

Research, Innovation and Teaching in Big Data Analytics

Challenges and Chances

Prof. Dr. Volker Markl

Chair of the Database Systems and Information Management Lab
Technische Universität Berlin

www.dima.tu-berlin.de



McKinsey & Co

McKinsey Global Institute



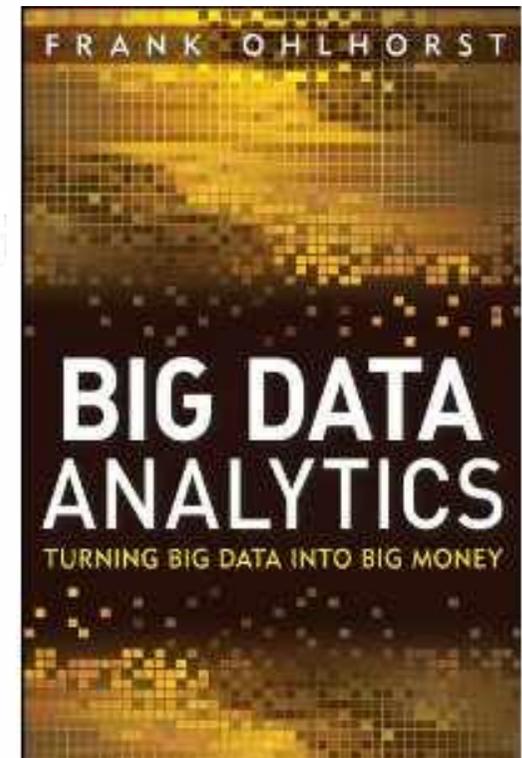
May 2011

Big data: The next frontier for innovation, competition and productivity

Source: McKinsey



Source: Information Week

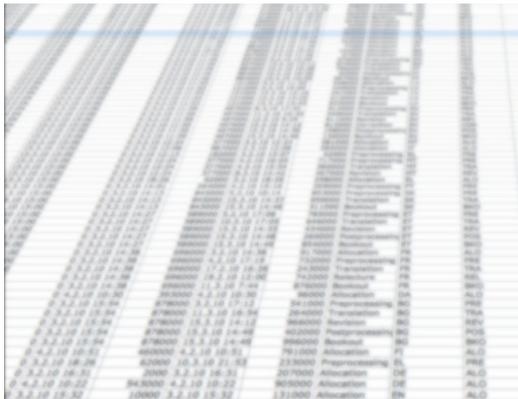


Source: Amazon

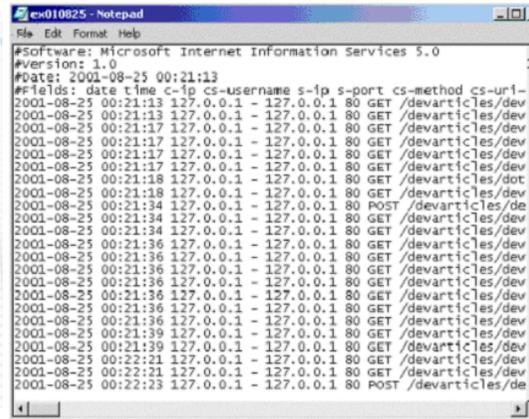


Source: NCSU

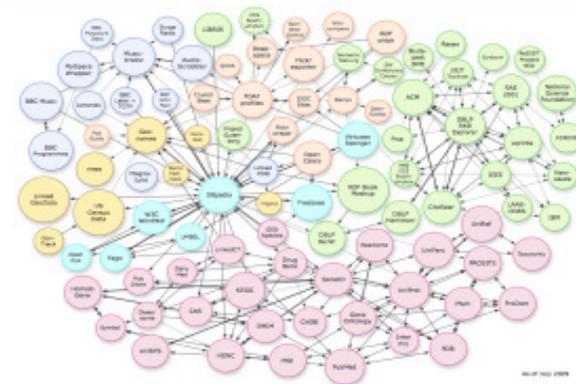
- What is Big Data?
- Examples
- Parallel Data Processing and Other Enabling Technologies
- Applications
- Challenges
- Call to Action



Transactional Data at „Webscale“

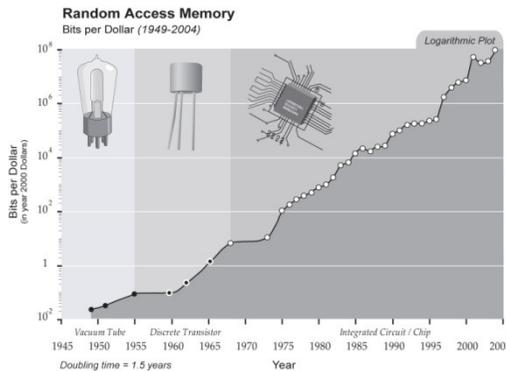


Web Logfiles

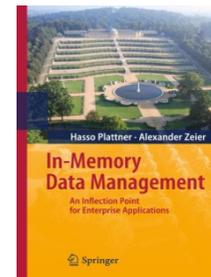


Linked Open Data and many other „human generated data sets“

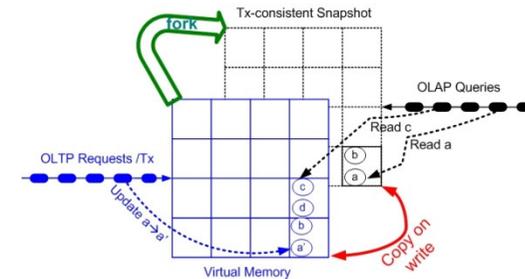
These data sets will fit into main memory soon!



IBM ISAO



SAP HANA



TUM Hyper

Many standard data management solutions exist!



It is cheaper to resequence than to store genome data!

Will Computers Crash Genomics?

New technologies are making sequencing DNA easier and cheaper than ever, but the ability to analyze and store all that data is lagging

Lincoln Stein is worried. For decades, computers have improved at rates that have boggled the mind. But Stein, a bioinformaticist at the Ontario Institute for Cancer Research (OICR) in Toronto, Canada, works in a field that is moving even faster: genomics.

The cost of sequencing DNA has taken a nosedive in the decade since the human genome was published—and it is now dropping by 50% every 5 months. The amount of sequence available to researchers has consequently skyrocketed, setting off warnings about a “data tsunami.” A single DNA sequencer can now generate in a day what it took 10 years to collect for the Human Genome Project. Computers are central to archiving and analyzing this information, notes Stein, but their processing power isn’t increasing fast enough, and their costs are decreasing too slowly to keep up with the deluge. The torrent of DNA data and the need to analyze it “will swamp our storage systems and crush our computer clusters,” Stein predicted last year in the journal *Genome Biology*.

Funding agencies have neglected bioinformatics needs, Stein and others argue. “Traditionally, the U.K. and the U.S. have not invested in analysis; instead, the focus has been investing in data generation,” says computational biologist Chris Ponting of the University of Oxford in the United Kingdom. “That’s got to change.”

Within a few years, Ponting predicts, analysis, not sequencing, will be the main expense hurdle to many genome projects. And that’s assuming there’s someone who can do it; bioinformaticists are in short supply everywhere. “I worry there won’t be enough people around to do the analysis,” says Ponting. Recent reviews, editorials, and scientists’ blogs have echoed these concerns (see Perspective on p. 728). They stress the need for new software and infrastructures to deal with computational and storage issues.

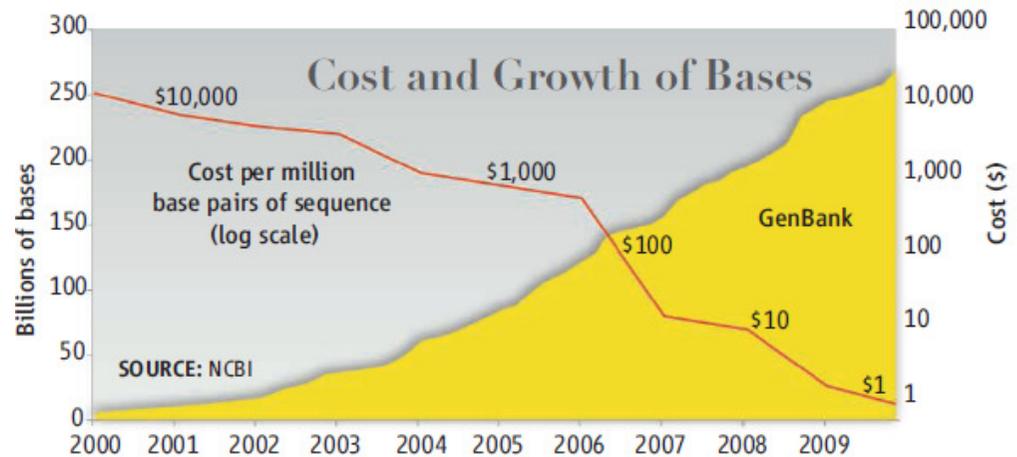
In the meantime, bioinformaticists are trying new approaches to handle the data onslaught. Some are heading for the clouds—cloud computing, that is, a pay-as-

you-go service, accessible from one’s own desktop, that provides rented time on a large cluster of machines that work together in parallel as fast as, or faster than, a single powerful computer. “Surviving the data deluge means computing in parallel,” says Michael Schutz, a bioinformaticist at Cold Spring Harbor Laboratory (CSHL) in New York.

Dizzy with data

The balance between sequence generation and the ability to handle the data began to shift after 2005. Until then, and even today, most DNA sequencing occurred in large centers, well equipped with the computer personnel and infrastructure to support the analysis of a genome’s data. DNA sequences churned out by these centers were deposited and stored in centralized public databases, such as those run by the European Bioinformatics Institute (EBI) in Hinxton, U.K., and the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Researchers elsewhere could then download the data for study. By 2007, NCBI had 150 billion bases of genetic information stored in its GenBank database.

Then several companies in quick succession introduced “next-generation” machines, faster sequencers that spit out data more cheaply. But the technologies behind these machines generate such short stretches of sequence—typically just 50



SCIENCE 11 FEBRUARY 2011
VOL 331, ISSUE 6018, PAGES 639-806



Video Streams



Smart Grids



Web Archive

EXABYTE
 (1,152,921,504,606,846,976 BYTES; 2^{60})
 approx. 1,000,000,000,000,000 or 10^{15}

5 EXABYTES: ALL WORDS EVER SPOKEN BY HUMAN BEINGS

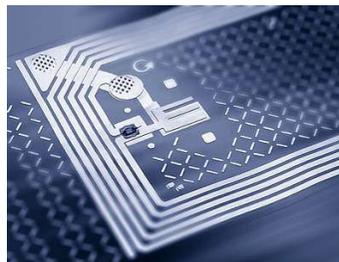


Sensor Data

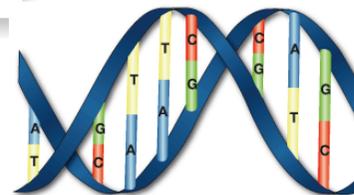


Audio Streams

ZETTABYTE
 (1,180,591,620,717,411,303,424 BYTES; 2^{70})
 approx. 1,000,000,000,000,000,000 or 10^{21}



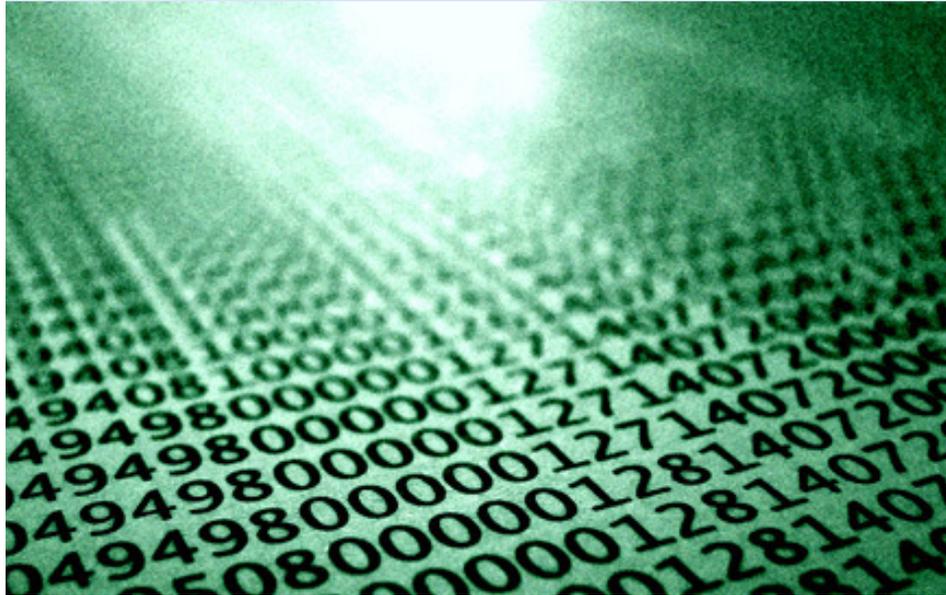
RFID Data



Thymine (Yellow) = T Guanine (Green) = G
 Adenine (Blue) = A Cytosine (Red) = C

Genome Data

**mostly machine generated
 or collaborative on web-scale**



Size
Format/Media Type
Uncertainty/Quality
Freshness
etc.

Data

Selection/Grouping
Relational Operators (Join)
Information Extraction & Integration
Data Mining
Predictive Models
etc.

Query

■ Volume

- Data size

■ Velocity

- Data ingestion speed
- Analysis time window

■ Variability

- Data formats
- Media types

■ Veracity

- Uncertainty
- Inconsistency

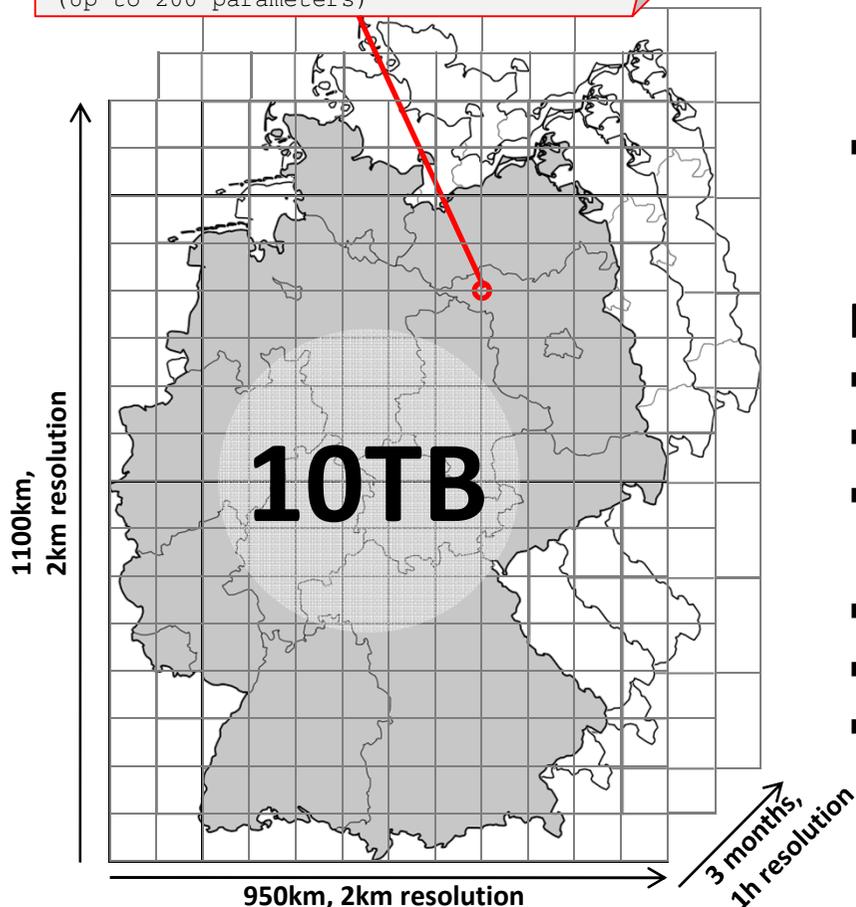
challenging for complex analysis questions

solutions exist for structured data and text
challenging for graph or audio/video data

challenging for automatic reasoning

Popular „definition“, but misses some other aspects of complexity

```
PS,1,1,0,Pa, surface pressure
T_2M,11,105,0,K,air_temperature
TMAX_2M,15,105,2,K,2m maximum temperature
TMIN_2M,16,105,2,K,2m minimum temperature
U,33,110,0,ms-1,U-component of wind
V,34,110,0,ms-1,V-component of wind
QV_2M,51,105,0,kgkg-1,2m specific humidity
CLCT,71,1,0,1,total cloud cover
...
(Up to 200 parameters)
```



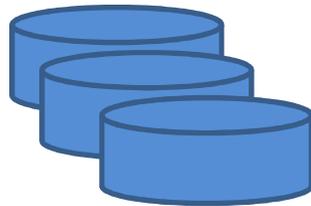
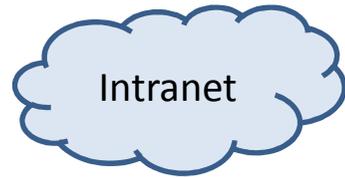
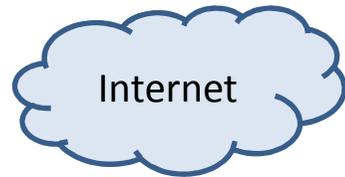
Analysis Tasks on Climate Data Sets

- Validate climate models
- Locate „hot-spots“ in climate models
 - Monsoon
 - Drought
 - Flooding
- Compare climate models
 - Based on different parameter settings

Necessary Data Processing Operations

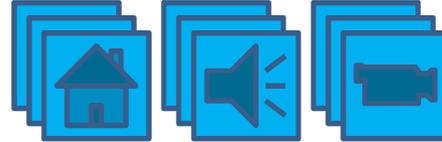
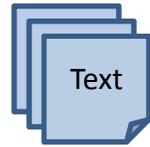
- Filter
- Aggregation (sliding window)
- Join

- Multi-dimensional sliding-window operations
- Geospatial/Temporal joins
- Uncertainty

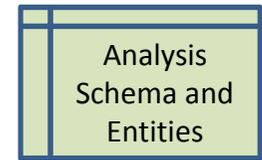
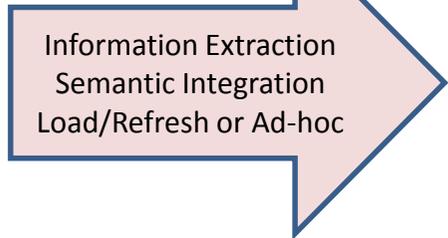


Data Warehouse
Data Marts

BI over Text



BI over Media Types



Situational BI

How did the sentiment towards my products change this year?

How is the customer satisfaction per region?

Who is leading in American Idol?



The *next generation of Business Intelligence (NGBI)* will correlate data warehouses with text and other modalities from web services of information providers, corporate Intranets and the Internet



Home Automation



Healthcare



Water Management



Lifecycle Management



Sales/Marketing



Traffic Management



Energy Management

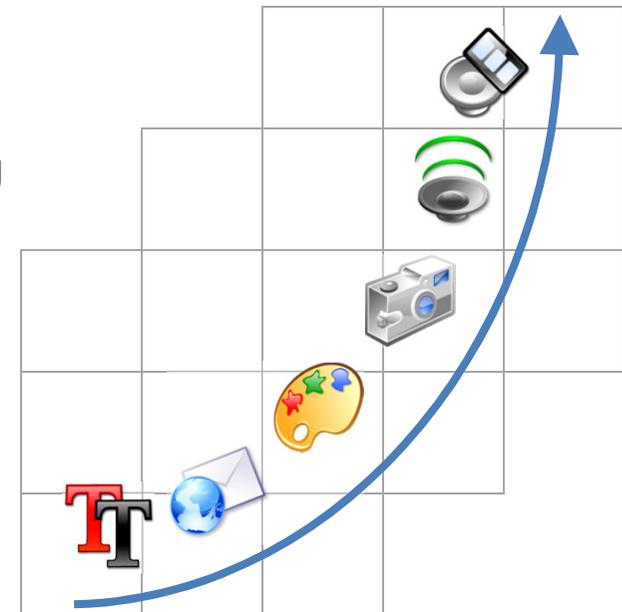
■ Need for data parallel processing

- Increase of data complexity
- Increase of query complexity
- Moore's Law: ManyCore and Cluster Computing
- Scale-up no longer possible

➔ **Scale-out is the name of the game**

■ Parallel programming is not easy

- (Network) Communication
- Concurrent programming: Divide & Conquer
- Synchronization as bottleneck (Amdahl's Law)
- Fault tolerance



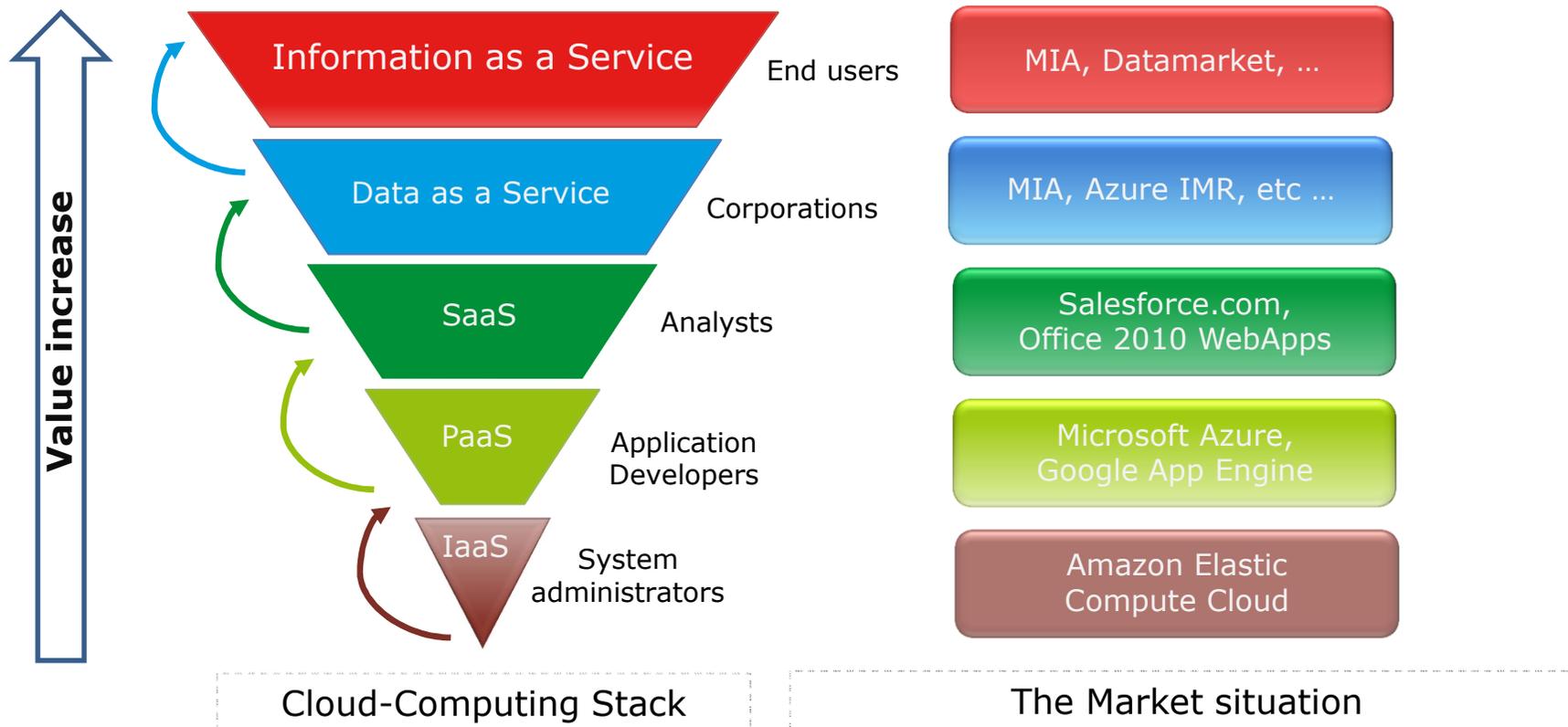
■ Data Programming models must ease parallel data processing

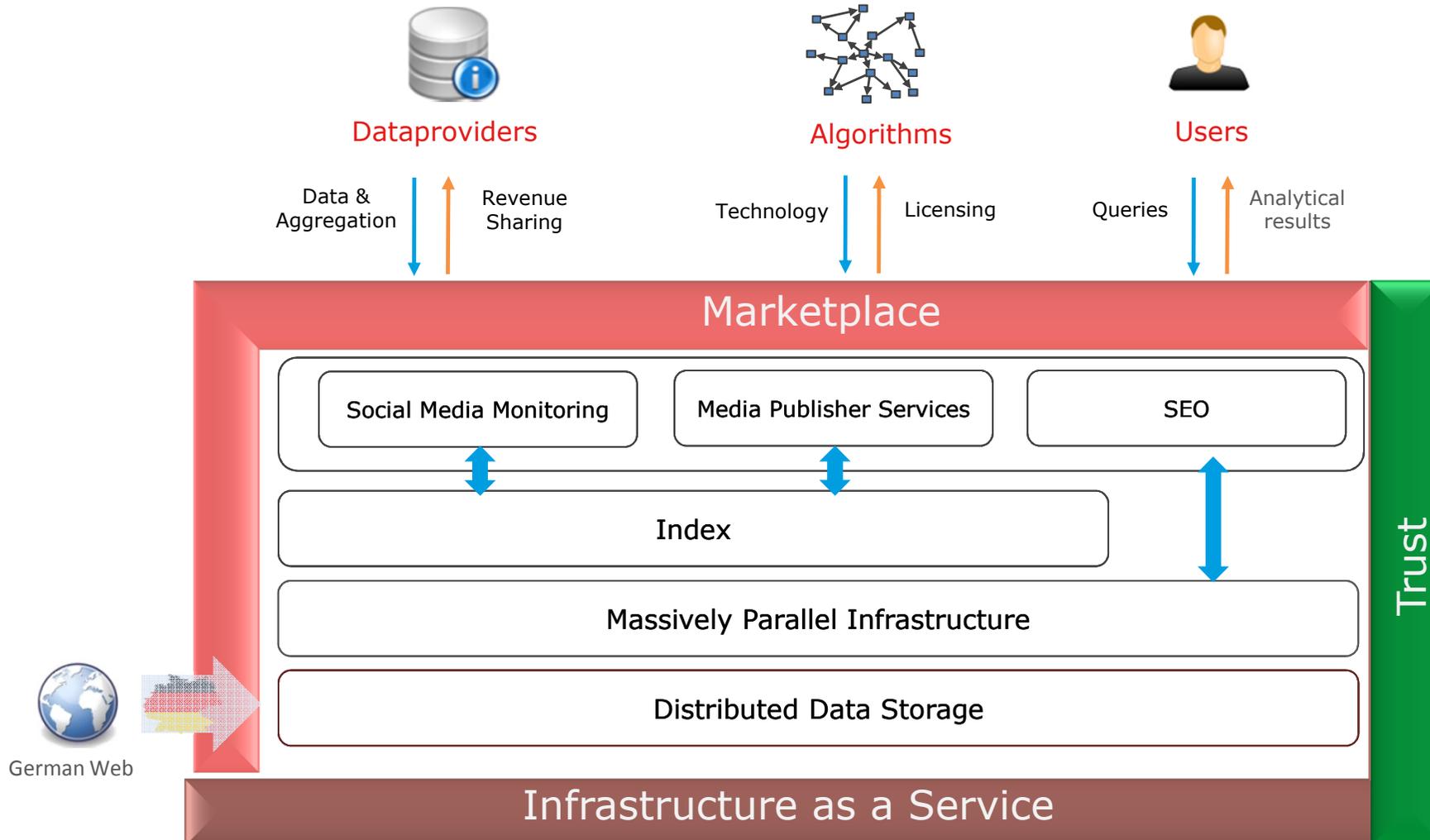
- Abstractions hide the gory details
- Automatic adaption to hardware
 - Parallelization and Optimization
- Beware: Data flow and control flow dependencies!
- Popular simple model: map/reduce (e.g., Hadoop)

A major new trend in information processing will be the trading of original and enriched data, effectively creating an information economy.

„When hardware became commoditized, software was valuable. Now that software is being commoditized, data is valuable.“ (TIM O'REILLY)

„The important question isn't who owns the data. Ultimately, we all do. A better question is, who owns the means of analysis?“ (A. CROLL, MASHABLE, 2011)





<http://www.trusted-cloud.de/de/778.php>

Use Cases

Scientific Data

Life Sciences

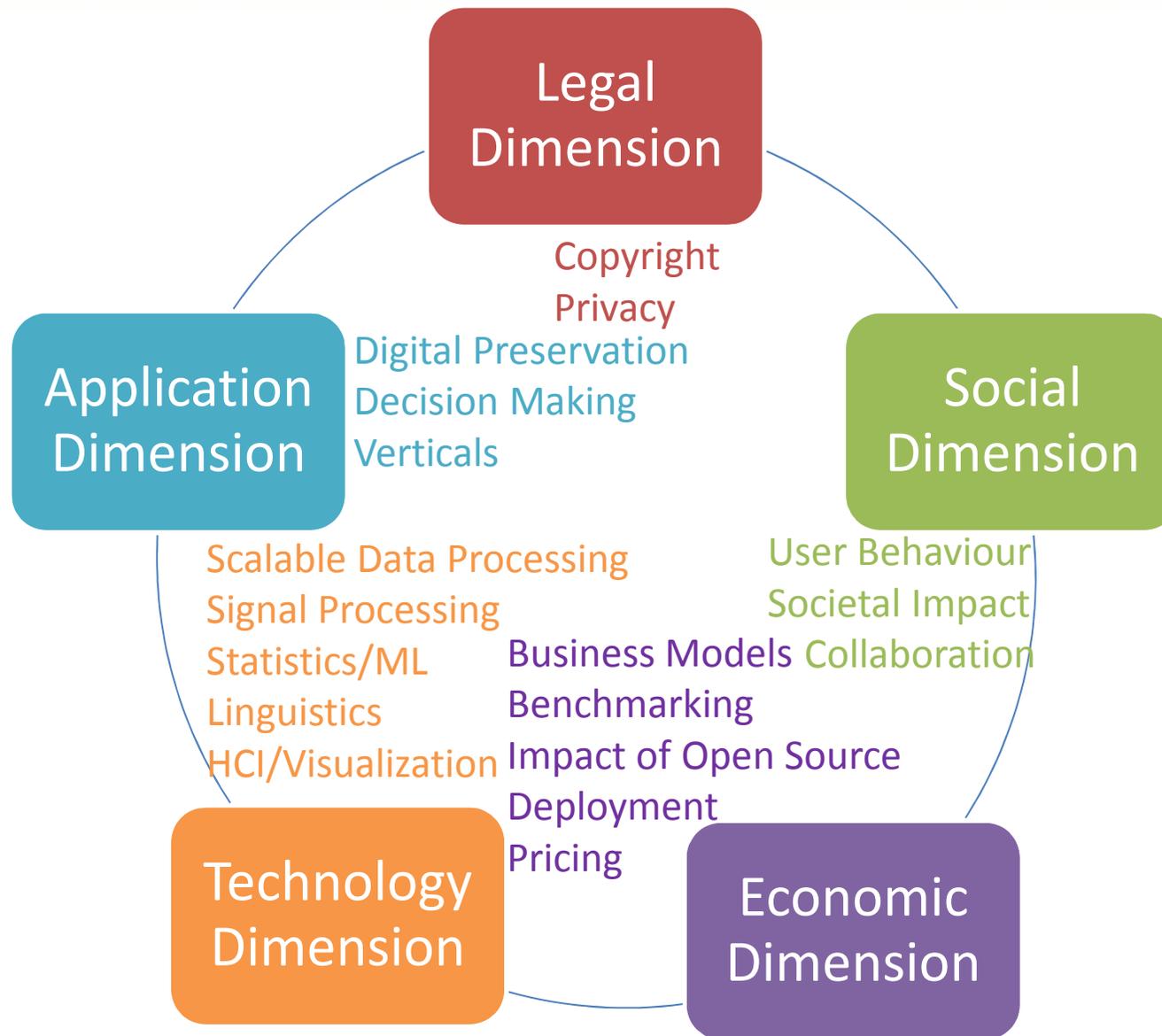
Linked Data

Parallel Analytics Query Processor

Infrastructure as a Service

- Explore the power of Cloud computing for complex analytical tasks
- Database-inspired approach
- Analyze, aggregate, and query
- Textual and (semi-) structured data
- Research and prototype a web-scale data analytics infrastructure

Stratosphere is publically funded by DFG and EIT. Open Source at www.stratosphere.eu



- **Educate Data Scientists** to create the required talent
 - T-shaped students
 - Information „literacy“
 - Data Analytics Curriculum

- **Research Big Data Analytics Technologies**
 - Data management (uncertainty, query processing under near real-time constraints, information extraction)
 - Programming models
 - Machine Learning and statistical methods
 - Systems architectures
 - Information Visualization

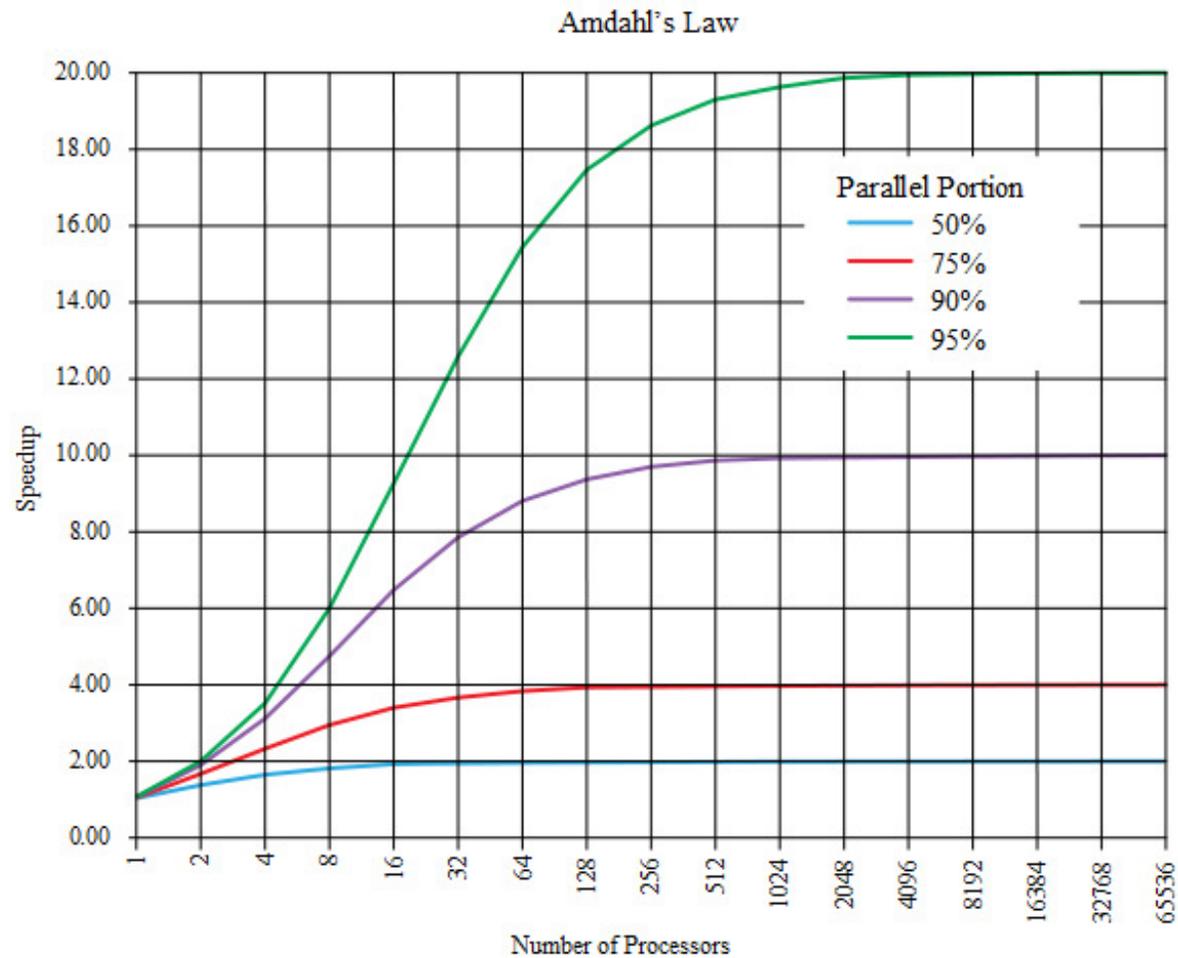
- **Innovate** to maintain competitiveness
 - Demonstrate flagship use-cases to raise awareness
 - Promote startups in the area of data analytics
 - Transfer technologies to German enterprises, in particular SMEs
 - Determine legal frameworks and business models

We need to ensure a German technological leadership role in „Big Data“

- Parallel Data Processing
- Parallel Data Management: ParStream et al.
- Map/Reduce: Hadoop, Stratosphere et al.

- The Speedup is defined as: $S_p = T_1 / T_p$
 - T_1 : runtime of the sequential program
 - T_p : runtime of the parallel program on p processors

- Ahmdal's Law: „The maximal speedup is determined by the non-parallelizable part of a program“ :
 - $S_{max} = \frac{1}{(1-f) + f/p}$ f: fraction of the program that can be parallelized
 - → Ideal speedup: $S=p$ for $f=1.0$ (linear speedup)
 - However - since usually $f < 1.0$ -, S is bound by a constant ! (e.g. ~ 10 for $f=0.9$)
 - → Fixed problems can only be parallelized to a certain degree!

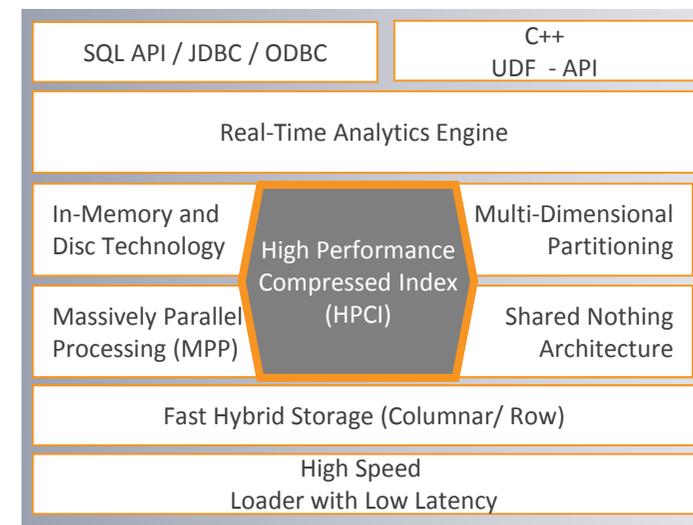


ParStream provides the unique combination of REAL-TIME – LOW-LATENCY – HIGH THROUGHPUT

ParStream's key advantages:

- Unique Highly compressed Bitmap Index which can be analyzed in compressed form (patent application filed)
- Real-time analytics through Massively Parallel Processing (MPP)
- Very high throughput due to massively reduced workload (no decompression, small index, efficient algorithms)
- Low-Latency through continuous import of new data without slowing down analytics
- Columnar data / index allows very flexible analytics (multi-column, multi-value)
- Specialized data / index types and algorithms
- Shared Nothing Architecture

PARSTREAM PARALLEL ARCHITECTURE



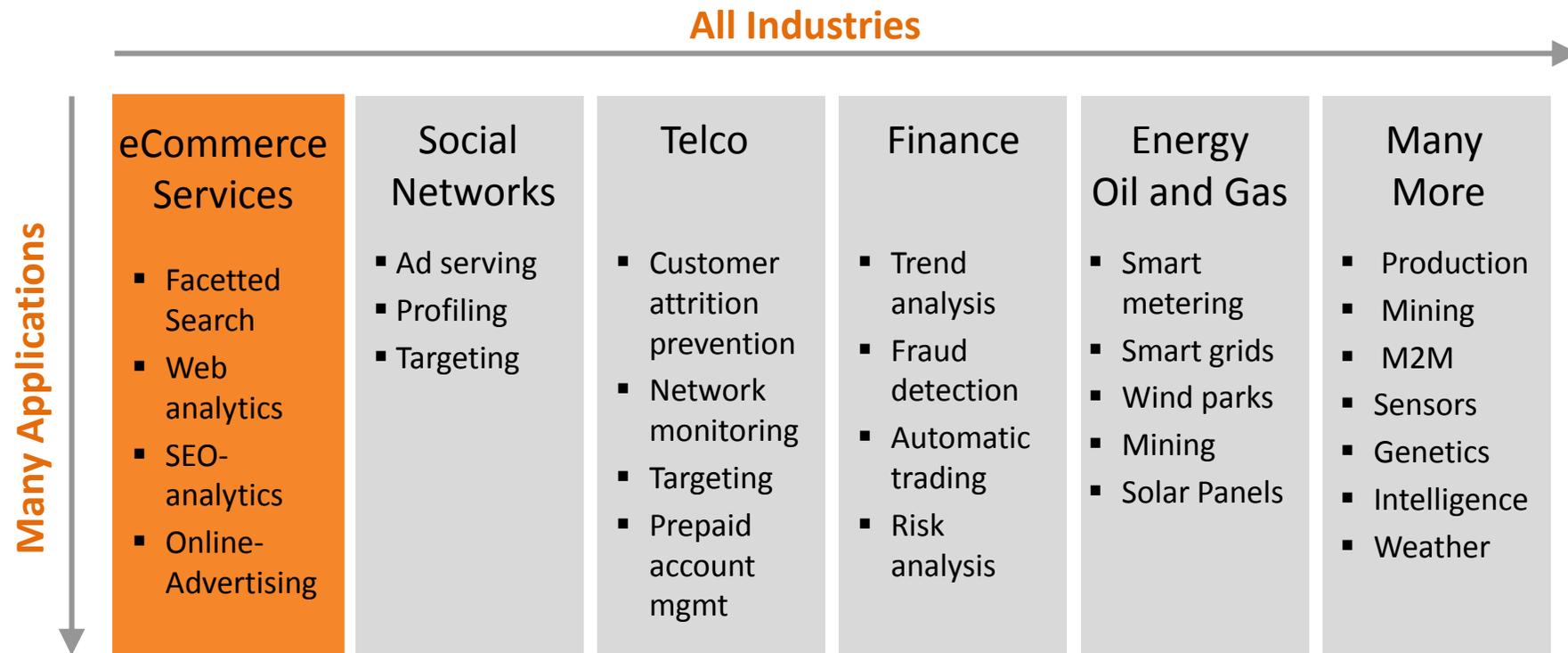
PARSTREAM INDEX ARCHITECTURE

Parallel Search within Compressed Index



© Parstream GmbH, used with permission, www.parstream.com

Big Data Analytics is a game changer in every industry
and is a huge market opportunity



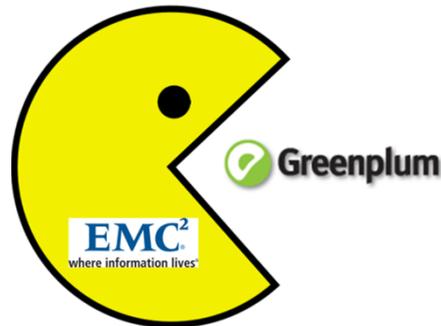
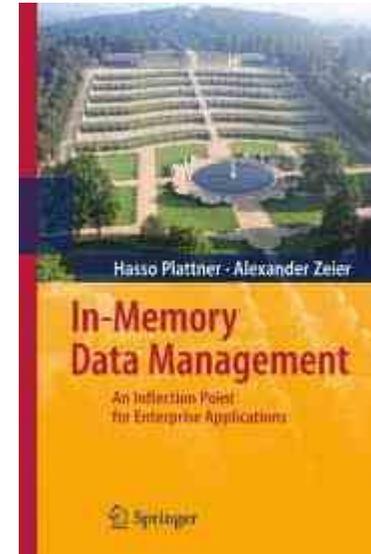
© Parstream GmbH, used with permission
www.parstream.com



JasperReports Server



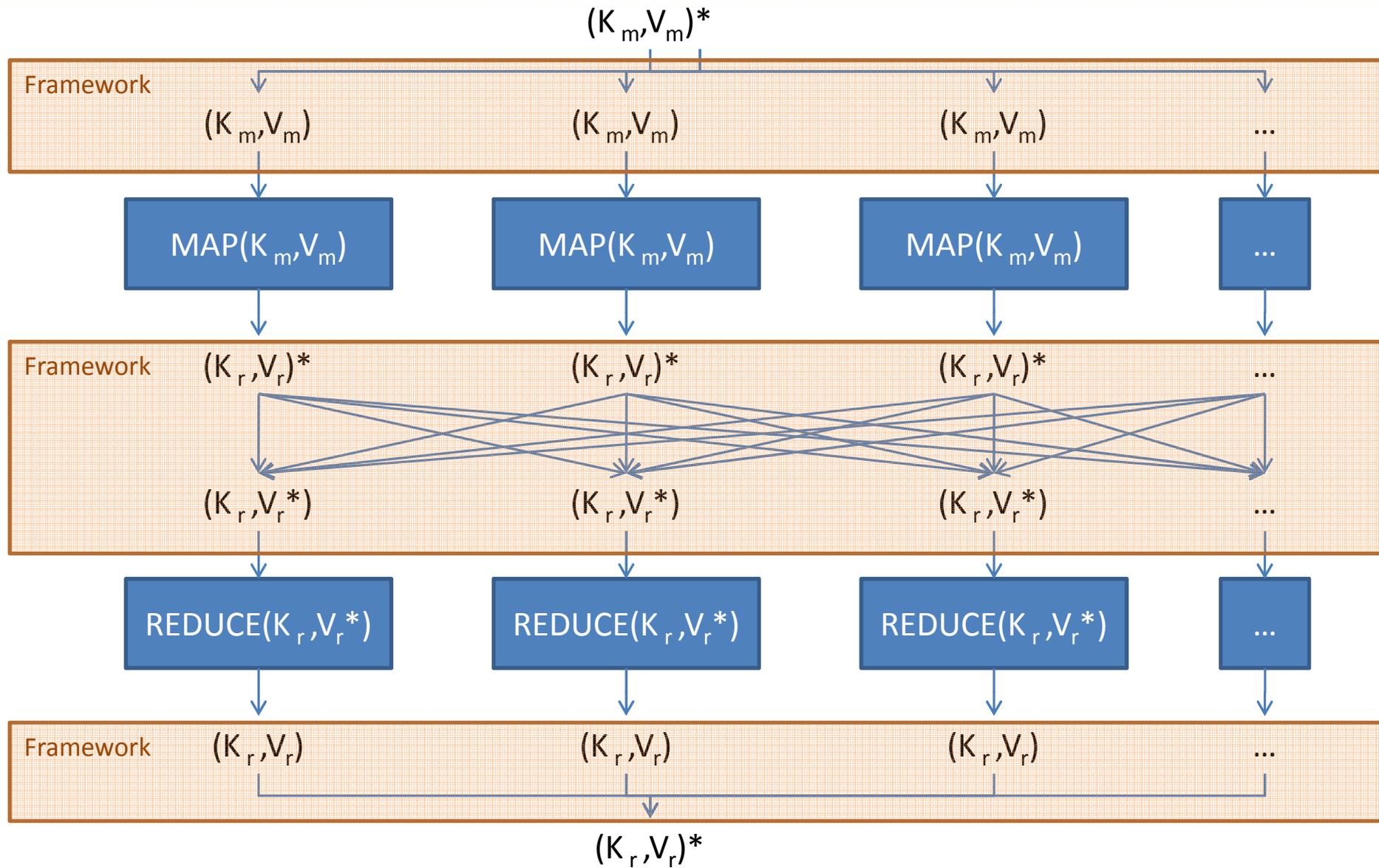
Sybase™
IQ



- Programming model for data-intensive programming
 - well-suited for large scale parallel execution
 - automatic parallelization & distribution of data and computational logic
 - easy to extend with fault tolerance schemes
 - clean abstraction for programmers

- Based on functional programming
 - treats computation as the evaluation of mathematical functions and avoids state and mutable data
 - no changes of states (no side effects)
 - output value of a function depends only on its arguments

- `Map` and `Reduce` are higher-order functions
 - take user-defined functions as argument
 - return a function as result
 - to define a map/reduce job, the user implements the two functions `m` and `r`

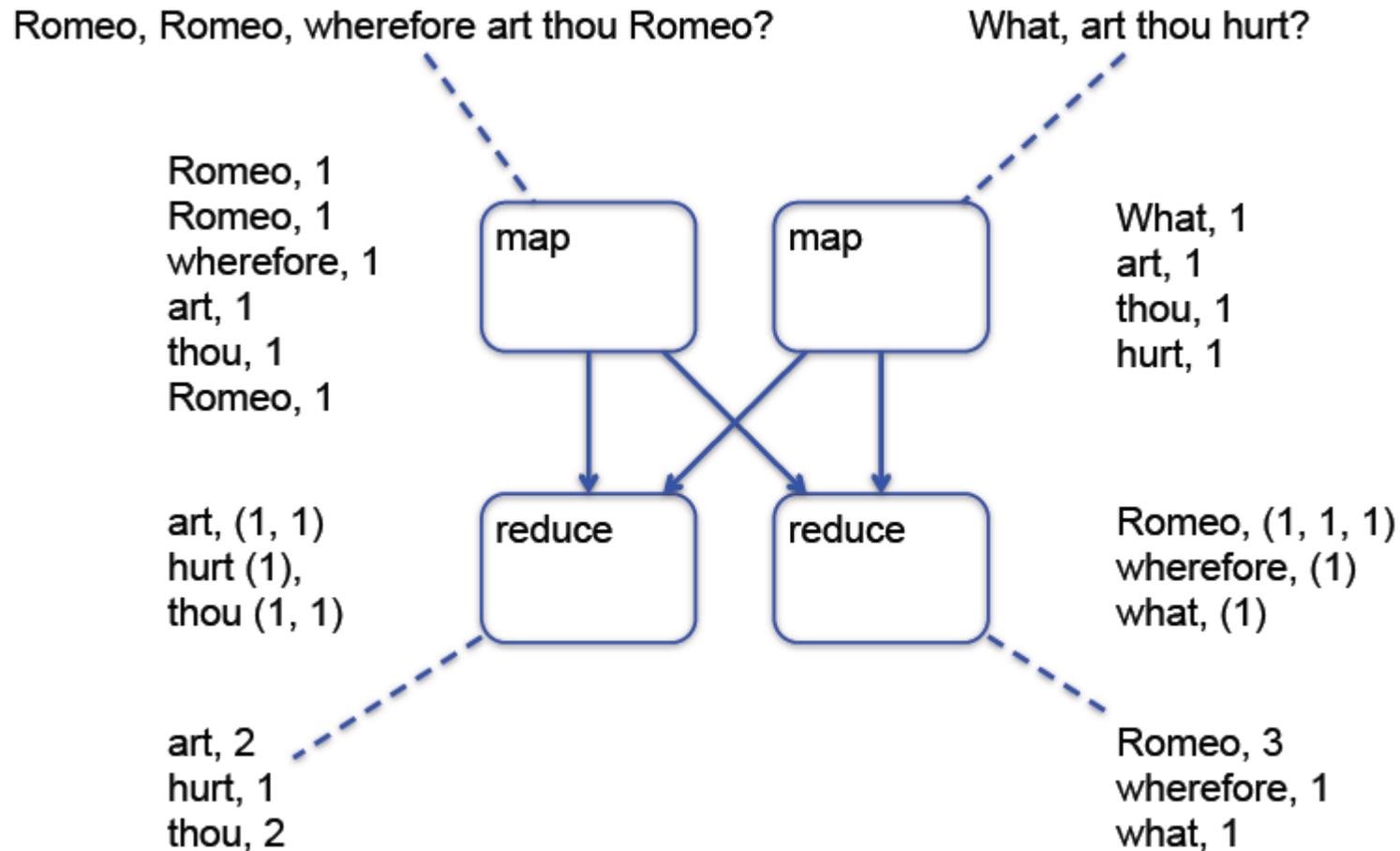


- *Problem*: Counting words in a parallel fashion
 - How many times different words appear in a set of files
 - **juliet.txt**: Romeo, Romeo, wherefore art thou Romeo?
 - **benvolio.txt**: What, art thou hurt?
 - Expected output: Romeo (3), art (2), thou (2), art (2), hurt (1), wherefore (1), what (1)

- *Solution*: Map-Reduce Job

```
m (filename, line) {
    foreach (word in line)
        emit(word, 1);
}

r (word, numbers) {
    int sum = 0;
    foreach (value in numbers) {
        sum += value;
    }
    emit(word, sum);
}
```



■ Selection / projection / aggregation

□ SQL Query:

```
SELECT year, SUM(price)
FROM sales
WHERE area_code = "US"
GROUP BY year
```

□ Map/Reduce job:

```
map(key, tuple) {
    int year = YEAR(tuple.date);
    if (tuple.area_code = "US")
        emit(year, {'year' => year, 'price' => tuple.price });
}

reduce(key, tuples) {
    double sum_price = 0;
    foreach (tuple in tuples) {
        sum_price += tuple.price;
    }
    emit(key, sum_price);
}
```

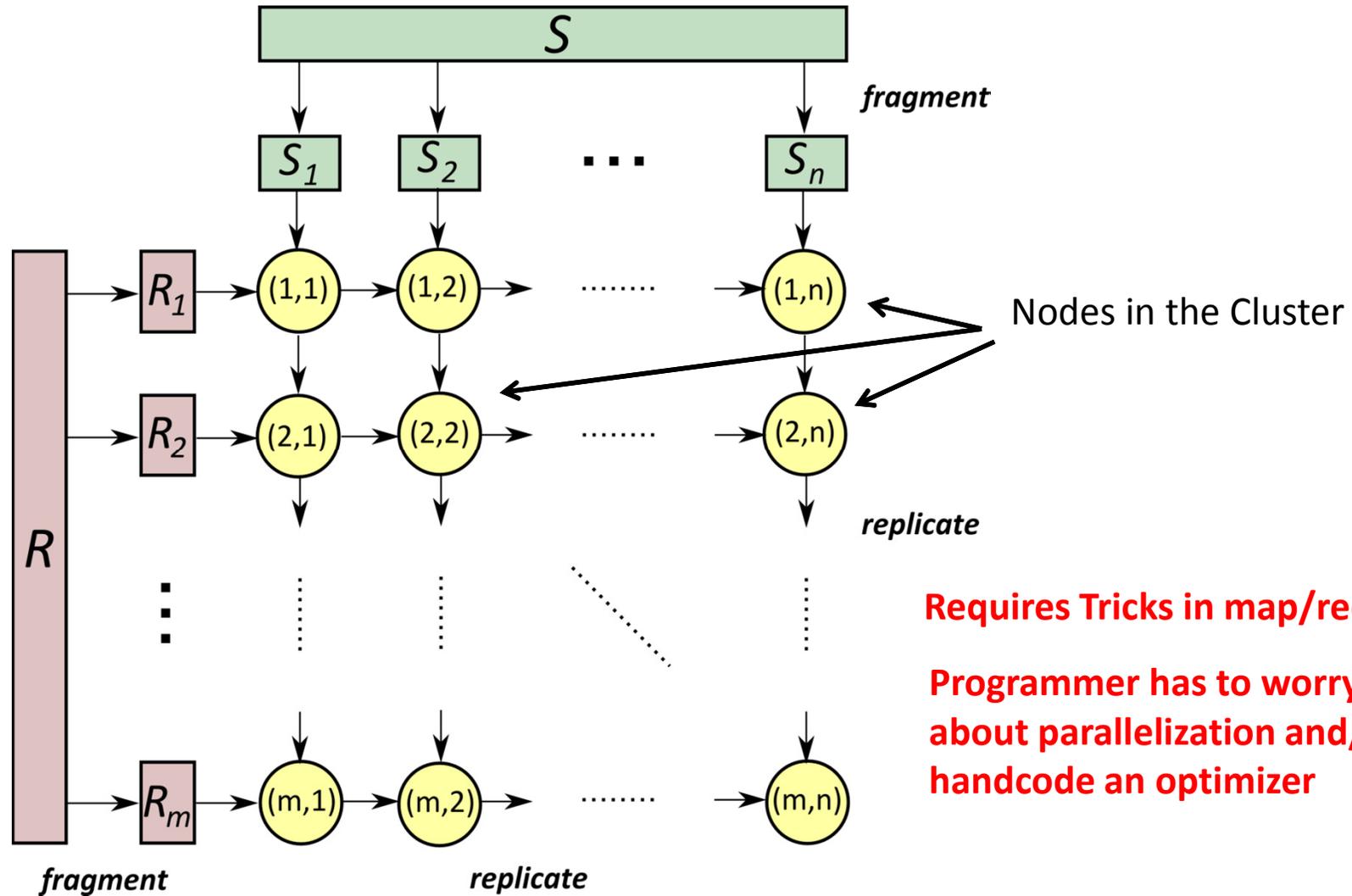
■ Sorting

□ SQL Query:

```
SELECT *  
FROM sales  
ORDER BY year
```

□ Map/Reduce job:

```
map(key, tuple) {  
    emit(YEAR(tuple.date) DIV 10, tuple);  
}  
  
reduce(key, tuples) {  
    emit(key, sort(tuples));  
}
```

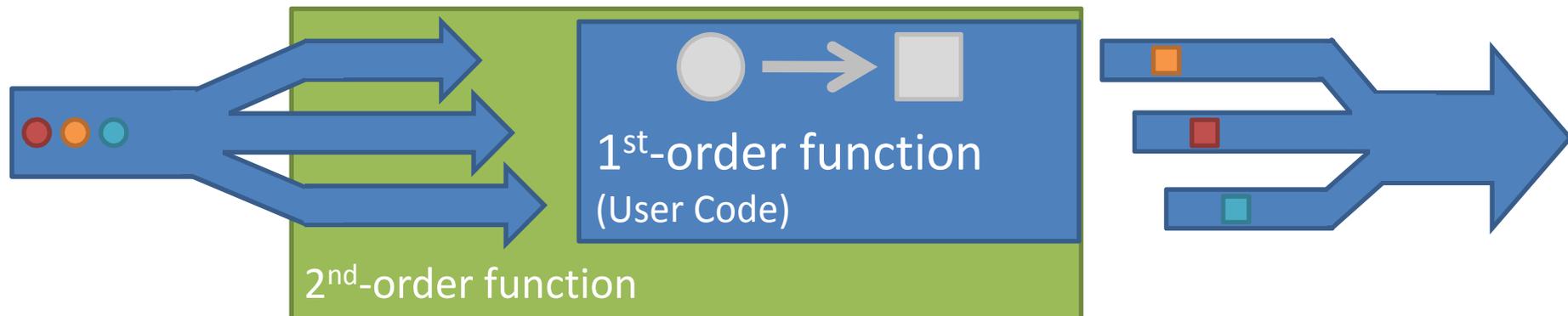


Requires Tricks in map/reduce
Programmer has to worry about parallelization and/or handcode an optimizer

Important generic example. Specialized parallel joins may exploit data locality

- PACT is a generalization and extension of MapReduce
 - PACT inherits many concepts from MapReduce

- Both are inspired by functional programming
 - Fundamental concept of programming model are 2nd-order functions
 - User writes 1st-order functions (user functions)
 - User code can be arbitrarily complex (albeit functional!)
 - 2nd-order function calls 1st-order function with independent data subsets
 - No common state should be held between calls of user function



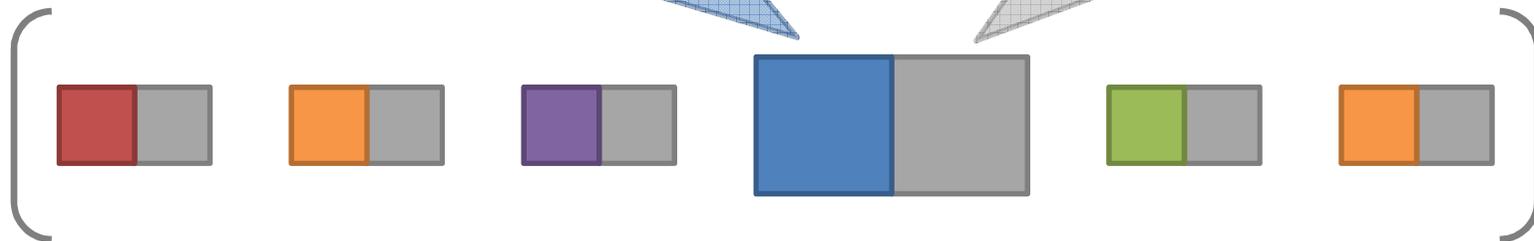
- Both use a common data format
 - Data is processed as pairs of keys and values
 - Keys and values can be arbitrary data structures

Key:

- Used to build independent subsets
- Must be comparable and hashable
- Does not need to be unique
 - no Primary Key semantic!
- Interpreted only by user code

Value:

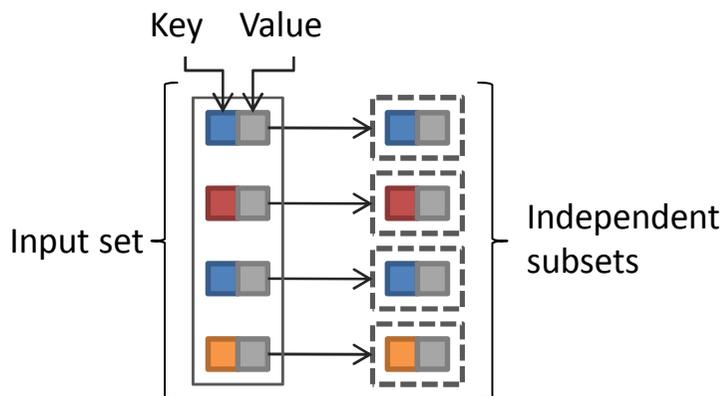
- Holds application data
- Interpreted only by user code



- MapReduce provides two 2nd-order functions

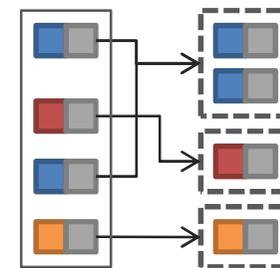
Map:

- All pairs are independently processed

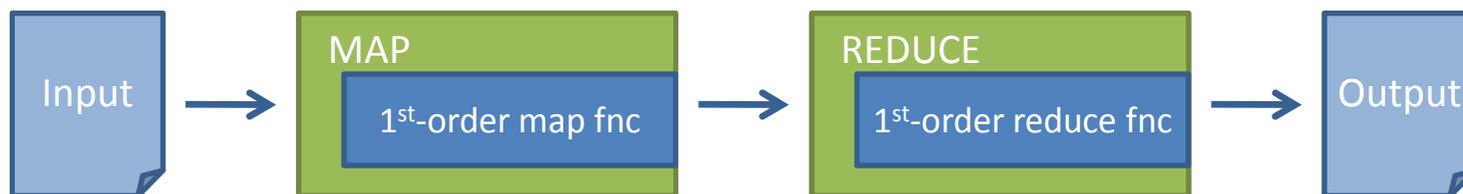


Reduce:

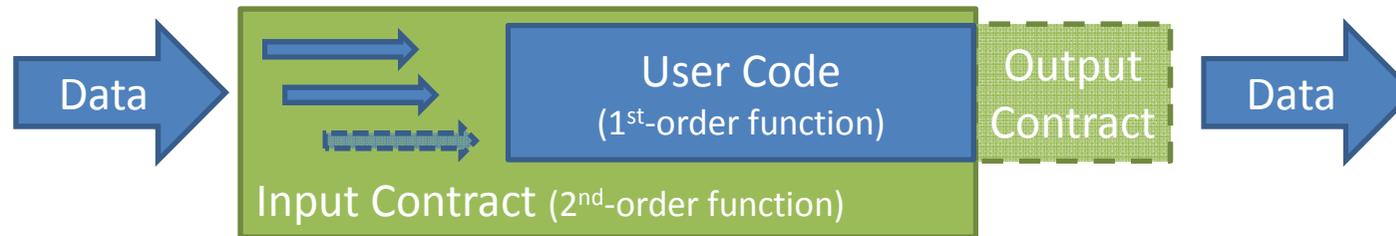
- Pairs with identical key are grouped
- Groups are independently processed



- MapReduce programs has fixed structure:



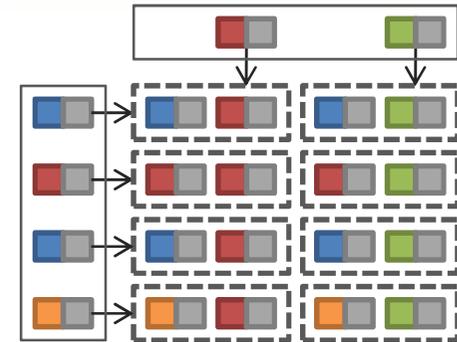
- Generalization and Extension of MapReduce Programming Model
- Based on Parallelization Contracts (PACTs)



- Input Contract
 - 2nd-order function; generalization of Map and Reduce
 - Generates independently processable subsets of data
- User Code
 - 1st-order function
 - For each subset independently called
- Output Contract
 - Describes properties of the output of the 1st-order function
 - Optional but enables certain optimizations
 - Think: "interesting properties"

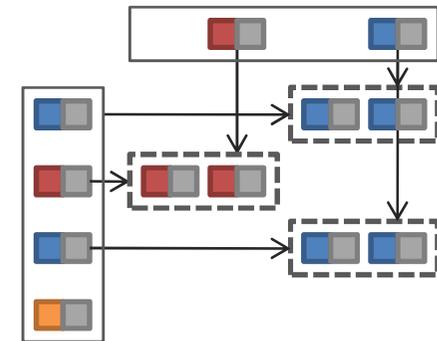
■ Cross

- Builds a Cartesian Product
- Elements of CP are independently processed



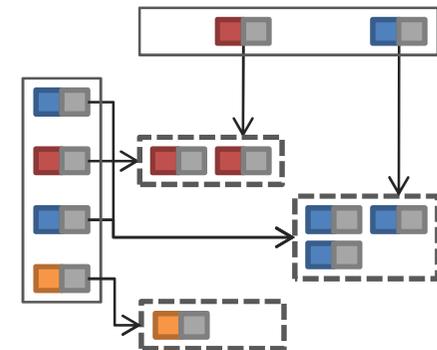
■ Match

- Performs an equi-join on the key
- Join candidates are independently processed



■ CoGroup

- Groups each input on key
- Groups with identical keys are processed together



■ Same-Key

- User Function does not alter the key



■ Super-Key

- Key generated by UF is a super-key of the input key



■ Unique-Key

- Data source or UF produces unique keys



- PACT Programs effectively are data flow graphs
 - Data comes from sources and flows to sinks
 - PACTs process data in-between sources and sinks
 - Multiple sources and sinks allowed
 - Arbitrary complex directed acyclic data flows can be composed

