HPI Hasso Plattner Institut

IT Systems Engineering | Universität Potsdam

**Analyze Genomes: An In-Memory Technology Use Case**
In-Memory Computing for Life Sciences
Dr. Matthieu-P. Schapranow, Hasso Plattner Institute

# Our Motivation
# Personalized Medicine

- **Motivation:** Can we analyze the entire data of a patient, incl. Electronic Health Records (EHR) and genome data, during a doctor's visit?

- Genome data analysis may add up to weeks,
i.e. biopsy, biological preparation, sequencing,
alignment, variant calling, full analysis, and evaluation

- Issue: Complex and time-consuming data processing tasks

- In-memory technology accelerates genome data processing

  - Highly parallel alignment / variant calling

  - Real-time analysis of individual patient or cohort data
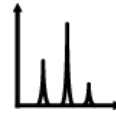
  - Combined search in structured / unstructured data

Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# Our Challenge
# Distributed Big Data Sources

**Human genome/biological data**
600GB per full genome
15PB+ in databases of leading institutes

**Human proteome**
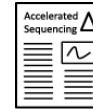160M data points (2.4GB) per sample
>3TB raw proteome data in ProteomicsDB

**Hospital information systems**
Often more than 50GB

**PubMed database**
>23M articles

Accelerated Sequencing

**Cancer patient records**
>160k records at NCT

**Medical sensor data**
Scan of a single organ in 1s creates 10GB of raw data

**Prescription data**
1.5B records from 10,000 doctors and 10M Patients (100 GB)

**Clinical trials**
Currently more than 30k recruiting on ClinicalTrials.gov

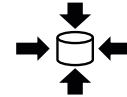Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014
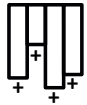
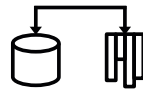Combined column and row store

Map/Reduce
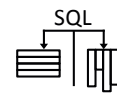
Single and multi-tenancy

Lightweight compression

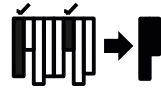Insert only for time travel

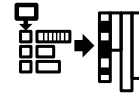Real-time replication

Working on integers

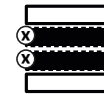SQL interface on columns and rows

Active/passive data store

Minimal projections

Group key

Reduction of software layers

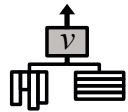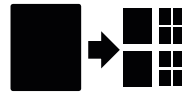Dynamic multi-threading

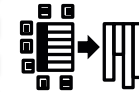Bulk load of data

Object-relational mapping

Text retrieval and extraction engine

No aggregate tables

Data partitioning

Any attribute as index

No disk

On-the-fly extensibility

Analytics on historical data

Multi-core/ parallelization
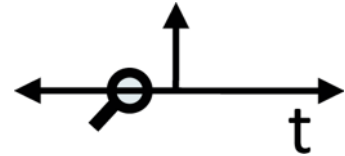
# Meaningful In-Memory Database Concepts
## Text Mining

- Full text indexing for any text attributes
- User-defined **dictionaries** to define entities
  http://scn.sap.com/community/developer-center/hana/blog/2013/12/27/hana-text-analysis-with-custom-dictionaries
- Custom Grouper User Language (**CGUL**) rules to create token-based regular expressions with linguistic attributes
  http://wiki.scn.sap.com/wiki/display/EIM/CGUL+Tips+and+Tricks+for+Entity+Extraction
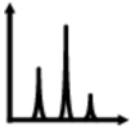
# Meaningful In-Memory Database Concepts
# Time Travel

- Process time series data by retrieving the complete database state at any period of history
- **History Database Tables**
- SELECT TEMPERATURE FROM "PATIENTS"."ICS_TEMP" WHERE PATIENT NAME = 'Matthieu Schapranow' AS OF COMMIT ID 209811
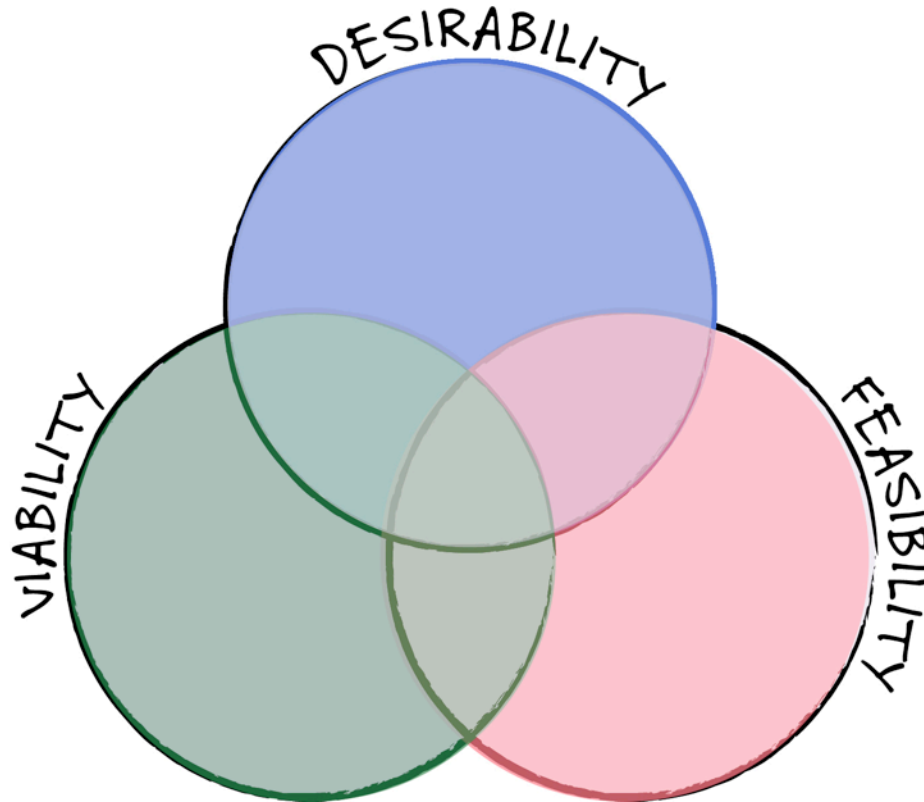- http://scn.sap.com/community/developer-center/hana/blog/2013/02/12/when-i-travelled-through-time-using-sap-hana

# Meaningful In-Memory Database Concepts
# Predictive Analysis Library (PAL)

- Provides specific **analysis functions** tightly integrated within the database, e.g. k-means or hierarchical clustering

- http://help.sap.com/hana/
  SAP_HANA_Predictive_Analysis_Library_PAL_en.pdf

# Our Vision
## Personalized Medicine



Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# Our Vision
## Personalized Medicine

**Desirability**

- Leveraging directed customer services
- Portfolio of integrated services for clinicians, researchers, and patients
- Include latest research results, e.g. most effective therapies

**Viability**

- Enable personalized medicine also in far-off regions and developing countries
- Share data via the Internet to get feedback from word-wide experts (cost-saving)
- Combine research data (publications, annotations, genome data) from international databases in a single knowledge base
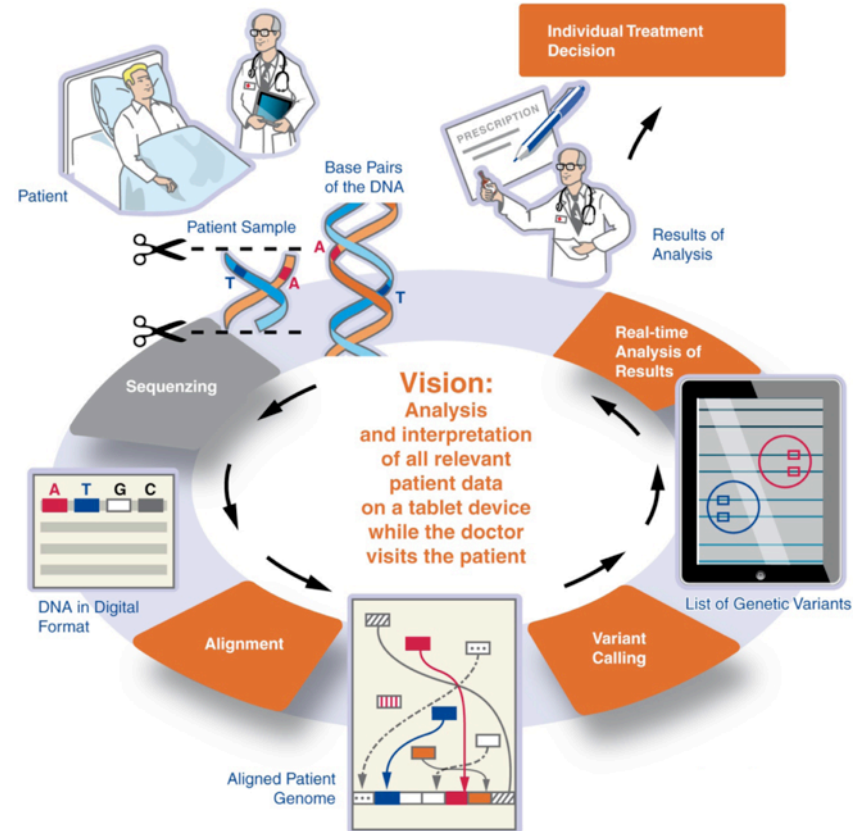
**Feasibility**

- HiSeq 2500 enables high-coverage whole genome sequencing in ≈1d
- IMDB enables allele frequency determination of 12B records within <1s
- 1 relevant out of 80M annotations <1s
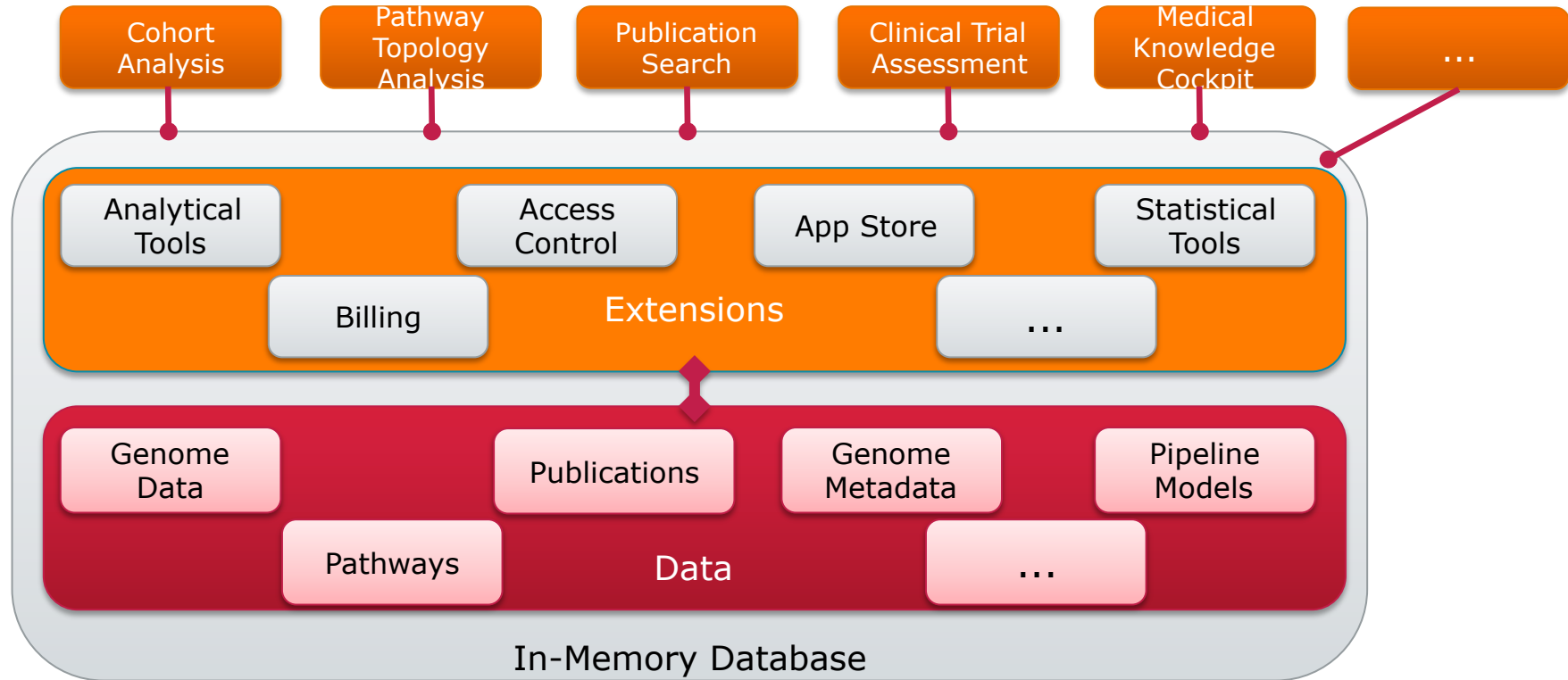- Data preparation as a service reduces TCO

# High-Performance In-Memory Genome Project
# Integration of Genomic Data

- Preprocessing of DNA (alignment, variant calling) can be modeled and is executed as integrated process

- Results are directly stored in in-memory databases, e.g. for
  - Statistical analyses, and
  - Links to latest research knowledge

- Real-time analysis of genome data enables completely new way of research and therapies, e.g. instant comparison with patient cohorts



Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# High-Performance In-Memory Genome Project Architectural Overview



Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# High-Performance In-Memory Genome Project
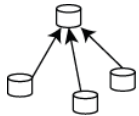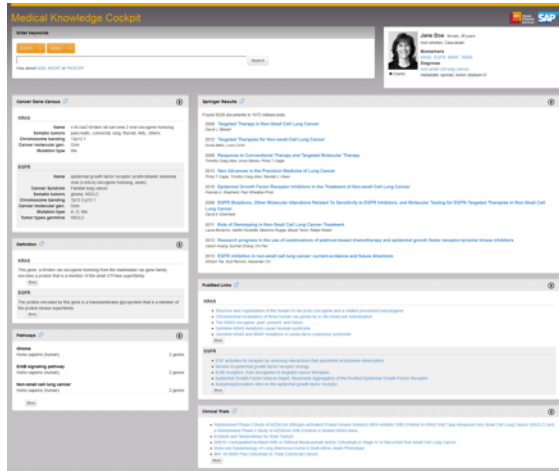# Selected Research Topics

**Improving Analyses:**

- Information combination, e.g. Medical Knowledge Cockpit, Oncolyzer
- Genome Browser enables deep dive into the genome
- Cohort Analysis, e.g. clustering of patient cohorts
- Combined Search, e.g. in clinical trials and side-effect databases
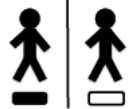- Pathways Topology Analysis, e.g. to identify cause/effect

**Improving Data Preparations:**

- Graphical modeling of Genome Data Processing (GDP) pipelines
- Scheduling and execution of multiple GPD pipelines in parallel
- App store for medical knowledge (bring algorithms to data)
- Exchange of sensitive data, e.g. history-based access control
- Billing processes for intellectual property and services

# High-Performance In-Memory Genome Project
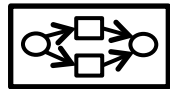# Medical Knowledge Cockpit



**Unified access** to structured and un-structured data sources



**Automatic clinical trial matching** build on text analysis features

- Search for affected genes in distributed and heterogeneous data sources

- Immediate exploration of relevant information, such as
  - Gene descriptions,
  - Molecular impact and related pathways,
  - Scientific publications, and
  - Suitable clinical trials.

- No manual searching for hours or days: In-memory technology translates searching into interactive finding!

# High-Performance In-Memory Genome Project
# Drug Response Testing



#Mutations by TASK_ID and RS_NUMBER

**Interactive analysis** of correlations between drugs and genetic variants
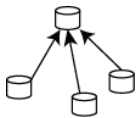
- Drug response depends on individual genetic variants of tumors

- Challenge: Identification of relevant genetic variants and their impact on drug response is a ongoing research activity, e.g. Xenograft models

- Exploration of experiment results is time-consuming and Excel-driven

- In-memory technology enables interactive exploration of experiment data to leverage new scientific insights

# High-Performance In-Memory Genome Project
## Search in Structured and Unstructured Medical Data



**Clinical Trials**

**Internal**
- Panitumumab Plus Pemetrexed and Cisplatin (PemCisP) Versus Pem...
- FLO +/- Pazopanib as First-line Treatment in Advanced Gastric Cance...

**External**
- Molecular Profiling and Targeted Therapy for Advanced Non-Small C...
- Erlotinib and Temsirolimus for Solid Tumors
- Molecular Epidemiology of Lung Adenocarcinoma in Multi-ethnic Asia...
- Safety and Efficacy Study of Neratinib and Cetuximab to Treat Patient...
- French National Observatory of the Patients With Non-small Cell Lung...

**Unified access** to structured and unstructured data sources

**Clinical trial matching** using text analysis features

- Extended text analysis feature by medical terminology
  - □ Genes (122,975 + 186,771 synonyms)
  - □ Medical terms and categories (98,886 diseases + 48,561 synonyms, 47 categories)
  - □ Pharmaceutical ingredients (7,099 + 5,561 synonyms)
- Indexed clinicaltrials.gov database (145k trials/ 30,138 recruiting)
- Extracted, e.g., 320k genes, 161k ingredients, 30k periods
- Select all studies based on multiple filters in less than 500ms

# High-Performance In-Memory Genome Project
# Analysis of Patient Cohorts



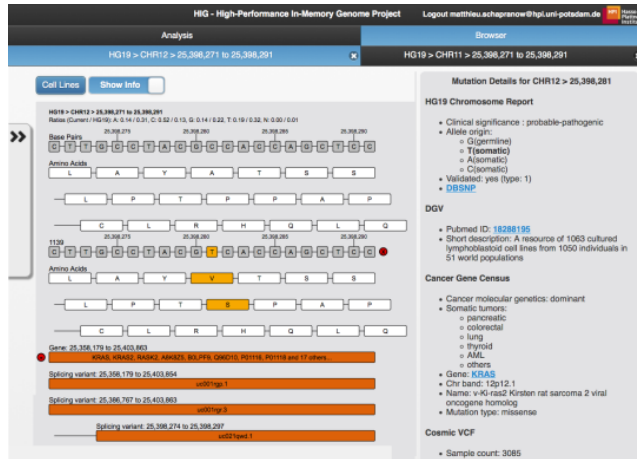**Fast clustering** directly performed within the in-memory database

- In a patient cohort, a subset does not respond to therapy – why?

- Clustering using various statistical algorithms, such as k-means or hierarchical clustering

- Calculation of all locus combinations in which at least 5% of all TCGA participants have mutations: 200ms for top 20 combinations

- Individual clusters are calculated in parallel directly within the database

- K-means algorithm: 50ms (PAL) vs. 500ms (R)
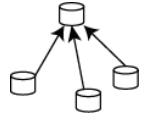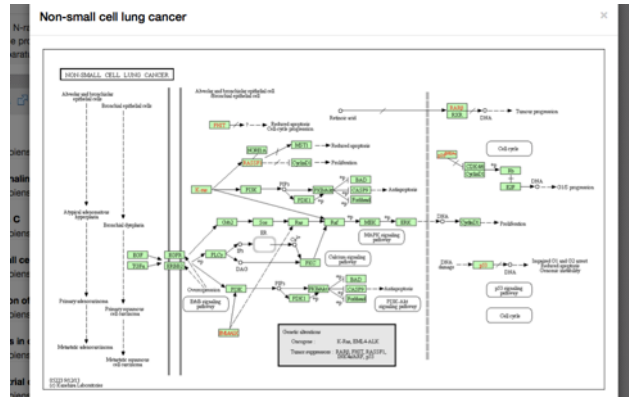
# High-Performance In-Memory Genome Project
# Genome Browser



**Unified access** to multiple formerly disjoint data sources

**Matching of genetic variants** and relevant annotations

- Genome Browser: Comparison of multiple genomes with reference
- Combined knowledge base: latest relevant annotations and literature, e.g. NCBI, dbSNP, and UCSC

- Detailed exploration of genome locations and existing associations
- Ranked variants, e.g. accordingly to known diseases

- Links to more detailed sources enable fast identification of relevant information while eliminating long-lasting searches.

Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# High-Performance In-Memory Genome Project
# Pathway Analysis



Non-small cell lung cancer



**Unified access** to multiple formerly disjoint data sources



**Pathway analysis** of genetic variants with graph engine

- Search in pathways is limited to "is a certain element contained" today

- Integrated >1,5k pathways from international sources, e.g. KEGG, HumanCyc, and WikiPathways, into HANA

- Implemented graph-based topology exploration and ranking based on patient specifics

- Enables interactive identification of possible dysfunctions affecting the course of a therapy before its start

# What to take home?
# Test it: http://we.AnalyzeGenomes.com

**For researchers**

- Enable real-time analysis of medical data
- Automatic assessment of data, e.g. scan of pathways to identify cellular impact of mutations
- Combined free-text search in publications, diagnosis, and EMR data, i.e. structured and unstructured data

**For clinicians**

- Preventive diagnostics to identify risk patients early
- Indicate pharmacokinetic correlations
- Scan for similar patient cases, e.g. to evaluate therapy success

**For patients**

- Identify relevant clinical trials and medical experts
- Start most appropriate therapy as early as possible

Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# Keep in contact with us!



Dr. Matthieu-P. Schapranow
schapranow@hpi.uni-potsdam.de
http://we.analyzegenomes.com/

Hasso Plattner Institute
Enterprise Platform & Integration Concepts
Dr. Matthieu-P. Schapranow
August-Bebel-Str. 88
14482 Potsdam, Germany

# BACKUP

Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# High-Performance In-Memory Genome Project
## Architectural Overview



Information and Feedback within the Window of Opportunity

Patients    Doctors    Insurers    Researchers

Real-Time Capturing and Data Analysis

In-Memory Database

*omics    Electronic Medical Records    Service Line Items    ...

All Relevant Medical Information

Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# High-Performance In-Memory Genome Project
# HANA Oncolyzer



- Research initiative for exchanging relevant tumor data to improve personalized treatment

- Honored 2012 by the Innovation Award of the German Capitol Region

- In-memory technology as key-enabler for real-time analysis of tumor data in seconds instead of hours

- Information available at your fingertips: In-memory technology on mobile devices, e.g. iPad

- Interdisciplinary cooperation between
  - Medical doctors,
  - Researchers, and
  - Software engineers.

Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# High-Performance In-Memory Genome Project
# HANA Oncolyzer
# Patient Details Screen

- Combines patient's time series data of specific patient and analysis results of patient cohort

- Real-time analysis across hospital-wide data whenever details screen is accessed

- http://epic.hpi.uni-potsdam.de/Home/HanaOncolyzer



Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# High-Performance In-Memory Genome Project
# HANA Oncolyzer
# Patient Analysis Screen

- Allows to real-time analysis on complete patient cohort
- Flexible filters and various chart types allow graphical exploration of data on mobile devices



Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# Medical Knowledge Cockpit
## Seamless Integration of Patient Specifics



- Google-like user interface for searching data
- Seamless integration of individual EMR data
- Search various sources for biomarkers, literature, and diseases

Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014

# Medical Knowledge Cockpit
## Publications



- In-place preview of relevant data, such as publications and publication meta data

- Incorporating individual filter settings, e.g. additional search terms

# Medical Knowledge Cockpit
## Publications



- Interactively explore relevant publications, e.g. PDFs
- Improved ease of exploration, e.g. by highlighted medical terms and relevant concepts

# Medical Knowledge Cockpit
## Relevant Scientific Findings at a Glance



- Personalized clinical trials, e.g. by incorporating patient specifics
- Classification of internal/external trials based on treating institute

Analyze Genomes, In-Memory Computing for Life Sciences, Dr. M.-P. Schapranow, Apr 8, 2014