# Enabling eScience

## Challenges of Supporting the Digital Scientist

# About Myself
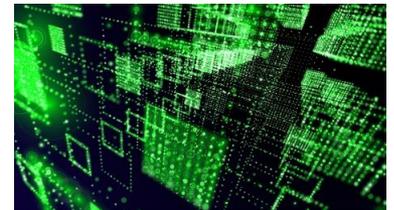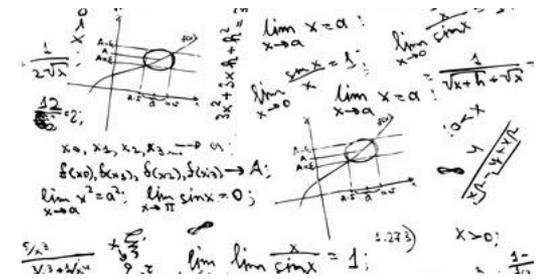
- PhD in theoretical physics – Bern
- Built large astrophysics science database – Sloan Digital Sky Survey, Johns Hopkins (3y)
- Headed data management tool development for LHC – CERN (5y)
- Built Swiss Tier2 at CSCS for LHC – CSCS (3y)
- Now leading SyBIT (5y)

# Goals of Talk

- eScience – characteristics
  - What is eScience
  - Who is an eScientist
- eScience – 2 examples
  - Astronomy / Astrophysics
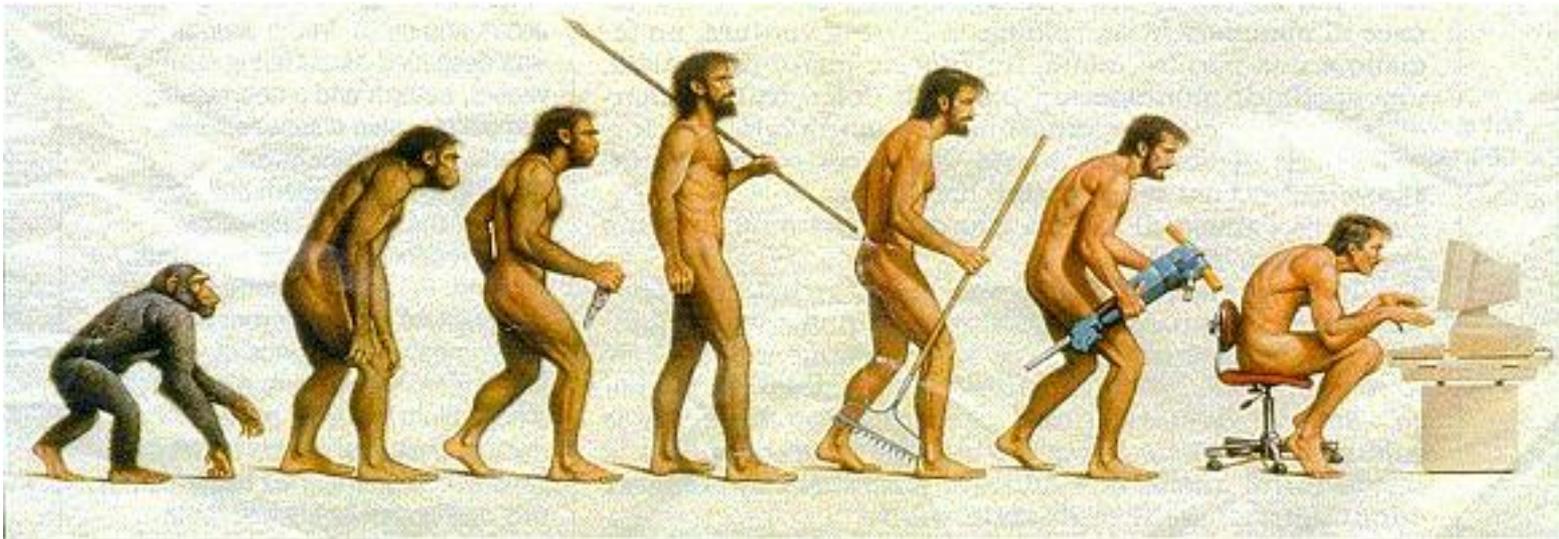  - Life Sciences
- Infrastructure for eScience
- Challenges

# Evolution of Science

- 1000 years ago:
  - **Empirical Observation**
- Few 100 years ago:
  - **Theory – Experiment Cycle**
- Since 50 years:
  - **Computer Simulation of Complex Phenomena**
- Now added:
  - **Data Exploration: eScience**

# Science Today

- Every measurement, experiment, simulation, etc produces **digital data**
- **eScience = Extract knowledge**

# eScience : Unify it all

- Theory, Experiment, Simulation
- Data captured by instruments
- Data generated by simulation
- Data Processed by Software Pipelines
- Previous Information/Knowledge from Digital Libraries

# eScience Methods

- Data Collection
  - Digital Instrumentation
- Simulation
  - HPC, distributed computing, etc
- Data Management
  - Collect data so that it is found again
  - Make sure the data is available for reuse
- Data Analysis
  - Statistical analysis, mining, visualization
- Data Curation
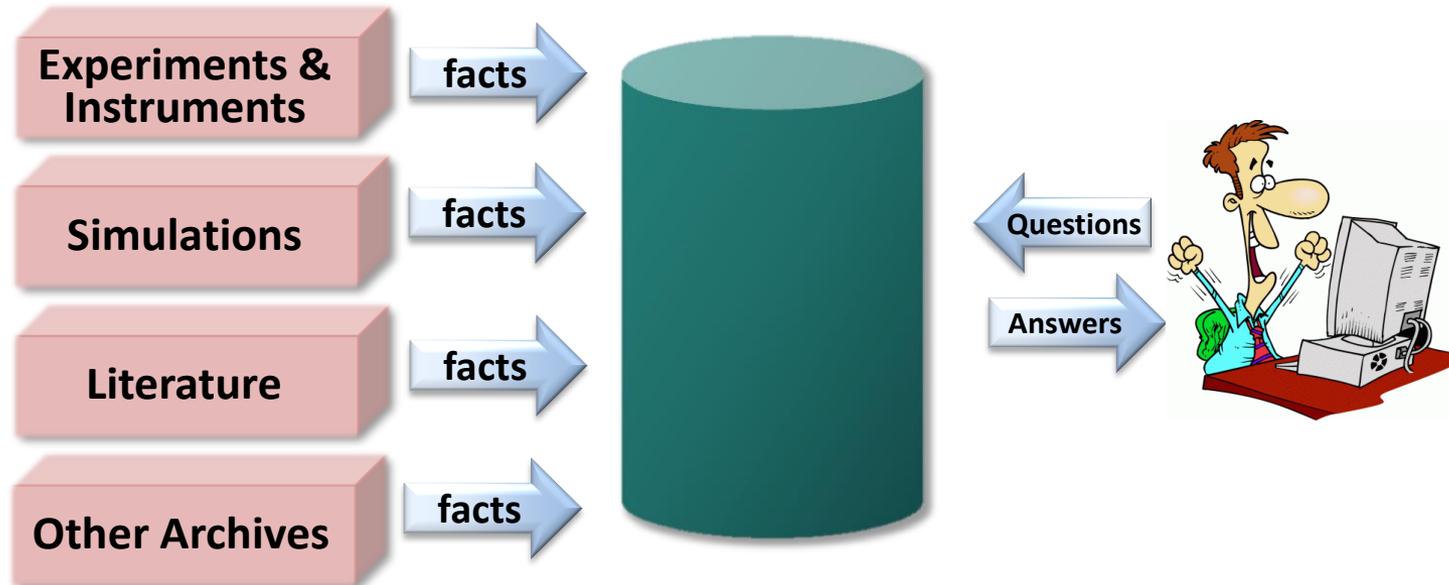  - Description, contextualization, annotation
  - Publication

# DOMAIN-informatics and computational DOMAIN

Each discipline is evolving its own computational research and informatics (tooling and engineering)

**How to codify and represent our knowledge is DOMAIN SPECIFIC**

# eScience People

**Data Engineer / Data Manager**
- Programmer / Database / NOSQL expert
- Managing very large data volumes
- Writing tools to collect data
- Making sure data is coherent and valid

**Algorithm Builder**
- Programmer of simulations, workflow engineering
- Models and algorithms
- HPC or Grid specialist

# eScience People

**× Data Analyst**

- × Build models based on data
- × Increasingly knowing the data
- × Ask right questions of the data

**× Data Steward**

- × Librarians, information specialists, ontologists
- × Make it discoverable, reusable, linked to other data
- × Well documented

# + **Large Collaborations**

✖ Increasing number of large / very large projects

✖ Dozens / hundreds of People

✖ Working together on the same project over years

# **Collaboration Examples**

- Astronomy : Large Sky Surveys
- High Energy Physics : CERN
- Systems Biology : SystemsX.ch
- Neurobiology : Human Brain Project

- Every discipline has examples
  - Geography, Earth Observation
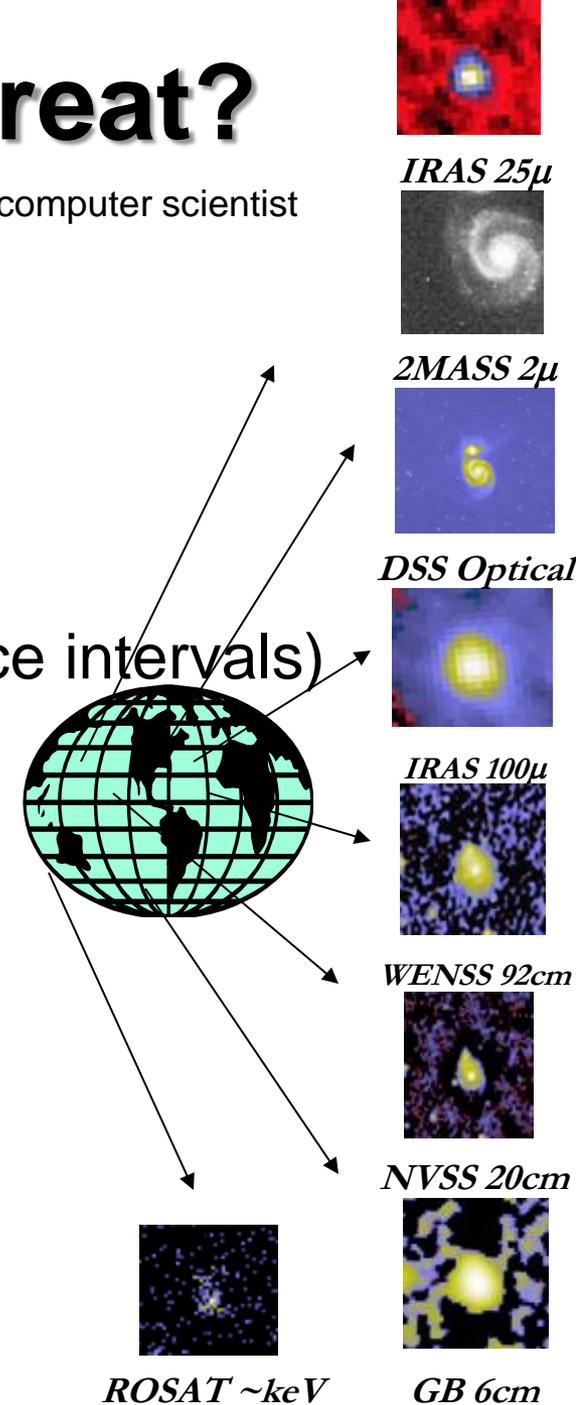  - Medicine
  - Human medicine
  - History
  - …

# Astronomy

- Used to be a single-person discipline
  - Data mining = the data is mine
  - No sharing of photographic plates

- But Large telescopes need a LOT of money
  - Large collaborations have emerged
  - Culture has completely changed

- Early engagement with computer science
  - Lots of method development
  - Lots of tooling
  - **Astronomers & Astrophysicists are good with computers**

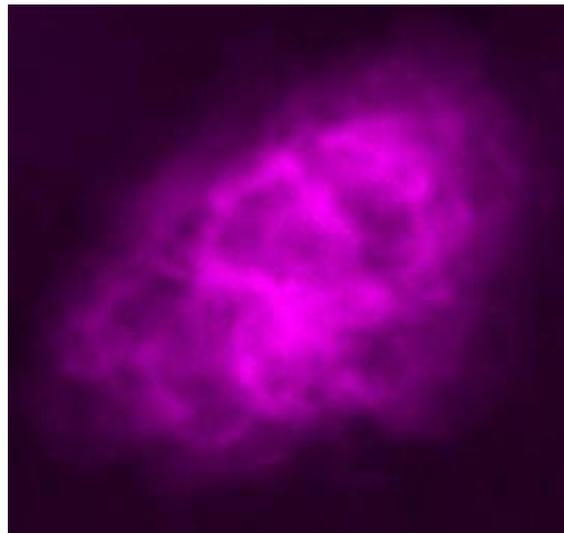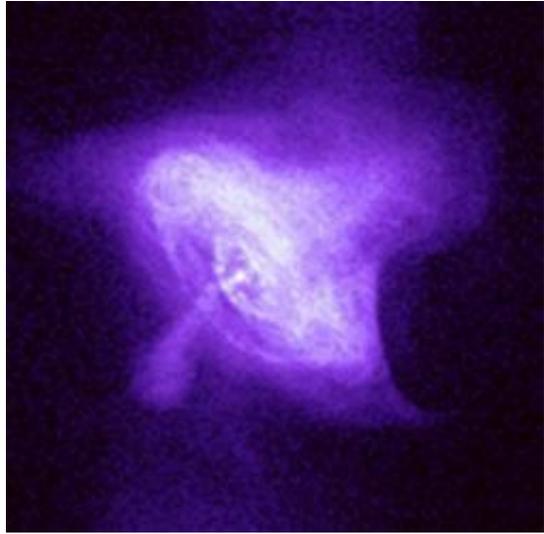# Why is Astronomy Data great?

For a computer scientist

- **It has no commercial value**
  - No privacy concerns
  - Can freely share results with others
  - Great for experimenting with algorithms

- **It is real and well documented**
  - **High-dimensional data** (with confidence intervals)
  - **Spatial** data
  - **Temporal** data

- Many **different instruments** from many **different places** and many **different times**

- **Federation is a goal**

- There is a lot of it (petabytes)

*IRAS 25μ*

*2MASS 2μ*

*DSS Optical*

*IRAS 100μ*

*WENSS 92cm*

*NVSS 20cm*

*ROSAT ~keV*

*GB 6cm*

Slide adapted from Jim Gray

# Time and Spectral Dimensions
# The Multiwavelength Crab Nebulae



Crab star 1053 AD

X-ray, optical, infrared, and radio views of the nearby Crab Nebula, which is now in a state of chaotic expansion after a supernova explosion first sighted in 1054 A.D. by Chinese Astronomers.

Slide courtesy of Robert Brunner @ CalTech.

# SkyServer.SDSS.org

- A modern archive
  - Access to Sloan Digital Sky Survey Spectroscopic and Optical surveys
  - Raw Pixel data lives in file servers
  - Catalog data (derived objects) lives in Database
  - Online query to any and all
- Also used for education
  - Hundreds of hours of online Astro
  - Implicitly teaches data analysis
- Interesting things
  - Spatial data search
  - Client query interface via Java Applet
  - Query from Emacs, Python, ….
  - Cloned by other surveys (a template design)
  - Web services are core of it.

Slide by Jim Gray

# World Wide Telescope Virtual Observatory

http://www.us-vo.org/        http://www.ivoa.net/

- Premise:  Most data is (or could be online)

- So, the Internet is the world's best telescope:
  - It has data on every part of the sky
  - In every measured spectral band: optical, x-ray, radio..
  - As deep as the best instruments (2 years ago).
  - It is up when you are up.
    The "seeing" is always great
    (no  working at night, no clouds no moons no..).
  - It's a smart telescope:
    links objects and data to literature on them..

Slide by Jim Gray

# Characteristics of Astronomy

- Data is well known
- Schema is defined and optimized
- Indices are given
- Community has learned SQL!
  - It is faster to query the archive than to write a computer program

# Astronomy is an eScience Discipline today

- Catalogs are omnipresent
- New projects are planned with data curation and integration in mind
- Communities are well organized
- Standards have been worked out and are in use
- People are being educated in these tools and methods, like SQL

# How to get there?

1.  Endorse domain as Data Intensive Science.
2.  Plan computing infrastructure to be inherently distributed: "scale-out" architecture.
3.  Bring computations to the data, rather than data to the computations wherever possible.
4.  Design system with concrete questions to the data. Do not go overboard with what-if scenarios.
5.  Plan for change, use an iterative, agile approach to everything.
6.  Automation of everything possible

# **New Challenges**

1. Needle in the Haystack, how to distinguish artifacts from real new objects / lost in automation

2. New instruments producing even more data, like LSST: full sky survey every 3 days

   - Monitor changing objects
   - Store light-curves for objects
   - Will start operations 2021
   - Hundreds of Petabytes of data

# Address Challenge 1:



- Citizen Science site to classify Galaxies by looking at them [www.galaxyzoo.org](www.galaxyzoo.org)
- Data from SDSS and Hubble
- Thousands of users
- Several interesting objects have been found, leading to new science!

# Addressing Challenge 2: LSST Database design

- Partitioning: distributed database architecture

- Shared-nothing data distribution

- Sharding of data and sharding of queries

# **LSST Database design**

https://dev.lsstcorp.org/trac/wiki/db

# **LSST Database design**

✕ Clustering of nearby objects using hierarchical triangular mesh

✕ Shared scanning – full table scans rather than indexed scans if more than just a few % of the data are being returned

  ✕ Takes care of special cases where every object needs to be touched

  ✕ Results to be stored in 'myDB'

# Systems Biology

# Systems Biology

As understood by a physicist

- OK, so we have the Gene sequence (Genome) of many organisms, but what does it say? How does it all work?

- Cell Biology and Molecular Biology are producing more and more high-resolution and high-quality data to answer this question

- Bottom-up approach, understanding the cell from first principles is very difficult.

- Systems Biology Approach: Understand the available data top-down. Study the complex interaction of many levels of biological information to understand how they work together.

# The Systems Biology Approach

**Systems Biology**:
- Biology
- Physics
- Chemistry
- Computer Science
- Mathematics
- Engineering
- Medicine

Biological System;
**Qualitative**,
Wet Lab Biology



**Quantitative**
Data sets of
Complete Systems;
Bio-Engineering

Network **Theory**;
Modeling,
Simulation

# **The Systems Biology Approach**

**Systems B...** Biological System;
* Biology **...tive**,
* Physics ...logy
* Chem...
* Compu...
* Mathematics
* Engineering
* Medicine

**This is new for Life Sciences!**

**Quantitative**
Data sets of
Complete Systems;
Bio-Engineering

**...heory**;
...g,
Sim...ion

# Characteristics of Systems Biology

- ✕ Data is NOT well known
- ✕ Schema is NOT defined and optimized
- ✕ A lot of data is being searched for connections to 'see' anything by chance
  - ✕ Clustering, PCA, networks
- ✕ What is the question again? Lots of exploration of the unknown.
- ✕ Community is less computer savvy
  - ✕ Programming is not a skill biologists learn
  - ✕ Disconnect between modelers and experimentalists

# SystemsX.ch

**Largest Swiss national research effort to date**



SCHWEIZERISCHE EIDGENOSSENSCHAFT
CONFÉDÉRATION SUISSE
CONFEDERAZIONE SVIZZERA
CONFEDERAZIUN SVIZRA

FNSNF
FONDS NATIONAL SUISSE
SCHWEIZERISCHER NATIONALFONDS
FONDO NAZIONALE SVIZZERO
SWISS NATIONAL SCIENCE FOUNDATION

ETH-RAT

SUK·CUS
Schweizerische Universitätskonferenz

**SystemsX.ch**
The Swiss Initiative in Systems Biology

# Interdisciplinary Research

Chemistry 5 · Physics 14 · Engineering 16 · Medicine 13 · Neurobio 2 · Semi-Bio 10 · Economics 3 · Biology 51 · Bioinformatics 6 · Comp Science 11 · Math 6

137 research groups

# SystemsX.ch Projects

- Research Technology Development projects: large collaborations
  - 8 large projects have concluded
  - 17 ongoing projects
  - 10 more to start in 2014
- Industry collaboration projects – 10+
- Pilot 'high risk' projects – 10+
- Student, PhD projects – dozens

# Challenge : Scale Up

✕ High Throughput Instruments
  - ✕ Much larger **data volumes**
  - ✕ Increased **data complexity**

✕ Large Collaborations
  - ✕ More people
  - ✕ More experiments and measurements

# SystemsX.ch is Data Intensive

Researchers

Data Intensive Life Science Research

Platforms

IT Resource providers

# SyBIT: Build the Distributed Infrastructure and Tools

SyBIT
SystemsX.ch
Biology IT

Researchers

Data Intensive Systems Biology Research

SyBIT Life Science Support

Platforms

IT Resource providers

# Fill the Gap

Researchers

**Support** | **Collaborate**

SyBIT

**provide**

**Enable Strengthen**

**Define requirements for team up with**

Platforms

IT Resource providers

Tools
Bioinformatics
Services
Data Management
Standards, Integration
Software engineering
Legacy support
Documentation
Education
Collaboration support
Setup, Configuration
Versioning
Commercial software

37

# Interactions in detail

# Also here:
# Data at the Center of Work

**Data**

- ✕ Almost everything is data driven

- ✕ Many formats and access patterns

# Data Production

| Mass Spec |
|-----------|
| Microarray |
| Microscope |
| HCS/HTS |
| Simulation |
| ... |

**O(TB)/Day** → **Data**

- ✕ Many different kinds of instruments
- ✕ Different data types
- ✕ New instruments produce much more data
- ✕ Data volume increasing exponentially

# Data Validation and Filtering

Mass Spec

Microarray

Microscope

HCS/HTS

Simulation

...

O(TB)

**Data**

✗ Validation, checks
✗ Conversion into standard formats
✗ Compression

✗ This can be very compute intensive
✗ This can produce a lot of new data
✗ Needs clusters to do it in a timely fashion

# Data Tracking and Metadata

Mass Spec

Microarray

Microscope

HCS/HTS

Simulation

...

**Data**

**O(GB)**

- Provenance metadata
- File catalogs
- Metadata on initial filtering

# Data Exploration

Mass Spec

Microarray

Microscope

HCS/HTS

Simulation

...

Data

- Interactive exploration of data
- Small-scale analysis
- Planning of large-scale data analysis

# Data Analysis and Modeling

- Large-scale analysis
  - On as much CPU power as possible
- Production of more data
  - Secondary datasets
  - Simulation, modeling
- Additional databases
  - Metadata
  - Result data

# Publication and Archiving



- ✕ Publication into public databases
- ✕ Curation
- ✕ Archiving for long-term storage

Mass Spec

Microarray

Microscope

HCS/HTS

Simulation

...

Data

O(100TB)/yr

O(TB)/yr

# Data Lifecycle

**All SystemsX.ch projects:**

• **O(PB)/yr kept data**

• **O(TB)/yr published data**

• **O(10^7) CPU Hours**

• **Several different DBs and formats**

✕ Some steps might be iterative

✕ Users are not interested in technology, it simply has to work

✕ Implementing ´Data´ such that all needs are met is a challenge

# **Repetitive problem**

EPFL

ETH

RTDs

UniXY

UZH

UniBas

✕ Same lifecycle everywhere

✕ Local Policies

✕ Local Infrastructures

✕ Local Services

✕ Local Access control

✕ Nontrivial coordination effort for RTDs to share data and services

# Mapping to Infrastructure



**SyBIT**
SystemsX.ch
Biology IT

Data Creation:
Instrument

Preprocessing:
Validation, zip, ..

Automated
Processing

Postprocessing:
Visualization, etc

Publication,
Archiving

Instruments

1. Instr. Store

2. Group Servers

5. Project Store

4. Cluster

3. User Home

User Desktop

6. Archive

7. Result Store

7. Result Servers

8. Metadata Catalog

# Mapping to Infrastructure

Data Creation:
Instrument

Preprocessing:
Validation, zip, ..

Automated
Processing

Postprocessing:
Visualization, etc

Publication,
Archiving

Instruments

1. Instr. Store
/local

2. Group Servers

5. Project Store
/project

4. Cluster
/scratch

3. User Home
/home

User Desktop

6. Archive
/store

7. Result Store
/www-data

7. Result Servers

8. Metadata Catalog

# Central vs. Local Resources

Data Analysis and
Modeling

Project Data
and Archive

Domain Resources
Collaboration Services
Result Data Servers

Central Cluster

Cluster FS

Hierarchical
Storage
And
Archive

Central
Servers
and VMs

Central Infrastructure

**NETWORK**

**Instruments**

Output
Storage

Group
Servers

Group laptops and desktops

Group Level Infrastructure

Data Production

Data Cache

Q/A and
Visualization

# So where are the problems?

✖ Well.. Where do we start

# **What is the data again?**

- Data is not very large but there are many **different kinds**
- People don´t understand it yet very well
  - How to evolve, schemas, versioning
  - How to be efficient about navigating in data – people still use excel sheets for that
  - CONTROLLED VOCABULARIES help
- **Big Data in life science = More, faster, more diverse than you can handle**

# Little motivation to annotate

- Just publish what is necessary and prescribed by the journals

- No recognition yet of producing ´good´ datasets

- Data quality as such is not yet very high, not much reuse (excel sheets limiting factor)

# Many formats

- **Not many standards exist, there are too many formats**
  - Instrument vendors often introduce new formats or conventions
- **No trust in central databanks**
- **People think their data is special and may have some value (in $)**
- **Not invented here effect is huge**
  - Everyone builds their own submission sysem, LIMS, etc

# **Data Loss**

- Student / postdoc leaves, nobody understands data anymore

- Who takes responsibility for data archives?

- No funding to build good archives

- Projects end, data is kept by chance only

# SyBIT Divide and Conquer Approach



**✖ Provide individual solutions for all projects**

  ✖ Allow for several solutions initially

  ✖ Consolidate over time and with new projects

# Approach: Divide and Conquer

## ✕ **Define Data Policies and User Access**

- ✕ Do not allow users to modify and move around large amounts of data themselves

- ✕ No Filesystem Access for large data – only through data catalogs

- ✕ Automation of data movement and storage

**Local – Central infrastrucuture splitup**

**Several metadata management systems are in place**

- **openBIS**
- **B-Fabric**
- **Dedicated databases**

# openBIS metadata store

**open**BIS
Biology Information System

Metadata

Knowledge Resources

Processed Data

Raw Data

Store

Organize

Annotate

Search

Share

Export

# Online Interface ...

# Automatic Data Management

# Approach: Divide and Conquer

✕ **Define Pipelines for each instrument type**

✕ Platform providers already exist and have a lot of experience we can build on

**A lot of work was invested into pipelines and workflows: we are in a reasonably good shape**

- **Proteomics**
- **Genomics**
- **Imaging**

**We are improving these continuously**

# Image Analysis Workflow

# iPortal Proteomics Portal

# iPortal Proteomics Portal

# Approach: Divide and Conquer

✕ **Define common storage types with storage providers**

  ✕ What kinds of storage services do we define?

  ✕ Which kind of service is offered where?

> **The process has started. Storage infrastructure is still an issue at most sites.**
>
> - **Local NAS, Central NAS solutions**
> - **IBM Remote NAS project**

# Approach: Divide and Conquer

**�֎ Define common data retention semantics**

   ✗ What kind of data is to be kept on what type of storage for how long

   ✗ How are the costs covered

**UNSOLVED PROBLEM TODAY**

**This is a policy issue. Few are willing to make and ENFORCE policy decisions.**

# Approach: Divide and Conquer

**✕ Define interfaces for data sharing**

- ✕ Keep responsibility for the data where it was produced (ie. the people who know what it is)

- ✕ Extract data into warehouses for each specific problem

**This is what we do by default.**

**Each project has their dedicated metadata catalog interface. Most project chose openBIS. We are working on interfaces between systems, most already exist.**

# Approach: Divide and Conquer

**✖ Define final data repositories for public access**

  - ✖ Some public repositories already exist also abroad, agree which one to publish to for each project

  - ✖ Where no such repository exists, define and build one for SystemsX.ch

# Approach: Divide and Conquer

**Define final data repositories for public**

**MOSTLY UNSOLVED.**

**The public published data hosters (EBI, EMBL, NCBI) also struggle with data volumes.**

**Sending data is difficult, a lot of data curation is necessary – very time consuming process.**

**For others we set up dedicated openBIS instances.**

**Projects cannot be responsible for long-term archiving!**

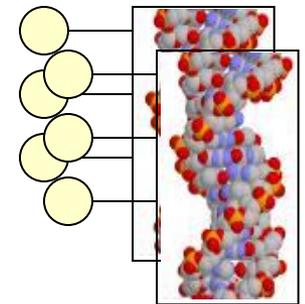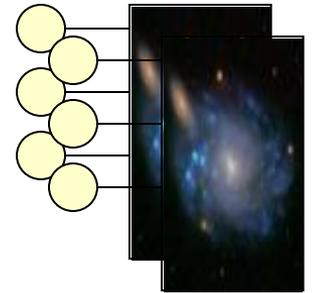# Long-term Data Archives : problem not only for LifeScience

✕ You have collected some data
   and want to publish science based on it.

✕ How do you publish the data
   so that others can read it and
   reproduce your results
   in 100 years?

  ✕ Document collection process?

  ✕ How document data processing
     (scrubbing & reducing the data)?

  ✕ Where do you put it?

# Generic Objectifying of Knowledge

✗ This requires agreement about

  ✗ **Units**: cgs

  ✗ **Measurements**:  who/what/when/where/how

  ✗ **CONCEPTS:**

    ✗ What's a planet, star, galaxy,…?

    ✗ What's a gene, protein, pathway…?

✗ **Need to objectify science:**

  ✗ what are the objects?

  ✗ what are the attributes?

  ✗ What are the methods (in the object sense)?

# Generic Objectifying of Knowledge

**SyBIT**
SystemsX.ch
Biology IT

✖ This requires agreement about

✖ **Un**

✖ **Me**

✖ **CO**

✖

✖

✖ **Nee**

✖ w

✖ what are the attributes?

✖ What are the methods (in the object sense)?
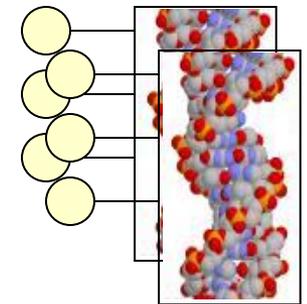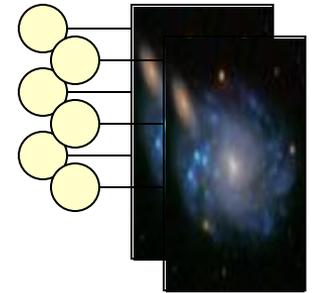
Agreements do not come easy. Need

Ontologies

Controlled Vocabularies

Schemas

**Domain Experts are still working it out!**

Slide adapted from Jim Gray

# Working Example: Entrez-GenBank
# http://www.ncbi.nlm.nih.gov/

- Sequence data deposited with Genbank
- Literature references Genbank ID
- BLAST searches Genbank
- Entrez integrates and searches
  - PubMedCentral
  - PubChem
  - Genbank
  - Proteins, SNP,
  - Structure,..
  - Taxonomy…
  - Many more

# Working Example: Entrez-GenBank
## http://www.ncbi.nlm.nih.gov/

- Sequence data deposited with Genbank
- Lite
- BLA
- Entr
  - Pu
  - Pu
  - Ge
  - Pr
  - Structure,..
  - Taxonomy…
  - Many more

**NCBI Funding is not secured**

Entrez Genomes

Complete Genomes

Genome Centers

3-D Structure

MMDB

Nucleotide sequences

Protein sequences

SyBIT
SystemsX.ch
Biology IT

# Data Sharing/Publishing Challenges

- How to keep programs that access and work with data? Cloud App Store for science? Who maintains those?

- What is the business model (reward/career benefit)?

- Journals starting to adapt, see Scientific Data http://www.nature.com/scientificdata

- But, many kinds of data are still orphaned: where to send Microscopy images?

# eScientists

## Researchers

- Domain specific expertise
- Measured by publications and impact
- Develops
  - New algorithms
  - New models
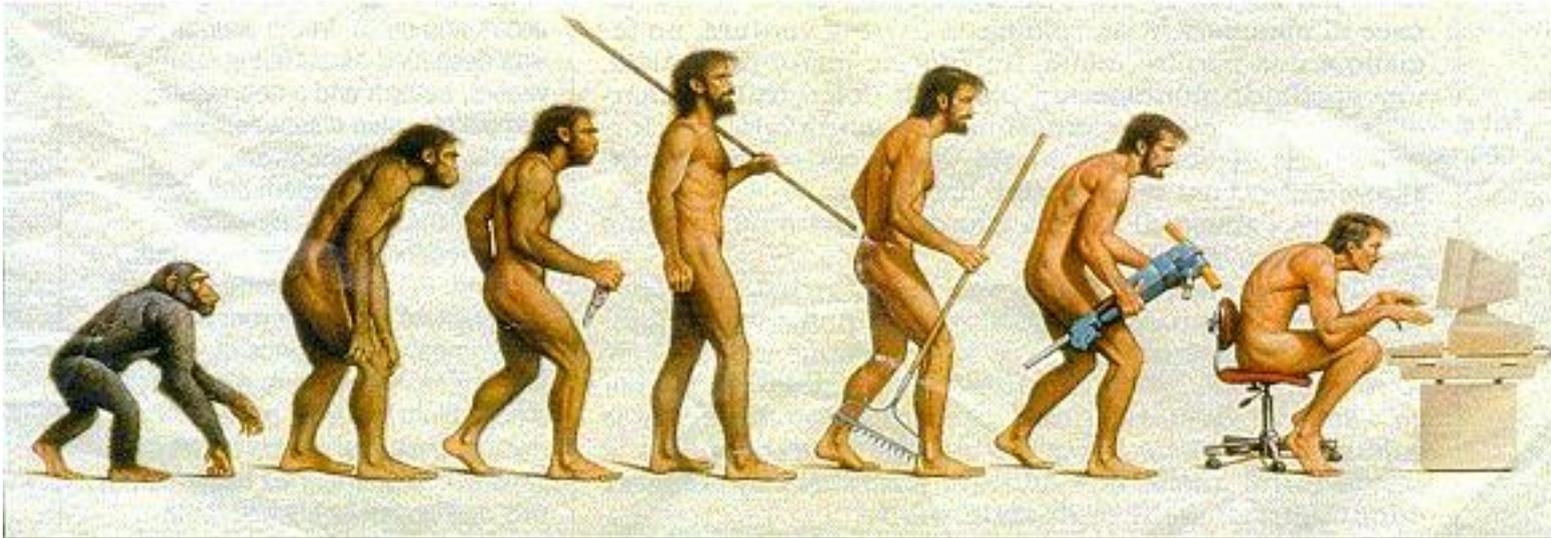- Find new connections
- Train new scientists

## Supporters / Stewards

- Domain specific expertise
- Measured by number of people supported
- Maintains existing tools and services, operations
- Trains scientists in usage of tools
- Curation of results, documentation

# Summary

- We need to manage the increased complexity (data and people)
- Need for specialized domain expert eScience supporters
- Scale-out distributed infrastructures enable us to grow: watch the clouds!
- Each scientific domain will address their challenges in their own way, no one size fits all.

# Coming Up?

# **Thanks to**

- Jim Gray and Tony Hey, Microsoft Research
- Alex Szalay, Johns Hopkins University
- Bernd Rinn, ETH Zurich

- [skyserver.sdss.org](skyserver.sdss.org)
- [www.lsst.org](www.lsst.org)
- [www.sybit.net](www.sybit.net)
- [http://www.nature.com/scientificdata/about/](http://www.nature.com/scientificdata/about/)

SyBIT
SystemsX.ch
Biology IT