

Databases & Data Mining

Erwin M. Bakker & Stefan Manegold

e.m.bakker

s.manegold

@liacs.leidenuniv.nl

<https://homepages.cwi.nl/~manegold/DBDM/>

<http://liacs.leidenuniv.nl/~bakkerem2/dbdm/>



DBDM: “Registration”

Please send an email

To: s.manegold@liacs.leidenuniv.nl

Subject: [DBDM-2018] Registration

containing the following information:

- *Your full name*
- *Your email address*
- *Your student ID*
- *Your affiliation (university)*
- *Your program / subject*

By Sunday 16 September 2018, 23:59 CEST.

DBDM: Overview

Period: September 11th - December 4th 2018 (Tuesdays)

Place: Room **312** (LIACS, Snellius building, Niels Bohrweg 1, 2333 CA Leiden)

Time: 15.30 - 17.15

ECTS: 6

Description:

The course Databases & Data Mining consists of a series of lectures in which advanced database and data mining techniques will be discussed, with applications to bioinformatics.

Grading:

There will be 2 database and 2 data mining assignments, i.e., 4 assignments in total, and a final exam (open book). The final grade will be based on a weighted average of the grades obtained for assignments P1, P2, P3, P4 and the Exam (E >5):

Final Grade = $(0.5 * P1 + P2 + 0.5 * P3 + P4 + 3 * E) / 6$.

DBDM: (tentative) Schedule

Date	Room	Subject (tentative)	Topic & Lecturer
11-09	312	Introduction	Databases and Data Management for Data Mining Stefan Manegold
18-09	312	Database Technology	
25-09	312	Database Technology	
02-10	312	Data Preprocessing	
09-10	312	<i>No class</i>	
16-10	312	Data Warehousing and OLAP	
23-10	312	Data Cube Technology	
30-10	312	Basic Data Mining Algorithms I	Data Mining Techniques and Applications Erwin Bakker
06-11	312	Basic Data Mining Algorithms II	
13-11	312	Advanced Data Mining Algorithms	
20-11	312	Mining in Bio-Data	
27-11	312	Graph Mining I	
04-12	312	Graph Mining II	

DBDM: Assignments

- 2 database assignments & 2 data mining assignments
- Will be announced individually during lectures and posted on website

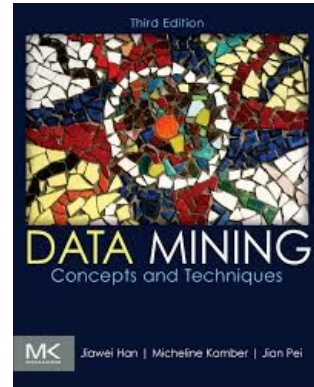
DBDM: Exam

- **open book exam:** you can take with you your book, and printed course notes (slides). *No electronic equipment is allowed, though.*
- **Materials to be studied:**
 - All content covered and discussed during lectures (slides will be shared).
 - More to be announced.
- **Date:** Monday, January 7, 2019
- **Time:** 14:00 - 17:00
- **Place:** Room F104, Van Steenisgebouw, Einsteinweg 2, 2333 CC Leiden

DBDM: Recommended Books

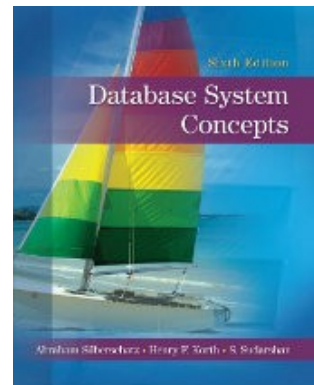
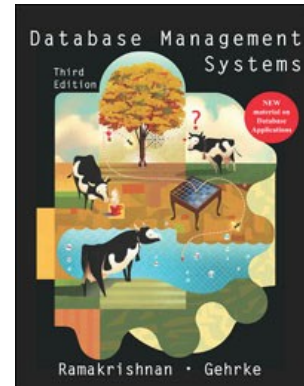
- **Data Mining:**

- *J. Han, M. Kamber, J. Pei. **Data Mining Concepts and Techniques (3rd Edition)**, Morgan Kaufman Publishers, July 2011 (ISBN 978-0123814791)*



- **Database systems (e.g.):**

- *Ramakrishnan, Gehrke: **Database Management Systems (3rd International Edition)**, McGraw-Hill, 2003 (ISBN 0-07-246563-8)*
- *A. Silberschatz, H. F. Korth, S. Sudarshan: **Database System Concepts (6th Edition)**, McGraw-Hill, 2010 (ISBN 0-07-352332-1)*



DBDM: “Registration”

Please send an email

To: s.manegold@liacs.leidenuniv.nl

Subject: [DBDM-2018] Registration

containing the following information:

- *Your full name*
- *Your email address*
- *Your student ID*
- *Your affiliation (university)*
- *Your program / subject*

By Sunday 16 September 2018, 23:59 CEST.

Databases & Data Mining

Stefan Manegold



Group leader Database Architectures
Centrum Wiskunde & Informatica (CWI)
Amsterdam

<http://homepages.cwi.nl/~manegold/>



<http://www.monetdb.org/>

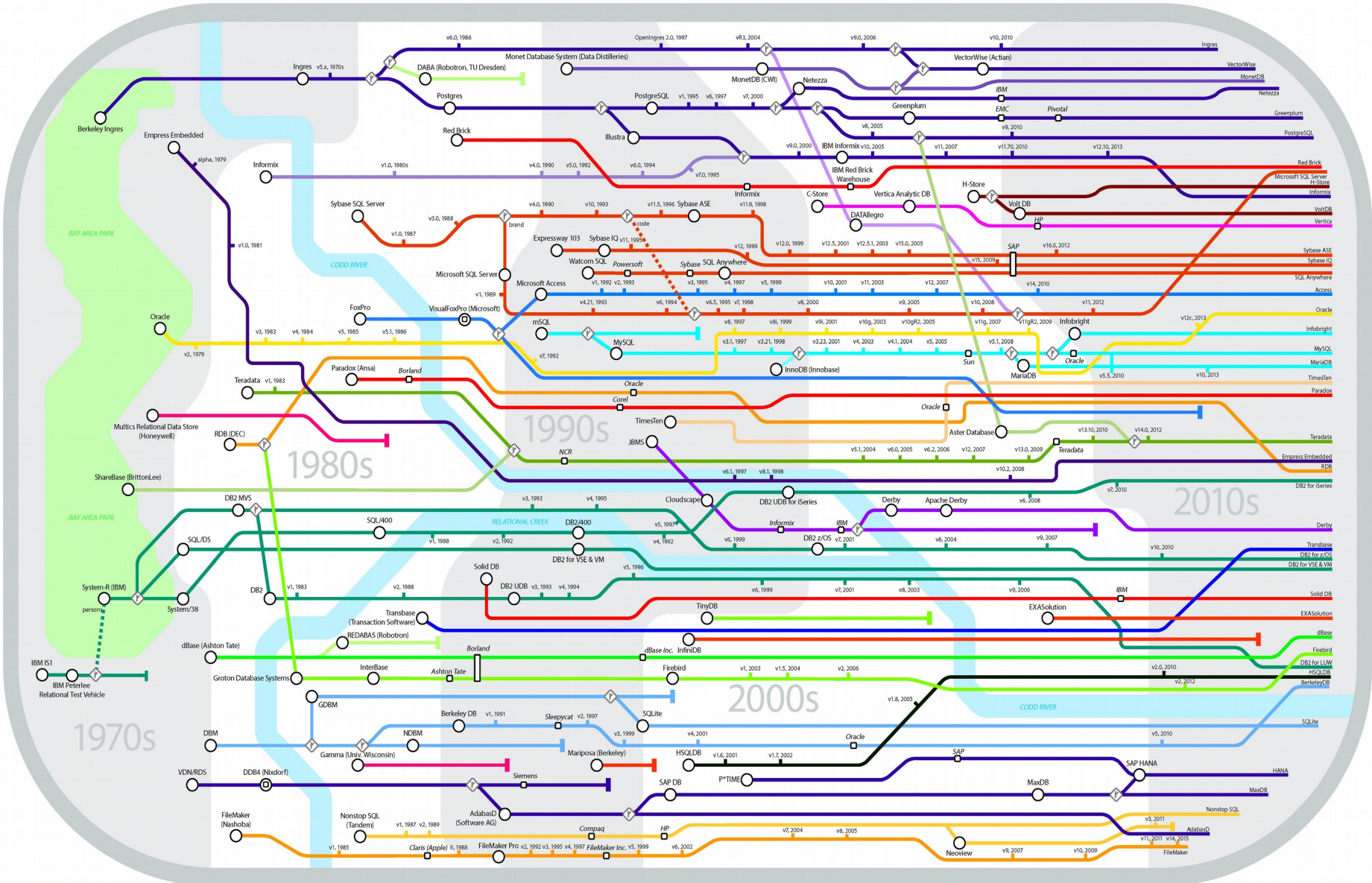


Prof. Data Management (0.2 fte)

LIACS & LCDS

Faculty of Science, Leiden University

Genealogy of Relational Database Management Systems



Key to lines and symbols

- Publishing Date
- ◻ Acquisition
- ↕ Versions
- ⊥ Discontinued
- ◇ Branch (intellectual and/or code)
- Crossing lines have no special semantics



Data



Data Management

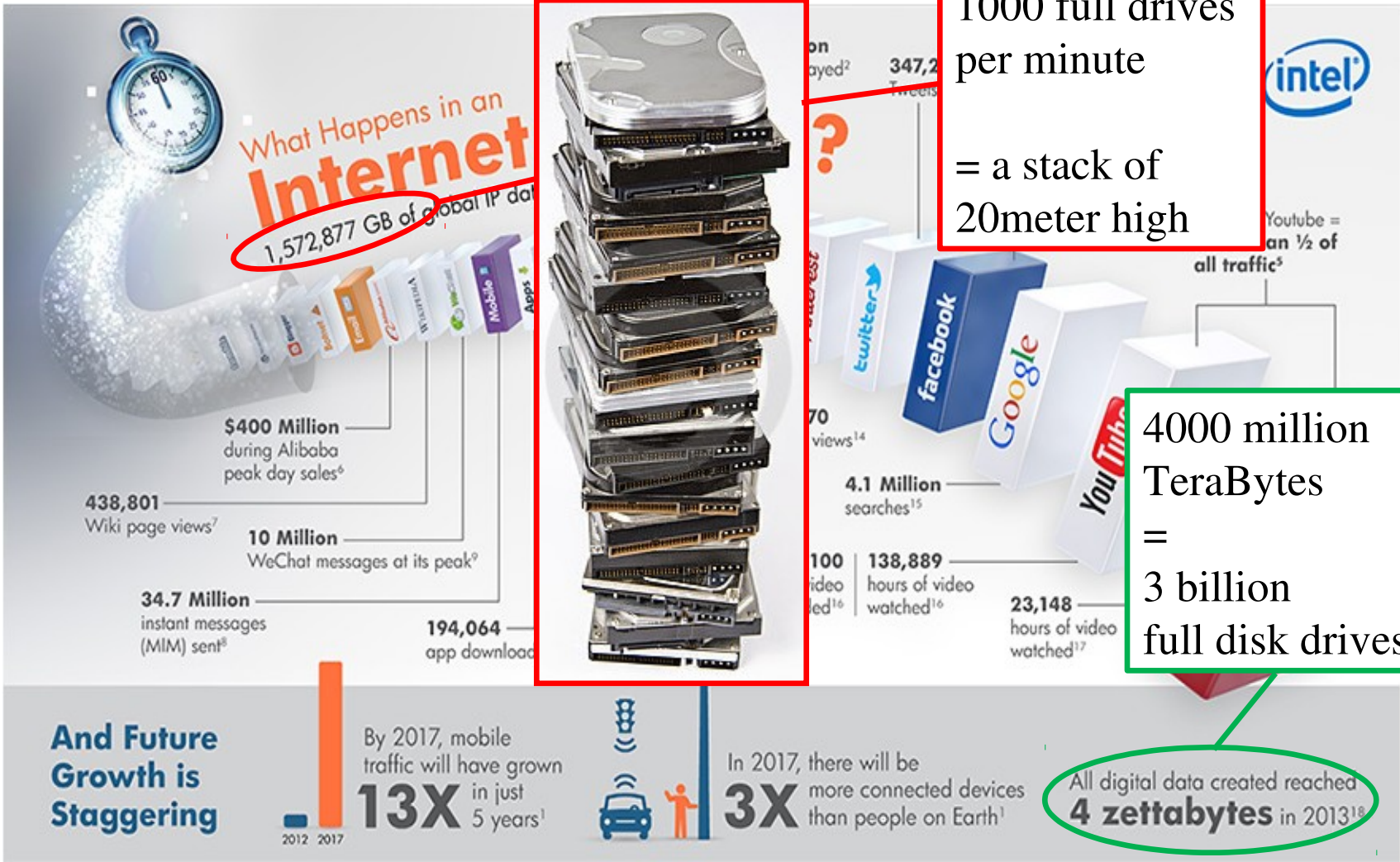


Database

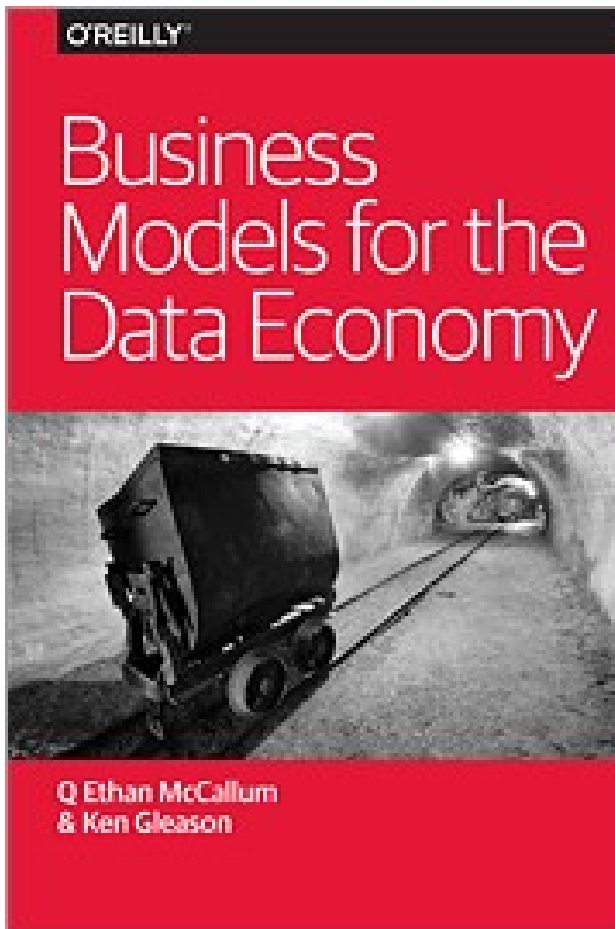


Data Mining

The age of Big Data



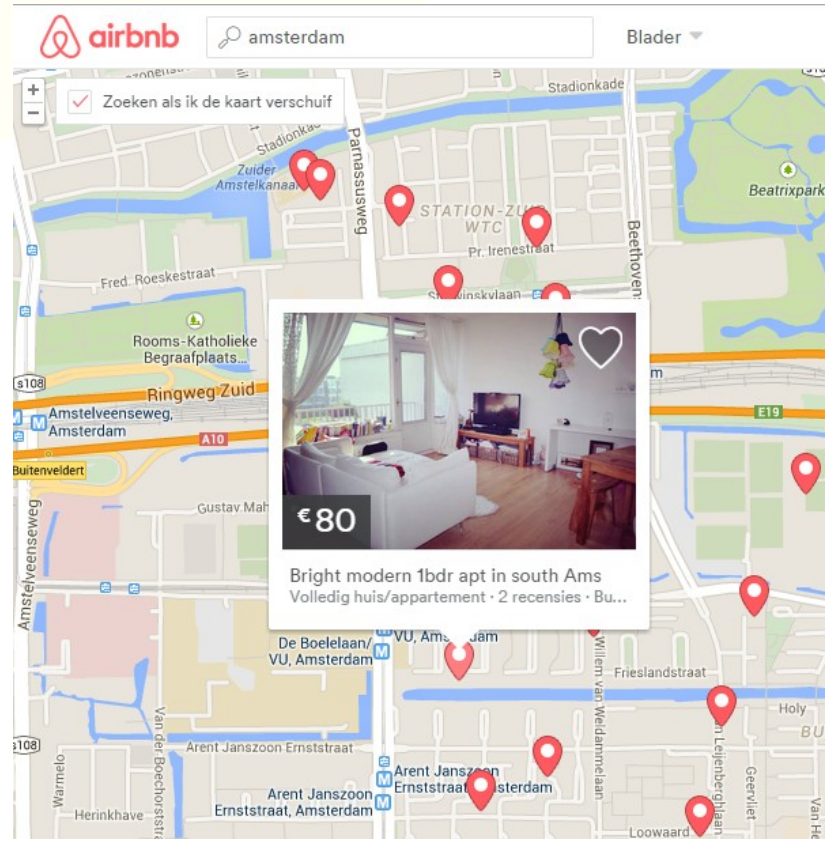
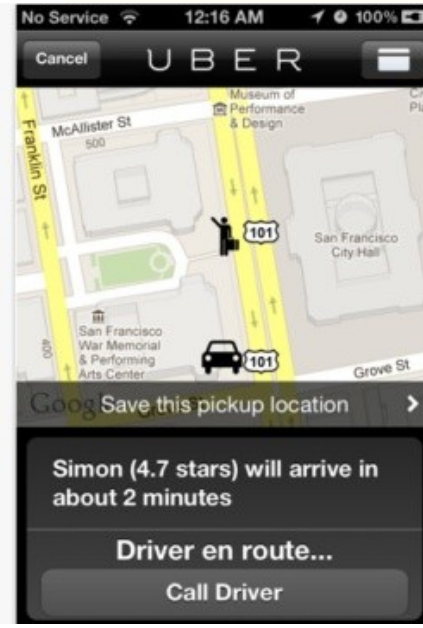
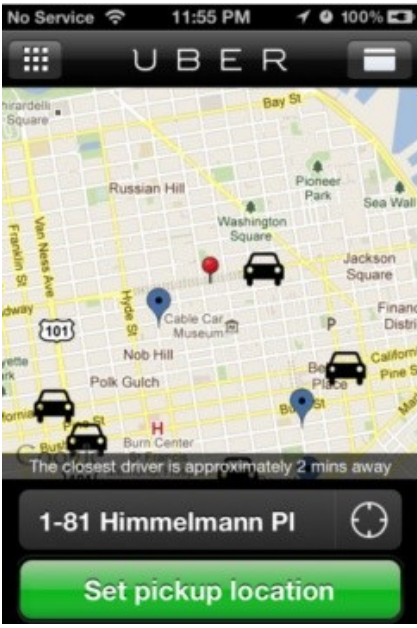






UBER

EVERYONE'S PRIVATE DRIVER™



DBDM: Selected Challenges

GIS (LIDAR):

Massive point clouds: 640 Billion (x,y,z) points / 15 TB
=> **spatial joins between point cloud and polygons**

netherlands

eScience

center

Logistics:



TOMTOM

> 5 trillion (10^{12}) GPS points (grows with >60k points/sec)

Seismology:



Koninklijk Nederlands
Meteorologisch Instituut
Ministerie van Infrastructuur en Milieu

~ 4 M files, ~ 500 GB (10x compressed)

=> **Transparent data ingestion: Data Vault**

COMMIT/

Remote sensing:



~2 PB satellite image data

=> **Array data processing: SciQL**



Astronomy:



LOFAR

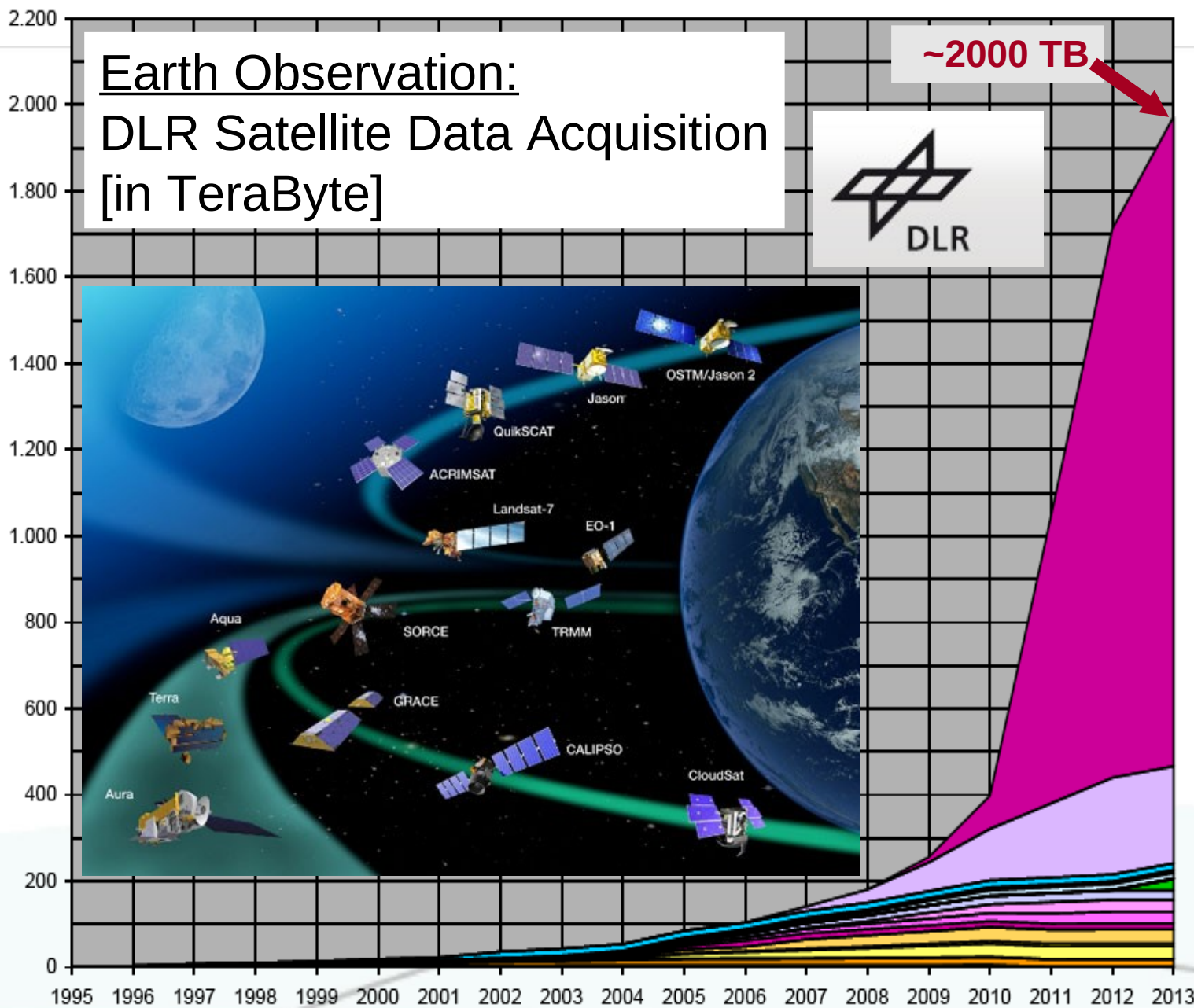
Raw data: 25 TB / hour; derived data: 100 TB / year

=> **Transient detection inside DBMS**

DBDM: Earth Observation

Earth Observation:
DLR Satellite Data Acquisition
[in TeraByte]

~2000 TB



- TanDEM-X
- TerraSAR-X
- ERS-SAR
- SRTM
- GEMOS
- XSAR-2
- XSAR-1
- AIR-RS
- ARES
- ROSIS
- DAIS
- HyMAP
- WDC
- EnMAP
- MODIS VA
- MODIS (Terra, Aqua)
- GOME (ERS und METOP VA)
- METOP-GOME2
- ENVISOLAR/MSG
- MSG
- METEOSAT
- ENVISAT-VA-Atmos
- ENVISAT SCIAMACHY
- ENVISAT MIPAS
- ENVISAT AATSR
- ERS-SAR
- NOAA





LOFAR Low Frequency Array for Radio Astronomy

International LOFAR Telescope (ILT)

Raw data:
~25 TB / hour

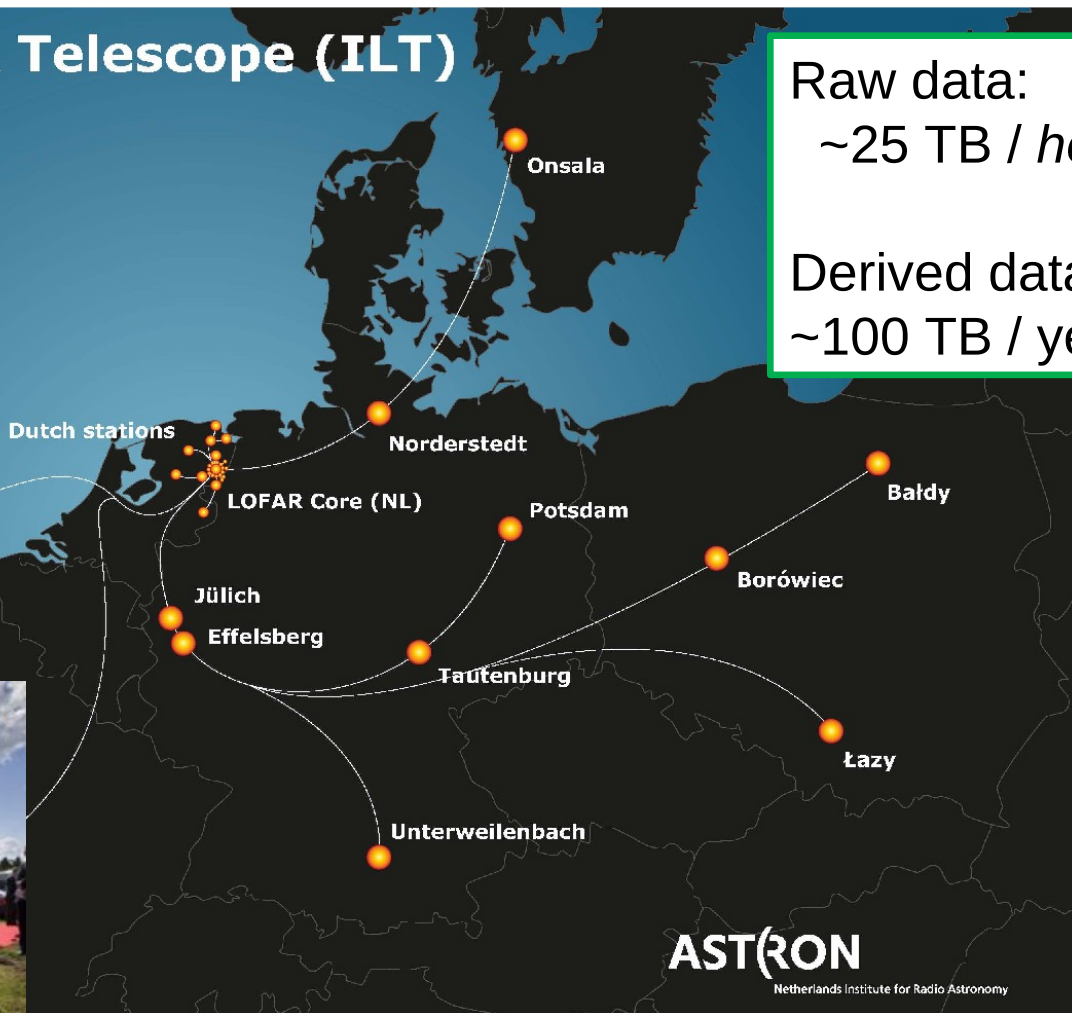
Derived data:
~100 TB / year



Chilbolton



Opening van de LOFAR telescoop door koningin Beatrix in juni 2010



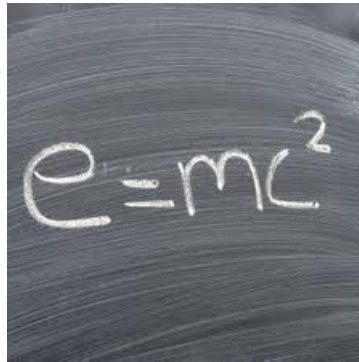
Data Disrupting Science: Paradigm Shift in Scientific Research



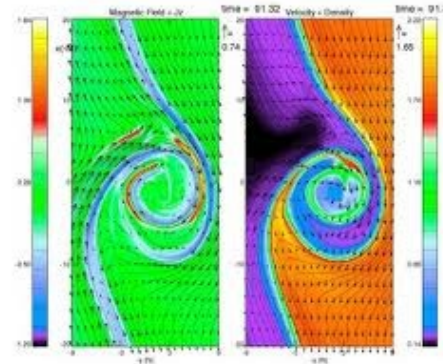
collecting &
analyzing data
data exploration
(eScience)



observing
empirical
1st



modeling
theoretical
2nd



simulating
computational
3rd



Jim Gray
(1944 - 2007)



The
**FOURTH
PARADIGM**

DATA-INTENSIVE SCIENTIFIC DISCOVERY

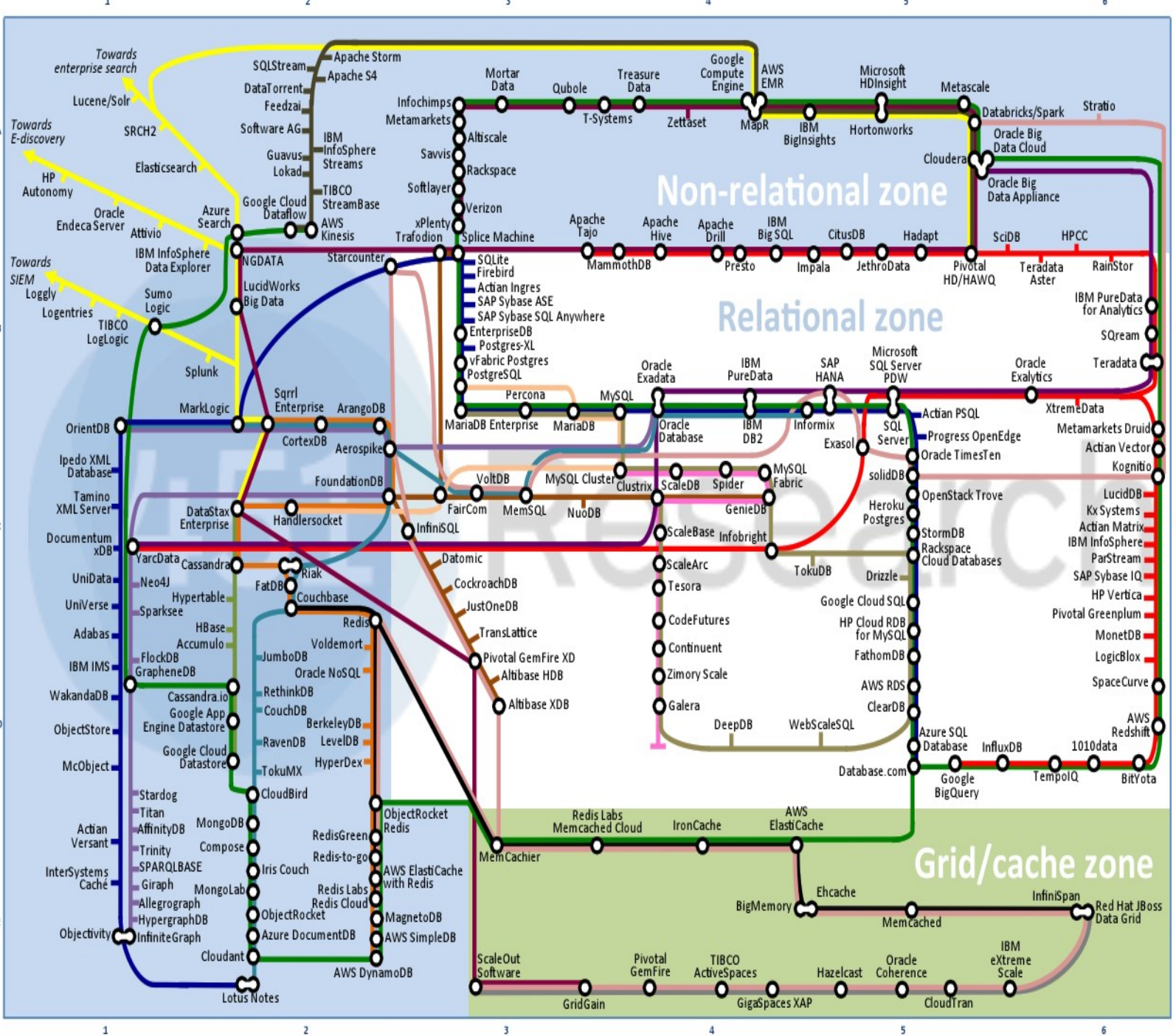
EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

Data Management & Data Mining



Data Platforms Map

October 2014



- Key:**
- General purpose
 - Specialist analytic
 - -as-a-Service
 - BigTables
 - Graph
 - Document
 - Key value stores
 - Key value direct access
 - Hadoop
 - MySQL ecosystem
 - Advanced clustering/sharding
 - New SQL databases
 - Data caching
 - Data grid
 - Search
 - Appliances
 - In-memory
 - Stream processing

<https://451research.com/dashboard/dpa>

BIG DATA LANDSCAPE, VERSION 3.0

Exited: Acquisition or IPO

Infrastructure NoSQL Databases: FOUNDATIONDB, DATASTAX, mongoDB, Couchbase, KERO SPIKE, HYPERTABLE, sqrrl, CLOUDANT, Amazon, Pivotal, Neo4j, OhmData, sonos Hadoop On Prem: HADAPT, cloudera, splice MACHINE, Zettaset, amazon, MAPR, Microsoft, Pivotal, IBM InfoSphere Big Data: MORTAR, infochimps, Qubole, JETHRO DATA, altiscale, amazon web services		Analytics Analytics Platforms: databricks, Quant Cell RESEARCH, PERSASIVE, GUAVUS, Datameer, KARMA SPHERE, collective IO, PRECISO, dataspota For Business Analysts: STAT WING, CIRRO, TREPAREL, OrigamiLogic, ClearStory, DataGravity Data Science Platforms/Tools: domino, nutonian, Alpine, Sense, MORTAR, CONTINUUM ANALYTICS, ploply, yhat, MODE Unstructured Data: BASIS, ATTIVO, GENERAL SENTIMENT, semantria, crimson hexagon, DIGITAL REASONING, ai, Quid, Narrative Science, Palantir BI Platforms: birst, bime, pentaho, GoodData, SiSense, boomi, platforma Machine Learning: SKYTREE, big ml, YOTTAMINE ANALYTICS, wise.io, context relevant		Applications Ad Optimization: aggregate knowledge, rocketfuel, TAPAD, ai Match, MediaMath, thetradedesk, across Marketing: LATTICE ENGINES, Sailthru, spinnaker, gainsight, Kontera, RelateIQ, TellApart, persado, bloomreach, CLICKFOX, Pursway Finance: Lenddo, BILL GUARD, wonga, cignifi, LendUp, KENSHO, OnDeck Security: mark43, enigma, SIGNIFYD, sift science, FORTSCALE, feedzai			
NewSQL Databases: MarkLogic, TRANSATTICE, RainStar, paradigm4, memsql, deepdb, nuODB, citusdata, skySQL, Clustrix, VoltDB, SQLFire Cluster Service: LexisNexis, HPCC Systems, mesosphere, Acunu Management/Monitoring: metafor, New Relic, StackIQ, tidemark, appromic, AppFirst, oceanSYNIC, AppFirst, DATADOG, bundy		Social Analytics: simple reach, bitly, synthesio, Dataminr, Statilizer, vicarious, DATA SIFT, tracx, bottlenose Analytics Services: THINK BIG, McKinsey & Company, OPERA, UO, VALANCE, DATA SCIENCE, Mu Sigma Statistical Computing: REVOLUTION, SAS, MATLAB, Log Analytics, splunk, loggly, sumologic, Kibana		Government/Regulator: Fiscal Note, FIRE SHOP, PREDPOL Education/Learning: KNEWTON, Panorama, Clever			
MPP Databases: TERADATA, ParStream, InfiniDB, Kognitio, NETEZZA, SQL Server, Pivotal, PARACCEL Graph Databases: Neo4j, Graph, aster data, InfiniteGraph Data Transformation: TRIFACTA, Palata, DataTamer, KALIDO, evelytix, TRAIKEIRON, syncsoft		Location/People/Events: RADIUS, Fliptop, LOCATE, Place IQ Big Data Search: hp, Autonomy, LucidWorks, ONTOLOGY Real-Time: kaggle, METAMARKETS, amiato, causata, DataKind		Security: TRACTIC, DATAGUISE, codefortytwo, Stormpath, VIPERA, Cleversafe, Panasas, nimblestorage, Compuverde App Dev: CONCURRENT, CONTINUITY, wibi: data		Health: Recombine, tubular, OPOWER, SIGHT MACHINE, THE CLIMATE CORPORATION, numberFire Industries: 23andMe, Ginger.io, FLATIRON, COUNSYL	
Crowd-sourcing: microtask, CROWD COMPUTING, CROWDPOWER, servio, mobileworks, mechanical Turk Storage: Cleversafe, Panasas, nimblestorage, Compuverde		Analytics: ANALYTIC SERVICES, THINK BIG, McKinsey & Company, OPERA, UO, VALANCE, DATA SCIENCE, Mu Sigma Statistical Computing: REVOLUTION, SAS, MATLAB, Log Analytics, splunk, loggly, sumologic, Kibana		Government/Regulator: Fiscal Note, FIRE SHOP, PREDPOL Education/Learning: KNEWTON, Panorama, Clever			

Cross Infrastructure / SAP, SAS, IBM, Google, Microsoft, vmware, amazon, 1010data, talend, hp, Autonomy, NetApp, TERADATA

Open Source

Framework: Spark, Hadoop YARN, HDFS	Query/Data Flow: HBASE, CouchDB, mongoDB, riak, Sqoop	Coordination/Work-flow: ZooKeeper, talend	Real-Time: Storm	Stat Tools: SciPy	Machine Learning: mLib, mahout	Cloud Deploy: Solr, LUCENE.net
-------------------------------------	---	---	------------------	-------------------	--------------------------------	--------------------------------

Data Sources

Data Mkts: Windows Azure Marketplace, bluekai, DataMarket, factual, knoema	Data Sources: DATA.GOV, premise, YODLEE, xignite, plaid, quandl, SPACE CURVE, STANDARD TREASURY, human/api	Sensor Data: kinsa, STREETLINE, SKYCATCH, fitbit, RunKeeper, JAWBONE, LUMA SENSE TECHNOLOGIES, Withings, BASIS, estimote	Incubators/Schools: zipfian, GA, INSIGHT, DataElite
--	--	--	---