



WWW.IARIA.ORG

**PANEL
ALLDATA/IMMM/KESA/MMEDIA**

**Computational Models for Big Data
Processing**

Moderator: Gary Weckman, Ohio University, USA

Getting lost

- **Open Data**
- **Linked Data**
- **Big Data**
- **Small Data**
- **Public Data**
- **Private Data**
- **Data sets**
- **Data accuracy**
- **5V+ (volume, velocity, veridicity,**

- **Computational Models**
 - Data Management**
 - Semantic Aspects**
 - Legal Aspects**

Computational Models for Big Data Processing

- **1 - In every sector, there is a growing need for companies to extract knowledge from Big Data in a fast and reliable way;**
- **2 – To do so, alternatively to stochastic analysis methods, companies could use rule-based NLP tools and environments, such as for example NooJ (www.nooj-association.org, <http://en.wikipedia.org/wiki/NooJ>). It would be interesting to practically see the advantages and disadvantages of such an adoption;**
- **3 - Specific professionals should be formed to perform Knowledge Extraction, starting from undergraduate courses.**

[Mario Monteleone]

Today's Panelists

Moderator

- **Gary Weckman, Ohio University, USA**

Panelists

- **Iryna Lishchuk, Institut für Rechtsinformatik/Leibniz Universität Hannover, Germany**
- **Venkat Gudivada, East Carolina University, USA**
- **Pedro Martins, Universidade de Coimbra, Portugal**
- **Gary Weckman, Ohio University, USA**
- **Jolon Faichney, Griffith University, Australia**
- **Jedrzej Rybicki, Forschungszentrum Juelich GmbH, Germany**
- **Mario Monteleone, Università degli Studi di Salerno, Italy *[input only]***

BIG DATA are here... their challenges, too

- **Iryna:** The legal side of big data, e.g., what legal issues may arise by processing of big data, e.g. privacy or copyright concerns
- **Venkat:** Though the parallel computing models such as MPI, OpenMP, CUDA, OpenACC, and OpenCL existed for long, they failed in effecting widespread adaptation. In contrast, Hadoop parallel computing framework quickly became a mainstream and widely used parallel computing model. Of late, Apache Spark is emerging as a replacement for Hadoop? Is Spark the right Hadoop replacement? What happened to Apache Storm? What is the future of MPI, OpenMP, CUDA, OpenACC, and OpenCL? shorter: Is Spark the right Hadoop replacement? What happened to Apache Storm? What is the future of MPI, OpenMP, CUDA, OpenACC, and OpenCL?
- **Pedro:** Some relevant discussion topics regarding parallelization tools and their advantages and disadvantages.
- **Gary:** Big data analytics, issues with finding patterns in Big Data and issues in modeling
- **Jolon:** Open Data allows organisations to share data. However, opening Data can provide a number of challenges. We propose Open Schemas as a step before Open Data to allow organisations to share schemas. Opening Schemas will motivate organisations to provide cleaner data models and better schema documentation
- **Jedrzej:** I would like to pledge for the (sometimes overlooked) interplay between the data management and resulting/possible computational models. This is an open question if it is actually possible to define the computation model first and then find out a possible data management strategy to fit it. Or is the order of things in case of Big Data pre-defined and one has to use (perhaps less sophisticated) computation model since nothing else works with distributed data.

Q & A

Qs & As



WWW.IARIA.ORG

Performance Challenges

over **Big** and small data

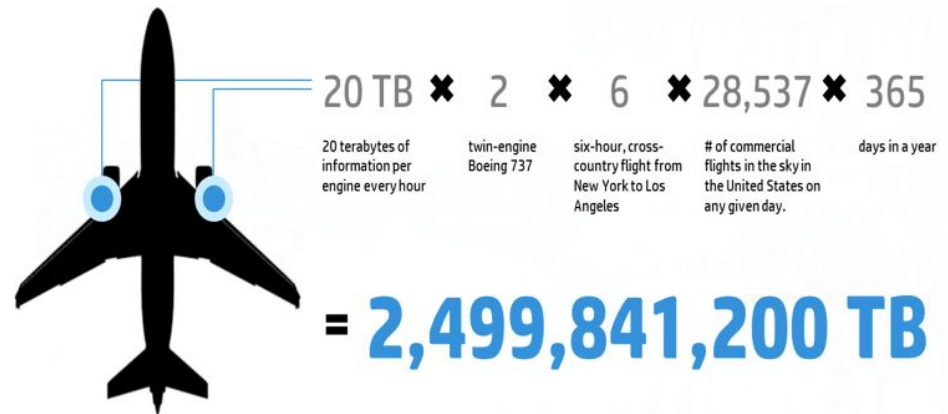
— Pedro Martins —

AllData 2016

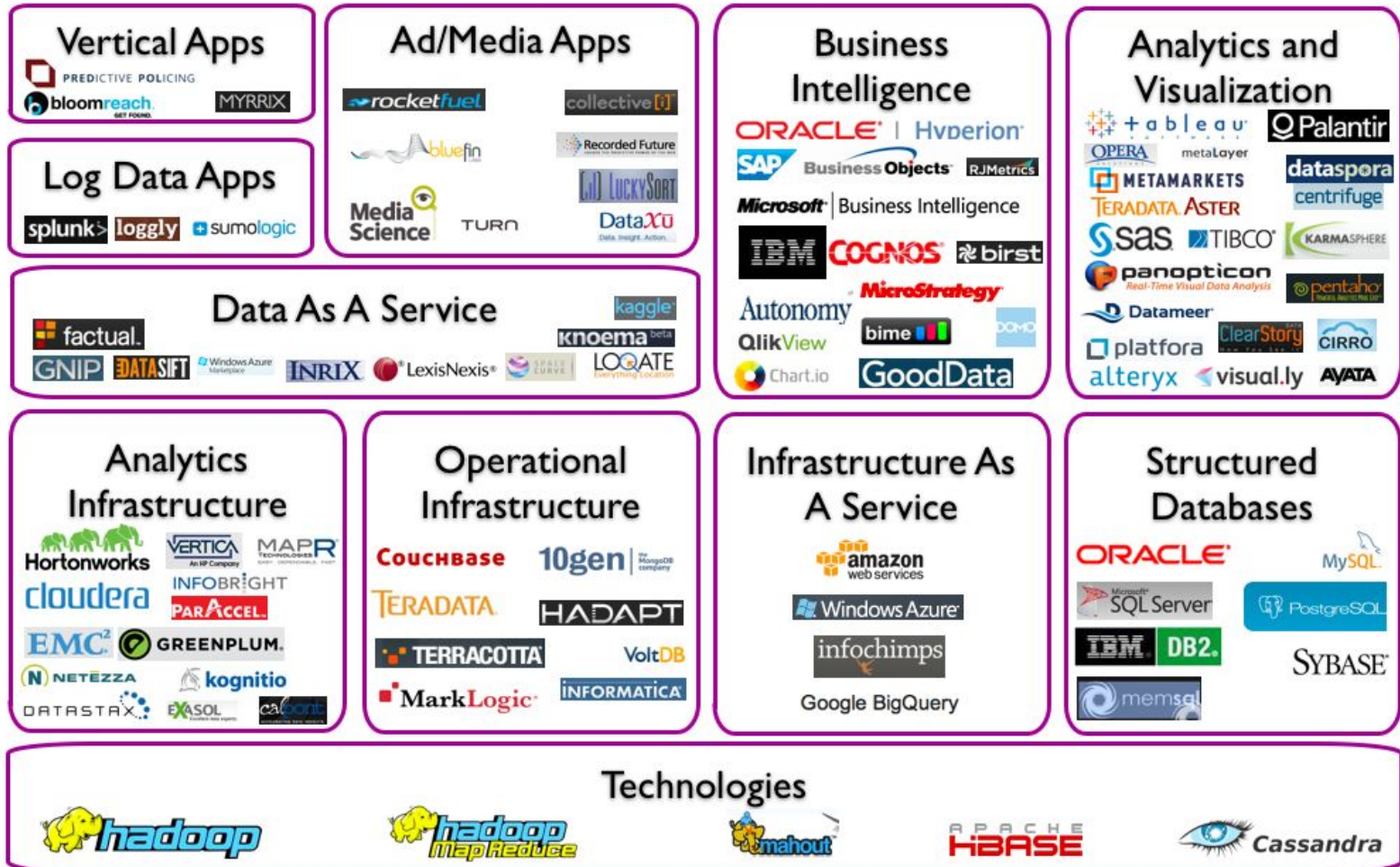
What is BigData ?



Sensor data from a cross-country flight



Big Data Landscape

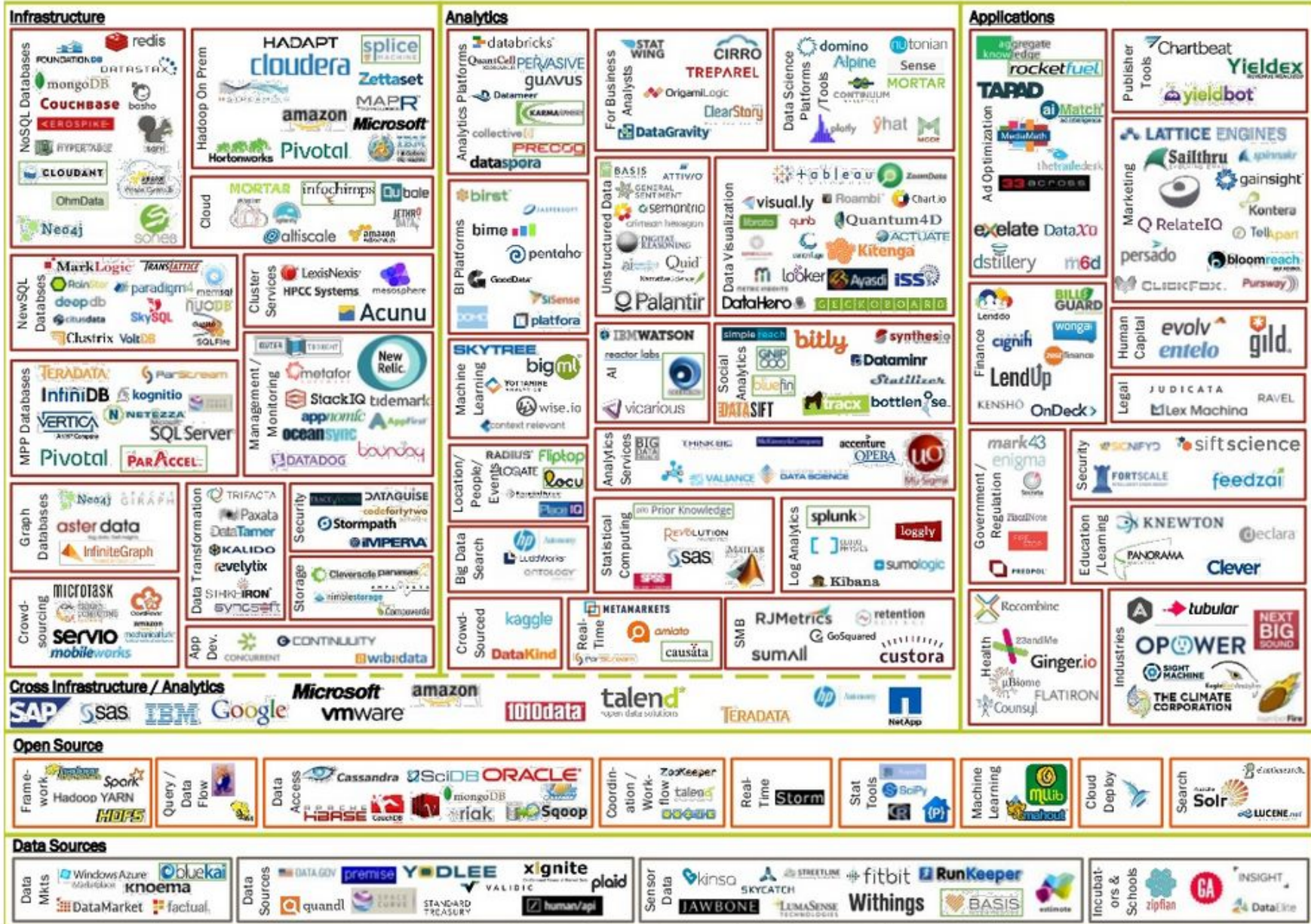


Big Data Landscape (Version 2.0)



BIG DATA LANDSCAPE, VERSION 3.0

Exited: Acquisition or IPO



| | | | | | | | | | | | |
|---|--|--|--|---|--|--|--|---|--|--|--|
| MARKETING EXPERIENCES Email Marketing Constant Contact, LiveIntent, Campaigner, StrongVoice, Bronco, dotMallor, yestmail, SMART FOCUS, Stream, Send Return Path, tower@ta zeta, JangoMail, Email-Kit, product, profusion, MailChimp, Mailer, Campaign Monitor, VerticalResponse, LeadSpend, epsilon, Experian, bluebonnet, CamSes, rapidmail, Movable, Customerio | | Creative & Design Adobe, balsamiq, offinnova, LiveReady, STEELHOUSE, AXURE, MockFlow, MATCH | | Calls & Call Centers keymetric, >ifybyphone, liveops, twilio, Five9, INVOKA, CallRail, Callmodo, CallFire, response, eGain, CallSource, tap, mongoosmetrics, Calltracking | | Marketing Apps wishpond, wizehive, SnapApp, offerpop, kontest, Woodbox, Gleam, voligo, CONTESTFACTORY, VIRALSWEEP, Wufoo, SocialAppsHQ, Rafflecopter, SurveyMonkey, PollDaddy, NORTH STACY, tagga, Strutta, conduit, SOCIAL, FluidSurveys, snap, surveygizmo | | Marketing Data epsilon, acxiom, NETPROSPER, ALLANT, InsideView, SalesLOX, leadspace, Infogroup, informatica, OneSource, factual, LiveRamp, AccuData, Gnip, VisualDNA, salesforce, Freshbase, IRI, Lead411, spoke, DATASIFT, STRICKIRON | | MARKETING OPERATIONS Marketing Analytics BECKON, collective, adomatrix, visualiq, nielsen, PivotLink, ThinkVme, converto, MARKETING EVOLUTION, Ijento, CONVERTMARKET, quantcast, distillery, COMSCORE, MarketShare, Anoss, Anamatrix, marketing, QED, Spot.tv, ACE, biScience, MOAT, measureful | |
| Mobile Marketing Tapjoy, Veeva, SponsorPay, ShopPad, hikerickit, AIRBURN, REVMOB, LeadBolt, Moba, mobilestorm, Trumba, Placed, INMOBI, OpenMarket, Prizmo, Kahuna, waterfall, media, VERVE | | Communities & Reviews jive, gigya, Ning, bazaarvoice, forumbe, satisfaction, ngage, socious, chute, communispace | | Customer Experience/VoC KANA, MEDALLIA, customerville, VERINT, perceptions, Confirm, CLARABRAND, Genesys, enkata, mindshare, gainsight, Kanyte, SATMETRIX, Qualaroo, ALLEGIANCE, USERVOICE | | SEO BRIGHT EDGE, altruiq, SYCARA, lotusjump, RANKABOVE, conductor, Wordtracker, AuthorityLabs, Shifts, rio, seo, seoClarity, calibri, UpCity, bloomreach, web ceo, traffic, yoast, SERPICO, yoast, ANALYTICS, SERPICO, linkdex, MOZ, SEOptics, searchmetrics, GinzaMetrics | | Channel/Local Mktg PTC, Zift-solutions, Gorilla, netafive, Channeltivity, brandmusic, nitromojo, Balihoo | | Dashboards Dundas, Logi, RJMETRICS, Dashboards, sweetspot, LEFTRONIC | |
| Search & Social Ads AdProof, KENSHOO, ACDHEMY, Marin, ACQUISIO, matchcraft, brandnetworks, Adaptly, TRIGGER, ADSPERT, nanigans, InsideVault, SearchForce, FINCH, sidetrack, MAGNETIC, SHIFT | | Social Media Marketing SOCIALFLOW, sprinklr, bottlen, ATTENTIFY, Socialware, Engagor, Visible, WILDFIRE, hootsuite, virah, shoutlet, awareness, trackur, newsie, sproutsocial, ArgyleSocial, EXPION, trackMaven, ArgyleSocial, moxie, Addvocate, Buzzint, ListenLogic, LITTLE BIRD, Lithium, UNIFIED, synthesio, infigy, @fluitive, mzingo, CAMPUSLY, hear, META, META, NETWORKED, socialbakers, Simply Measured, CURALATE, conversal, thismoment, EngageSciences, SocialMetrics, Zuberance, BuzzBundles, meitwater, Analytica | | Loyalty & Gamification Badgeville, CROWDTWIST, TIBCO, SOCIALTWIST, loyaltygator, PunchTab, vermt, 500Friends, AIMIA, REWARDSTREAM, Actionable, BUNCHBALL | | Personalization Evergage, Baynote, RapLeaf, DEMANDBASE, LIVEPERSON, monetate, MONOLOOP, PREDICTA, enecto, MYBUYS, richrelevance, APPTUS, SAILTHRU, choice, stream, PERSONYZE, CERTONA, PREDICTIVE EDGE, MAGIQ | | Content Marketing kapost, curata, cadence, Zerys, Springcm, DIVVY, Percolate, Skyword, COMPENDIUM, Contently, CMM, GatherContent, Zemanta, MARKETING.AI, collective, PAPERSHARE, Kontera, eassely, utobrain, NewsCred, atomic reach, RebelMuso, oero, Scoop, CrowdSource, interwet, publishthis, cloudwords, Corporate, trapit, uberflip, vorak, copylogger, magnify, Storfory, scripted, tribrer | | Web & Mobile Analytics Google, CLIC, TALET, @KISSmetrics, crazyegg, IBM, webtrends, compete, ServiceTick, EVIDON, spring, VisiStat, mixpanel, mouseflow, Kontagant, inspectlet, evolution, bitly, Clickdensity, Localytics, GoSquared, celebrus | |
| Display Advertising doubleclick, DataXu, bizo, sitescout, BRANDSREVENUE, neustar, DATA, netmining, AdRoll, CRITEO, TruSignal, OpenX, elite, Admant, rocketfuel, kwanzoo, bluecava, OPTIMIZE, Simplifi, Chango, Taykey | | Events & Webinars Eventbrite, CITRIX, Brandscopic, Acteva, ONZ4, INXPO, xert, FUZE BOX, InfoNeedle, InterCall, tmi, implex, etouches, MeetingBurner, TalkPoint, Webinars, saba, cornex, LENOS, Webinars, EventKloud, ReadyTalk | | Testing & Optimization on, webtrends, maxymiser, Google, Wingify, SITESPECT, Optimizely, APT, accenture, Adobe, unbounce, abmio, Useit, Testing.com, Lander, Pagella, Experiment.ly, GROUPTAG, pluralr, INSTA PAGE, SYNPERENCE, SeeWiz, CONVERSION, VERTster, NELIO, A/B Tasty, zen, AVENSED, SparkPage, convert, KAMELEON, Marketizator, VIBETRACE, PAGEWIZ, GlobalMaxer, IMPREGO, LeadPages | | Sales Enablement elastic, postwire, cloze, innovation, Stride, Bloomfire, Salespod, Qvidian, Yesware, SCEDOS, SKURA, MindMatrix, clearslide, KnowledgeTree, shoupod, Allmean, WittyParrot, Aligned, pipedrive, SAVO, Primary Intelligence, Velocity, UpSync, contactually, pipeliner, TappCtrl | | Digital Asset Mgmt WIDEN, EMC, BYNDER, Adgistics, CELLIUM, MOSAIC, ADAM, Cognizant, mediaevault, DigiEyeZ, AssetBank, WEDAM, MEDIA, BEACON, brandwatch, Camio, ThirdLight | | Business Intelligence pentaho, alteryx, ORACLE, Information, QlikView, Data, SAP, Builders, METRIC INSIGHTS, GoodData, PROGNOS, tableau, BITAM, ACTUATE, SAS, Predixion, ASANA, TIBCO, entrinsti, Microsoft, PANORAMA, spogobi, REVOLUTION, THINKMAP, LANAWASTORA, IASPERSOFT, board, Entalysis, Yellowfin | |
| Video Ads & Marketing OOYALA, vimeo, brightcove, eyeview, Jivox, BrightRoll, WIZ, WISTIA, onscreen, ramp, SpotMixr, pixability, viduad, ramp, SpotMixr, vid, caster, spot, change, Optiomatic, piksel, ustudio, SundaySky, YuMe, Kaltura, MEASURES, videology, mixpo, iddler | | Events & Webinars | | Testing & Optimization | | Sales Enablement | | Digital Asset Mgmt | | Business Intelligence | |
| MIDDLEWARE Data Management Platforms/Customer Data Platforms blueai, exelate, optimove, LOTAME, PARIO, LYRICS, turn, Fabric, krux, RedPoint, infer, MINTIGO, IGNITIONONE, XAXIS, AGILONE, Knofice, core, audience | | Tag Management DC Storm, enlighten, ValueClick, TEALIUM, Qubit | | User Mgmt GIGYA, janrain, LoginRadius, OpenID, oneall | | Cloud Connectors snaplogix, Jitterbit, key, elastic, Zapier | | APIs Layer7, apigee, 3scale, Mashery | | | |
| BACKBONE PLATFORMS CRM Microsoft, Nutshell, insightly, ORACLE, salesforce.com, netsuite, nimbix, RelatedQ, BUDGARCRM, SAS, action, CAPILLARY, Pega, fullcircle, SW, Steelwedge, OnePageCRM, saleslogix, salesforce, Sage, SAP, Highrise, suite, LOGICBOX | | Marketing Automation / Integrated Marketing TERADATA, IBM, VIZ, SharpSpring, Spectate, whatsthexx, ORBR, VOCUS, pardot, salesfusion, lyris, Leadberry, hatcheduck, Marketo, SAS, action, GENGO, SIGNAL, target360, interact, HubSpot, ExactTarget, SALESIGN, LeadLife, integrate, etrigo, NBT-Results, leadsquared, Intercom, RIGHTON, Bislr, genius, ORACLE, ONTRAPORT, GreenRope, LEADLANDER, Venntive, Infusionsoft, RightWave, ClickSquared, INBOX25, CallidusCloud, SIGNPOST | | Web Site / WCM / WEM dynamicWEB, Aektron, SDL, Acquia, OpenCms, Wordpress, CRAFT, Adobe, OPENTEXT, e-Spirit, Joomla!, atex, Dupeal, BRIDGELINE, IBM, COREMEDIA, eZ, INGENUIX, WIX, EPISERVER, elcom, hp, Kentico, GX, Paper Thin, weebly, sitecore, Kentico, GX, Limelight, ORACLE, censhare, GILTY, dotCMS, percussion, HFFPO | | E-commerce INTERSHOP, ZOOPY, demandware, Bigcommerce, Commerce, go-cart, avangate, hybrid, Node, volusion, spree, mozu, clearbridge, Goodsie, mozu, ULTRACART, ebay, enterprise, nimblecommerce, commerce, Digital River, ORACLE, vinda | | | | | |
| INFRASTRUCTURE Databases ORACLE, Microsoft, IBM, MySQL, HYPERBASE, MarkLogic, Cassandra, mongoDB | | Big Data Hadoop, splunk, Zettaset, Hortonworks, data, BitYota, IBM, Pivotal, DATASTAX, MAER | | Cloud SCALAR, heroku, flexscale, max.com, rackspace, CloudPulse, Atamai, DigitalOcean, amazon, Microsoft, GIGASPACE, verizon | | Mobile App Dev Google, Microsoft, Apsalar, TapStream, Parse, Kinvey, PIMMO, DISTIMO, Xamarin, Proga, netbiscuits, IBM, kony, swrve, VESSEL | | Web Dev django, python, JavaScript, node, #.net, stackoverflow, Bootstrap | | Marketing Environment Google, Microsoft, ebay, facebook, LinkedIn, twitter, YAHOO!, Pinterest, YouTube, Quora, reddit, slideshare, ampyspace | |

How to achieve more performance ?

E (extraction)

- from data sources
- cell phone towers
- supermarket network

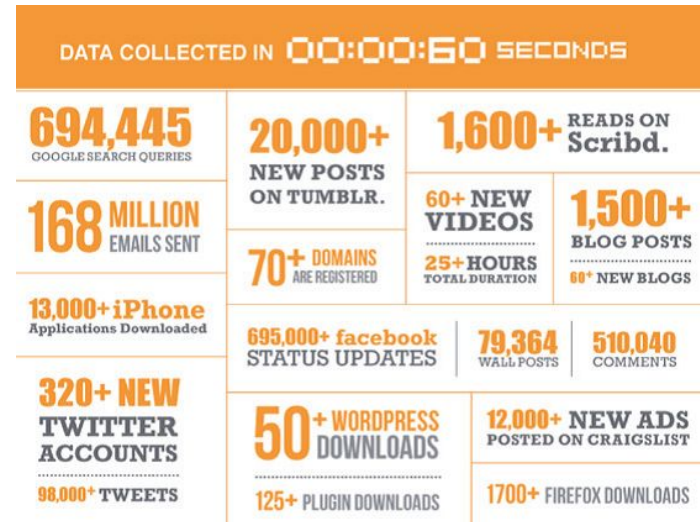
T (transformation)

- filter, data cleansing, validation, sort...

L (load)

- into any DW for later analysis

Q (query)



How to achieve more performance ?

Application solution?

Oracle, MySQL, Vertica, TeraData, MapReduce Architecture, etc

Is MapReduce the “holy grail”?



How to achieve more performance ?

Application solution?

(e.g. Oracle, MySQL, Vertica, TeraData, MapReduce Architecture)

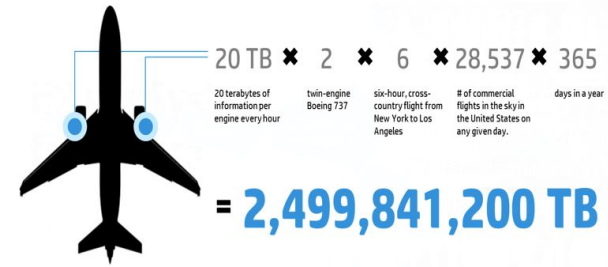
OR

Divide to conquer?

And this?

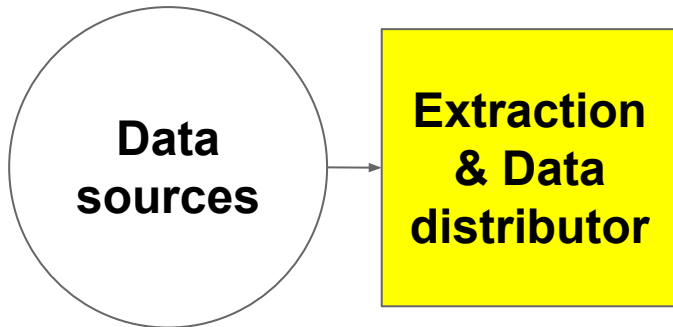
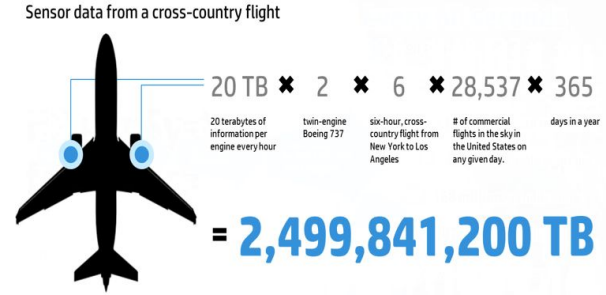


Sensor data from a cross-country flight

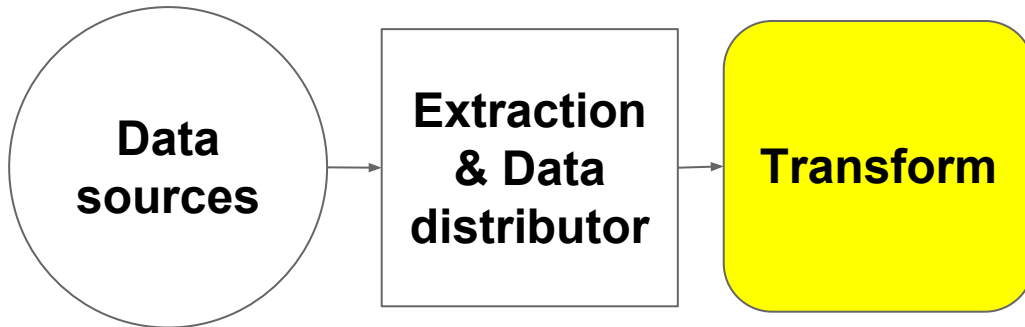
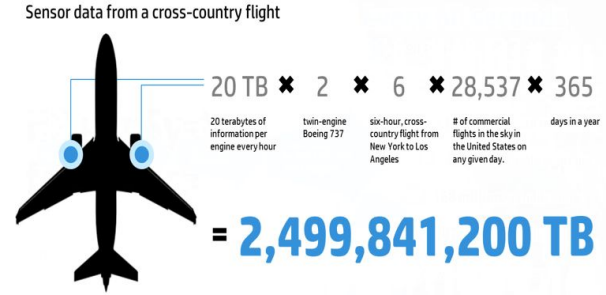


Data sources

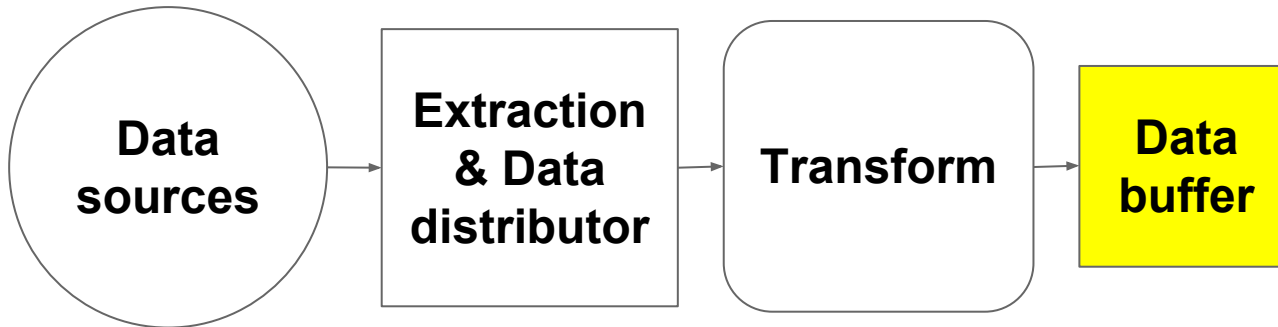
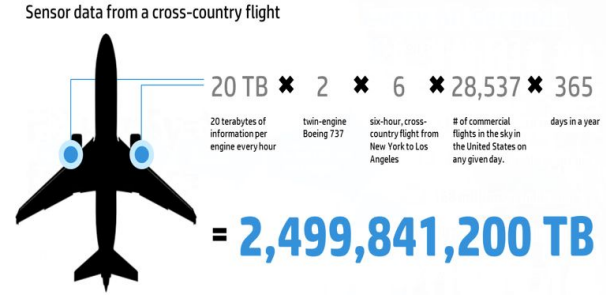
And this?



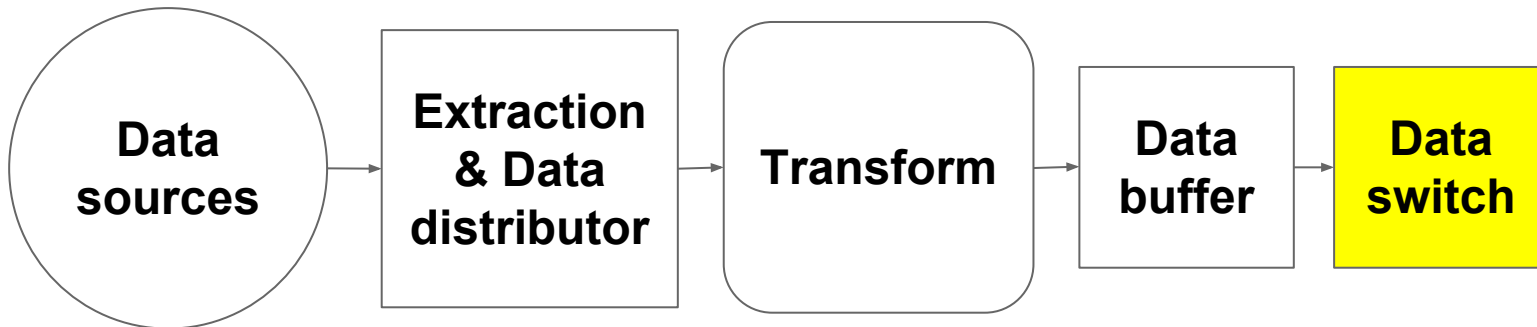
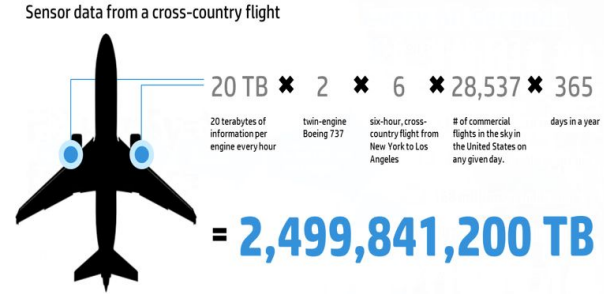
And this?



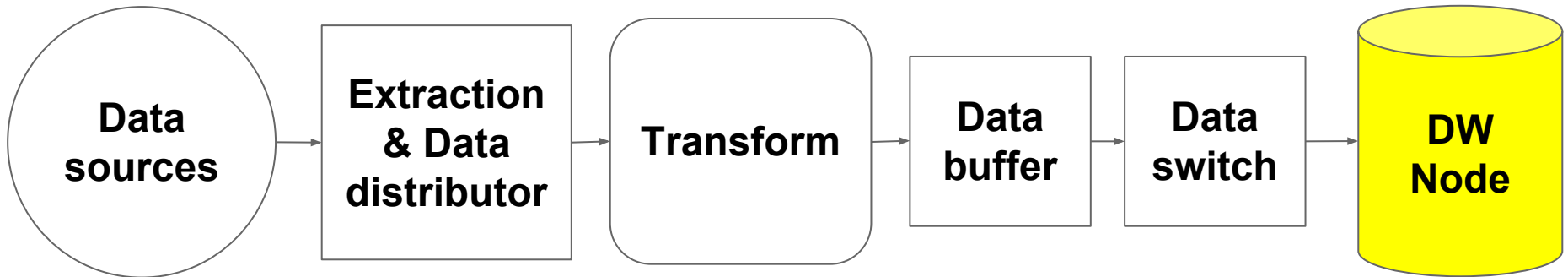
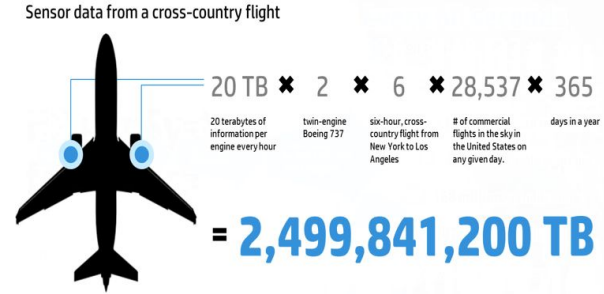
And this?



And this?

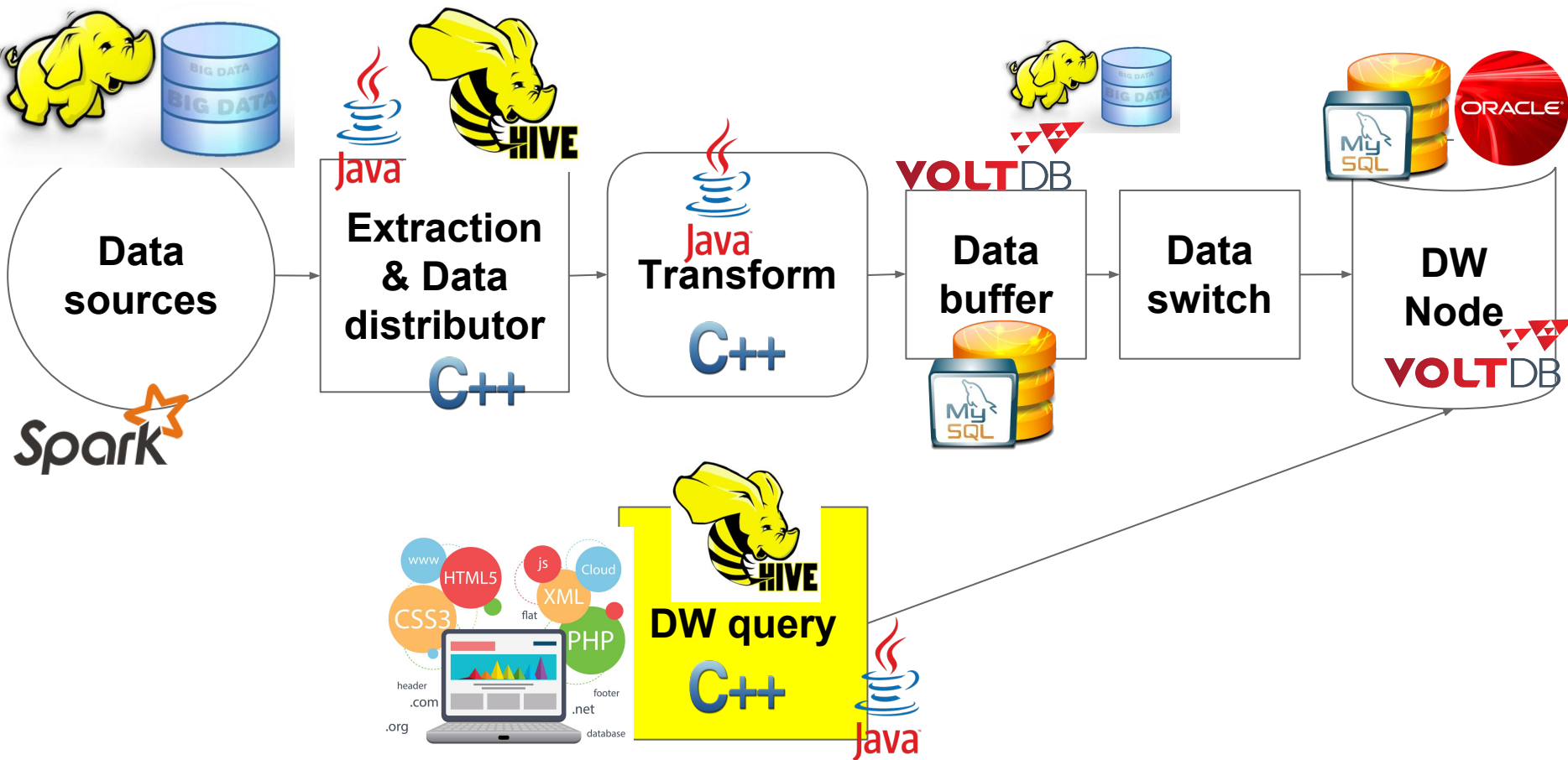


And this?



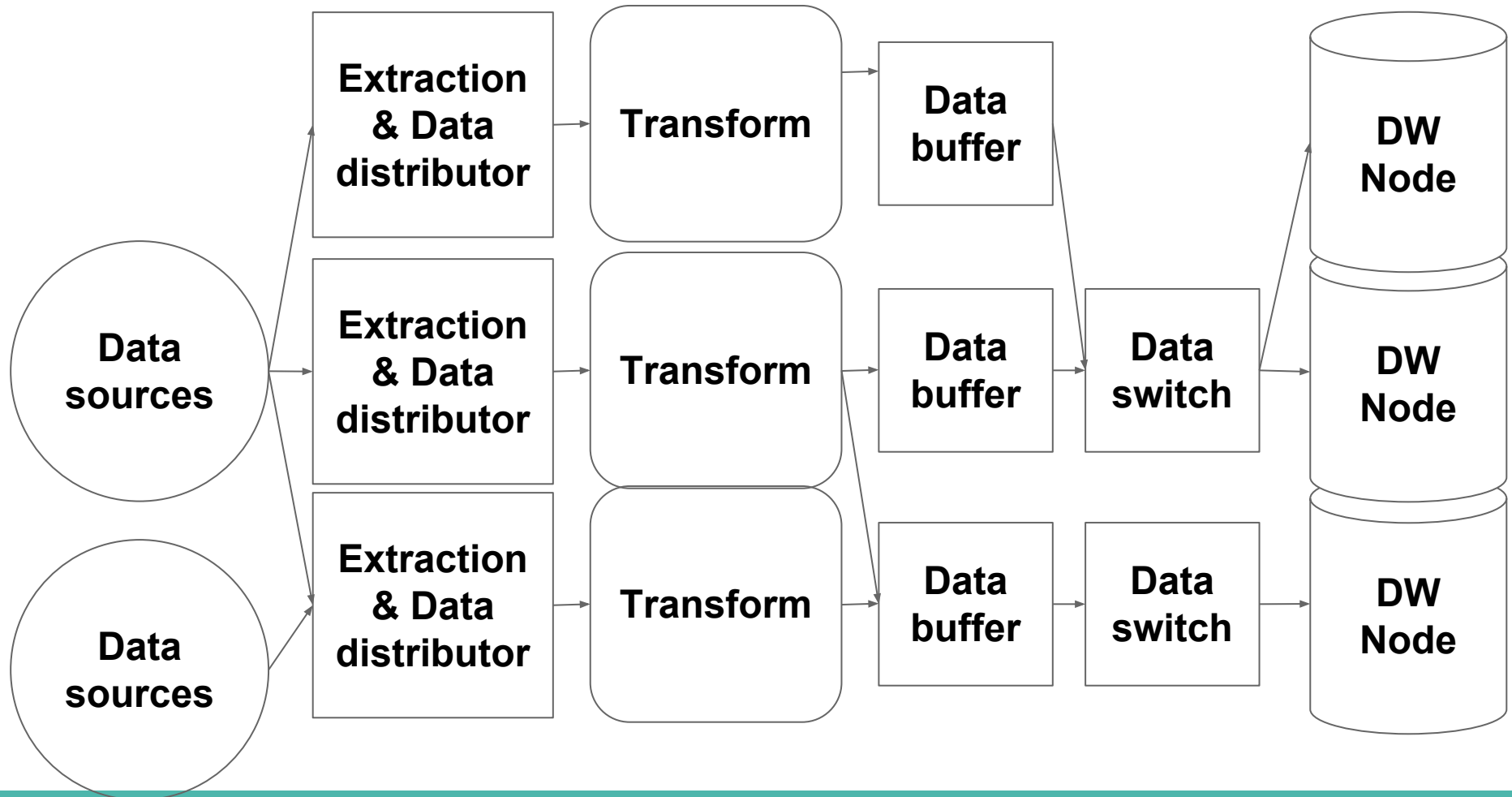
And this? - We call it AScale

(Automatic scaler, developed at University of Coimbra - Pedro Martins)



And this?

- A: scale-out (horizontally);



“The simplest explanation is usually the correct one”

William Ockham

Questions & Comments ?

Is MapReduce the “holy grail”?

Why not Map-Reduce architectures? (VLDB 2015)

Re-use permitted when acknowledging the original © Daniel Abadi, Shivnath Babu, Fatma Ozcan, and Ippokratis Pandis (2015)

Presenters

2



Fatma Özcan
IBM Research
IBM Big SQL



Daniel Abadi
Yale University and
Teradata
HadoopDB/Hadapt



Ippokratis Pandis
Cloudera
Cloudera Impala



Shivnath Babu
Duke University
Starfish

Is MapReduce the “holy grail”?

Why not Map-Reduce architectures? (VLDB 2015)

- **Little control of data flow**
- **Fault tolerance guarantees not always necessary**
- **Does not interface with existing analysis software**

Is MapReduce the “holy grail”?

Why not Map-Reduce architectures? (VLDB 2015)

- **Little control over storage (Append only file system)**
- **Little control over resource management**
- **Often used for “data dump” (irregular and unreliable)**

Computational Models for Big Data Processing

Venkat N Gudivada
East Carolina University
Greenville, North Carolina
USA

High Performance Computing Models

- Availability of multiprocessor and multi-core chips and GPU accelerators at commodity prices
- Personal Supercomputers
- Cloud-hosted Cluster Computers for the masses
- Parallel computing for the elite (<http://top500.org/>)

High Performance Computing Models - Shared-memory Paradigm

- OpenMP (compiler directives)
- OpenACC (targets acceleration devices)
- CUDA (architecture and a programming model)
- OpenCL (heterogeneous platforms)
- Haskell Concurrent Programming Model

High Performance Computing Models - Distributed Memory Paradigm

- MPI (message passing)
- MapReduce (Hadoop Ecosystem – HDFS, MapReduce, Pig, Pig Latin, Hive, Cascading, Scalding, Cascalog, Storm, and Spark; Google Cloud DataFlow)
- Erlang Message Based Concurrent Programming Model

Future Trend

- Personal supercomputers (packaged and pre-installed applications)
- Parallel computing for the masses
- Exascale computing for the elite

Open Schemas

A step towards Open Data

Dr Jolon Faichney

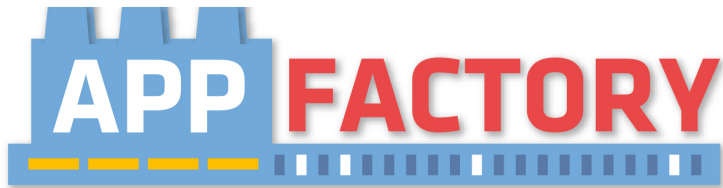
School of Information and Communication
Technology

Griffith University, Australia

j.faichney@griffith.edu.au



- Lecturer, Griffith University, Australia
- Founder, App Factory Student Enterprise
- Member of ODIQ Certificate Localisation Working Group



What is Open Data?

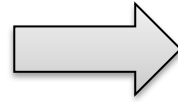
- Governments in democratic countries serve the public:
 - Therefore public should have access to data
- Transparency
- Accountability
- Productivity

- Founded by Tim Berners-Lee and Nigel Shadbolt in 2012
- UK-based, worldwide
- Promote the concept of:

"Open by Default"



Traditional Government Data



Pre-Web Government Data



Previously, access to Government data required person-to-person interaction

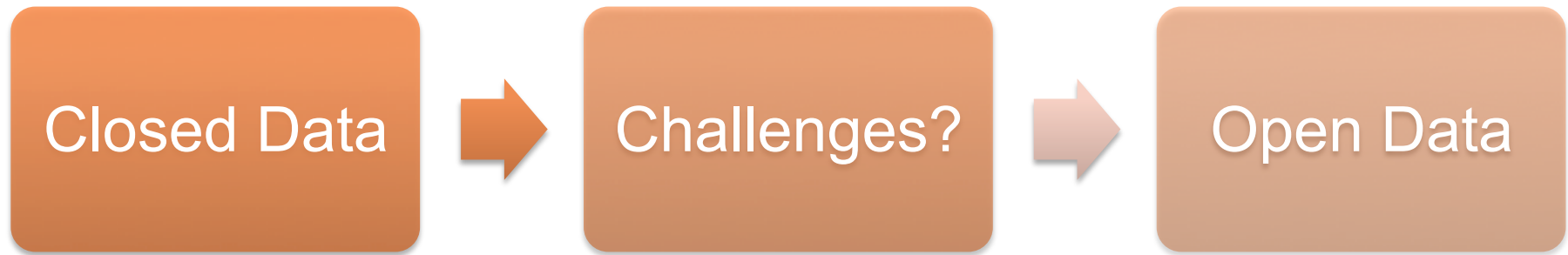
Post-Web Government Data



Open Data Challenges

- Cultural/policy challenges
 - Policy changes required
 - Resistance to data sharing within organisation
 - Exposure of poor/inaccurate quality data
- Technical challenges
 - Writing new software/APIs
 - Realtime data
- Privacy issues (de-identifying data)
- Resources/costs required

Open Data Challenges



Open Schema Approach



- Encourage departments to Open Schemas as a first step towards Open Data
- Advantages:
 - Opening Data has policy challenges such as privacy
 - Can also be applied to data that won't be opened
 - Improves interoperation within organisation
 - Improves interoperation external to organisation
 - Work towards consistent schemas
 - Provides the public with an indication of the availability of data leading to demand-driven open data (DDOD)

- Challenges:
 - Schemas can be messy
 - Not intended for public viewing
 - Evolved over many years and many people
 - Schemas may be poorly documented
 - Exposing schemas may be perceived as a security risk
- But
 - May motivate organisations to maintain well designed and documented schemas

- ODI's Open Banking Standard:
 - <http://theodi.org/open-banking-standard>



- DDI:

- The Data Documentation Initiative (DDI) is an international standard for describing statistical and social science data.
- Documenting data with DDI facilitates interpretation and understanding -- both by humans and computers.
- The freely available international DDI standard describes data that result from observational methods in the social, behavioral, economic, and health sciences.
- Use DDI to **D**ocument, **D**iscover, and **I**nteroperate



- Google, Microsoft, Yahoo, Yandex
- Entities, Relationships, Actions

MedicalCondition

[Thing](#) > [MedicalEntity](#) > [MedicalCondition](#)

Any condition of the human body that affects the normal functioning of a person, whether physically or mentally. Includes diseases, injuries, disabilities, disorders, syndromes, etc.

Usage: Between 100 and 1000 domains

[\[more...\]](#)

| Property | Expected Type | Description |
|--|---|---|
| Properties from MedicalCondition | | |
| associatedAnatomy | AnatomicalSystem or SuperficialAnatomy or AnatomicalStructure | The anatomy of the underlying organ system or structures associated with this entity. |

Topics for Discussion

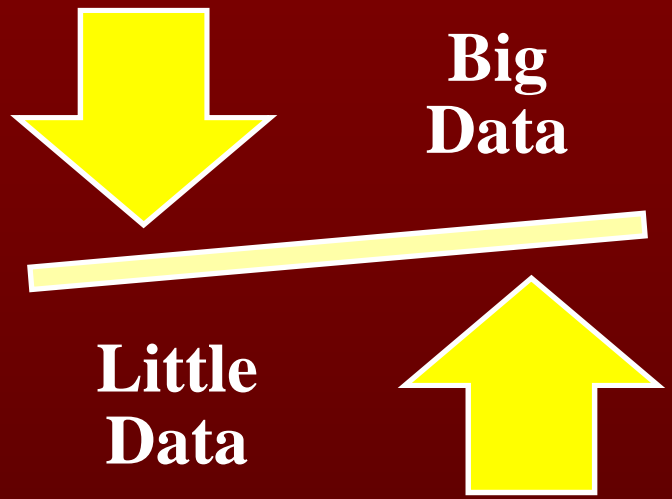
- Is "Open by Default" too ambitious/impractical?
- What challenges do organisations face in Opening Schemas?
- What benefits do Open Schemas provide?
- How can we work towards consistent schemas?

Topic: Computational Models for Big Data Processing

Gary R. Weckman

Ohio University

- **Means in which data can be collected is more readily available than ever**
- **Big Data more relevant than ever because it can be used to improve decisions and insights within the domains it is used**
- **The term Big Data can be loosely defined as data that is too large for traditional analysis methods and techniques.**



Big: Random reduction
Little: Synthetic (SMOTE)
Imbalance Data
Missing Data

Big data analytics

issues with finding patterns in Big Data and issues in modeling

- **Due to the volume of Big Data**
 - **Traditional methods of analysis can function poorly**
 - **finding patterns that do not exist (Ratner, 2003)**
- **As datasets grow in size (leaning towards Big Data)**
- **spurious structures tend to also be discovered in**
- **datasets and interpreted as meaningful when they**
- **have no true meaning (Ratner, 2003).**

Table 5: Pattern Recognition Methods Applicable to Big Data

| Topic | Reference and Examples |
|-----------------------------------|--|
| Anomaly Detection | A form of unsupervised classification where statistically rare events are of interest. Useful for removing background from data. Methods: PCA reconstruction error/residual analysis (Jackson & Mudholkar, 1979) (Jablonski, Bihl, & Bauer, 2015), Mahalanobis distance (De Maesschalck, Jouan-Rimbaud, & Massart, 2000) , RX detector (Reed & Yu, 1990) |
| Artificial Neural Networks | A variety of nonlinear classifiers that employ gradient descents and a variety of algorithms to classify and predict data Methods: Feedforward, recurrent, and self-organizing map methods (Young, Bihl, & Weckman, 2014) (Ward, Bihl, & Bauer, 2014) |
| Classification | Applying methods to create a model that accurately represents the data with respect to known classes. Methods: ANNs (Duda, Hart, & Stork, 2001), Classification Trees (Duda, Hart, & Stork, 2001), Discriminant Analysis (Dillon & Goldstein, 1984), Logistic Regression (Hosmer & Lemeshow, 2000), |
| Class Imbalance | Dealing with imbalances in classes, which can create bias in analytical methods Methods: Over-sampling, boosting, bootstrapping (Zhang J. , 2004), artificial sampling (Bui, 2004) |
| Clustering | Clustering, or unsupervised classification, refers to methods that search for underlying patterns in data Methods: Hierarchical (distance and linkage based) (Gordon, 1987) (Johnson, 1967) (Milligan & Cooper, 1987) (Milligan & Cooper, 1985), k-means (Jain, 2010), affinity propagation (Frey & Dueck, 2007), density based (Ester, Kriegel, Sander, & Xu, 1996), and other methods (Jain, 2010) |
| Crowd Sourcing | Employing multiple users to analyze data, contribute to a solution, or leverage their computer power Methods/Examples: Crowd sourced games to find patterns (Mavandadi, et al., 2012) (Martin, et al., 2013), multiple opinions for clinical data (Celi, Mark, Stone, & Montgomery, 2013), distributed projects (Schreiner, 2001) |
| Dimensionality Reduction Analysis | Dimensionality reduction through transforming data into a new space (feature extraction) or selecting subsets of original data features (feature selection). Methods: Principal Component Analysis (PCA) (Dillon & Goldstein, 1984), Factor Analysis (FA) (Dillon & Goldstein, 1984), Independent Component Analysis (Jain, Duin, & Mao, 2000), Kernel methods (Jain, Duin, & Mao, 2000), Stepwise, forward, and backward selection methods (Jain, Duin, & Mao, 2000), ANN signal to noise ratio feature screening (Bauer, Alsing, & Greene, 2000), input reduction (Young W. , Weckman, Thompson, & Brown, 2008), Wilk's Lambda (Dillon & Goldstein, 1984) (Eisenbeis, 1977), F-test (Bihl, Temple, Bauer, & Ramsey, 2015), and other methods (Jain, Duin, & Mao, 2000).methods (Jain, Duin, & Mao, 2000) |
| Imputation | Filling in missing observation through various methods, missing data can appear randomly through a dataset, in rows (e.g. survey responses), or in columns (possibly from a bad sensor or attribute) |

MyHealthAvatar: Your Lifetime Companion for Healthcare



MyHealthAvatar

Computational Models for Big Data Processing Legal Perspective

Iryna Lishchuk, LL.M.
Institut für Rechtsinformatik
Leibniz Universität Hannover
Hannover, Germany

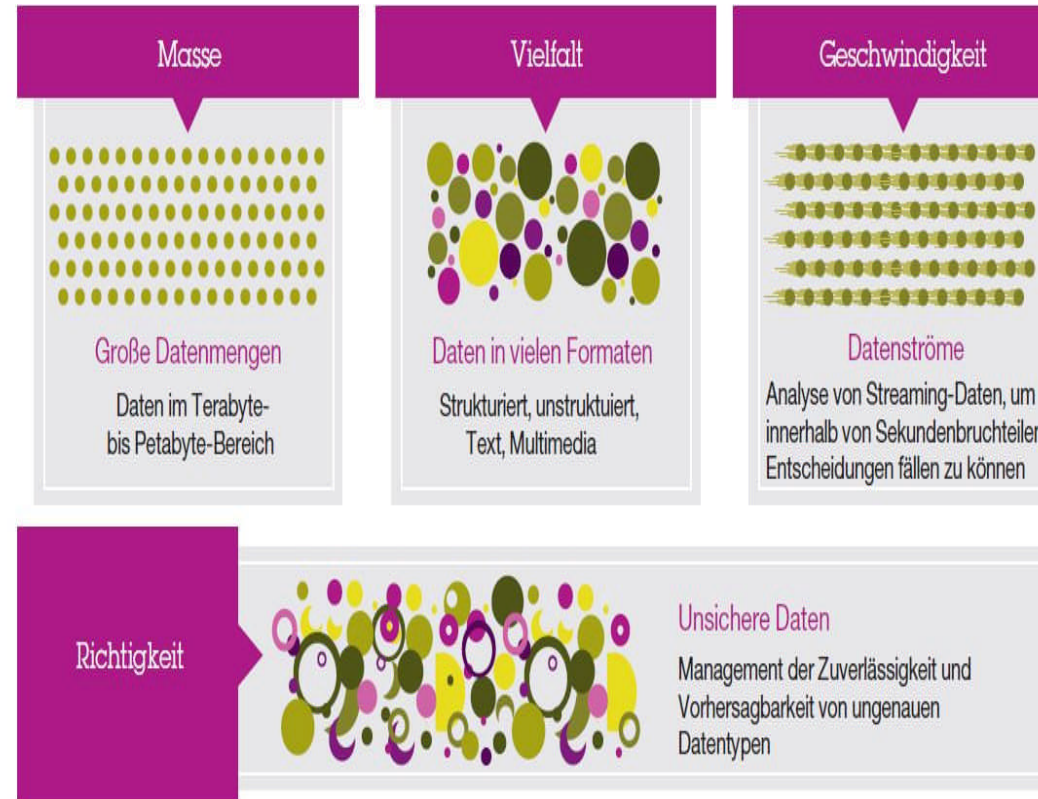
ALLDATA 2016, February 21 - 25, 2016 - Lisbon,
Portugal





Big Data - Concept

- Volume
- Velocity
- Variety
- Veracity



<http://www-935.ibm.com/services/de/gbs/thoughtleadership/GBE03519-DEDE-00.pdf>

ALLDATA 2016, February 21 - 25, 2016 - Lisbon, Portugal

Personal Data - Concept

- Article 2, Data Protection Directive 95/46/EC

“'personal data' shall mean any information relating to an identified or identifiable natural person ('data subject');

an identifiable person is one who can be identified, directly or indirectly, in particular by reference to an identification number or to one or more factors specific to his physical, physiological, mental, economic, cultural or social identity”

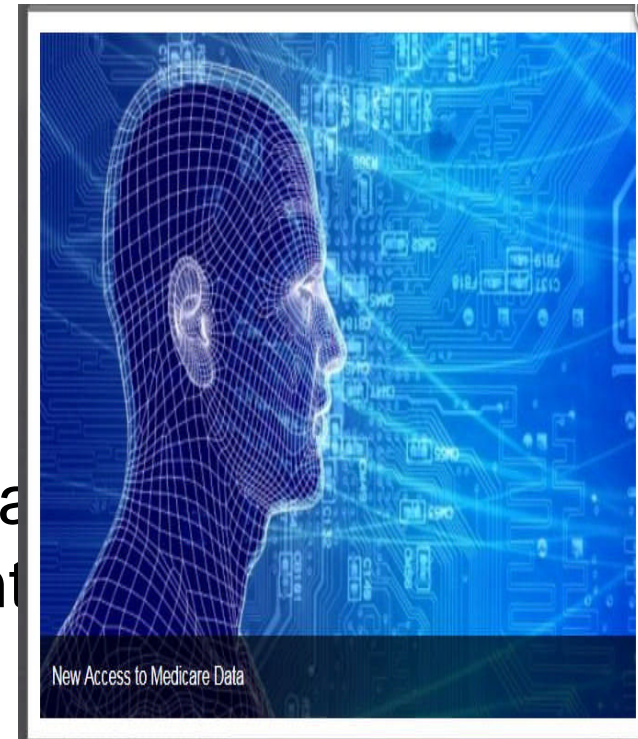


Data Linking v. Data Protection

→ Linking isolated data sets can make a person behind identifiable

→ Pure data procession can turn into procession of personal data

→ Procession of personal data subject to legal requirements





Data Linking v. Data Protection

Article 2 Data Protection Directive 95/46/EC

Procession, collection and use of personal data subject to consent of the data subject, which is legitimate, when:

- Free
- Express
- Purpose related
- Informed



Data Linking v. Data Protection

- Linking of already available and for specific purpose collected data may constitute change of purpose
- Change of purpose requires for legitimation consent of the data subject
- Data linking for a specific purpose may be justified if qualifies as compatible uses of data

Computational Models for Big Data Processing

Jędrzej Rybicki

23rd February 2016

Computer Science Archeology

Computer Science Archeology

```
shasum * | awk {'print $1'} | sort | uniq -c | grep -v " 1 "
```

Computer Science Archeology

```
shasum * | awk {'print $1'} | sort | uniq -c | grep -v " 1 "
```

- ▶ Unix pipelines

Computer Science Archeology

```
shasum * | awk {'print $1'} | sort | uniq -c | grep -v " 1 "
```

- ▶ Unix pipelines
- ▶ Ken Thomson et al. in PDP-11 (1972)

Computer Science Archeology

```
shasum * | awk {'print $1'} | sort | uniq -c | grep -v " 1 "
```

- ▶ Unix pipelines
- ▶ Ken Thomson et al. in PDP-11 (1972)
- ▶ limits on how big a single problem could be

Computer Science Archeology

```
shasum * | awk {'print $1'} | sort | uniq -c | grep -v " 1 "
```

- ▶ Unix pipelines
- ▶ Ken Thomson et al. in PDP-11 (1972)
- ▶ limits on how big a single problem could be
- ▶ no disks (or very small)

Data Management Models for Big Data Processing

Jędrzej Rybicki

23rd February 2016

Don't store everything

Don't store everything

“[...] in theoretical work it is often possible to cherry pick assumptions to produce a given result.”

— “Chameleons: The Misuse of Theoretical Models in Finance and Economics,” by Paul Pfleiderer

Don't store everything

“[...] in theoretical work it is often possible to cherry pick assumptions to produce a given result.”

— “Chameleons: The Misuse of Theoretical Models in Finance and Economics,” by Paul Pfleiderer

- ▶ no problem was solved just by putting more data on it

Don't store everything

“[...] in theoretical work it is often possible to cherry pick assumptions to produce a given result.”

— “Chameleons: The Misuse of Theoretical Models in Finance and Economics,” by Paul Pfleiderer

- ▶ no problem was solved just by putting more data on it
- ▶ monolith store vs. micro stores

Don't store everything

“[...] in theoretical work it is often possible to cherry pick assumptions to produce a given result.”

— “Chameleons: The Misuse of Theoretical Models in Finance and Economics,” by Paul Pfleiderer

- ▶ no problem was solved just by putting more data on it
- ▶ monolith store vs. micro stores
- ▶ point of highest ignorance

Avoid web scale envy

- ▶ perhaps the big data companies/scientist are where they are because they started small?

Avoid web scale envy

- ▶ perhaps the big data companies/scientist are where they are because they started small?
- ▶ do you think that MapRed is good?

Avoid web scale envy

- ▶ perhaps the big data companies/scientist are where they are because they started small?
- ▶ do you think that MapRed is good?
- ▶ *no!* nothing else is possible with the selected data management solution

Avoid web scale envy

- ▶ perhaps the big data companies/scientist are where they are because they started small?
- ▶ do you think that MapRed is good?
- ▶ *no!* nothing else is possible with the selected data management solution
- ▶ ... but perhaps you can use something else

Avoid web scale envy

- ▶ perhaps the big data companies/scientist are where they are because they started small?
 - ▶ do you think that MapRed is good?
 - ▶ *no!* nothing else is possible with the selected data management solution
 - ▶ ... but perhaps you can use something else
 - ▶ graph algorithms: example of how data management can influence computational models
- ⇒ vertex-oriented vs. pattern matching

Beware of difference between academia and private sector

- ▶ data sharing: crazy idea!

Beware of difference between academia and private sector

- ▶ data sharing: crazy idea!
- ⇒ sharing infrastructure?

Beware of difference between academia and private sector

- ▶ data sharing: crazy idea!
- ⇒ sharing infrastructure?
- ▶ data management for sharing might be different?

Beware of difference between academia and private sector

- ▶ data sharing: crazy idea!
- ⇒ sharing infrastructure?
 - ▶ data management for sharing might be different?
 - ▶ sharing with whom?
- ⇒ limit sophistication

Data Management Models influence Data Processing Models

Data Management Models for Big Data Processing:

- ▶ Don't store everything
- ▶ Avoid web scale envy
- ▶ Academia vs. private sector
- ▶ Don't be too sophisticated