# MOLGENIS
# bioinformatics toolkit
# & XGAP
# eXtensible Genotype And Phenotype platform

**GMOD meeting Europe**
**Cambridge, Sept 13, 2010**

Morris A. Swertz, K Joeri van der Velde, Alexandros Kanterakis, Juha Muilu, Tomasz Adamusiak, Martijn Dijkstra, Gudmundur A. Thorisson, George Byelas, Danny Arends, Members of EU-GEN2PHEN, NL-NBIC, EU-CASIMIR, BBMRI-NL, EU-PANACEA, Anthony J. Brookes, Ritsert C. Jansen and Helen Parkinson

# Outline

- MOLGENIS?
  - Flexible bioinformatics application toolkit
  - Demo: Model -> Generate -> Use

- XGAP?
  - eXtensible Genotype And Phenotype model
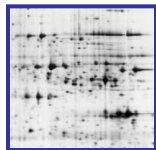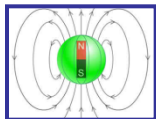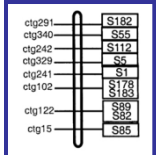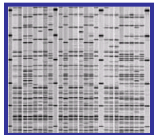  - MOLGENIS  generated xQTL & GWAS software
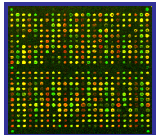
- Link to GMOD?
  - Chado? DAS? BioMART? Intermine? Gbrowse?

# MOLGENIS
*Flexible bioinformatics application toolkit for data management and interfacing*
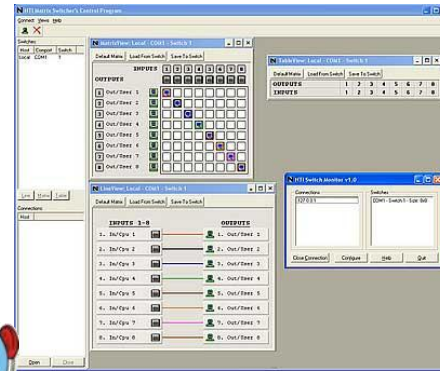
*etc.*

*etc.*

# Challenge



**DB**

**Logic**

```
static void main(String[] args) throws Exceptio
String path = args[0];
final String expr = args[1];

List l = new ArrayList();
findFile(new File(path), new P() {
        public boolean accept(String t) {
            return t.matches(expr) || isZi
        }
    });

List r = new ArrayList();
for (Iterator it = l.iterator(); it.hasNext(););
        File f = (File) it.next();
        String fn = f + "";
        if (fn.matches(expr)) r.add(fn);
        if (isZip(f.getName())) {
            findZip(fn, new FileInputStrea
                public boolean accept(
                    return n.match
```

**GUI**

*biologist*

**Exchange services**

**APIs**

$f(x)$ → SOAP → Service Provider

Service Requester

*bioinformatician*

*Etc*

# Challenge multiplied by project

# Needed alternative method

nature REVIEWS GENETICS

## Beyond standardization: dynamic software infrastructures for systems biology

*Morris A. Swertz and Ritsert C. Jansen*

Abstract | Progress in systems biology is seriously hindered by slow production of suitable software infrastructures. Biologists need infrastructure that easily connects to work that is done in other laboratories, for which standardization is helpful. However, the infrastructure must also accommodate the specifics of their biological system, but appropriate mechanisms to support variation are currently lacking. We argue that a minimal computer language, and a software tool called a generator, can be used to quickly produce customized software infrastructures that 'systems biologists really want to have'.

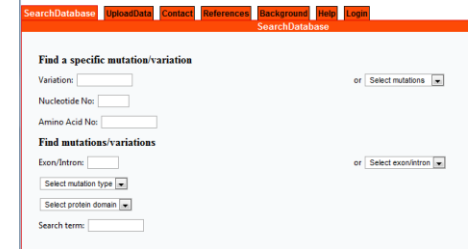# MOLGENIS



**Model in DSL**

NextGenSeq

Mutation database

Model organisms

**Generator**

**Use generated software**

Solexa Sequencer LIMS

*database of COL7A1 mutations*

*Animal Observatory*

*repeat often*

# MOLGENIS: Reuse in light of large variation



**Model in DSL**

NextGenSeq

Mutation database

Model organisms

**Generator**

NEW

**Use generated software**

Solexa Sequencer LIMS

*database of COL7A1 mutations*

*Animal Observator*

*repeat often*

http://www.molgenis.org
Swertz & Jansen (2007) *Nature Reviews Genetics* 8, 235-243
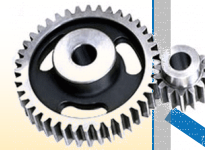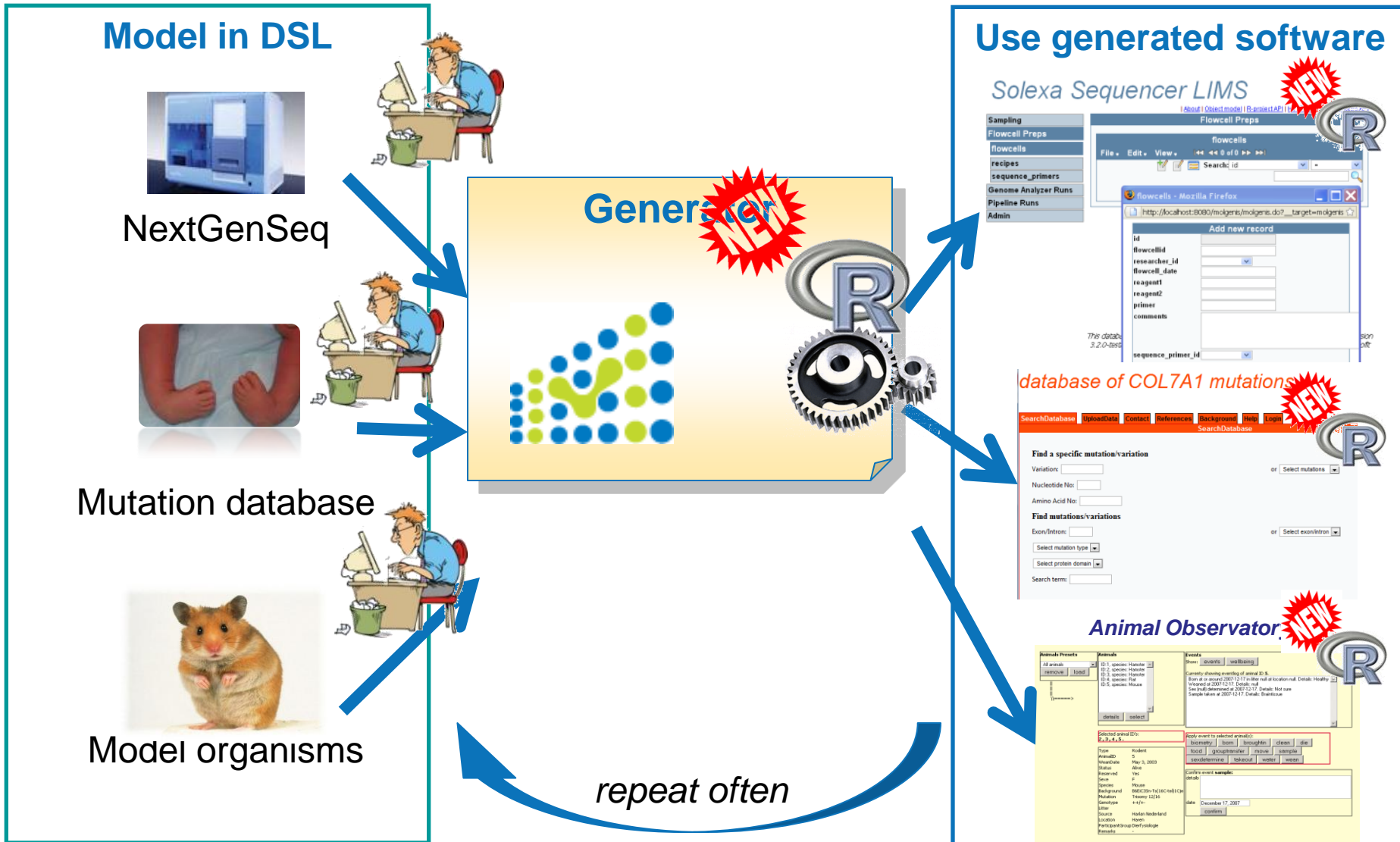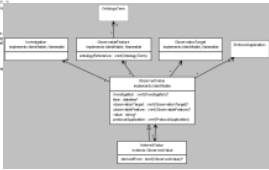Swertz et al (2004) *Bioinformatics* 20(13), 2075-83

# Example output

**❶ UML documentation of your model**

*Investigation*
*implements Identifiable, Nameable*

*The Investigation class defines well-contained units of study, each having a unique name and a group of actions (protocol applications) and/or results (or ObservedValues). For instance, Framingen study. Maps to ISA/MAGE Investigation and MAGE-TAB experiment. Discussion: should be about MAGE-TAB IDF type of minimal information about an investigation?*

*ObservableFeature*
*implements Identifiable, Nameable*

*The ObservableFeature class defines anything that can be observed (there may be many alternative protocols to measure them). For instance systolic blood pressure, Diastolic blood pressure. Treatment for hypertension. These values are unique within a data set. Preferably each ObservableFeature should be named according to a well-defined ontology. Risk class maps to ISA/MAGE biosource and PodD9. ObservableFeature. Multi-value features can be grouped by protocol. For instance, blood pressure consists of observations for features systolic and diastolic blood pressure.*

*ObservedValue*
*implements Identifiable*

*The ObservedValue class defines the individual value observed. To attach to the trial source...*

**❷ Edit & trace your data**

| Report |
| **Investigations** |
| Observable Features |
| Observation Targets |
| Protocols |
| Ontology Terms |

**❸ Import/export to Excel**

File ▾   Edit ▾   View ▾

⬇ Download visible
⬇ Download selected
⬇ D...
◆ Ad...

protocol_observableFeature...
Text Document 15,2 KB
protocol.txt
Text Document 1,51 KB
panel_individuals.txt
Text Document 106 KB
panel.txt
Text Document 150 bytes
ontologyterm.txt
Text Document 1,25 KB
ontologysource.txt
Text Document 48 bytes
observedvalue.txt
Text Document 4,20 MB
observablefeature.txt
Text Document 26,2 KB
investigation.txt
Text Document 606 bytes
individual.txt
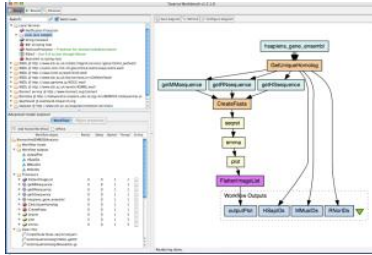Text Document 141 KB

**❹ Connect to statistics** ®

```
find.investigation()
102 downloaded

obs<-find.observedvalue(
43,920 downloaded

#some calculation
add.inferredvalue(res)
36 added
```
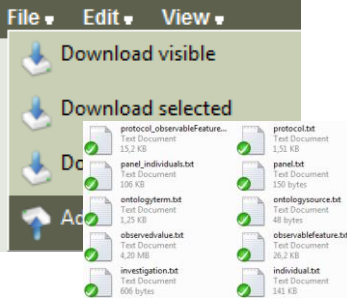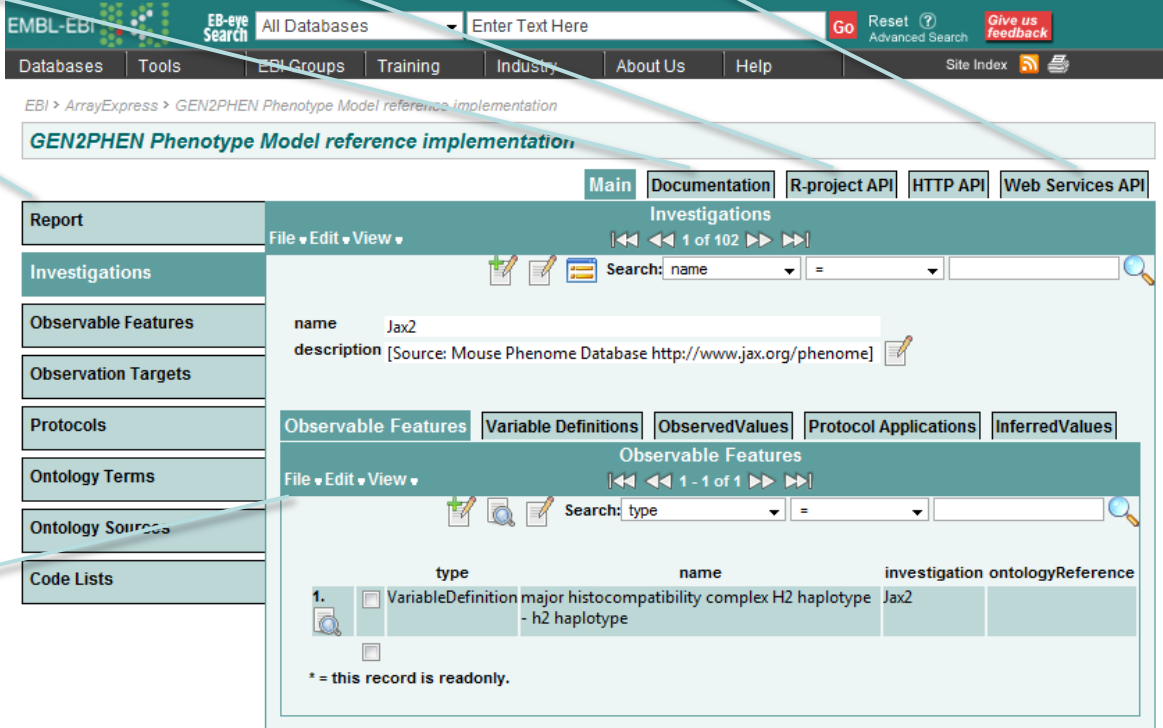
**❺ Workflow ready web-services**

**❻ plugin your own scripts (eg OntologyBrowser)**

EMBL-EBI
EB-eye Search   All Databases ▾   Enter Text Here   **Go**   Reset ⓘ Advanced Search   **Give us feedback**

| Databases | Tools | EBI Groups | Training | Industry | About Us | Help | Site Index 🔊 🖶 |

EBI > ArrayExpress > GEN2PHEN Phenotype Model reference implementation

**GEN2PHEN Phenotype Model reference implementation**

**Main** | Documentation | R-project API | HTTP API | Web Services API

| Report |
| **Investigations** |
| Observable Features |
| Observation Targets |
| Protocols |
| Ontology Terms |
| Ontology Sources |
| Code Lists |

**Investigations**
File ▾ Edit ▾ View ▾      |◀◀ ◀◀ 1 of 102 ▶▶ ▶▶|

📝 📝 📄   Search: name ▾  = ▾   🔍

name   Jax2
description [Source: Mouse Phenome Database http://www.jax.org/phenome] 📝

**Observable Features** | Variable Definitions | ObservedValues | Protocol Applications | InferredValues

**Observable Features**
File ▾ Edit ▾ View ▾      |◀◀ ◀◀ 1 - 1 of 1 ▶▶ ▶▶|

📝 🔍 📝   Search: type ▾  = ▾   🔍

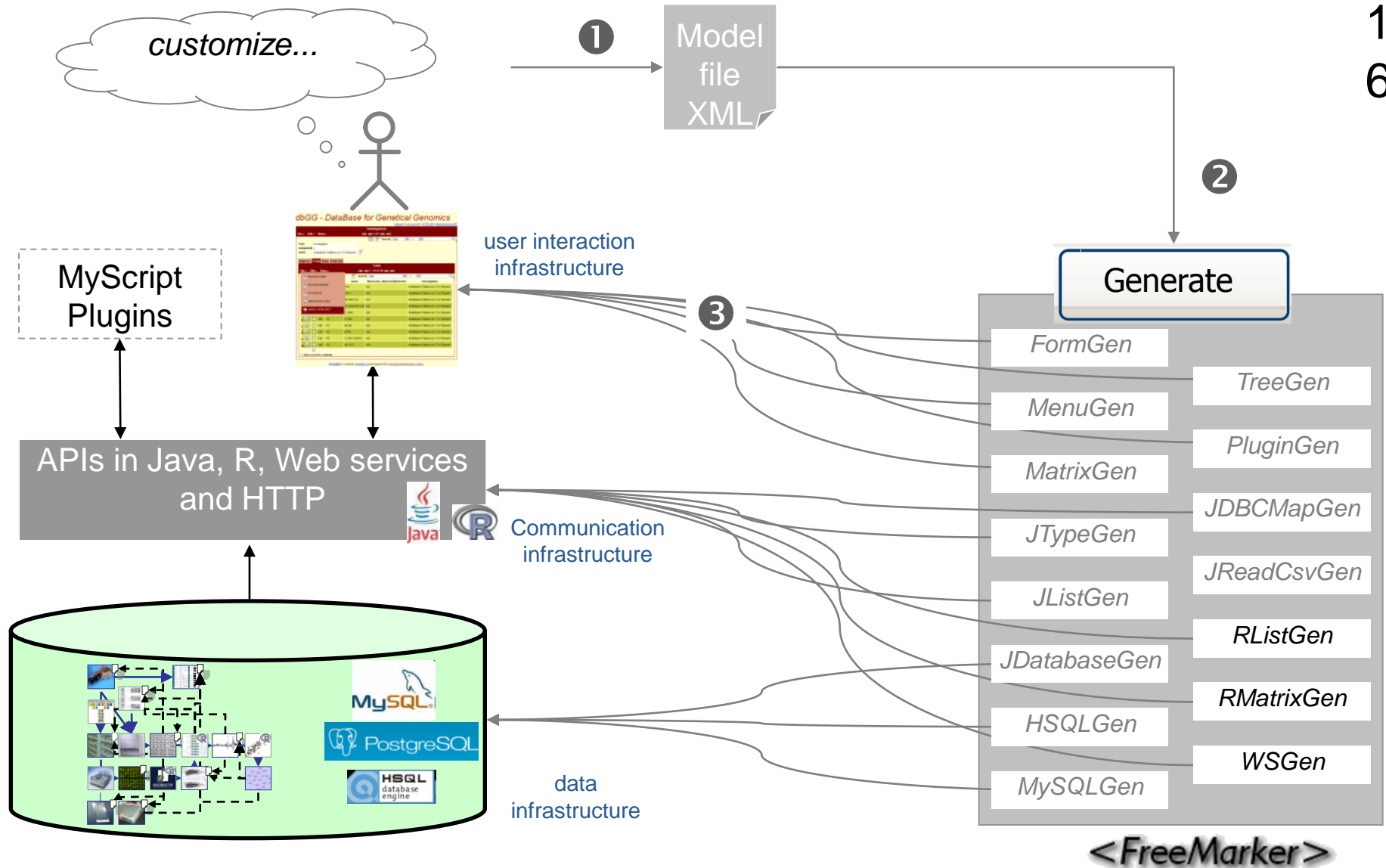| | | type | name | investigation | ontologyReference |
|---|---|---|---|---|---|
| 1. 🔍 | ☐ | VariableDefinition | major histocompatibility complex H2 haplotype - h2 haplotype | Jax2 | |
| | ☐ | | | | |

* = this record is readonly.

This database was generated using the open source MOLGENIS database generator version 3.3.0-testing.
Please cite Swertz et al (2004) or Swertz & Jansen (2007) on use.

# MOLGENIS demo
### *Model -> Generate -> Use*

# Model -> Generate -> Use



*customize...*

❶ Model file XML

16

❷

Generate

MyScript Plugins

user interaction infrastructure

❸

APIs in Java, R, Web services and HTTP

Communication infrastructure

data infrastructure

*FormGen*

*TreeGen*

*MenuGen*

*PluginGen*

*MatrixGen*

*JDBCMapGen*

*JTypeGen*

*JReadCsvGen*

*JListGen*

*RListGen*

*JDatabaseGen*

*RMatrixGen*

*HSQLGen*

*WSGen*

*MySQLGen*

<FreeMarker>

# A generator = template

## e.g. `${Name(entity)}` -> `ExperimentMapper`

### (A) Generator Template

```
public class ${Name(entity)}Mapper
  extends DataMapper<${ Name(entity)}> {
  public String addSql(${Name(entity)} e) {
   return String.format(
    "insert into ${ Name(entity)} ( "
    +"${csv(entity.Fields, "name($i)")}"
    +") values ("
    +"${csv(entity.Fields, "'%s'")}"
    + ")",
    ${csv(entity.Fields,
"e.get${Name(i)}()")}
    );
  } ...
```

*generates*

### (B) Generated source file

```
public class ExperimentMapper
  extends DataMapper<Experiment> {
  public String addSql(Experiment e) {
    return String.format(
    "insert into Experiment ( "
    +"ID,Name,Medium,Stress,Log,
visibleToGroup"
    +") values ("
    +"'%s','%s','%s','%s','%s','%s'
    +")",
    e.getID(), e.getName(),
    e.getMedium(),e.getStress(),
    e.getLog(),e.getVisibleToGroup()
    );
  } ...
```

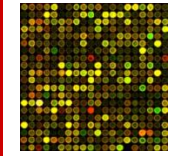# Usage examples in Life Sciences

Mutation

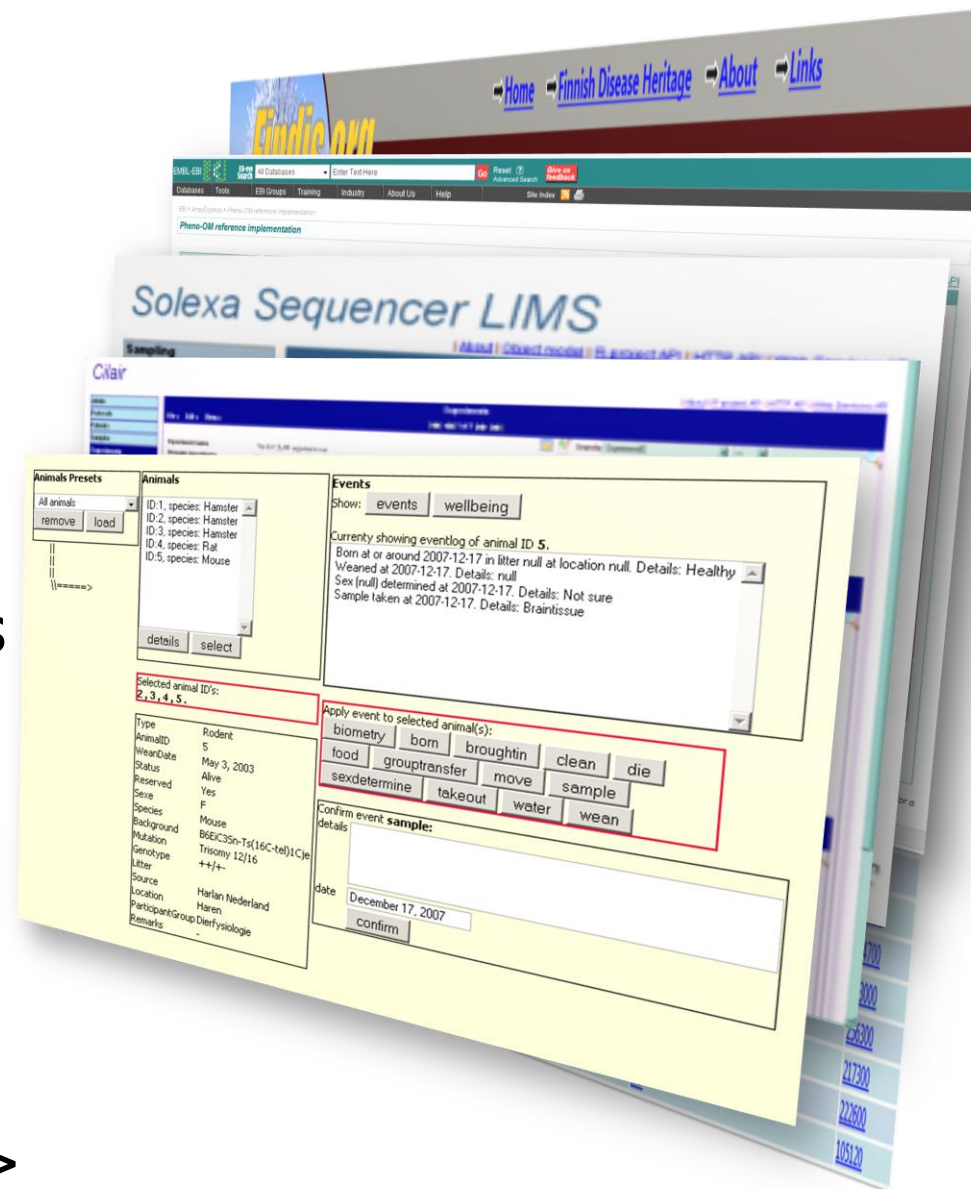Phenotype

Sequencing LIMS

Proteo/Metabolomics

Animal LIMS

GWAS / GWLS

*<add your project here>*

**XGAP**
*extensible genotype and phenotype data model for xQTL*

**Problem domain:
xQTLs
GWAS**

10 — strains
10.000 — markers
genome
inbreed
100 — individuals
genotype
1,000,000 — genotypes
map
10,000 — QTL
correlate
hybridize
100,000 — expression
preprocess
10,000,00 — norm exprs.
network
100 — microarray
100,000 — probes

http://www.xgap.org
Swertz, van der Velde et al (2010) Genome Biology 9;11(3): R27.

# Data in practice

Data in matrices

| Subjects: PANELS |
|---|
| T M r A a R i K t E s: R S | DATA ELEMENTS |

*TRAIT × SUBJECT*

# Annotations in practice

- Annotations in tables, e.g. Marker

| Name | Symbol | Chr | cM | bpStart | mb |
|------|--------|-----|-----|---------|-----|
| C1M1 | I_1_pkP1050 | 1 | -18.2603 | 168807 | 0.168807 |
| C1M2 | I_2_pkP1101 | 1 | -17.2825 | 992188 | 0.992188 |
| C1M3 | I_3_pkP1103 | 1 | -11.959 | 1884415 | 1.884415 |
| C1M4 | I_4_pkP1052 | 1 | -6.1004 | 2818973 | 2.818973 |
| C1M5 | I_5_egPE107 | 1 | -3.5488 | 3502476 | 3.502476 |
| C1M6 | I_6_egPF101 | 1 | -1.4887 | 4338254 | 4.338254 |
| C1M7 | I_7_pkP1054 | 1 | -0.6162 | 4845515 | 4.845515 |
| C1M8 | I_8_egPH102 | 1 | 0.4597 | 5893622 | 5.893622 |
| C1M9 | I_9_pkP1057 | 1 | 0.9366 | 6359867 | 6.359867 |
| C1M10 | I_10_pkP1116 | 1 | 2.1576 | 7589863 | 7.589863 |
| C1M11 | I_11_egPK103 | 1 | 2.4087 | 7894081 | 7.894081 |
| C1M12 | I_12_pkP1059 | 1 | 2.9456 | 8654360 | 8.65436 |
| C1M13 | I_13_pkP1122 | 1 | 3.7959 | 9569914 | 9.569914 |
| C1M14 | I_14_egPN104 | 1 | 4.7801 | 10259909 | 10.259909 |
| C1M15 | I_15_egPO105 | 1 | 6.0193 | 11085295 | 11.085295 |
| C1M16 | I_16_pkP1068 | 1 | 7.5226 | 11760182 | 11.760182 |

# Model: try 1

# But…

# XGAP model: <any trait> X <any subject>

# Extending on FuGE



Jones et al (2007) Functional Genomics Experiment model.
Nature Biotechnology

**XGAP**

 *extensible genotype and
 phenotype <u>software platform</u> for xQTL*

# Generated: user interfaces



XGAP - eXtensible Genotype and Phenotype platform

http://www.xgap.org

Swertz, van der Velde et al (2010) Genome Biology 9;11(3): R27.

# Data exploration

**Phenotypes**

|◄◄ ◄◄ 1 - 8 of 8 ►► ►►|

| name | Description | |
|------|-------------|---|
| PCTT10 | Percent time spent in center of arena (interval of 10 min) | Identifi |
| TOTDIST | Total distance | Identifi |
| TOTREAR | Total rearing | Identifi |
| AMBEPIS | Ambulatory episodes | Identifi |
| AVGVELO | Average velocity | Identifi |
| PCTREST | Percent resting | Identifi |
| ACTFACT | Activity factor | Identifi |
| ANXFACT | Anxiety factor | Identifi |

**Individuals**

|◄◄ ◄◄ 1 - 10 of 362 ►► ►►|

| name | Strain | |
|------|--------|---|
| 138422 | C57BL/6J (B6) + C58/J | Identifi |
| 138423 | C57BL/6J (B6) + C58/J | Identifi |
| 138424 | C57BL/6J (B6) + C58/J | Identifi |
| 138425 | C57BL/6J (B6) + C58/J | Identifi |
| 140942 | C57BL/6J (B6) + C58/J | Identifi |
| 140943 | C57BL/6J (B6) + C58/J | Identifi |
| 140944 | C57BL/6J (B6) + C58/J | Identifi |
| 141427 | C57BL/6J (B6) + C58/J | Identifi |
| 141428 | C57BL/6J (B6) + C58/J | Identifi |
| 141429 | C57BL/6J (B6) + C58/J | Identifi |

**File** ▾    |◄◄ ◄◄ Phenotype 1-8 of 8 ►► ►►|

Individual 1-10 of 362

Panspeed 10
Width 15
Height 10

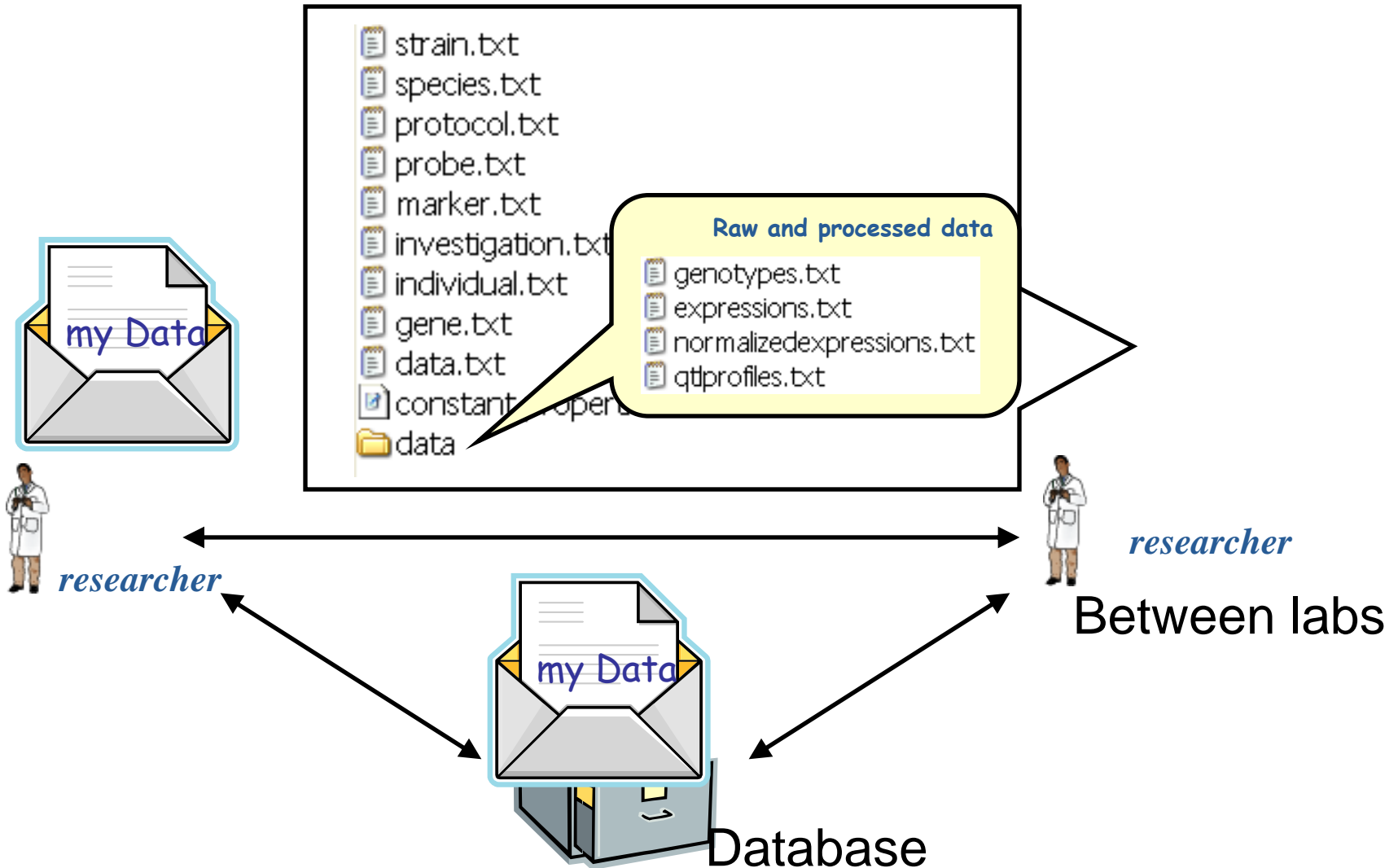| | PCTT10 | TOTDIST | TOTREAR | AMBEPIS | AVGVELO | PCTREST | ACTFACT | ANXFACT |
|------|--------|---------|---------|---------|---------|---------|---------|---------|
| 138422 | 35.35 | 3,818.8 | 57 | 138 | 43.16 | 54.47 | -0.01 | 1.92 |
| 138423 | 18.82 | 3,741 | 67 | 115 | 37.48 | 48.83 | 0.18 | 0.05 |
| 138424 | 17.2 | 3,569 | 108 | 117 | 33.49 | 53.63 | 0.11 | -0.49 |
| 138425 | 19.93 | 3,466.4 | 70 | 113 | 35.45 | 52.67 | -0.12 | -0.07 |
| 140942 | 20.38 | 5,296.4 | 123 | 136 | 35.46 | 38.47 | 1.78 | -0.29 |
| 140943 | 17.57 | 2,689.8 | 91 | 91 | 37.29 | 56.25 | -0.79 | 0.08 |
| 140944 | 30.27 | 4,108.2 | 63 | 141 | 41.64 | 46.38 | 0.61 | 1.29 |
| 141427 | 28.97 | 3,466.5 | 112 | 127 | 36.14 | 48.17 | 0.33 | 0.59 |
| 141428 | 13.25 | 2,391.7 | 83 | 76 | 29.55 | 60.18 | -1.2 | -1.1 |
| 141429 | 22.12 | 3,140.5 | 62 | 107 | 35.48 | 52.55 | -0.37 | 0.11 |

http://www.xgap.org

Swertz, van der Velde et al (2010) Genome Biology 9;11(3): R27.

# Generated: common database/format

Simple text based format

# Generated: common database/format

# Plugin: import wizard

# Generated: rich user documentation

## XGAP 1.4 distro prototype documentation.

### Table of contents

### Investigation
*extends FugeInvestigation*

Inherited atttributes:
annotations, id, name, description, start, end,

Constraints:

**unique(id)**:
Field id is unique within an Investigation.

**unique(name)**:
Name is unique.

### ProtocolApplication
*extends FugeProtocolApplication*

Inherited atttributes:
annotations, id, name, description, Investigation, activityDate, inputData, protocol, protocolDeviation, outputMaterials, outputData,

http://www.xgap.org
Swertz, van der Velde et al (

# Generated: connection to R statistics



**XGAP**

*markers*  *subjects*
*phenotype*  *genotypes*

```
#connect to my XGAP database
source("http://aserver/xgap/api/R")

#upload my 'metanetwork' investigation
add.investigation(name="metanetwork")

#use 'metanetwork' investigation
use.investigation(name="metanetwork")

#upload subjects and traits
add.marker(name=rownames(markers),
          chr =markers$chr,
          cm  =markers$cM)
add.metabolite(name=rownames(metabolites))
add.subject(name=colnames(genotypes))

#upload genotype and phenotype data matrices
add.datamatrix(geno,
    name="geno" rowtype="marker"
    coltype="subject" valuetype="text")
add.datamatrix(mpheno,
    name="mexpr" rowtype="metabolite"
    coltype="subject" valuetype="double")
```
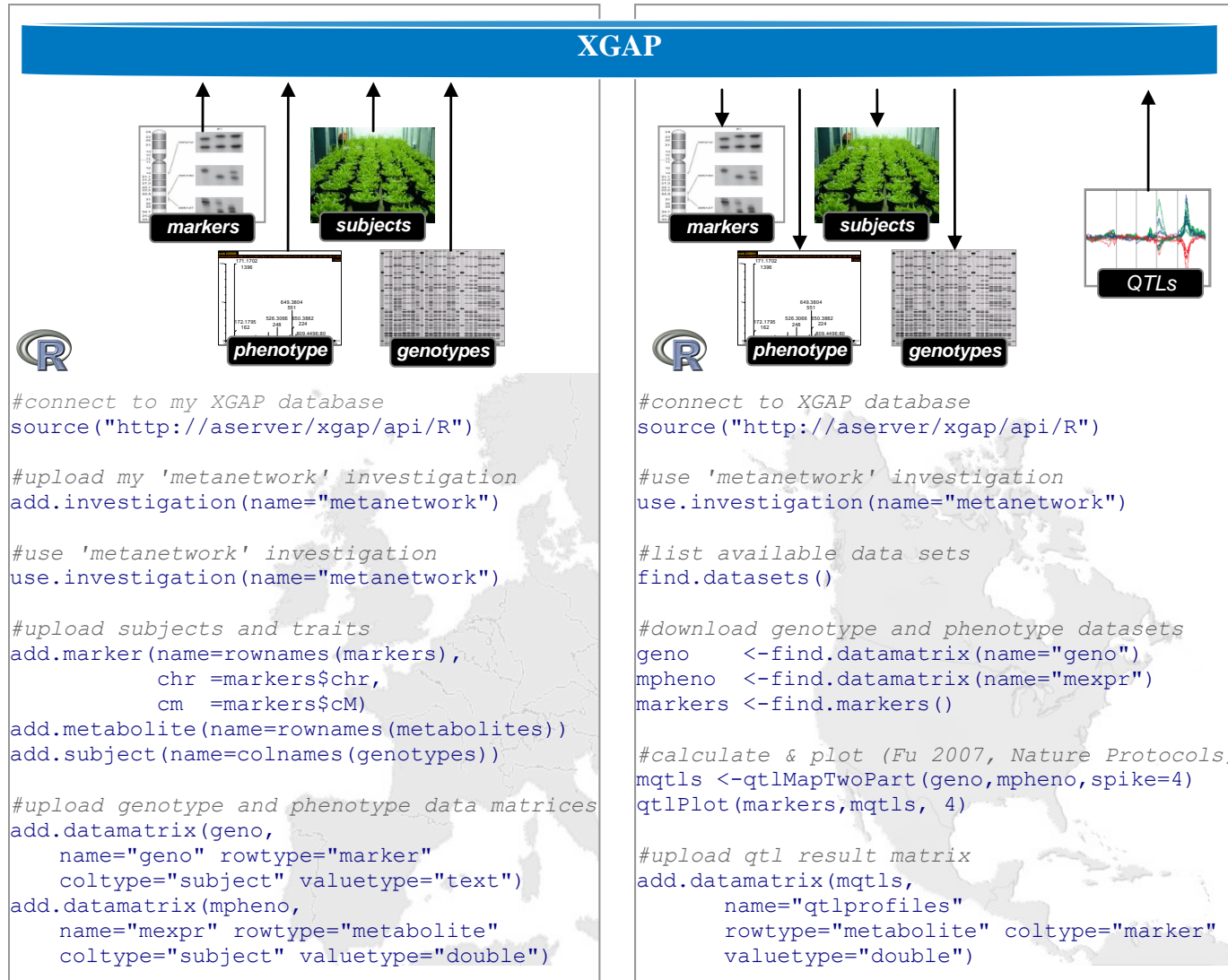
Scientist **A** uploads raw data

*markers*  *subjects*  *QTLs*
*phenotype*  *genotypes*

```
#connect to XGAP database
source("http://aserver/xgap/api/R")

#use 'metanetwork' investigation
use.investigation(name="metanetwork")

#list available data sets
find.datasets()

#download genotype and phenotype datasets
geno    <-find.datamatrix(name="geno")
mpheno  <-find.datamatrix(name="mexpr")
markers <-find.markers()

#calculate & plot (Fu 2007, Nature Protocols)
mqtls <-qtlMapTwoPart(geno,mpheno,spike=4)
qtlPlot(markers,mqtls, 4)

#upload qtl result matrix
add.datamatrix(mqtls,
      name="qtlprofiles"
      rowtype="metabolite" coltype="marker"
      valuetype="double")
```
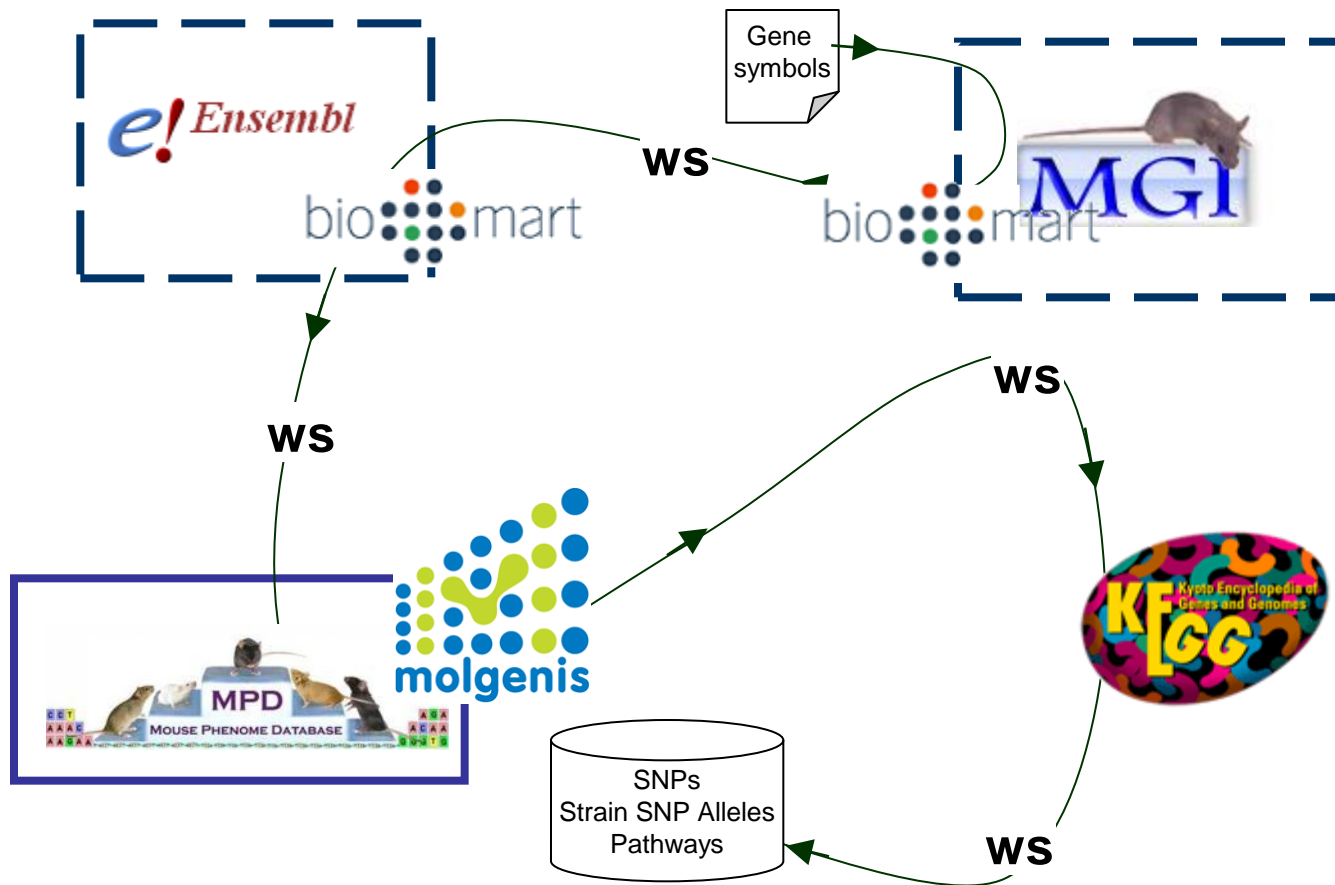
Scientist **B** uploads analysis results

Swertz et al (2010) XGAP: a uniform and extensible data model and software platform for genotype and phenotype experiments. Genome Biology 11(3).

# Generated: tool integration interfaces

- REST, SOAP, RDF



Smedley, Swertz, Wolstencroft et al (2008) Solutions for data integration in functional genomics: a critical assessment and case study. Briefings in Bioinformatics 9(6)

# Plugin: Data analysis using cloud/cluster

# Discussion & Conclusion
*GMOD links?*

# GMOD link ideas

- Chado
  - XGAP harmonization towards Chado?
  - MOLGENIS 4 Chado?

- Gbrowse & DAS
  - Have XGAP data projected on genome browser?
  - Serve XGAP data as custom tracks?

- BioMART/InterMine
  - Consume BioMART data to auto-annotate experimental data?
  - Export XGAP experiments into MART/MINE query environments?

# Ontologizing....



HPO:
Abnormally shaped ears
Auricular malformation
Deformed auricles

MP:
Malformed auricles
Malformed ears
Malformed external ears
etc

Deformed ears?
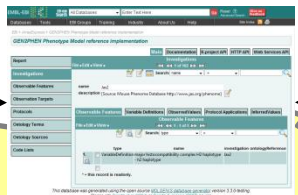
Local ontologies (OLW or OBO)

BioPortal

OLS

*query*

*expansion*

RDF + SPARQL

Panacea

GEN2PHEN

LifeLines

IOP

OntoCAT – Ontology common API tasks
http://www.ontocat.org  and  http://precedings.nature.com/documents/4666

# Getting started

*http;//www.molgenis.org*

## molgenis

Search

Login    Preferences    Register

| Wiki | News | Documentation | Download | ChangeLog | Tickets | Browse Source | Roadmap |

★★★★★    Start Page    Index    History    Last Change

### MOLGENIS development manual

- Software needed
  - Java
  - Tomcat
  - MySql/Postgresql
  - Eclipse
  - MOLGENIS (svn or zip)

- Model + Generate
  - New database
  - Existing databases

- Use
  - Web interrace
  - R, REST, SOAP, JAVA interfaces

LGPL3
Free Software

# Acknowledgements

Morris Swertz

Joeri van der Velde

Joris Lops

Danny Arends

Alex Kanterakis

Erik Roos

Richard Scheltema

Martijn Dijkstra

Rudi Alberts

Bruno M. Tesson

Gonzalo Vera Rodriguez

Tomasz Adamuziak

Juha Muilu

Gudmundur Thorisson

Damian Smedley

Katy Wolstencroft

Ritsert C. Jansen

Cisca Wijmenga

Carole Goble

John M. Hancock

Andrew R. Jones

Klaus Schughart

Paul Schofield

Anthony Brookes

Helen E. Parkinson

BBMRI-NL biobanking (Hs)

EU-GEN2PHEN consortium (Hs)

EU-PANACEA consortium (Ce)

NL Brassica Nutr. consortium (At)

EU-CASIMIR consortium (Mm)

NBIC/BioAssist consortium (bioinfo)

# Thank you! Questions?

m.a.swertz@rug.nl
k.j.van.der.velde@rug.nl

## Web

- MOLGENIS: http://www.molgenis.org
- XGAP: http://www.xgap.org
- OntoCAT: http://www.ontocat.org

molgenis
.org
*Your database at the push of a button*

## Pubmed

- Swertz et al (2010) *Genome Biology* 9;11(3): R27.
- Smedley et al (2008) *Briefings in Bioinformatics* 9(6): 532-544
- Swertz & Jansen (2007) *Nature Reviews Genetics* 8, 235-243
- Swertz et al (2004) *Bioinformatics* 20(13), 2075-83