


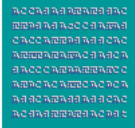


gdw Genome Data Warehouse		
		

Processing Genome Data using Scalable Database Technology

Johann Christoph Freytag, Ph.D.
freytag@dbis.informatik.hu-berlin.de
<http://www.dbis.informatik.hu-berlin.de>

Stanford University, February 2004

© 2004 DBIS, Humboldt Universität zu Berlin
Prof. Johann Christoph Freytag

Processing Genome Data
using Scalable Database Technology

My Background

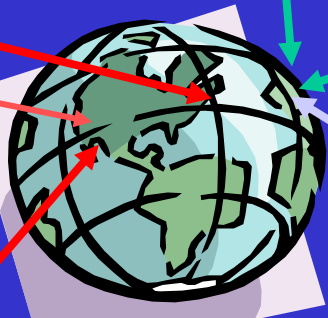
Professor of DBIS
@HUB

PHD @ Harvard Univ.

ERCR (European Computer Industry Research Centre), München (87-89)

DEC's Database Technology Center, München (90-93)

Visiting Scientist, Microsoft Res. (2002)



- Starburst project, IBM Almaden Research Center (85-87)
- Visiting Scientist, Almaden Research Center (97/98)
- Visiting Scientist, IBM SVL (2001)

© 2004 DBIS, Humboldt Universität zu Berlin
Prof. Johann Christoph Freytag

Processing Genome Data
using Scalable Database Technology

“What’s the Meaning of Life”

DNA → **RNA** → **Protein**

Transcription Translation

Genomic „Transmitter“ Messenger Gene product

Replication

© 2004 DBIS, Humboldt Universität zu Berlin
Prof. Johann Christoph Freytag

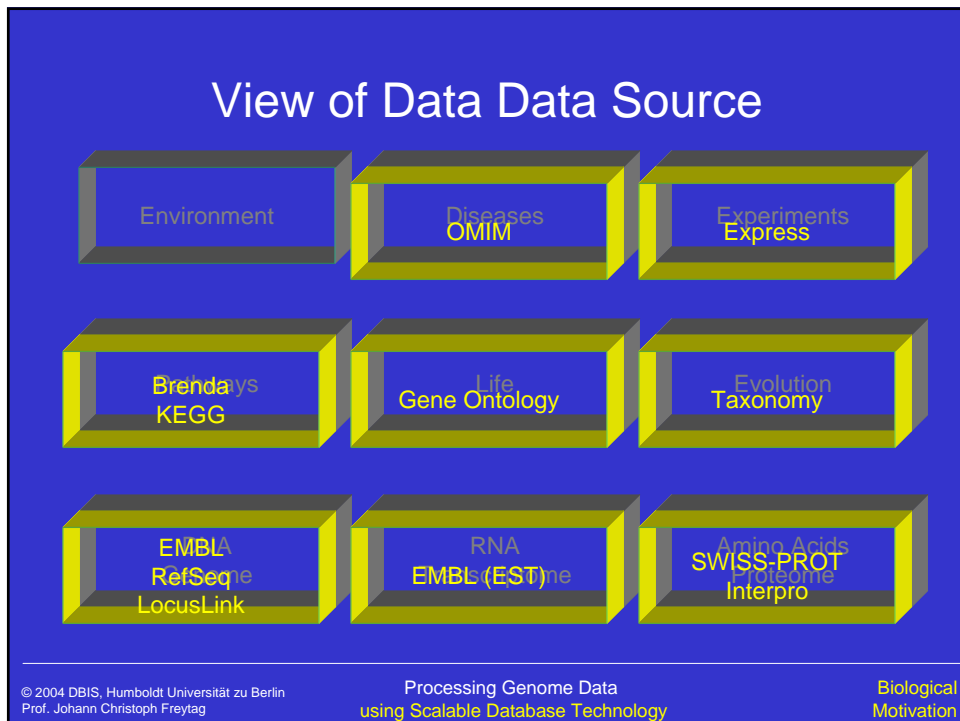
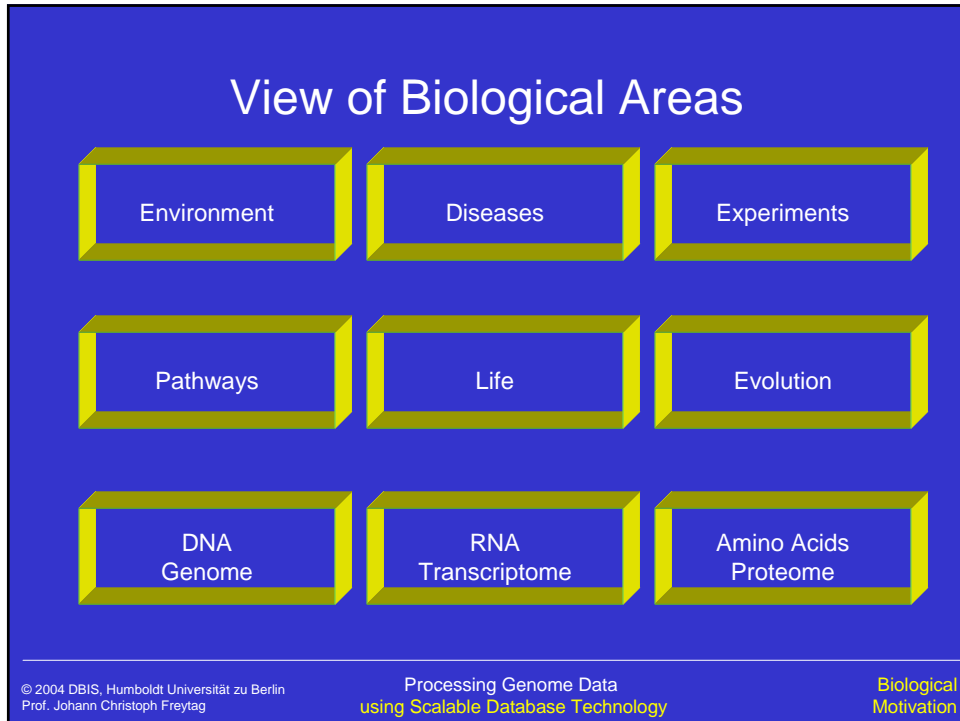
Processing Genome Data
using Scalable Database Technology

Overview

- (Biological) Motivation/Problems
- Using Database Technology
 - Gene-EYe Integration-Plattform
 - Data Cleansing
 - BLAST-Integration into GDB
 - In-and-Out-the-Database: Using Workflow for “dry” Experiments
- Summary

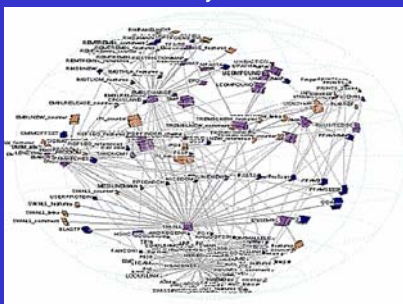
© 2004 DBIS, Humboldt Universität zu Berlin
Prof. Johann Christoph Freytag

Processing Genome Data
using Scalable Database Technology



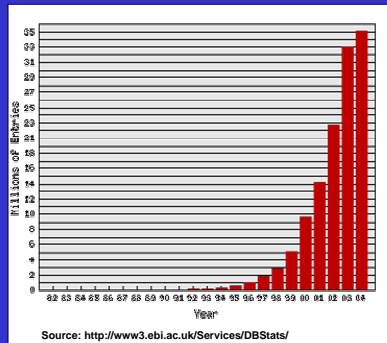
Complex Relationships

A graph depicting the relationships between 400+ biological data sources served by the EBI via SRS



More than 400 Data Sources on the WEB

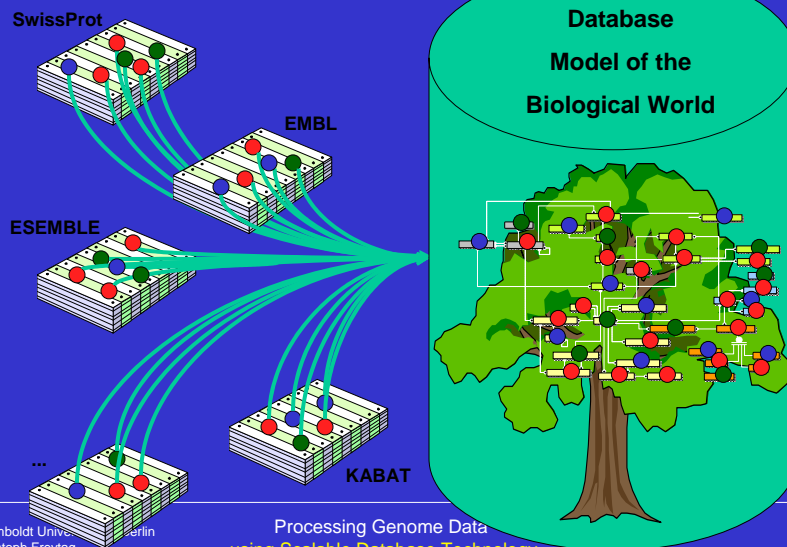
Database Growth of EMBL (# of records)



© 2004 DBIS, Humboldt Universität zu Berlin
Prof. Johann Christoph Freytag

Processing Genome Data
using Scalable Database Technology

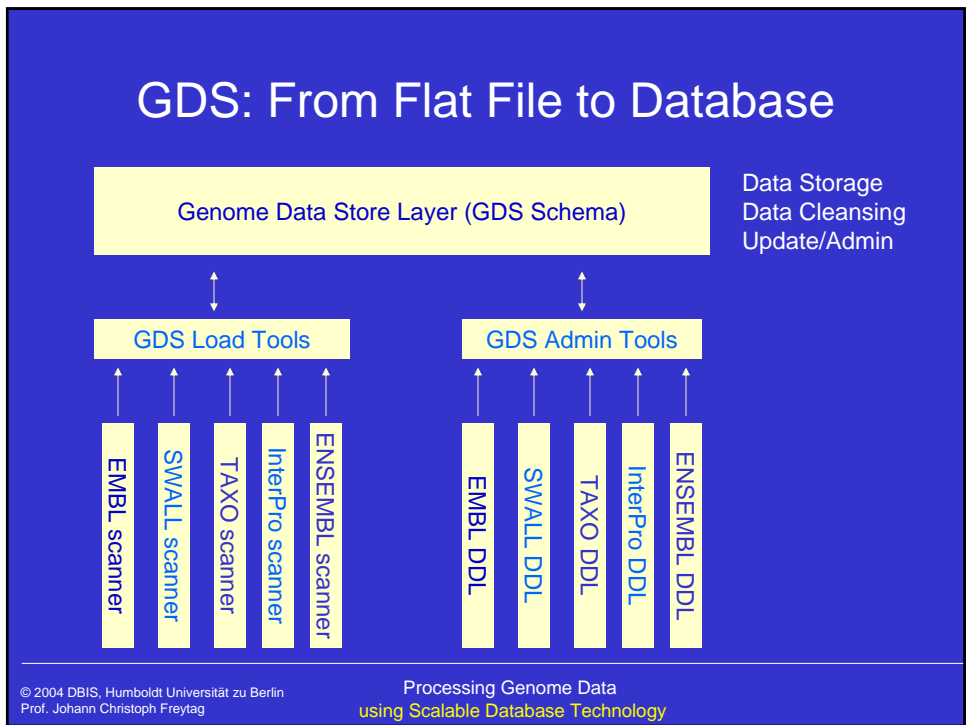
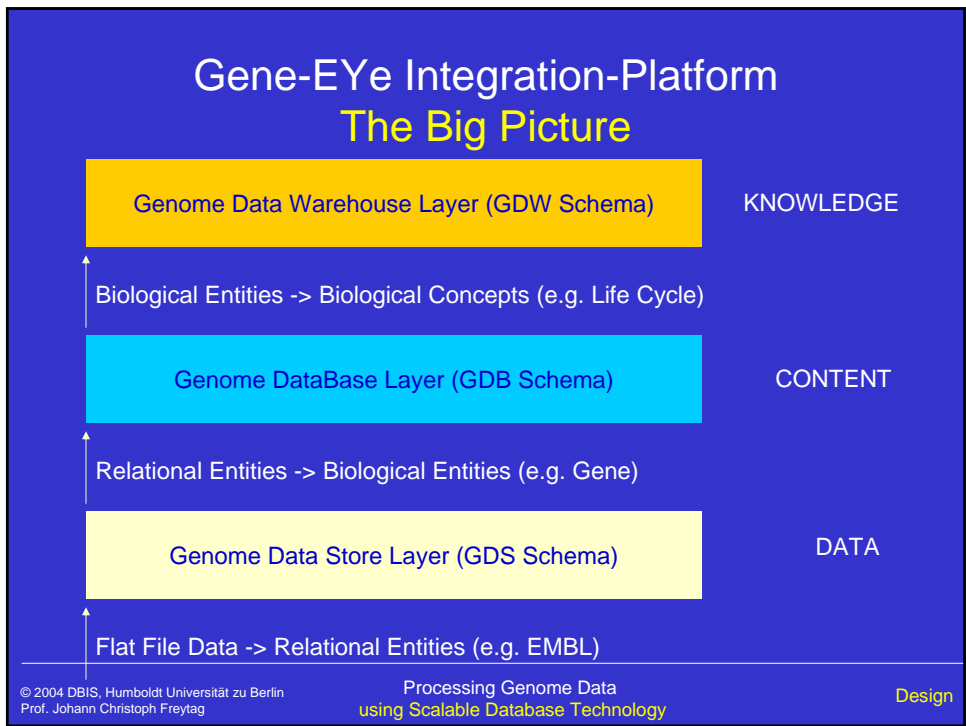
DBIS (our) Approach

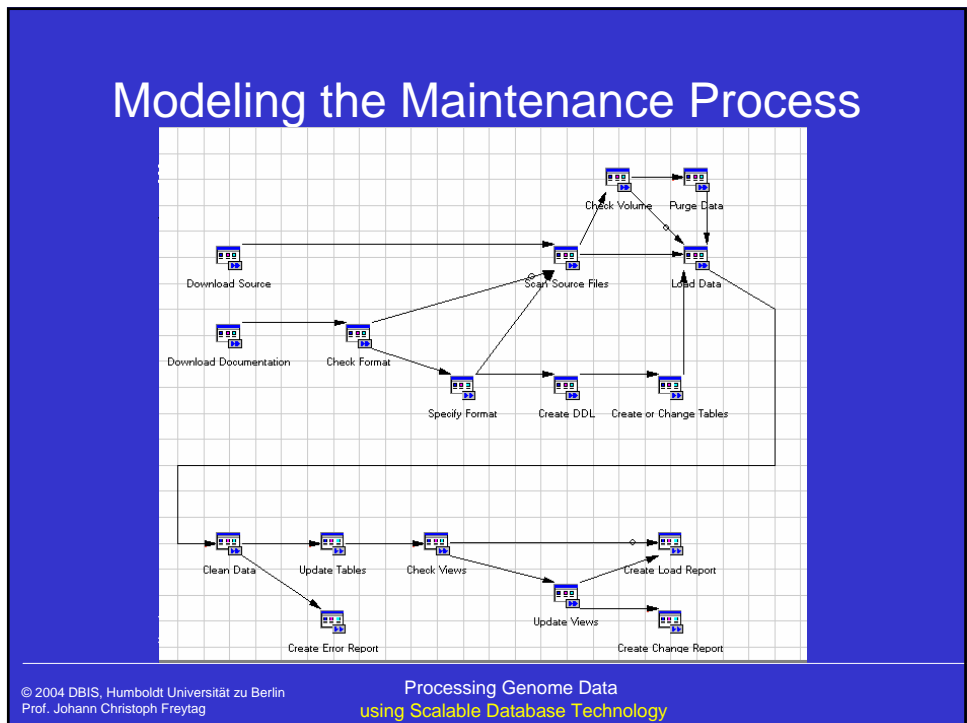
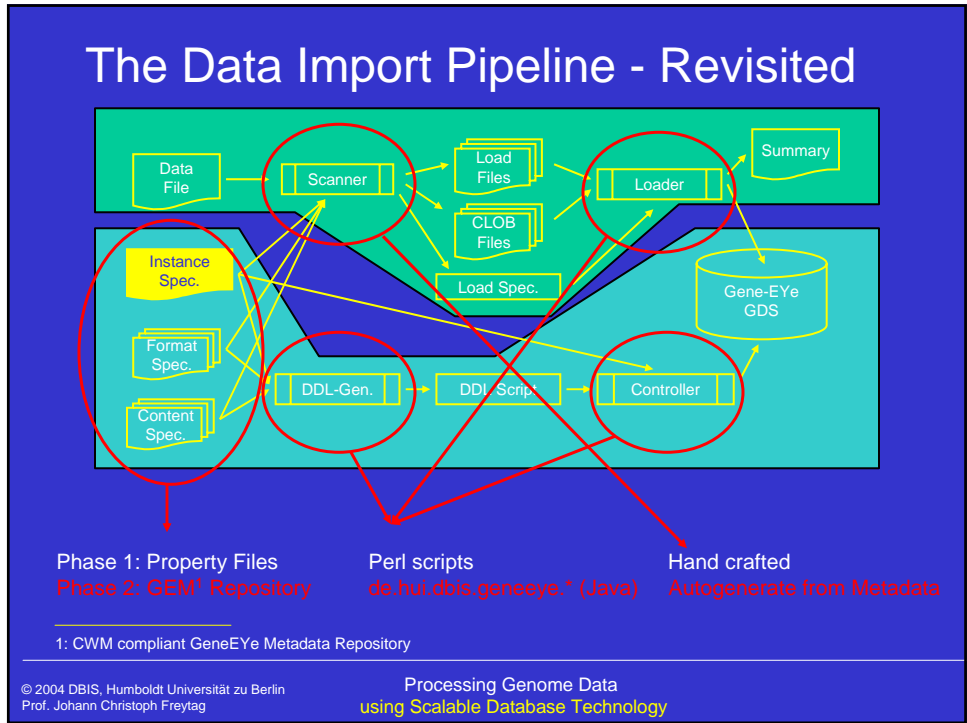


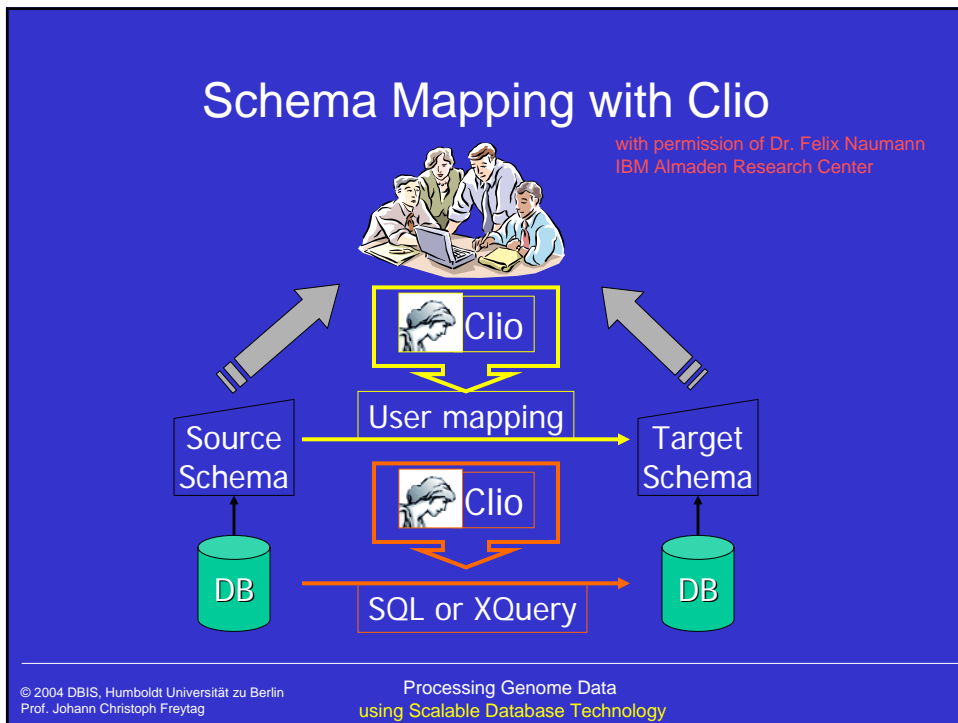
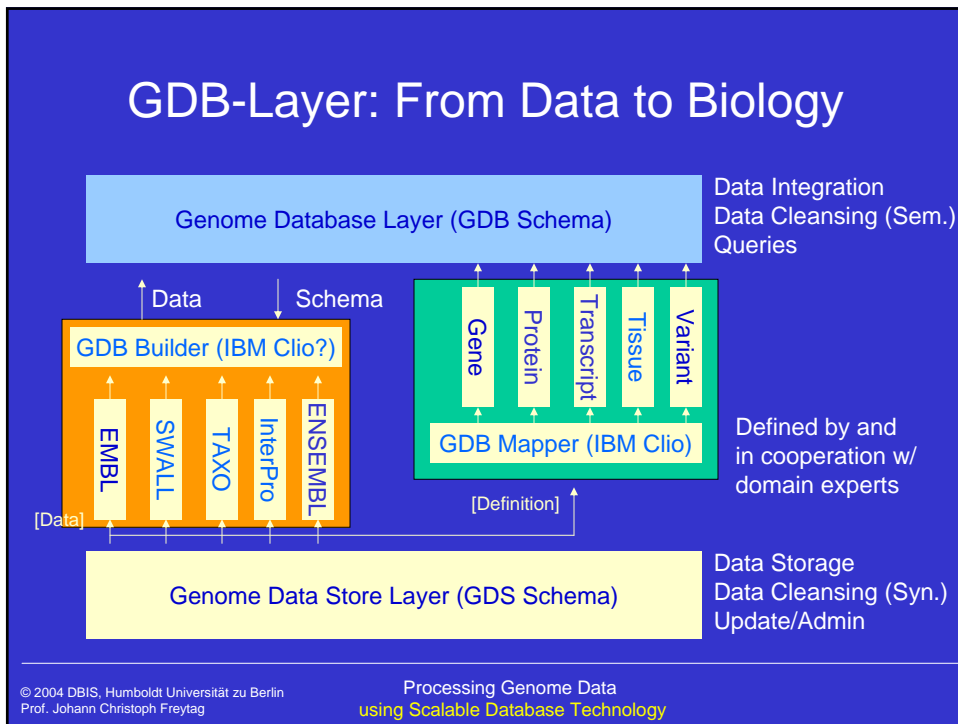
- (Biological) Motivation/Problems
- Using Scalable Database Technology
 - Gene-EYe Integration-Platform
 - Data Cleansing
 - BLAST-Integration into GDB
 - In-and-Out-the-Database: Using Workflow Concepts for “dry” Experiments
- Summary

Gene-EYe Integration-Platform Vision

- Provide mechanisms for
 - unified handling of different data sources
 - data source integration
 - change management
 - user defined data preparation
- Provide
 - relevant tools for sequence manipulation and retrieval
 - work flow support for operation and administration







Clio Features

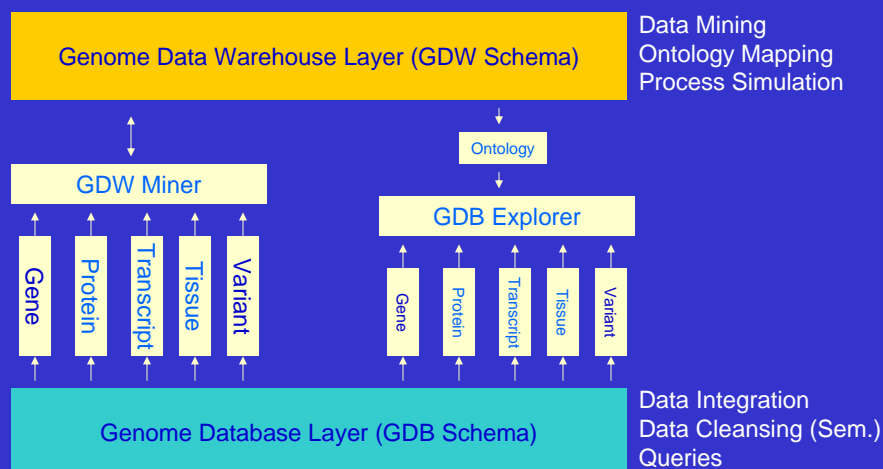
with permission of Dr. Felix Naumann
IBM Almaden Research Center

- Schema Viewer
 - Visual mapping between schema elements
- Attribute Matcher
 - Intelligent suggestions of likely mappings
- Data Viewer
 - Data examples for mapping queries
- Queries
 - SQL, XSLT, Xquery
 - Use and adhere to source and target schema constraints

© 2004 DBIS, Humboldt Universität zu Berlin
Prof. Johann Christoph Freytag

Processing Genome Data
using Scalable Database Technology

GDW: Providing Facts for Research





© 2004 DBIS, Humboldt Universität zu Berlin
Prof. Johann Christoph Freytag

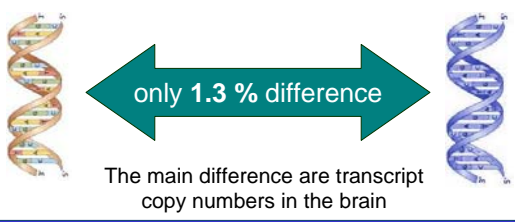
Processing Genome Data
using Scalable Database Technology

- (Biological) Motivation/Problems
- Using Scalable Database Technology
 - Gene-EYe Integration-Platform
 - Data Cleansing
 - BLAST-Integration into GDB
 - In-and-Out-the-Database: Using Workflow Concepts for “dry” Experiments
- Summary

Errors in Genome Data







only 1.3 % difference

The main difference are transcript copy numbers in the brain

A Sequence	
agatcc	0,23%
atagcc	–
agcagc	2,58%

Annotation	
←	5% - 30%

Sequence	
QAERYDEMVESMKKVD	
SVAYKNVIGARRASWY	
EDKLMIREYRQMVVER	

Annotation	
CT AS IN	
LING COM	
DIRECTL	
VED IN A	

© 2004 DBIS, Humboldt Universität zu Berlin
Prof. Johann Christoph Freytag

Processing Genome Data
using Scalable Database Technology

Reliability-based Merging (cont.)

- Domain expert identifies reliable parts for merging
- Definition of a set of views for integration

r ₁	ID	A ₁	A ₂	A ₃
	1	A	1	2000
	2	A	2	2000
	3	B	3	2000
	4	B	3	2000
	5	C	5	2002

r ₂	ID	A ₁	A ₂	A ₃
	1	B		
	2	B		
	3	B	3.1	2000
	4	B	3.4	2000
	5	D	5.6	2002

Current work:

- Which are the relevant *mismatch patterns*?
- How to assess their *relevance & importance*?

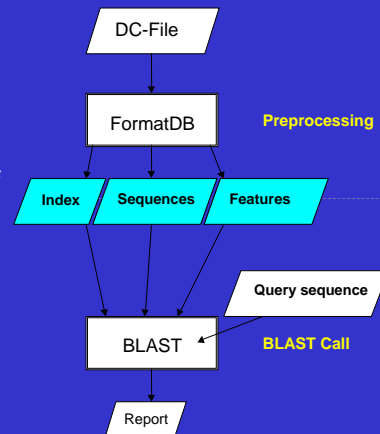
A ₂	A ₃
1.2	2000
2.3	2000
3.1	2000
3.4	2000
5.6	2002

e.g. MIN()

- (Biological) Motivation/Problems
- Using Scalable Database Technology
 - Gene-EYE Integration-Platform
 - **BLAST-Integration into GDB**
 - In-and-Out-the-Database: Using Workflow Concepts for “dry” Experiments
- Summary

BLAST: General Introduction

- Algorithm/Package: Similarity Search
- Developed by Altschul et al. (1990)
- Three Steps:
 1. Search for Word Pairs (Iseq, DSeq) of Length L on the Data Collection of Sequences above Threshold T
 2. Expansion of each Word Pair until the Value V of their Alignment is Δ away from the local maximum
 3. Output of complete alignment (*High-scoring Segment Pair, HSP*), if $\text{Value}(\text{Alignment}) > S$



Output: Powerset of Alignments

© 2004 DBIS, Humboldt Universität zu Berlin
Prof. Johann Christoph Freytag

Processing Genome Data
using Scalable Database Technology

BLAST UDF Implementation

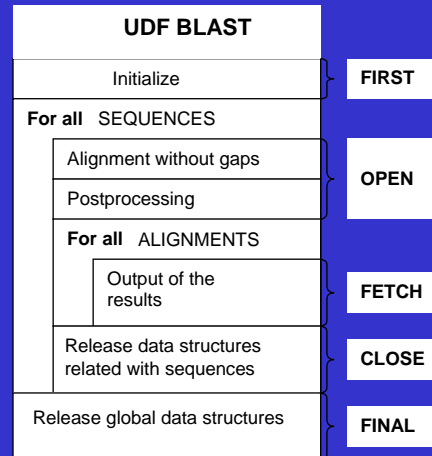
- Goal: Using BLAST in SQL-statements
 - How?
 - BLAST-UDF implemented as Table Function
 - Use in SQL Query
- ```
SELECT *
FROM TABLE(BLAST(<Parameter>, <Query Sequence>,
 <Comparison Sequence>))
```
- Each call returns a set of alignments over Sequences in the Database

© 2004 DBIS, Humboldt Universität zu Berlin  
Prof. Johann Christoph Freytag

Processing Genome Data  
using Scalable Database Technology

## Structure of UDB Table Function

- Implementation:
  - Mapping of program into calling structure for table functions
  - Communication between the different calls via *scratchpad*
  - scratchpad: Storage area which remains intact and unchanged between UDF calls
    - Storage of data structures for different steps
    - especially for output from postprocessing: `SeqAlign`



© 2004 DBIS, Humboldt Universität zu Berlin  
Prof. Johann Christoph Freytag

Processing Genome Data  
using Scalable Database Technology

- (Biological) Motivation/Problems
- Using Scalable Database Technology
  - Gene-EYE Integration-Platform
  - BLAST-Integration into GDB
  - **In-and-Out-the-Database: Using Workflow Concepts for “dry” Experiments**
- Summary

© 2004 DBIS, Humboldt Universität zu Berlin  
Prof. Johann Christoph Freytag

Processing Genome Data  
using Scalable Database Technology

Overview

## The Challenge: Exon Skipping

Gene

Protein

Gen

Intron Exon Intron Exon Intron Exon

**One Gen with 100 Exons  $\Rightarrow 2^{100} \sim 10^{30}$  Variations**

n Exons within one Gene  $\rightarrow$  linearly combined (splicing)  $\rightarrow$  Used as Pattern for Protein Generation

© 2004 DBIS, Humboldt Universität zu Berlin  
Prof. Johann Christoph Freytag

Processing Genome Data  
using Scalable Database Technology

## Challenge: Exon Skipping

Exon Exon Exon Exon

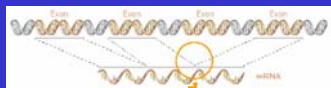
mRNA

Do **alternative fusion points** new funktional (i.e. biologigagl meaningful) „patterns“?

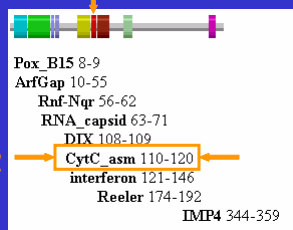
© 2004 DBIS, Humboldt Universität zu Berlin  
Prof. Johann Christoph Freytag

Processing Genome Data  
using Scalable Database Technology

## Functional Genomics: Gain of New Insight First Horizon: Simple Exon Skipping



110110000111111

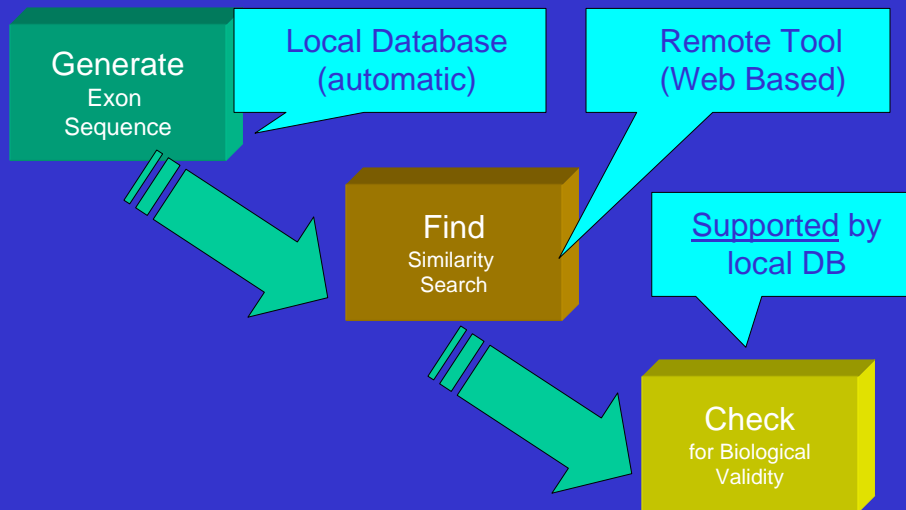


New Functionality!

© 2004 DBIS, Humboldt Universität zu Berlin  
Prof. Johann Christoph Freytag

Processing Genome Data  
using Scalable Database Technology

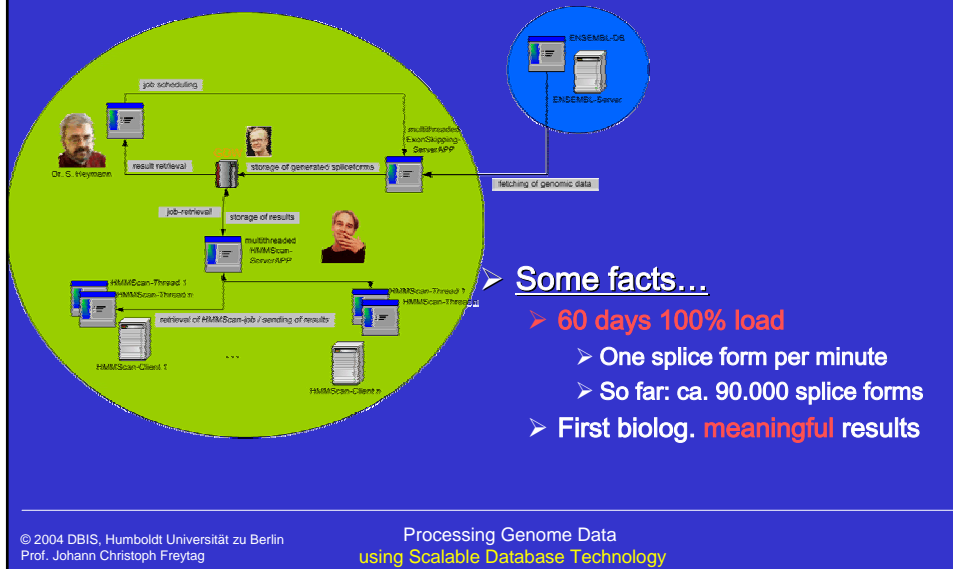
## Flow of Processing Steps



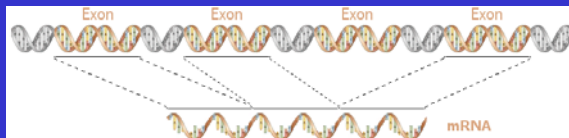
© 2004 DBIS, Humboldt Universität zu Berlin  
Prof. Johann Christoph Freytag

Processing Genome Data  
using Scalable Database Technology

# Implementation



# Cooperations



- Cooperation with
  - Univ. of Jena (Rolf Backofen)
  - Berlin Center of Bioinformatics (BCB)
    - Charite, FU, Max-Planck-Institut (M. Vingron)
  - Industry: IBM, small companies, ...

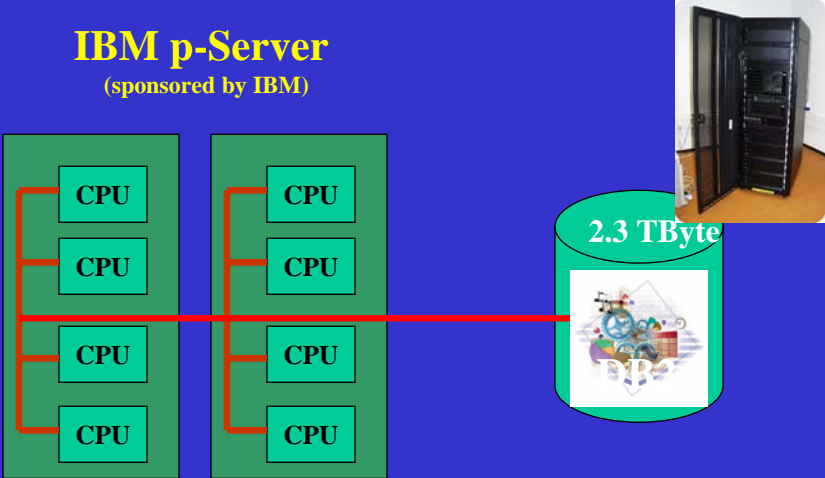


Patrick Chappatte, Switzerland



## Database Environment

**IBM p-Server**  
(sponsored by IBM)




2.3 TByte

© 2004 DBIS, Humboldt Universität zu Berlin  
Prof. Johann Christoph Freytag

Processing Genome Data  
using Scalable Database Technology

## Summary

- Lesson learnt
  - Highly Dynamic Environment
    - Data: changes frequently
    - User: changes frequently
  - Provide a framework for
    - Date integration
    - Data processing
    - Data changes
    - Data dependencies.....
    - Meta data management
- Future Work
  - Query processing
    - Include domain knowledge
  - Data cleansing
  - Set of UDFs for biological data processing
  - Visualization of Data



© 2004 DBIS, Humboldt Universität zu Berlin  
Prof. Johann Christoph Freytag

Processing Genome Data  
using Scalable Database Technology

Summary