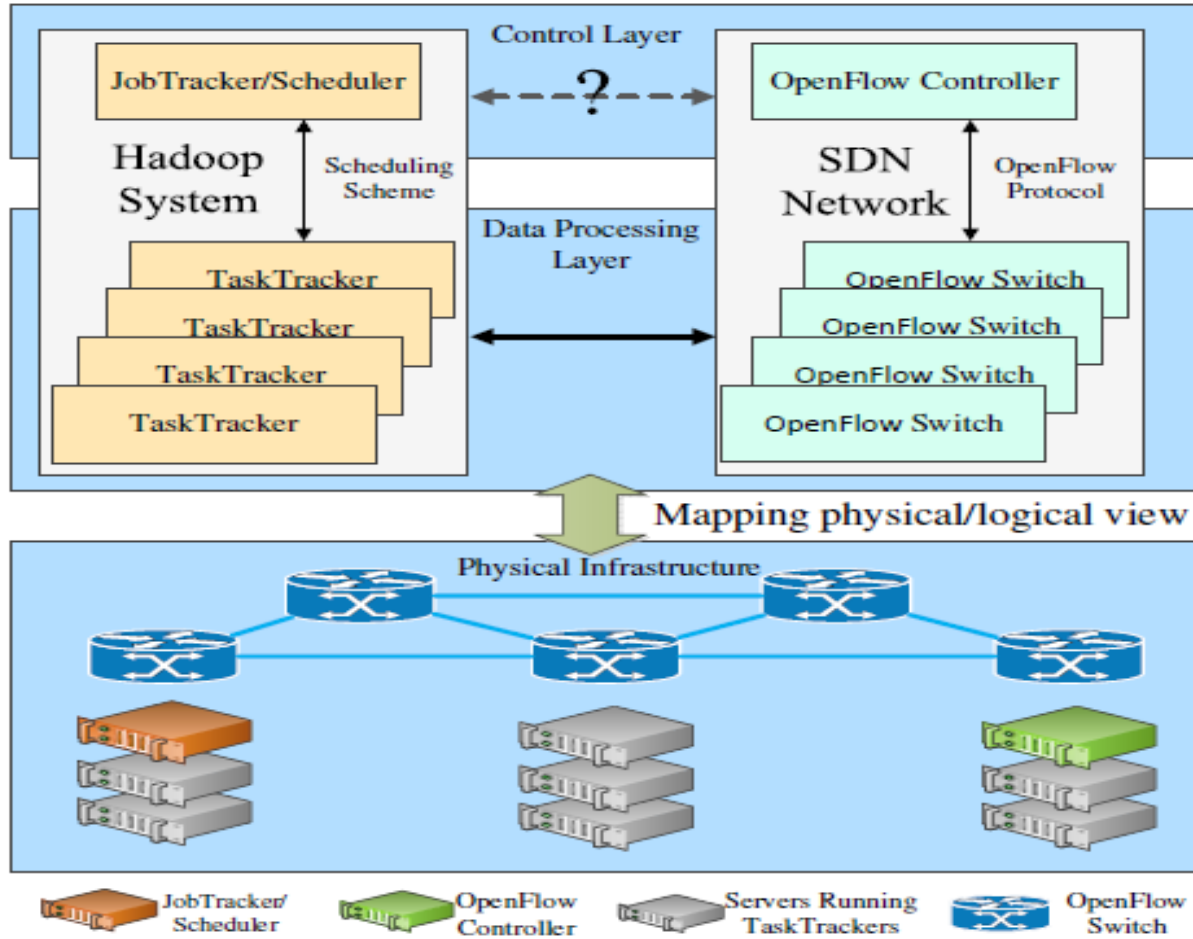


- Inherent performance limitations of application-layer overlays
- Difficult to perform TCP-friendly congestion control in network multicast traffic.
- Blast uses optical transmission to realize a physical-layer broadcast medium, via passive optical power splitting, to connect a data source to its receivers
- tailor-made control plane, capable of collaborating with data analytics applications interactively, making resource allocation decisions nearly optimally, and directing the data flows in optical and electrical components in the network

Optical Multicast



SDN & Hadoop



Big Data for Networking

Big data for Networking

- Well investigated problems:
 - Traffic classification
 - Intrusion detection
 - Anomaly detection
 - Cognitive networks

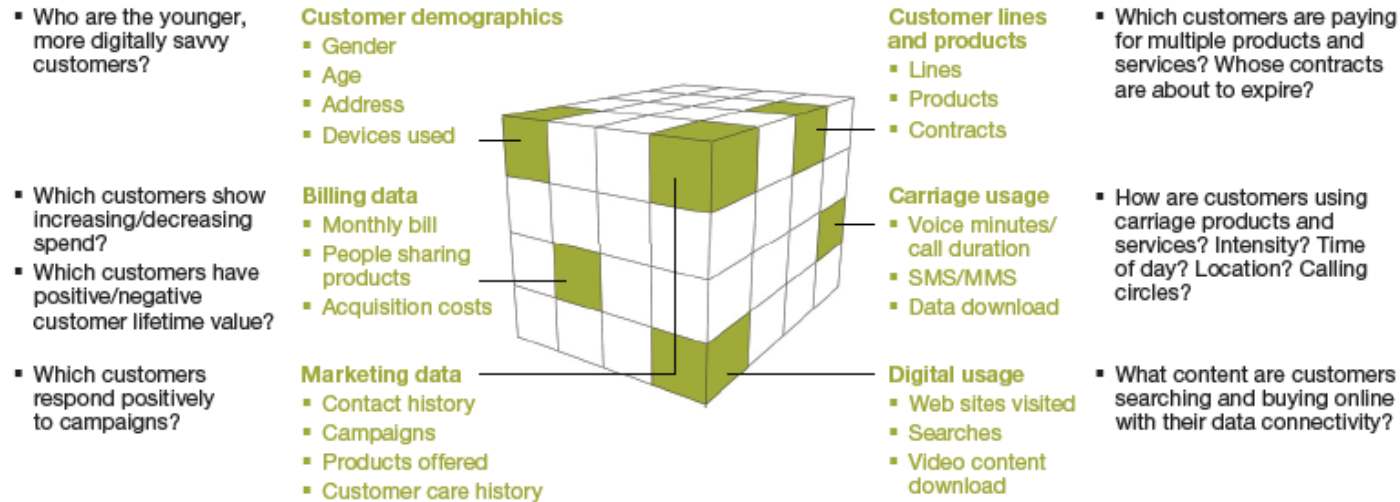
Big data for Networking

Recent Work

- Definition of new services
- Service level prediction
 - ✓ VoD
 - ✓ Scheduling in hybrid clouds
- Network management and configuration
- Anomaly in forwarding tables
- Fault detection and location
- Construction of radio environment maps

Definition of new Services and Class of Services

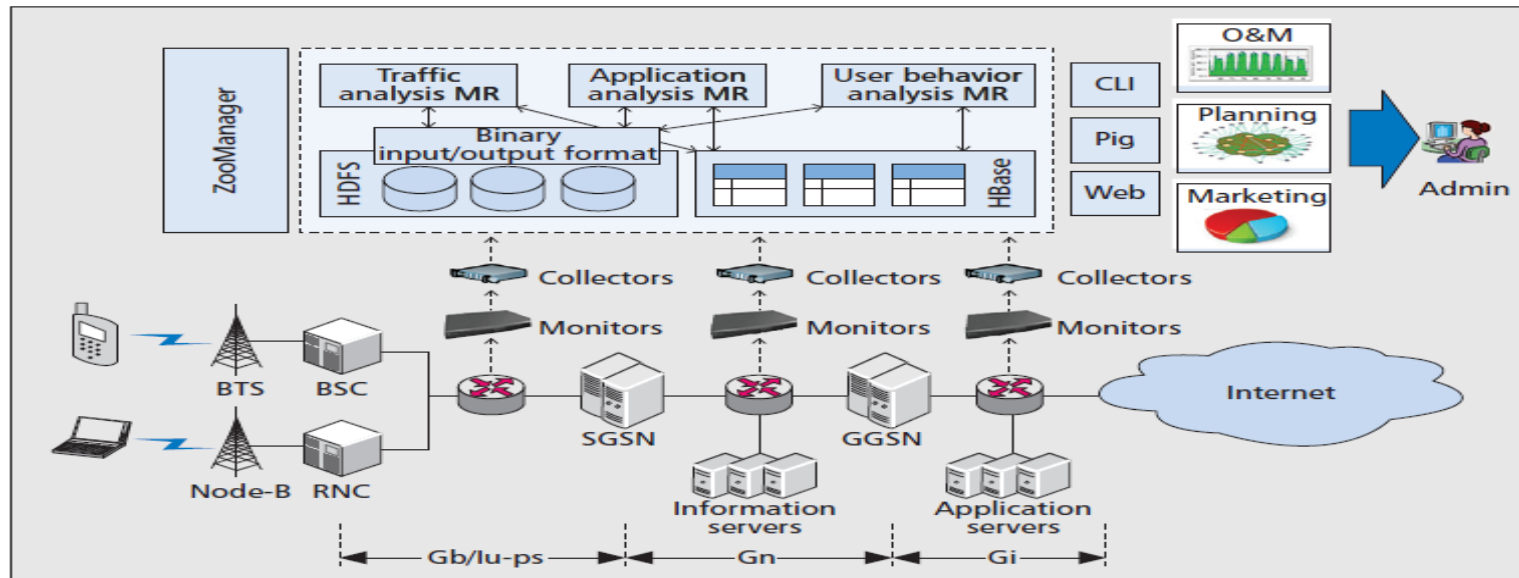
Based on detailed customer profiles, telcos can differentiate customer service models and develop individualized offer recommendations



SOURCE: McKinsey

Identifying Communities of Website users

- Mobile operator in China
- Objective: to find website communities of users and identify their usage behavior



Identifying Communities of Website users

- Traffic Monitor:
 - Line-speed packet parsing (PPPoE, GRE at various interface Gb, Gn, Gi)
 - Real-time traffic classification (19 sub-service classes)
 - Multi-level traffic statistics (packet, flow and aggregate level)
- Hbase in HDFS Hadoop, key/value stored in a columnar manner
- Mining logs of HTTP:
 - Affinity graphs models website usage
 - Sparsified affinity graph constructed by employing a scale free fitting index (node in-degree larger than a threshold)
 - Nodes are ranked by an influence score

Identifying Communities of Website users

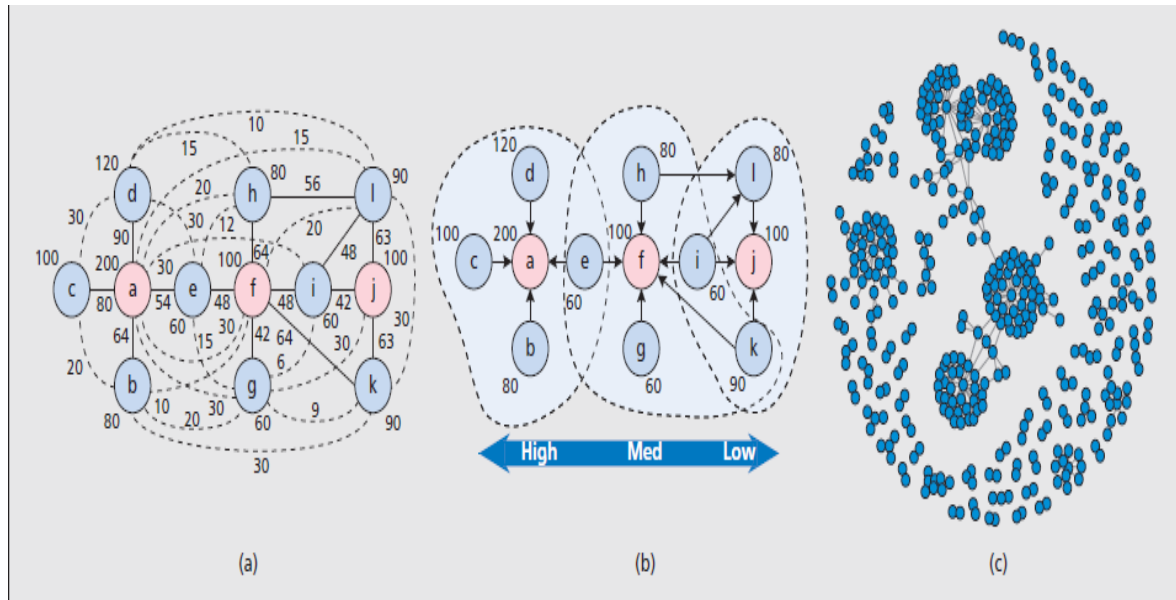


Figure 3. Website community identification: a) example of user distribution among websites; b) example sparsified affinity graph; c) identified website communities.

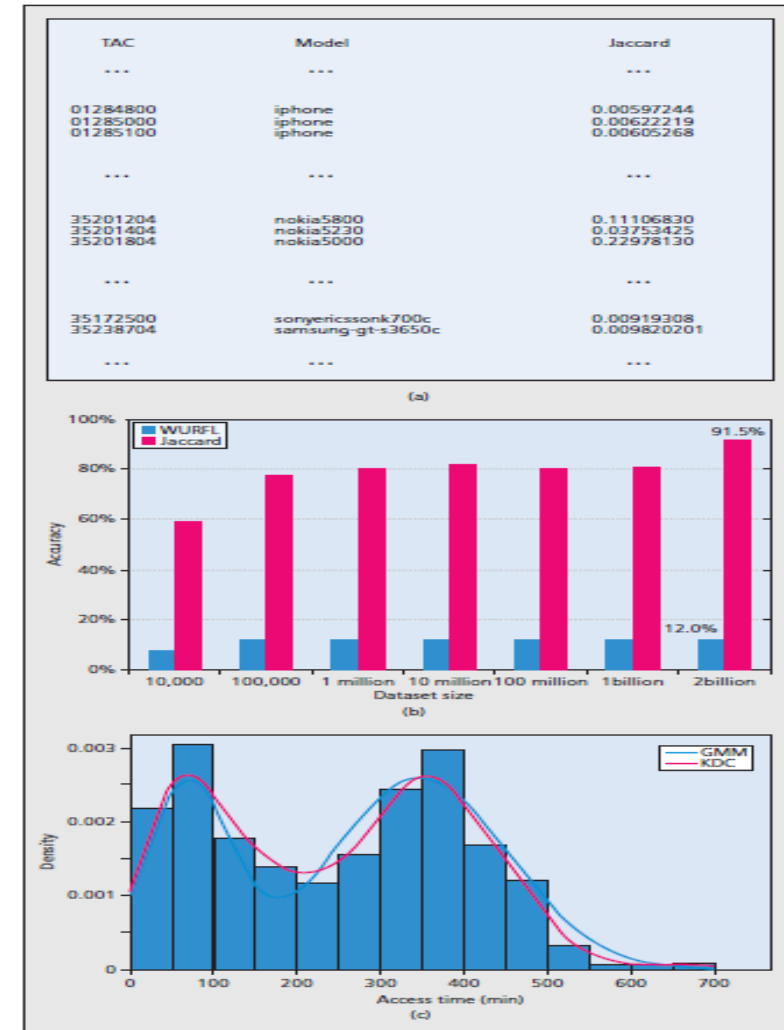
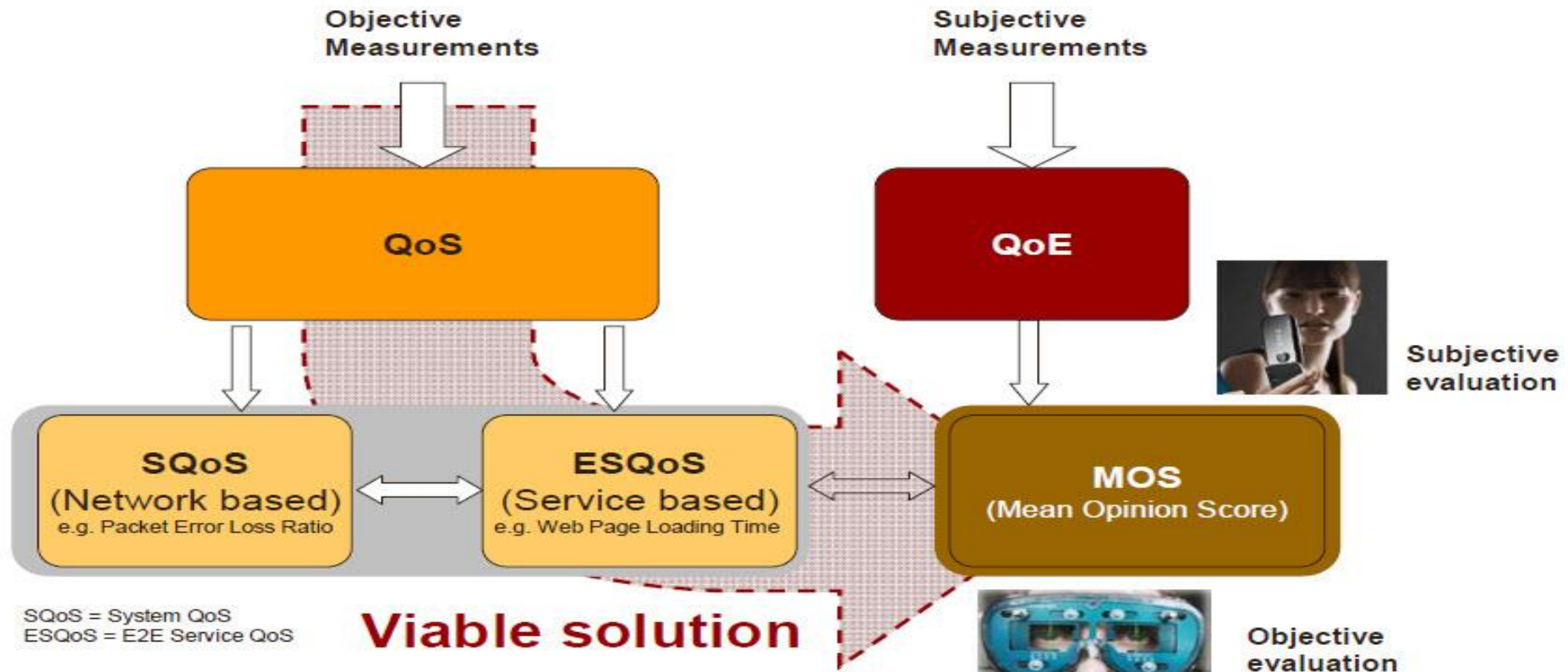


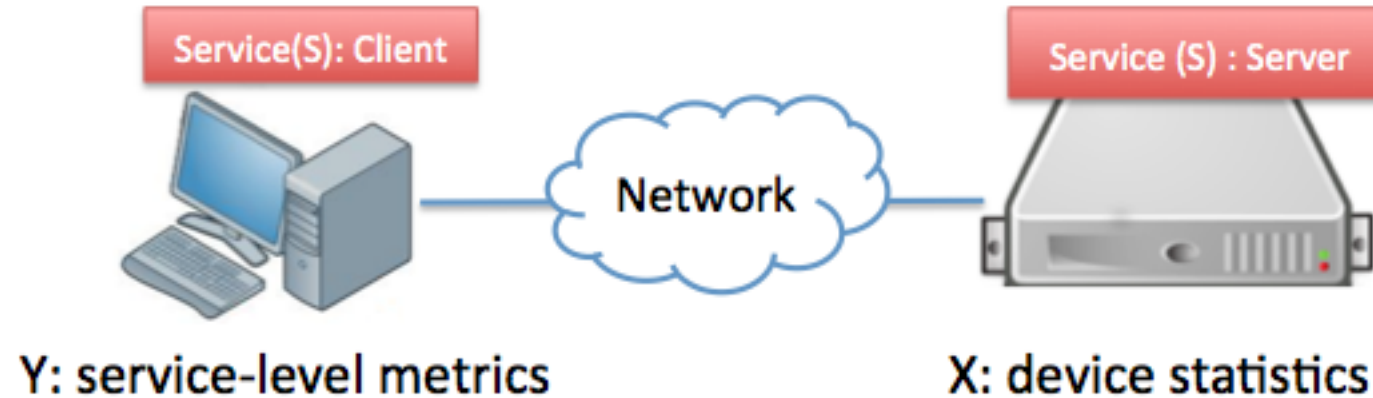
Figure 4. Mobile client model identification and user behavior analysis: a) example of identification results; b) accuracy of identification; c) density of network access time.

Service Level Prediction

QoS and QoE parameters – Mapping Model



Service Level Prediction - VoD



- **Video streaming:** video frame rate, audio buffer rate, RTP packet rate

- CPU load, memory load, #network active sockets, #context switching, #processes, etc..
- raw data from /proc provided by Linux kernel

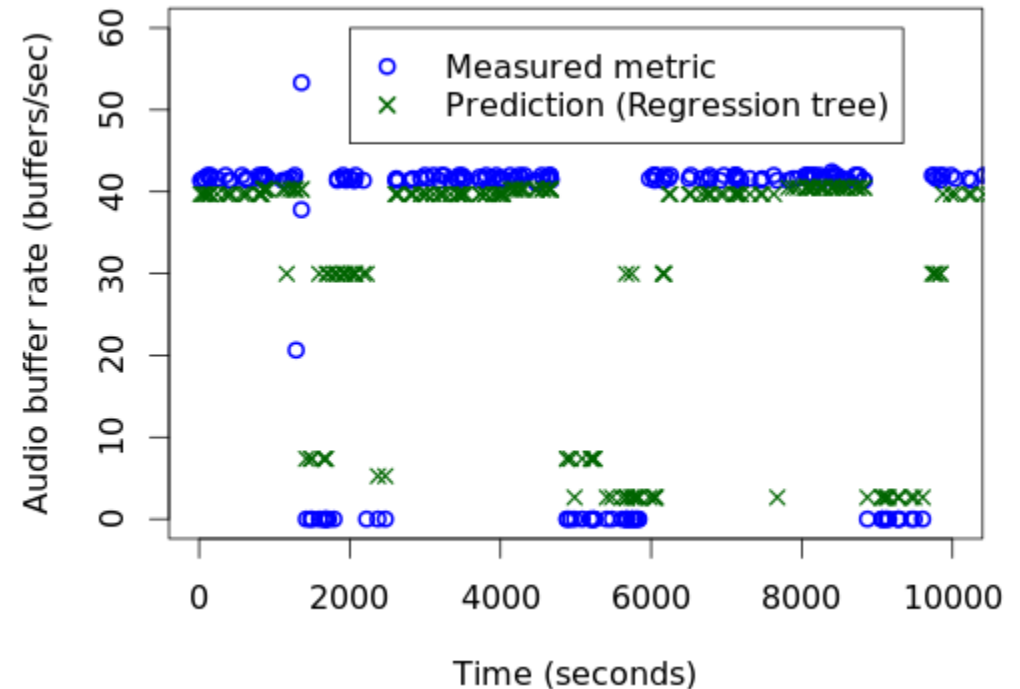
- Building block for real-time service assurance for a telecom cloud

Service Level Prediction - VoD

- Statistical learning on low-level (OS-level) metrics, taking a large number of features (> 4000) rather than few service-specific features (≤ 10)
- Measured metrics
 - Video frame rate (frames/sec)
 - Audio buffer rate (buffers/sec)
 - RTP packet rate (packets/sec)
- Linear regression, Regression tree, Random forest, Lasso regression
- Network statistics and client low-level metrics not considered
- Network and client machine are lightly loaded

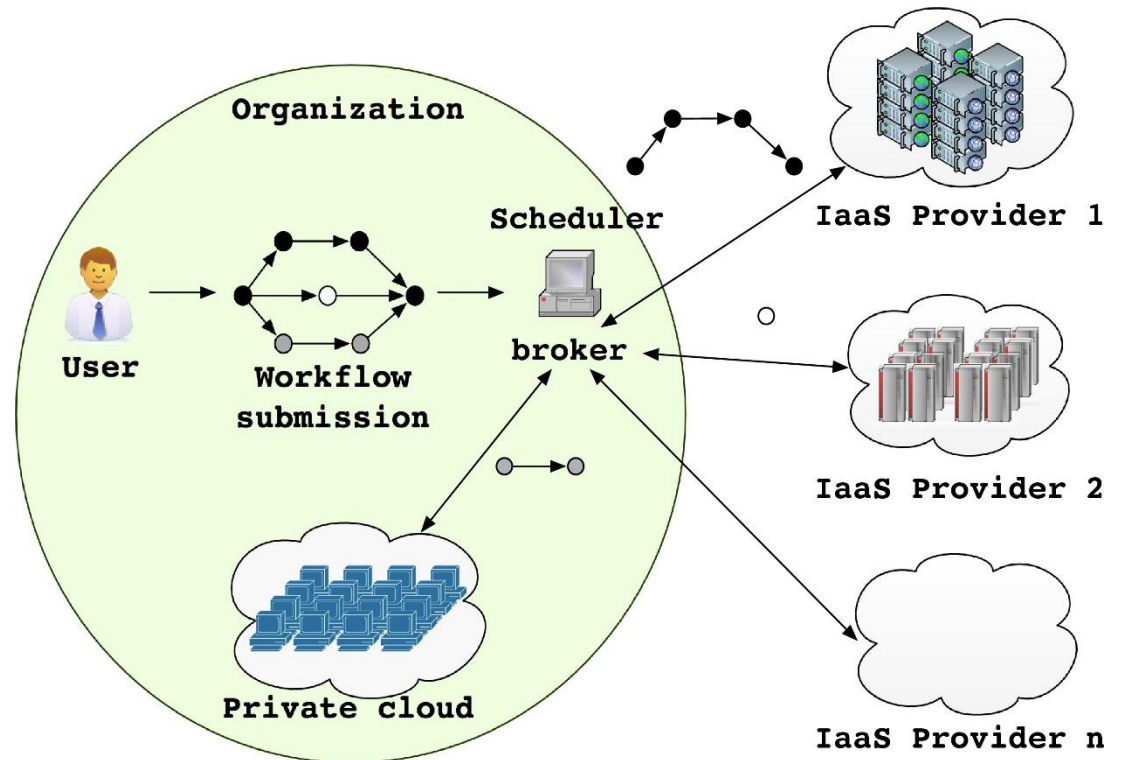
Service Level Prediction-VoD

- It is feasible to accurately predict client-side metrics based on low-level device statistics
 - Normalized Mean Absolute Error below 15% across service-level metrics and traces
- Preprocessing of X is critical
- Trade-off between computational resources vs. prediction accuracy
- No time dependence considered



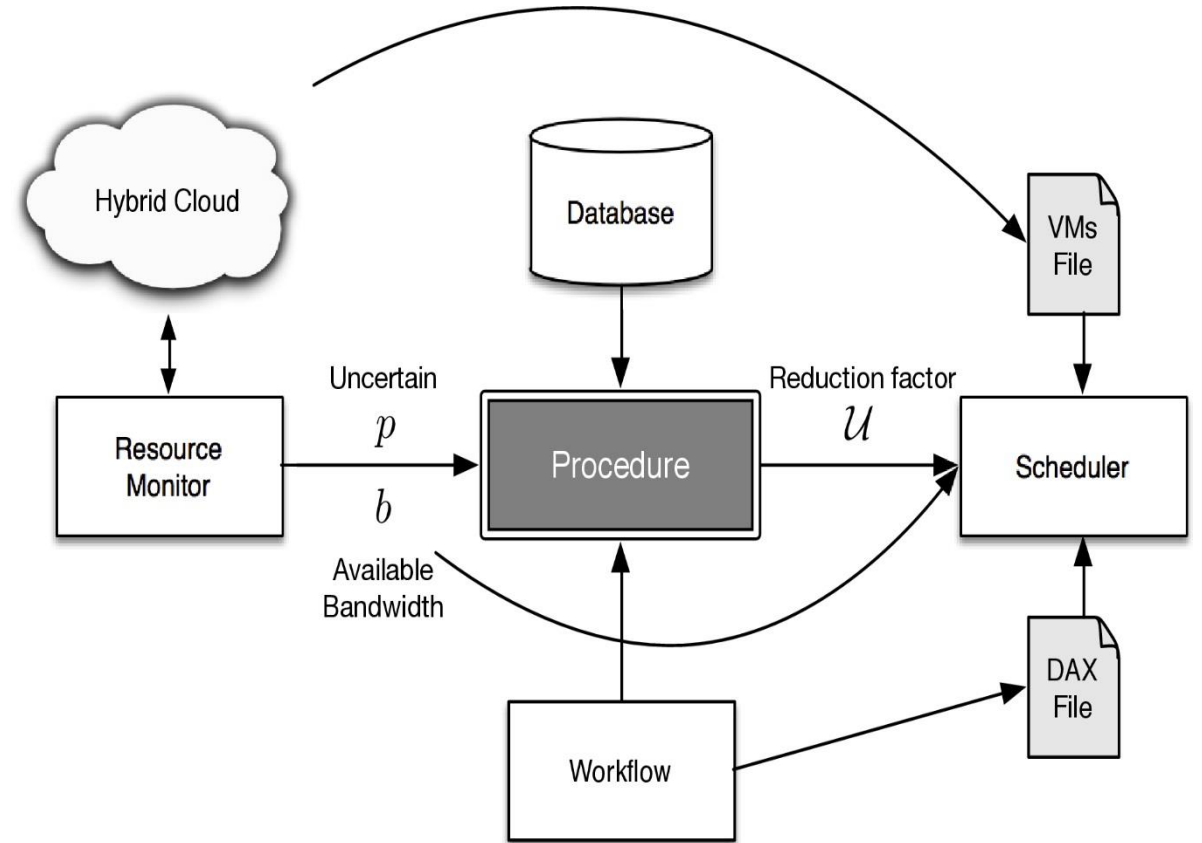
Service Level Prediction- hybrid cloud

- Impact of intercloud link bandwidth on makespan of workflow execution
- Overestimation of available bandwidth can lead to increased makespan and higher cost
- Underestimation of the bandwidth leads to unnecessary leasing of resources
- Available bandwidth varies during execution of workflow

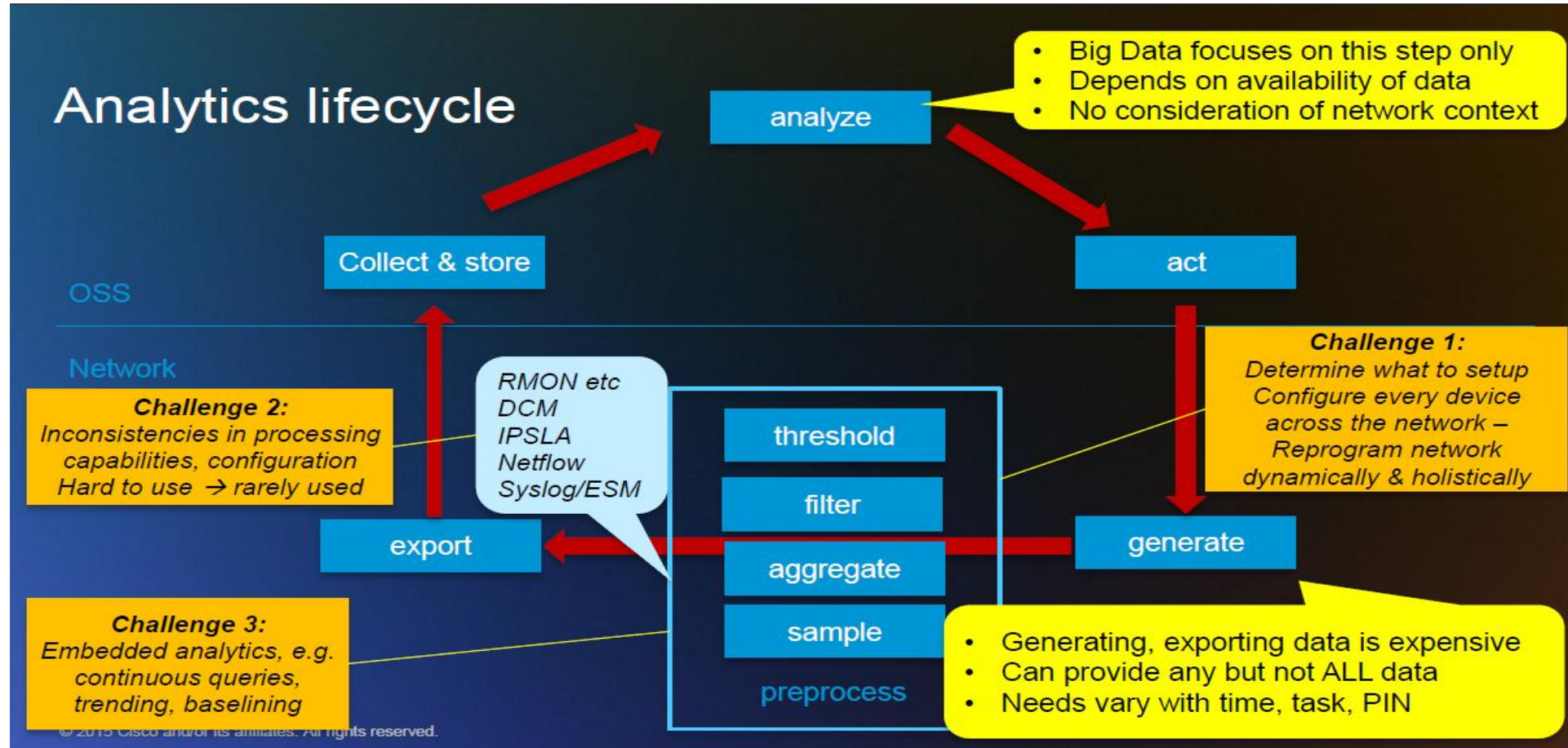


Service Level Prediction- hybrid cloud

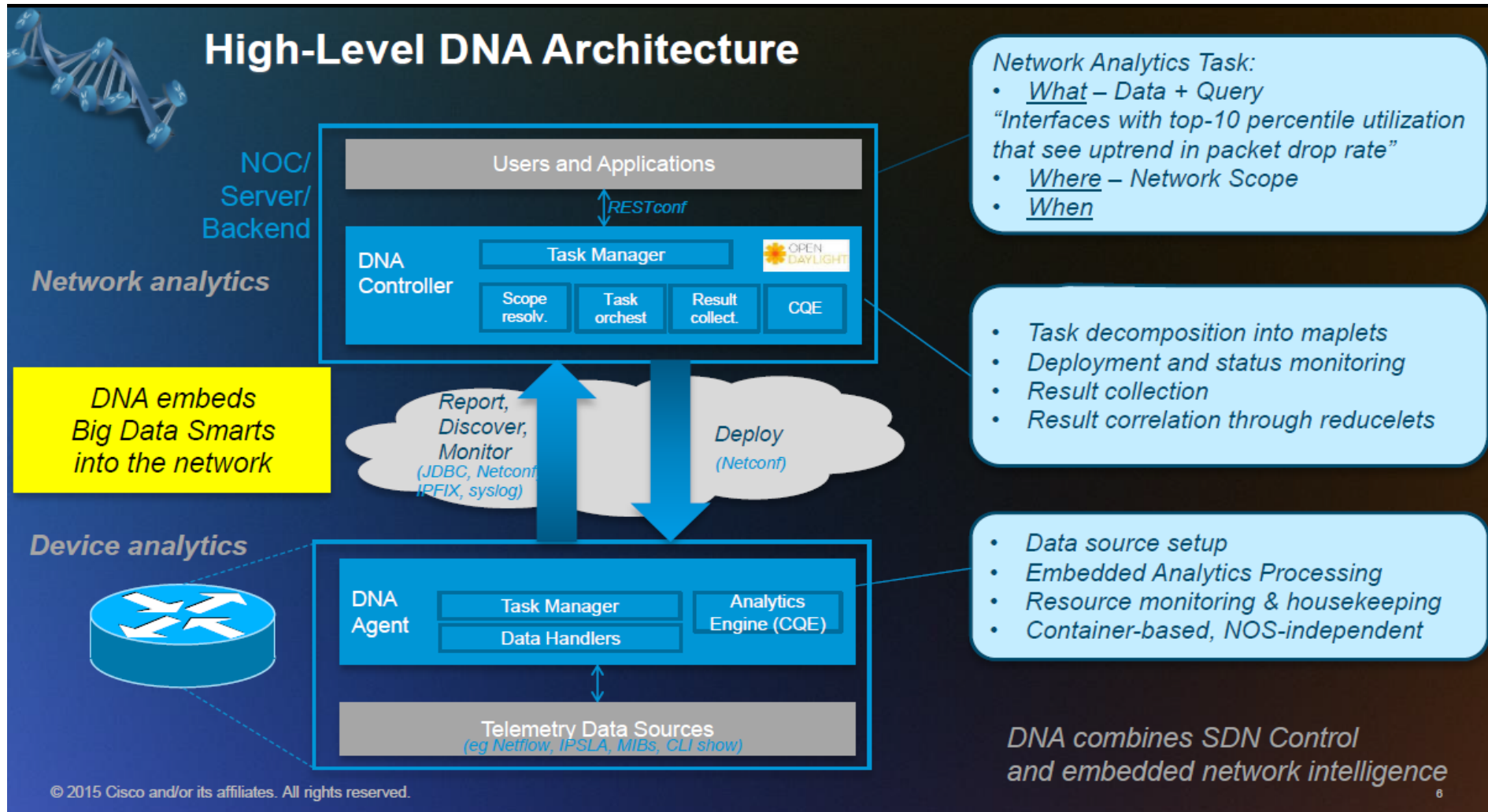
- Deflating fator to measured available bandwidth
- Database with historical data on performance relating available bandwidth, type of workflow, deflating fator, quality of information and estimation erros
- Multiple linear regression
- Reduced the number of disqualified solutions and reduced costs



Big Data for Network Management and Operation



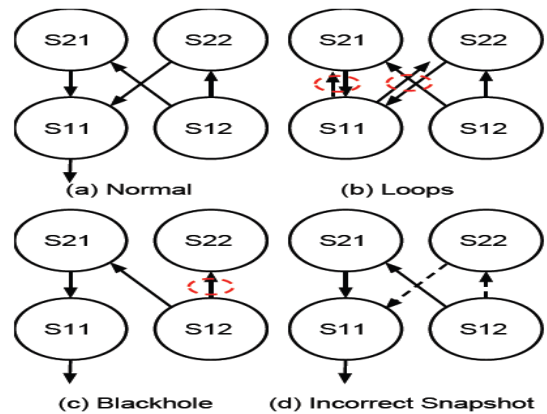
Big Data for Network Management and Operation



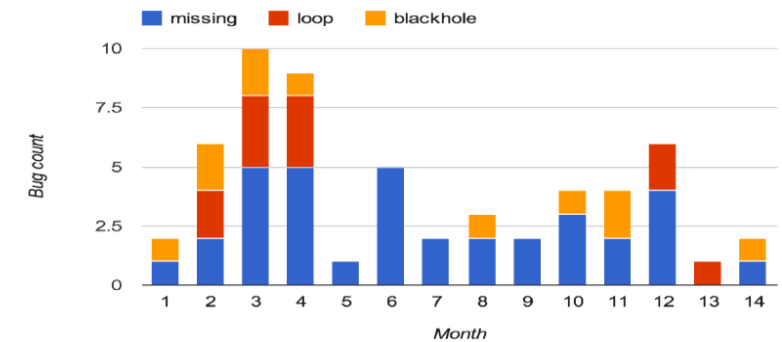
Verification of Forwarding Tables

- 10,000s switches in modern data centers
- Traditional approach of scaling-out and redundant design assumes correct reaction to errors
- Dormant bugs in routing systems triggers rare boundary conditions, some times perceived when benign event happens
- Two major deficiencies of available tools:
 - Assumption of consistent snapshot of forwarding tables
 - Not sufficiently fast to meet requirements of modern datacenters

Verification of Forwarding Tables



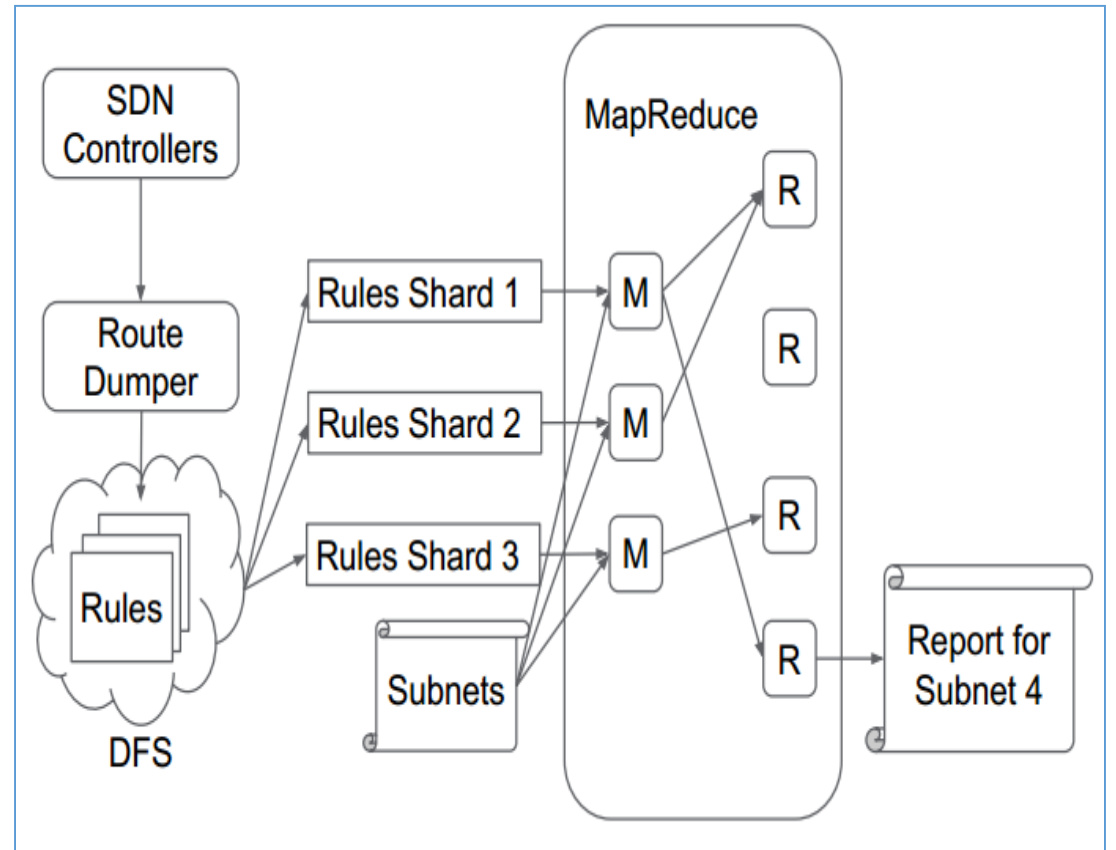
Data Plane Tickets in a Google Data Center



- Loops - usually caused by prefix aggregation
- Black-holes - lost BGP updating information
- Inconsistent snapshot - in an event of failure, changes to forwarding table may conflict with previous forwarding information

Verification of Forwarding Tables

- Reachability - subnet can be reached from any switch; depth-first-search from the subnet switch, verify if all other switches are reachable
- Loop detection - strongly connected component
- Black-role - if a switch does not have a matching entry for the subnet



Verification of Forwarding Tables

- Libra: Dive-and-Conquer approach

Data set	Switches	Rules	Subnets
DCN	11,260	2,657,422	11,136
DCN-G	1,126,001	265,742,626	1,113,600
INET	316	151,649,486	482,966

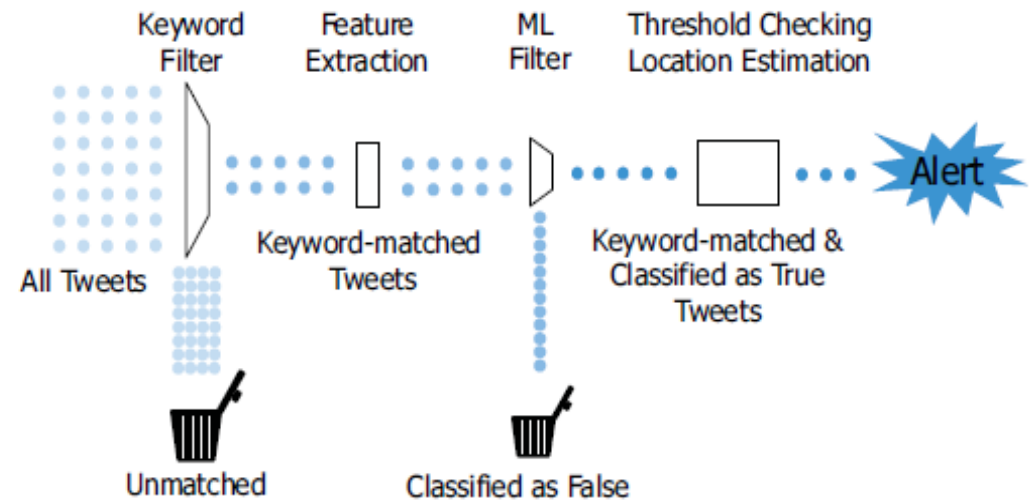
	DCN	DCN-G	INET
Machines	50	20,000	50
Map Input/Byte	844M	52.41G	12.04G
Shuffle Input/Byte	1.61G	16.95T	5.72G
Reduce Input/Byte	15.65G	132T	15.71G
Map Time/s	31	258	76.8
Shuffle Time/s	32	768	76.2
Reduce Time/s	25	672	16
Total Time/s	57	906	93

Fault Identification and Location



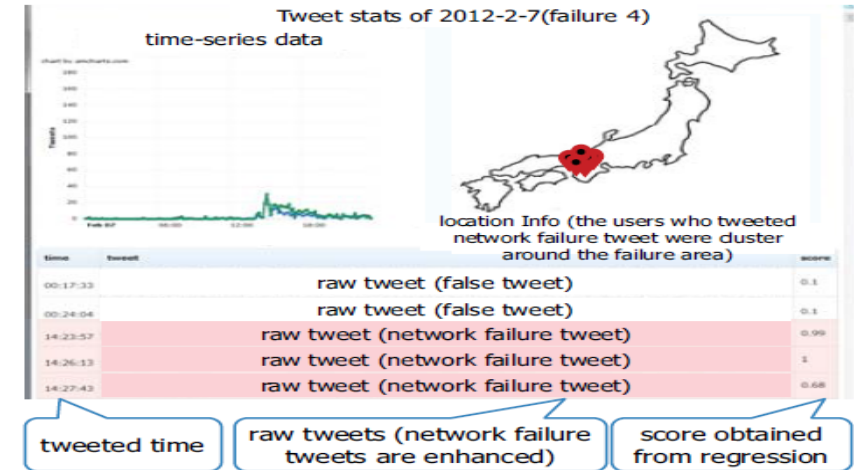
- Some network failures become silent failures to mobile operators; difficult to establish rules for failure detection
- Most failures are not reported to call centers
- Monitoring Tweeter to detect failure in mobile services
 - ✓ "Why can't I text messages?"
- Most tweets are not related to network failure
 - ✓ I dropped my phone in toilet when I called my friend

Fault Identification and Location



- Two requirements:
 - ✓ reduction of false positives and
 - ✓ location and when incident happened
- Imbalance between network failure tweets and others
- Closeness of network failure tweets and false positives
- Three stages: keyword filtering, machine learn filter and alert system

Fault Identification and Location



- Keyword filtering: selection of features related to network failure
 - ✓ 4398 network failure and 6924 false positives
 - ✓ 10K tweets/sec for 40 keywords
- Four different techniques in machine learning filtering: SVM with and without Gaussian Kernel, Naive Bayes and Adaptive Regularization of weights
 - ✓ 1% of raw tweets,
 - ✓ Capacity of 1 K tweets per second
 - ✓ 1 hour to construct a model
- Location: name of city, station or landmark; use of gazetteer and kernel density estimation to estimate location of fault; known GPS location of tweets
 - ✓ Heavy processing load; 0.2 users/second

Radio Environment Maps



- Spatial maps of received signal strengths can be used in dynamic spectrum access to, for example, discover coverage holes in cellular networks
- Perform drive test to collect data and perform spatial interpolation
- Methods from spatial statistics, such as the Kriging method, are accurate and robust. However, its complexity is $O(n^3)$ where n is the number of measurement points
- Current used approaches can only give a rough estimation since propagation simulations have inherent inaccuracies due to limited information on landscape data. Moreover, drive tests are too expensive

Radio Environment Maps

- Minimization of drive tests developed by 3GPP makes every mobile phone a spectrum measurement device, making available a large number of path loss or received signal strength and GPS location information
- Operators can harvest unprecedented amount of data
- Employment of fixed rank kriging techniques with linear computational complexity to process hundreds of thousands of measurements, spatial estimate
- Model fitting takes 20 seconds of computation in a desktop and individual predictions less than milliseconds

Radio Environment Maps

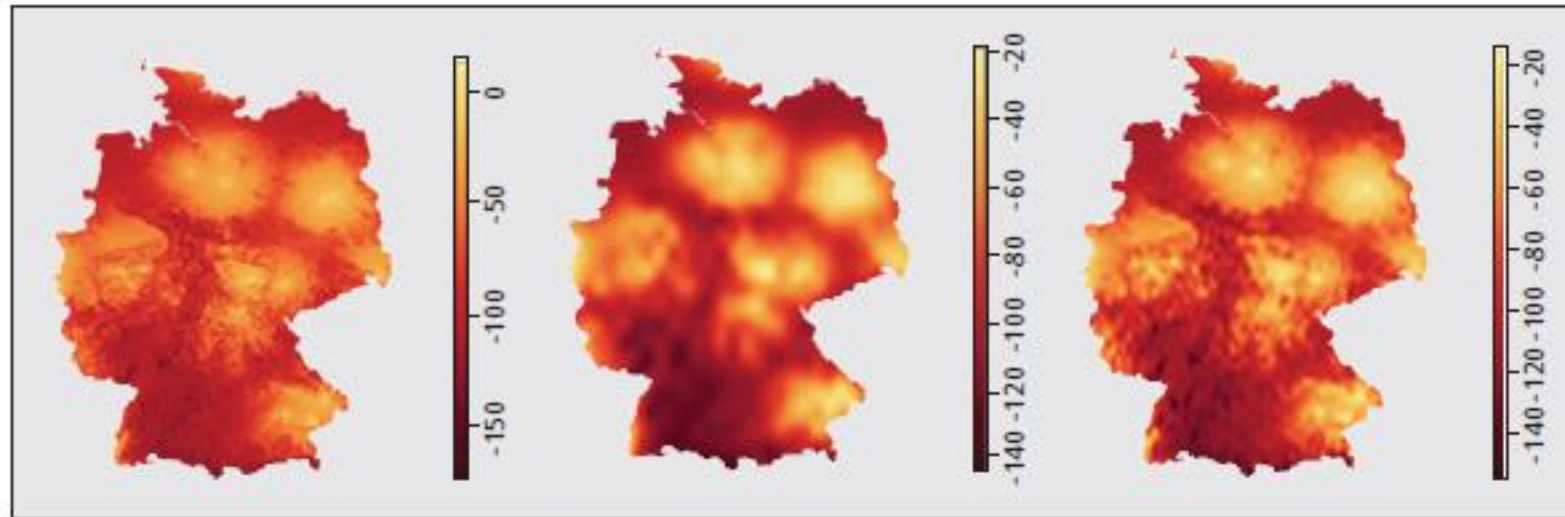


Figure 1. Example of a typical spatial estimation task in management of wireless networks. The map on the left shows the actual coverage of a digital TV network over Germany (total area of approximately 350,000 km²), whereas the middle and right figures show spatial estimates based on 10,000 and 27,000 distributed measurements, respectively.