# The Berlin Big Data Center

**Jonas Traub**

**Technische Universität Berlin / DFKI IAM**

www.dima.tu-berlin.de | bbdc.berlin | jonas.traub@tu-berlin.de

with materials from Tilmann Rabl and Volker Markl

Beuth Hochschule

ZIB Berlin

Christof
Schütte
Information-
based Medicine

Alexander
Reinefeld
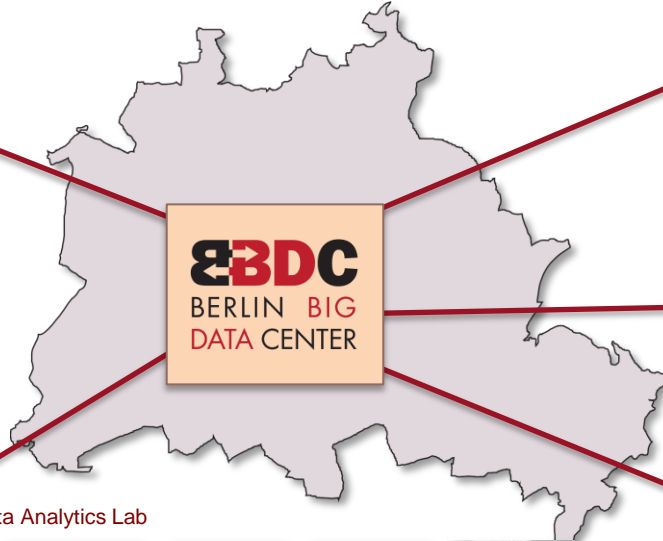File Systems,
Supercomputing

Tim
Conrad
Bioinformatics

Stefan Edlich
Software Engineering

Fritz-Haber-Institut,
Max-Planck-Gesellschaft

Matthias Scheffler

Material Science

Technische Universität Berlin, Data Analytics Lab

Volker
Markl
Data Management

Klaus R.
Müller
Machine
Learning

Anja
Feldmann
Computer
Networks

Odej
Kao
Distributed
Systems

Thomas
Wiegand
Video Mining

Ziawasch
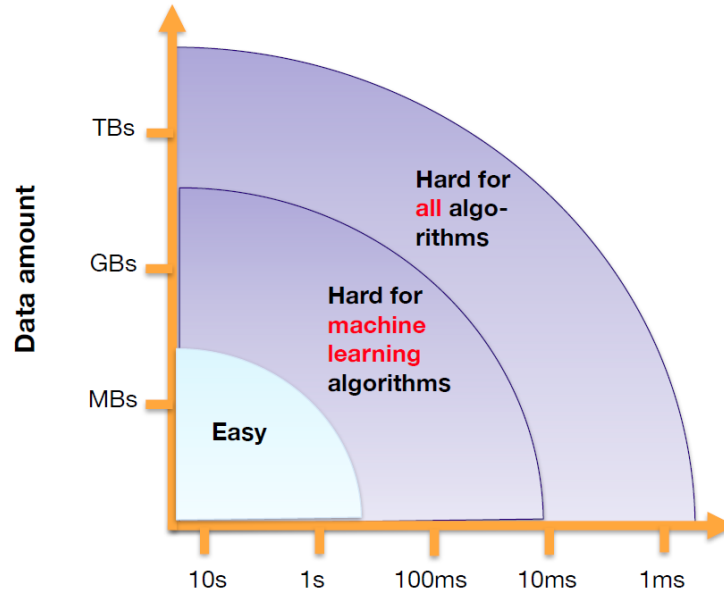Abedjan
Big Data
Management

DFKI
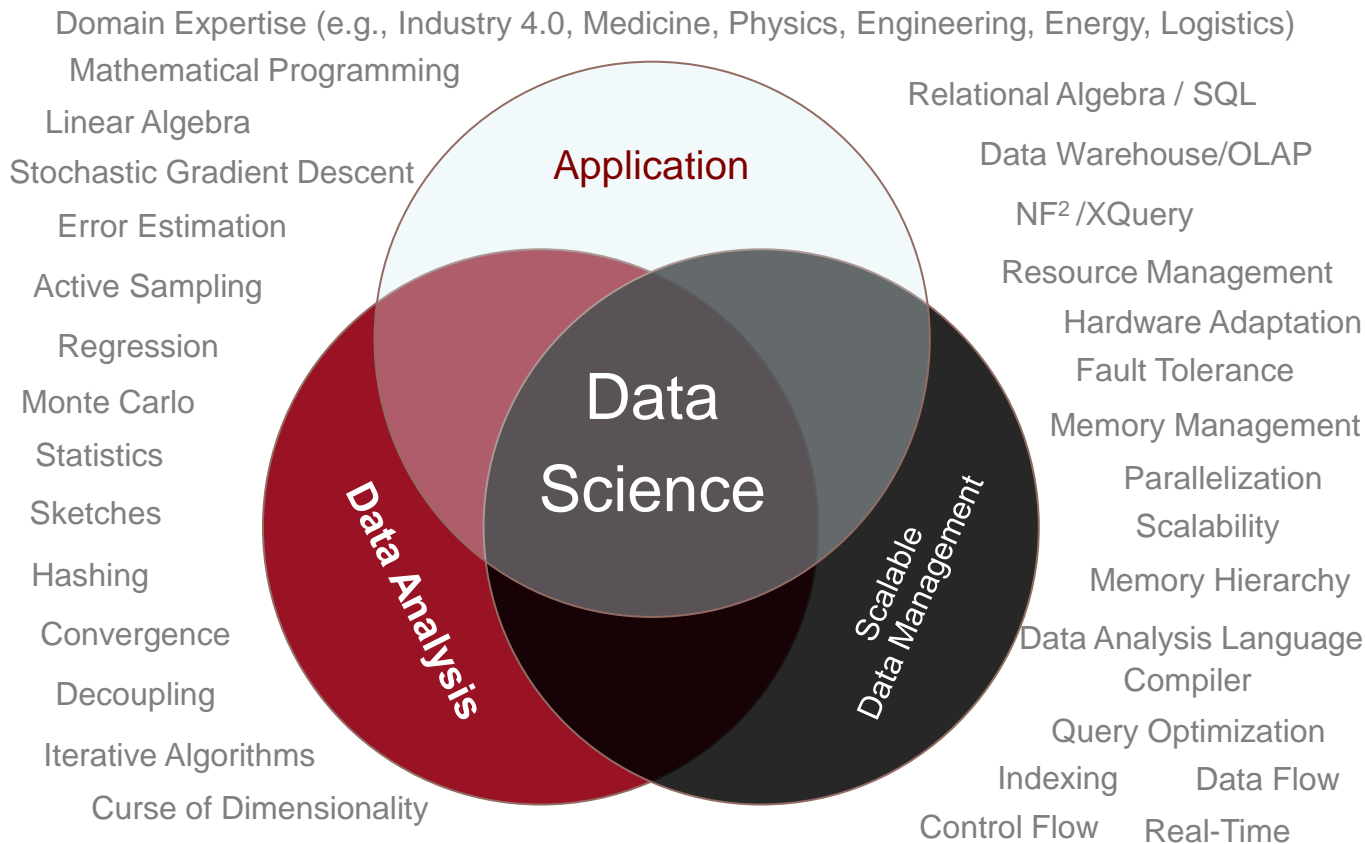
Hans Uszkoreit
Language Technology

# BBDC Goals

- **Pooling expertise** in scalable data management, data analytics, and big data applications.

- **Conducting fundamental research** to develop novel and automatically scalable technologies capable of performing "deep analysis" of big data.

- **Developing an integrated, declarative, highly scalable, open-source system** that enables the specification, automatic optimization, parallelization, hardware adaptation, fault-tolerance, and efficient execution of advanced data analysis problems using varying methods that leverage our work on Apache Flink.

- **Transferring technology and know-how** to support innovation in companies and startups.

- **Educating data scientists** with respect to the five big data dimensions via leading educational programs.

- **Empowering people to leverage** "Smart Data."

- **Enabling the general public to conduct sound data-driven decision-making.**
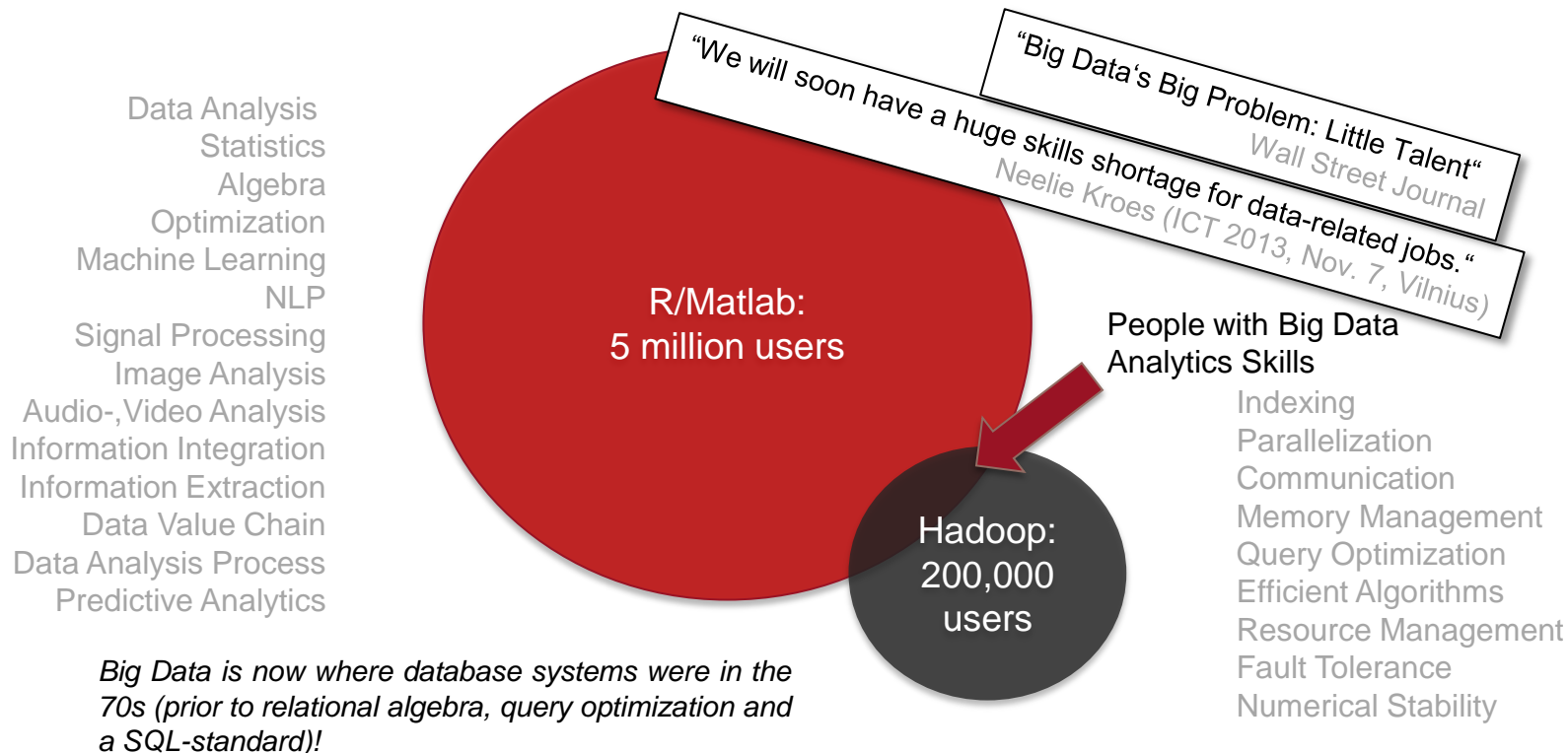
# Big Data – System View



Tension between *performance* and *algorithmic expressiveness*

# "Data Scientist" – "Jack of All Trades!"

Domain Expertise (e.g., Industry 4.0, Medicine, Physics, Engineering, Energy, Logistics)

Mathematical Programming

Linear Algebra

Stochastic Gradient Descent

Error Estimation

Active Sampling

Regression

Monte Carlo

Statistics

Sketches

Hashing

Convergence

Decoupling

Iterative Algorithms

Curse of Dimensionality

Application

Data Analysis

Data Science

Scalable Data Management

Relational Algebra / SQL

Data Warehouse/OLAP

NF$^2$ /XQuery

Resource Management

Hardware Adaptation

Fault Tolerance

Memory Management

Parallelization

Scalability

Memory Hierarchy

Data Analysis Language Compiler

Query Optimization

Indexing        Data Flow

Control Flow    Real-Time

# Big Data Analytics Requires Systems Programming

Data Analysis
Statistics
Algebra
Optimization
Machine Learning
NLP
Signal Processing
Image Analysis
Audio-,Video Analysis
Information Integration
Information Extraction
Data Value Chain
Data Analysis Process
Predictive Analytics

R/Matlab:
5 million users

Hadoop:
200,000
users

"We will soon have a huge skills shortage for data-related jobs."
Neelie Kroes (ICT 2013, Nov. 7, Vilnius)

"Big Data's Big Problem: Little Talent"
Wall Street Journal

People with Big Data
Analytics Skills

Indexing
Parallelization
Communication
Memory Management
Query Optimization
Efficient Algorithms
Resource Management
Fault Tolerance
Numerical Stability

*Big Data is now where database systems were in the 70s (prior to relational algebra, query optimization and a SQL-standard)!*

**Declarative languages to the rescue!**

# "Data Scientist" – "Jack of All Trades!"

Domain Expertise (e.g., Industry 4.0, Medicine, Physics, Engineering, Energy, Logistics)

Mathematical Programming

Linear Algebra

Stochastic Gradient Descent

Error Estimation

Active Sampling

Regression

Monte Carlo

Statistics

Sketches

Hashing

Convergence

Decoupling

Iterative Algorithms

Curse of Dimensionality

Relational Algebra / SQL

Data Warehouse/OLAP

$NF^2$ /XQuery

Resource Management

Hardware Adaptation

Fault Tolerance

Memory Management

Parallelization

Scalability

Memory Hierarchy

Data Analysis Language Compiler

Query Optimization

Indexing        Data Flow

Control Flow     Real-Time

**Domain Expertise**

Traditional Research

Danger Zone!

Statistics

Data Science

Hacking Skills

Machine Learning

**New Technology to the Rescue!**

# Danger Zone!



**Total revenue generated by arcades**
correlates with
**Computer science doctorates awarded in the US**

Arcade revenue
Computer science doctorates

● Computer science doctorates ◆ Arcade revenue

tylervigen.com

# Danger Zone! Contd.



**Divorce rate in Maine**
correlates with
**Per capita consumption of margarine**

# Apache Flink–
# A Success Story
# Born in Berlin

http://flink.apache.org

# Timeline

# What can I do with it?

Batch
processing

Machine Learning at scale

Stream
processing

Graph Analysis

Flink

A big data processing system that can **natively** support all these workloads.

# What is Apache Flink?

Apache Flink® is an open-source stream processing framework for distributed, high-performing, always-available, and accurate data streaming applications.
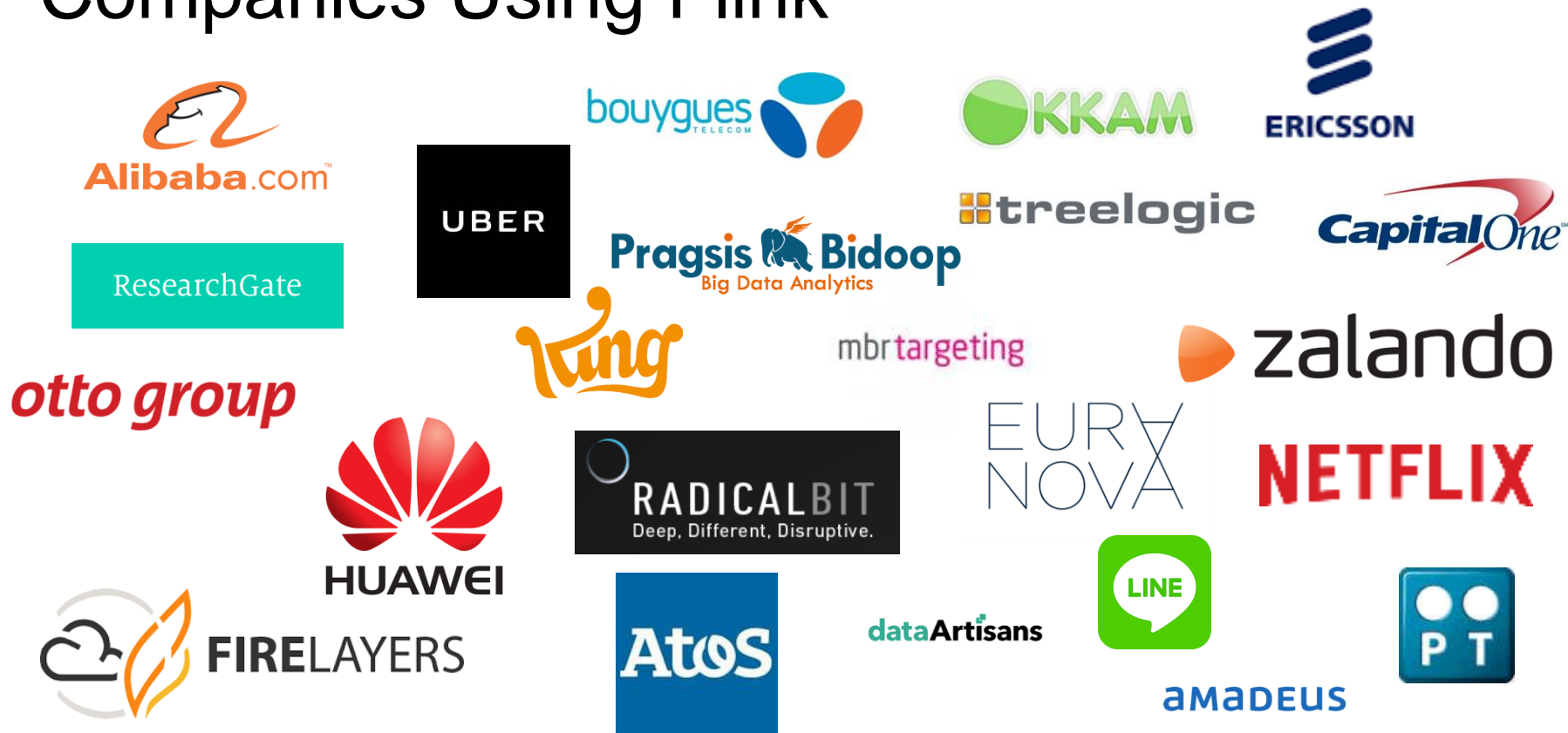
**Key Features**

- Bounded and unbounded data
- Event time semantics
- Stateful and fault-tolerant
- Running on thousands of nodes with very good throughput and latency
- Exactly-once semantics for stateful computations
- Flexible windowing based on time, count, or sessions in addition to data-driven windows
- **DataSet** and **DataStream** programming abstractions are the foundation for user programs and higher layers



P. Carbone, A. Katsifodimos, S. Ewen, V. Markl, S. Haridi, K. Tzoumas: Apache Flink™: Stream and Batch Processing in a Single Engine. IEEE Data Eng. Bull. 38(4): 28-38 (2015)
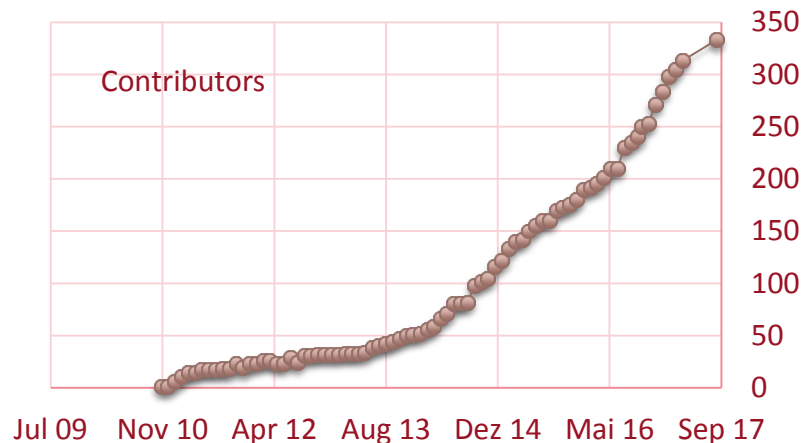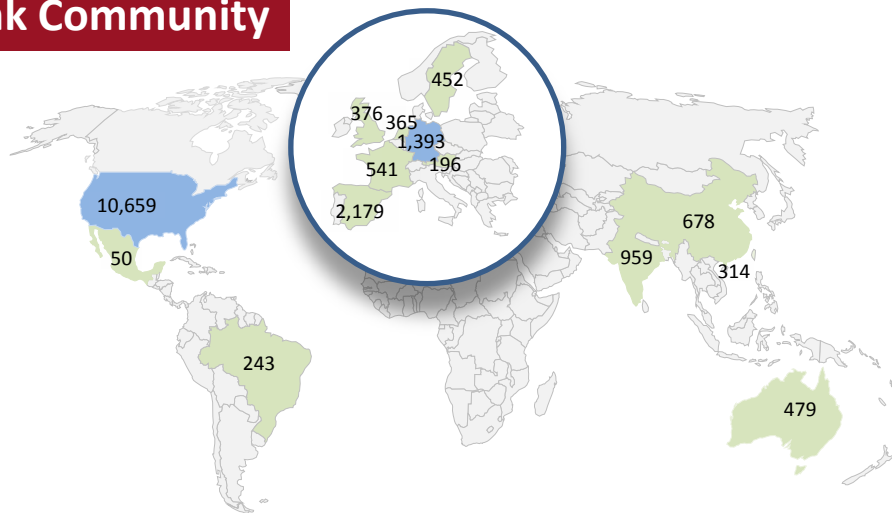
# Companies Using Flink
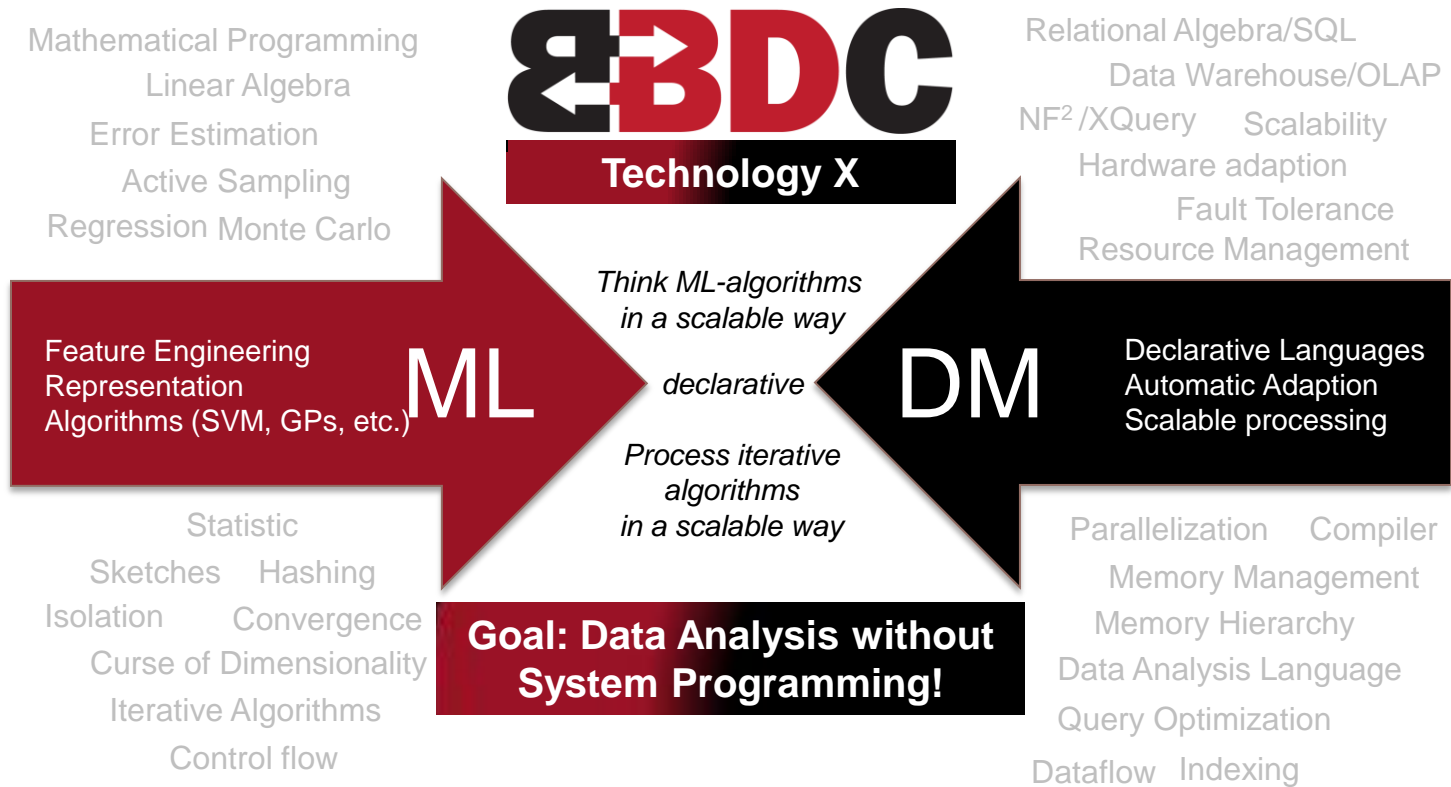
# Innovation und Transfer am Beispiel Apache Flink

## Flink Community



Contributors

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Jul 09 | Nov 10 | Apr 12 | Aug 13 | Dez 14 | Mai 16 | Sep 17 |

Map values:
452
376 365
1,393
541 196
2,179
10,659
50
243
678
959
314
479

| | | | |
|---|---|---|---|
| 17.820+ | Meetup-Mitglieder weltweit | 14+ | Länder mit regelmäßigen Meetup-Events |
| 328+ | open-source Entwickler (contributors) | 30+ | Anwenderunternehmen |
| 40+ | Meetupgruppen weltweit | | Firmengründung dataArtisans |

# Machine Learning + Data Management = X



**BBDC**

**Technology X**

Mathematical Programming
Linear Algebra
Error Estimation
Active Sampling
Regression Monte Carlo

Relational Algebra/SQL
Data Warehouse/OLAP
$NF^2$ /XQuery    Scalability
Hardware adaption
Fault Tolerance
Resource Management

**ML**

Feature Engineering
Representation
Algorithms (SVM, GPs, etc.)

*Think ML-algorithms
in a scalable way*

*declarative*

*Process iterative
algorithms
in a scalable way*

**DM**

Declarative Languages
Automatic Adaption
Scalable processing

Statistic
Sketches    Hashing
Isolation    Convergence
Curse of Dimensionality
Iterative Algorithms
Control flow

**Goal: Data Analysis without
System Programming!**

Parallelization    Compiler
Memory Management
Memory Hierarchy
Data Analysis Language
Query Optimization
Dataflow  Indexing

# What, Not How! Consider K-means Clustering.



"What"
(Apache Flink)
(Scala frontend)

65 lines of code
short development time
robust runtime

Declarative data analysis program with automatic optimization, parallelization and hardware adaption
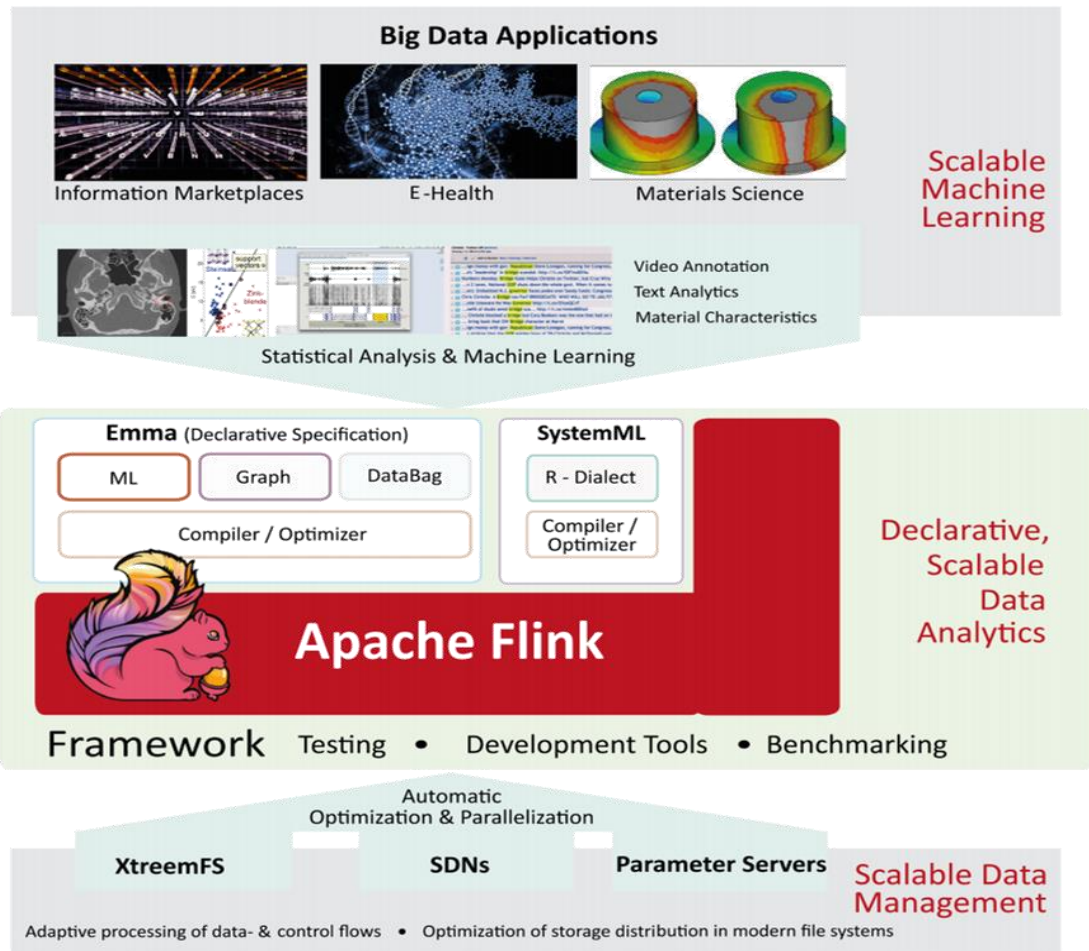
"How"
(Hadoop)

486 lines of code
long development time
non-robust runtime

Hand-optimized code
(data-, load- and system dependent)

# Big Data Analytics Without Systems Programming!
# (What, Not How!)

**Data Analyst**

Larger human base of "data scientists"
Reduction of "human" latencies
Cost reduction

**Machine**

Description of "What?"
(declarative specification)

Description of "How?"
(state of the art in scalable data analysis)
Map/Reduce, MPI

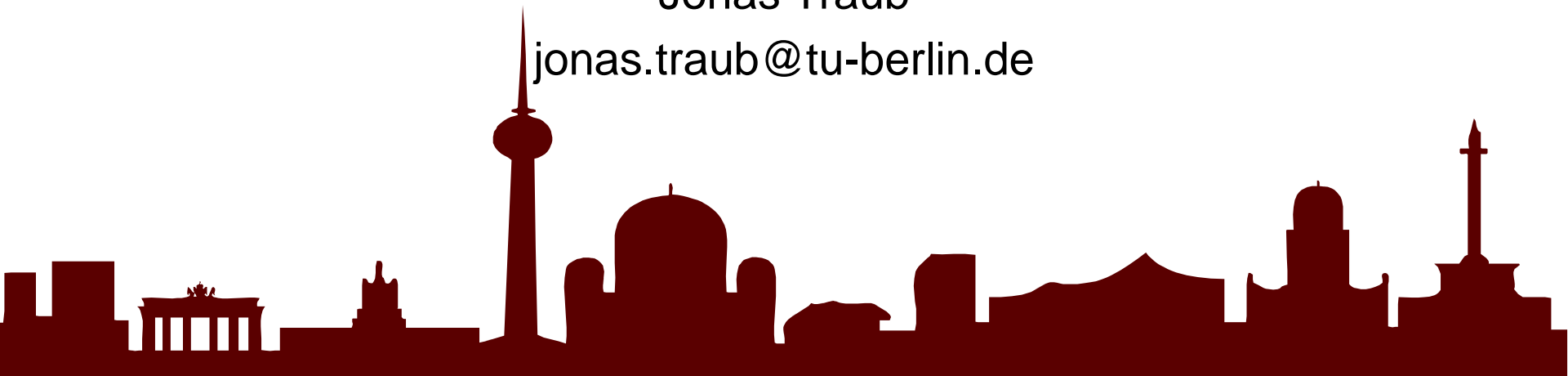# Flink in the BBDC Stack

# BBDC Lessons Learned & Challenges

- International competitors have far more funding and visibility
  - Berkeley AmpLab: $30 Mio for 6 years, recently grown into Berkeley Institute for Data Science with another $40 Mio funding
  - UK Turing Institute: GBP 67 Mio funding
  - And many others (e.g., across the US, China, Korea, and Japan)

- German companies are followers, not leaders in big data
  - Many of Germany's large companies have not yet developed a big data strategy and are risk-averse, or focus too much on short-term and established solutions.
  - It is far easier to work with "new" companies to transfer novel technologies
    - ResearchGate, Zalando, King, IMR, and Spotify, among others
  - Open source solutions and/or establishing new companies are the best route to turn research into innovation
  - US and large international companies are easier to collaborate with, often times via their respective German subsidiary.

# Thank You

Contact:

Jonas Traub

jonas.traub@tu-berlin.de

# The Berlin Big Data Center

## Jonas Traub

## Technische Universität Berlin / DFKI IAM

www.dima.tu-berlin.de | bbdc.berlin | jonas.traub@tu-berlin.de

with materials from Tilmann Rabl and Volker Markl