

TARGET Conference 2013
Probing Big Data for answers
3 - 5 April, 2013
Groningen, Netherlands

Flexible datamodels for life sciences

K. Joeri van der Velde, Martijn Dijkstra, Morris Swertz,
members of the Genomics Coordination Center



university of
groningen

genomics coordination
center

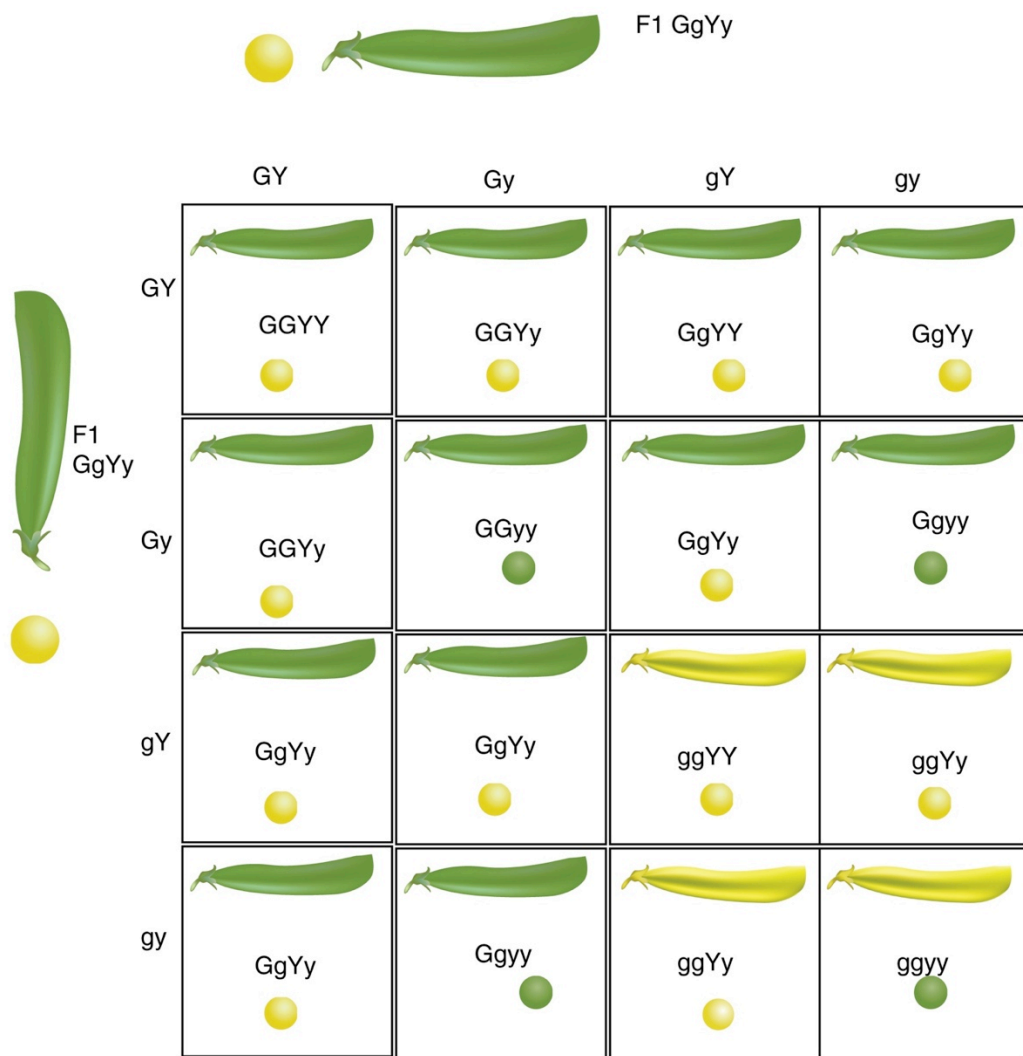
Outline

- Introduction
 - Motivation
 - Dealing with variation in research needs
 - Design-time vs runtime configuration
- Results
 - XGAP model for homogeneous data ('molecules')
 - Observ-OM model for heterogenous data ('phenotypes')
- Showcase
 - EB Registry
 - WormQTL
- Current work

Introduction

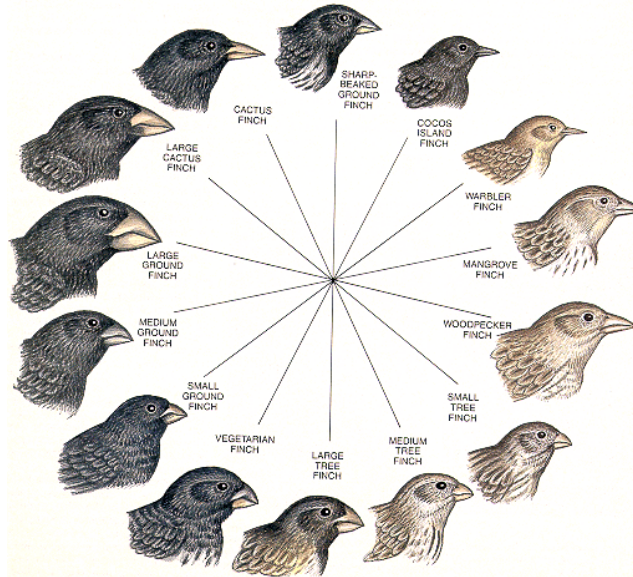
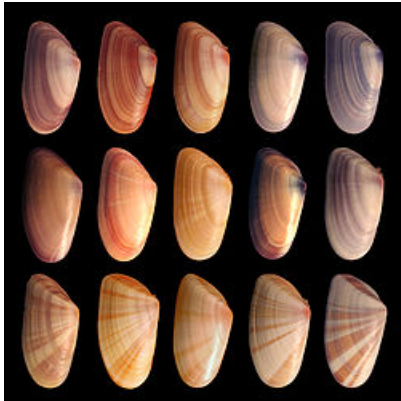
Motivation

Biological variation first explained by genetics



*Mendel got lucky..
single-locus traits*

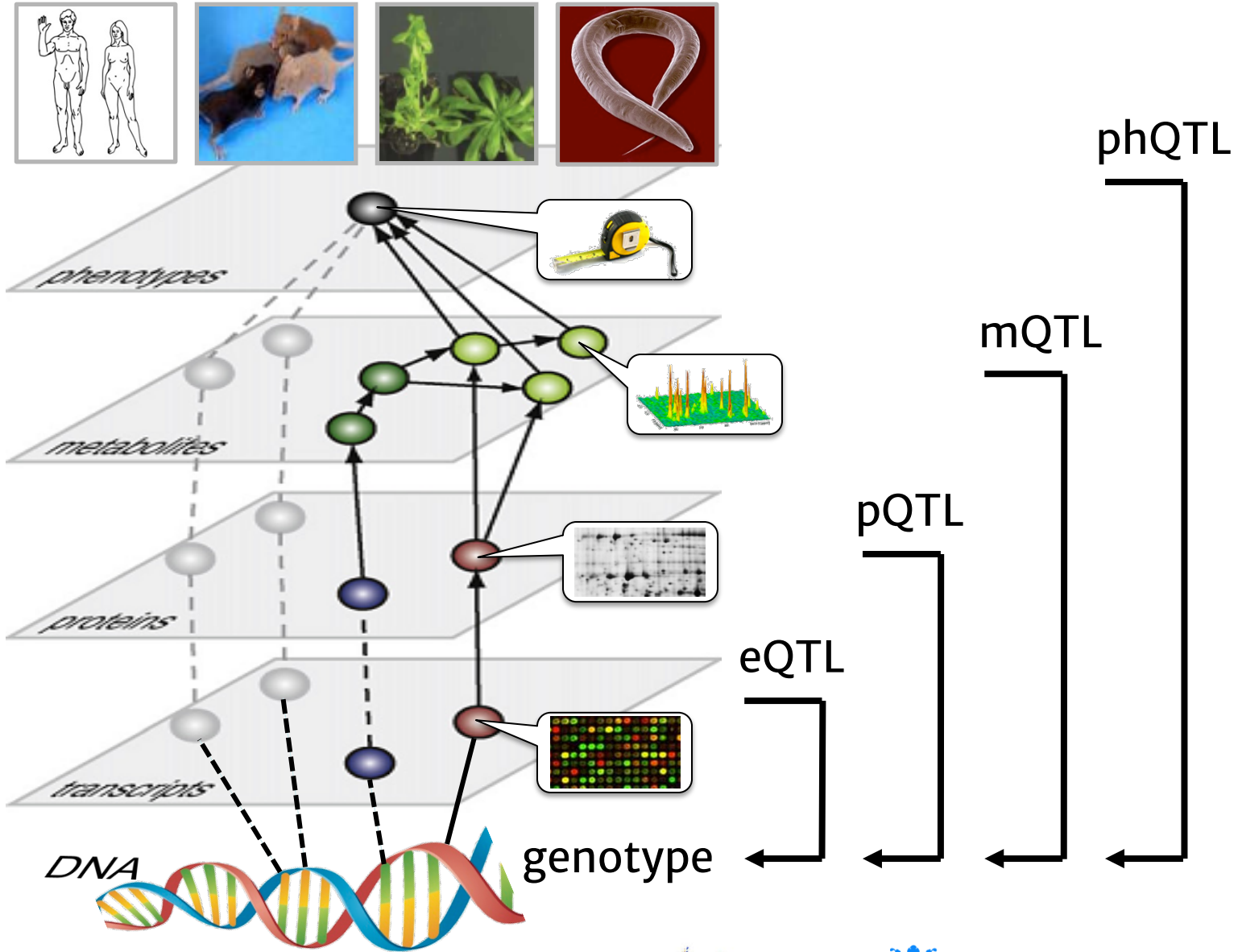
Motivation: Understanding complex variation



	wild-type	wild-type	<i>C. elegans</i> (RNAi)
Polar body location (5)			
	<i>C. elegans</i>	<i>Rhabditella axei</i>	<i>zyg-11</i> (RNAi)
Detached centrosome (7)			
	<i>C. elegans</i>	<i>Bursilla</i> sp.	<i>sun-1</i> (RNAi)
Pseudo-cleavage (8)			
	<i>C. elegans</i>	<i>Protorhabditis</i> sp.	<i>pph-6</i> (RNAi)
Centrosome shape (25)			
	<i>C. elegans</i>	<i>R. inermis</i>	<i>lin-5</i> (RNAi)
Asynchrony (34)			
	<i>C. elegans</i>	<i>O. myriophila</i>	<i>par-1</i> (RNAi)
Spindle orientation (36)			
	<i>C. elegans</i>	<i>Protorhabditis</i> sp.	<i>par-6</i> (RNAi)



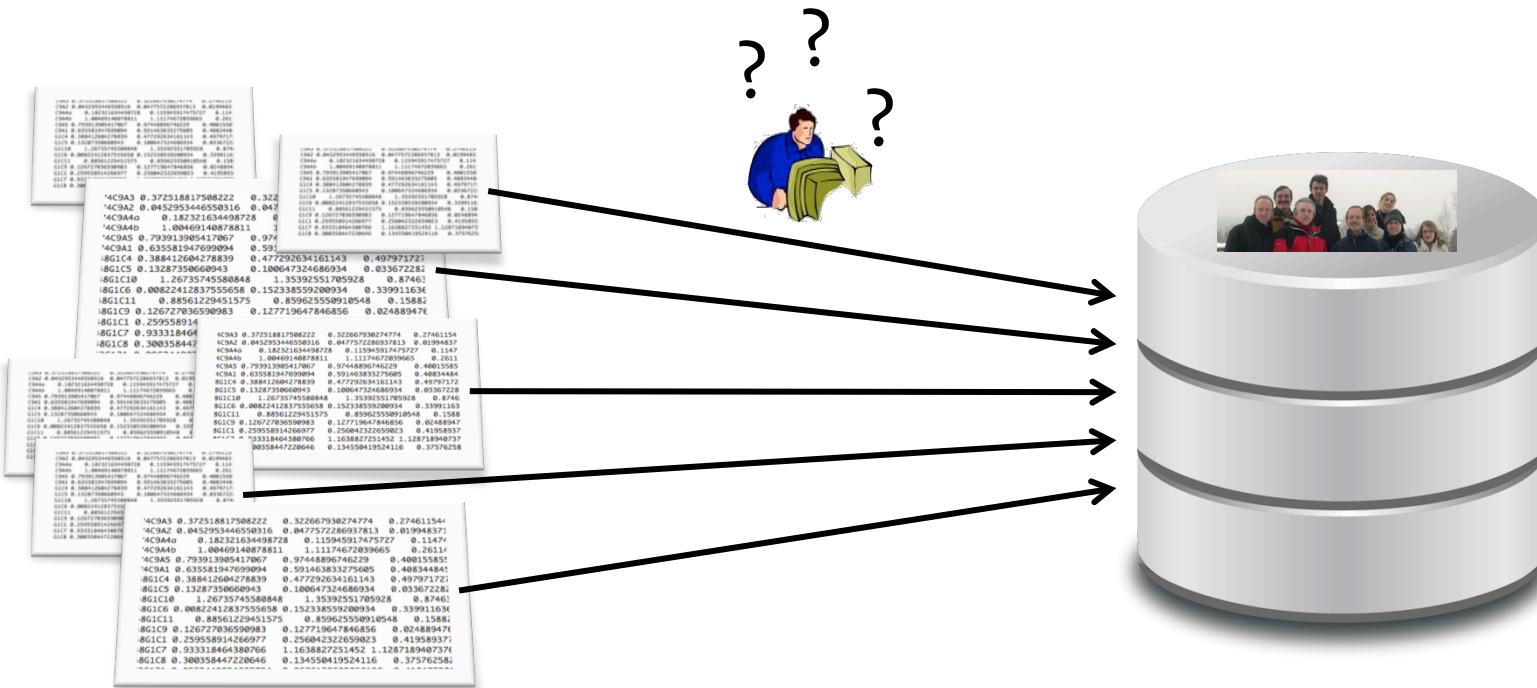
How? Understanding geno-to-pheno



Introduction

Dealing with variation in research needs

Challenge: Building a 'team' database



Introduction

Design-time vs runtime configuration

Autogenerate the software

Model in DSL



NextGenSeq



Gen

What models to use?
Can we have a model that rules them all?

Mutation c

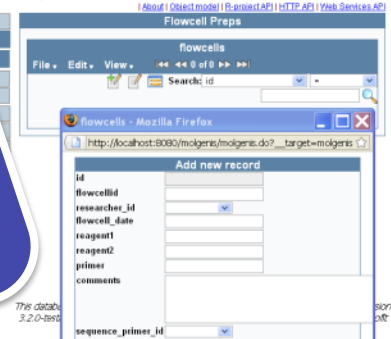


Model organisms

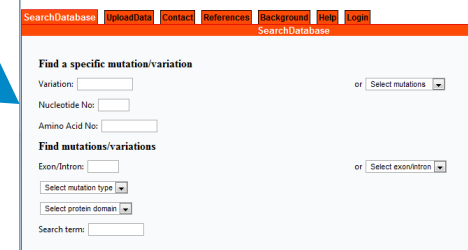


Use generated software

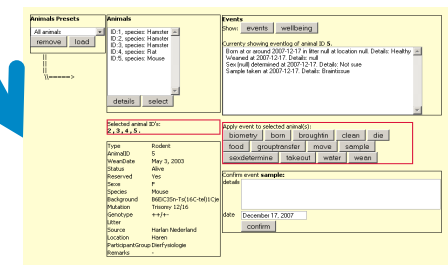
Solexa Sequencer LIMS



database of COL7A1 mutations



Animal Observatory



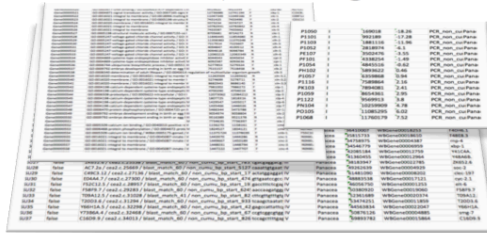
<http://www.molgenis.org>

Swertz & Jansen (2007) *Nature Reviews Genetics* 8, 235-243

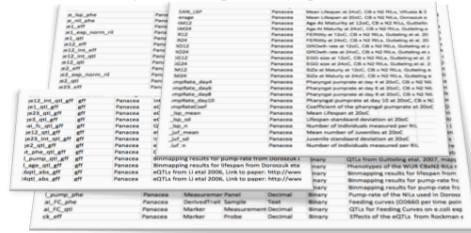
Swertz et al (2004) *Bioinformatics* 20(13), 2075-83

What are we dealing with?

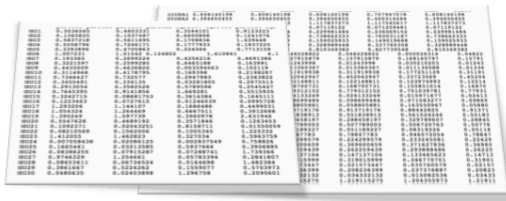
Genomic features,
individuals, ontologies ..



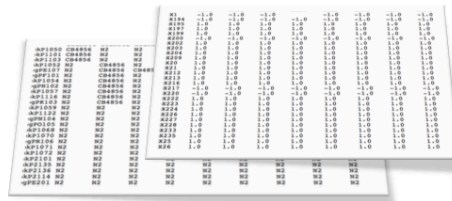
Metadata for phenotypes,
datasets, samples, panels ...



Biomolecular measurements,
association results ...



Genotypes,
conditions ...



Stable?

Annotations of concepts
used in data sets,
mostly static content

Dynamic?

Experimental data sets,
usually flexible and
volatile content

Apparently we need something **stable** AND **dynamic**
..without becoming exhaustive

Example: eQTL data

Stable?

Stable?

Probe (annotation)

name	mismatch	description
WSU1	true	NA / SpotReport / blast_match_NA / n
WSU2	false	C25A1.8 / cea2.c.00914 / blast_match
WSU3	false	F21F3.6 / cea2.c.02677 / blast_match
WSU4	false	F25H2.9 / cea2.c.02801 / blast_match
WSU5	false	F56H1.4 / cea2.c.04344 / blast_match

Marker (annotation)

name	chromosome	bpstart	cm	description
pkP1050	I	169018	-18.26	PCR_non_cu
pkP1101	I	992189	-17.28	PCR_non_cu
pkP1103	I	1881116	-11.96	PCR_non_cu
pkP1052	I	2818974	-6.1	PCR_non_cu
egPE107	I	3502476	-3.55	PCR_non_cu

	pkP1050	pkP1101	pkP1103	pkP1052	egPE107
WSU1	0.5036565	0.4603331	0.3544101	0.9123223	0.4157701
WSU2	0.1365825	0.1037487	0.6600898	0.1241076	0.1672705
WSU3	0.5837218	0.5611695	0.1708836	1.439448	1.94431
WSU4	0.5558796	0.7246171	0.1777933	0.1937225	0.4413371
WSU5	0.3393896	0.4705863	0.224066	0.7713159	0.01334126

Dynamic?

eQTL profiles (data set)

Stable = good for code generation

Annotations: Column-oriented data

name	chromosome	bpstart	cm	description
pkP1050	I	169018	-18.26	PCR_non_cu
pkP1101	I	992189	-17.28	PCR_non_cu
pkP1103	I	1881116	-11.96	PCR_non_cu
pkP1052	I	2818974	-6.1	PCR_non_cu
egPE107	I	3502476	-3.55	PCR_non_cu

← *Attributes*

3. import

↓ 1. model

```
<entity name="Locus" abstract="true">
  <description> position. Typical examples of such traits are genes,
  probes and markers. Common structure for entities that have a
  genomic</description>
  <field name="Chromosome" label="Chromosome" type="xref"
  xref_entity="Chromosome" xref_field="id" xref_label="name" nillable="1"
  description="Reference to the chromosome this
  position belongs to." />
  <field name="cm" label="cmPosition" type="decimal" nillable="true"
  description="genetic map position in centi_morgan (cM)." />
  <field name="bpStart" label="Start (5')" type="long" nillable="true"
  description="numeric basepair posion (5') on the chromosome" />
  <field name="bpEnd" label="End" type="long" nillable="true"
  description="numeric basepair posion (3') on the chromosome" />
  <field name="Seq" type="text" nillable="true"
  description="The FASTA text representation of the sequence." />
  <field name="Symbol" type="varchar" nillable="true"
  description="todo" />
</entity>
<entity name="Chromosome" extends="ObservableFeature">
  <field name="orderNr" type="int" />
  <field name="isAutosomal" type="bool" description="Is 'yes' when number of chr
  <field name="bpLength" type="int" nillable="true" description="Lenght of the c
  <field name="Species" label="Species" type="xref" xref_entity="Species"
  xref_field="id" xref_label="name" nillable="true"
  description="Reference to the species this
  chromosome belongs to." />
</entity>
```

2. generate



id	name	description	investigation ontology	reference	alternative identifiers	label	Chromosome cM
105050	WSU1	NA / SpotReport / blast_match_NA / non_cumulative_bp_start_0	Public				
105051	WSU2	C25A1.8 / cea2.c.00914 / blast_match_60 / non_cumulative_bp_start_10184580	Public			dec-87	I
105052	WSU3	F21F3.6 / cea2.c.02677 / blast_match_60 / non_cumulative_bp_start_4912043	Public			F21F3.6	I
105053	WSU4	F25H2.9 / cea2.c.02801 / blast_match_60 / non_cumulative_bp_start_10567120	Public			pas-5	I
105054	WSU5	F56H1.4 / cea2.c.04344 / blast_match_60 / non_cumulative_bp_start_5741975	Public			rpt-5	I
105055	WSU6	H06001.1 / cea2.c.04508 / blast_match_60 / non_cumulative_bp_start_7011070	Public			pdl-3	I

Challenge: Modeling 'dynamic' data sets

Can we use 'entity-attribute-value' modeling?

	BF221L	FD85C	C6L9	T7M24	BF151L
RIL1	A	A	A	B	B
RIL2	B	A	B	B	A
RIL3	B	B	A	B	A
RIL4	A	B	B	A	B
RIL5	A	B	A	A	B

← *Not stable columns!
Different per data set...*

↓ 1. model

*Values all comparable:
Columns all of the same type*

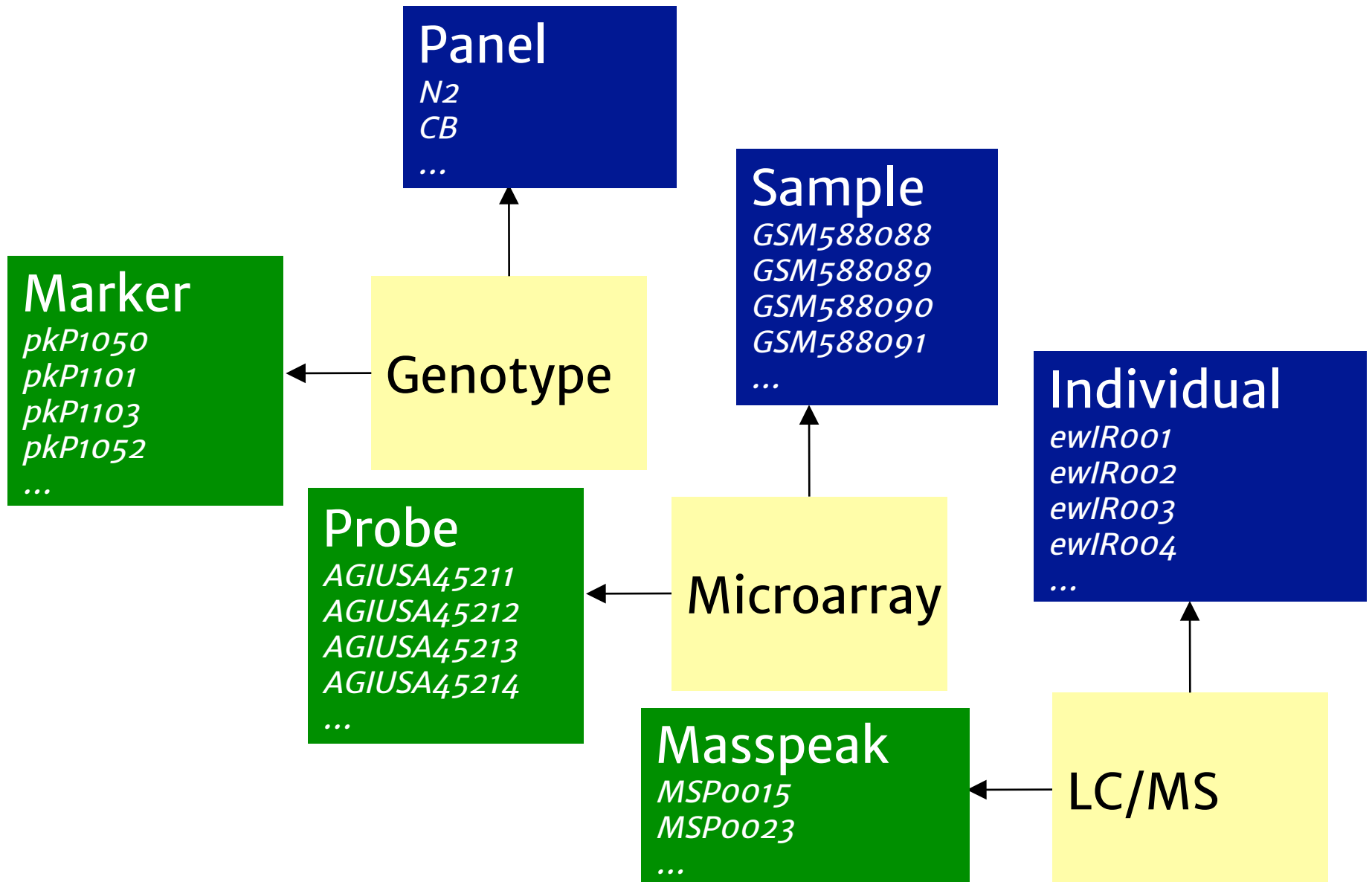
```
<entity name="Locus" abstract="true">
  <description> position. Typical examples of such traits are genes,
  probes and markers. Common structure for entities that have a
  genomic</description>
  <field name="Chromosome" label="Chromosome" type="xref"
    xref_entity="Chromosome" xref_field="id" xref_label="name" nillable="true"
    description="Reference to the chromosome this
    position belongs to." />
  <field name="cM" label="cMPosition" type="decimal" nillable="true"
    description="genetic map position in centi morgan (cM)." />
  <field name="bpStart" label="Start (5')" type="long" nillable="true"
    description="numeric basepair postion (5') on the chromosome" />
  <field name="bpEnd" label="End" type="long" nillable="true"
    description="numeric basepair postion (3') on the chromosome" />
  <field name="Seq" type="text" nillable="true"
    description="The FASTA text representation of the sequence." />
  <field name="Symbol" type="varchar" nillable="true"
    description="todo" />
</entity>
<entity name="Chromosome" extends="ObservableFeature">
  <field name="orderNr" type="int" />
  <field name="isAutosomal" type="bool" description="Is 'yes' when number of chr
  <field name="bpLength" type="int" nillable="true" description="Lenght of the c
  <field name="Species" label="Species" type="xref" xref_entity="Species"
    xref_field="id" xref_label="name" nillable="true"
    description="Reference to the species this
    chromosome belongs to." />
</entity>
```



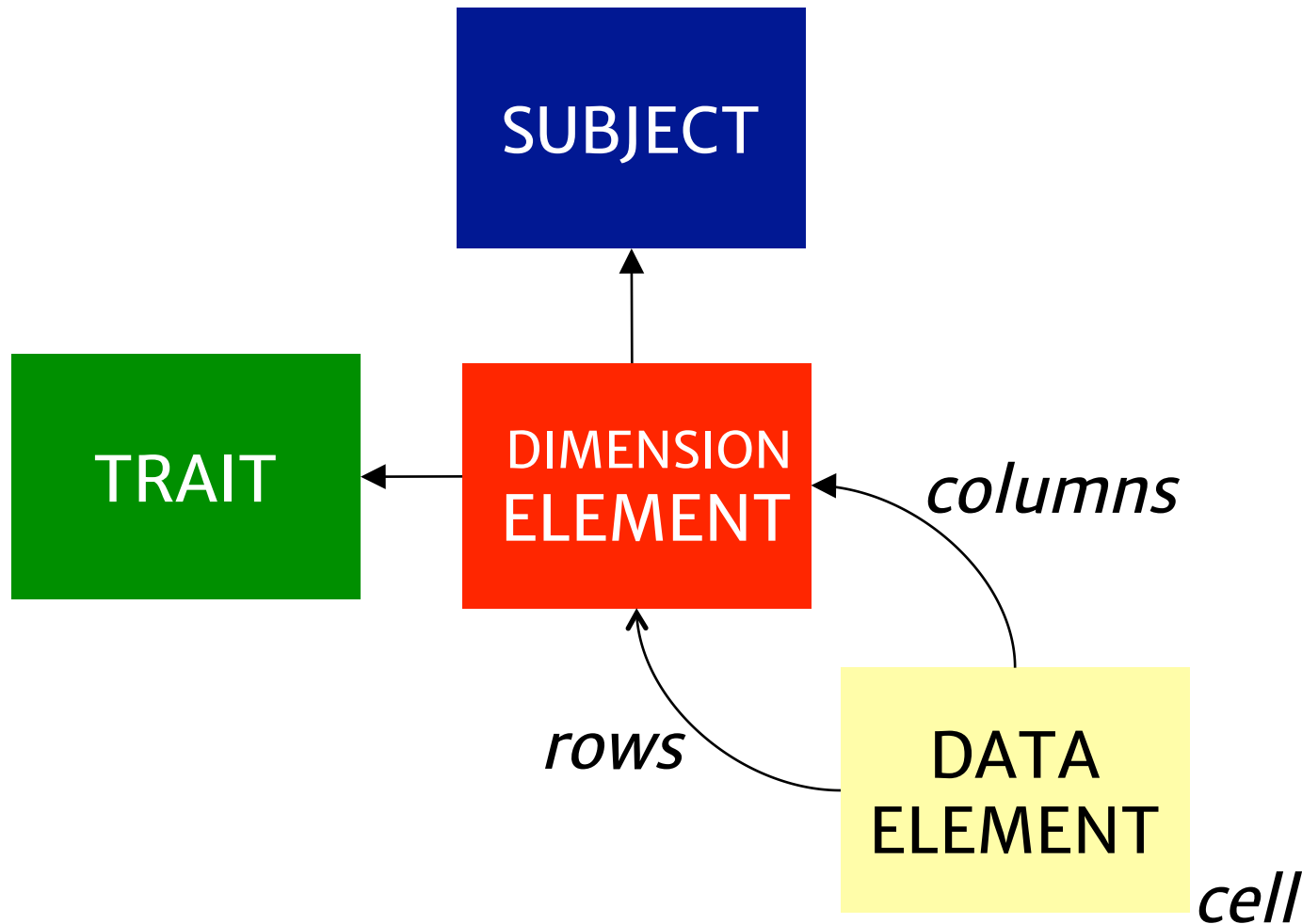
Results

XGAP model

Challenge: Data sets can be variable combinations

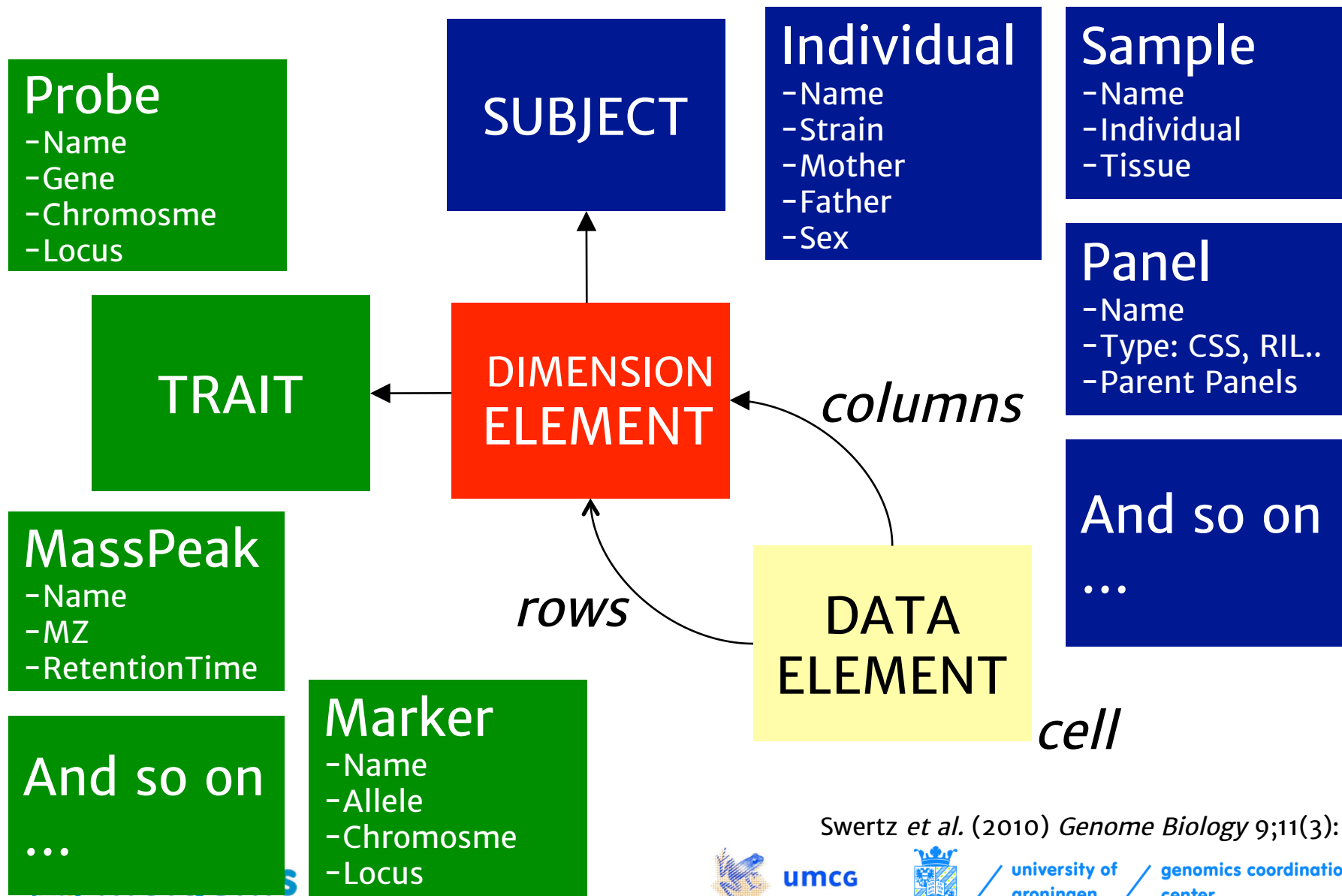


XGAP model: <any trait> X <any subject>



Swertz et al. (2010) *Genome Biology* 9;11(3): R27.

Extensible core model



Swertz et al. (2010) *Genome Biology* 9;11(3): R27.



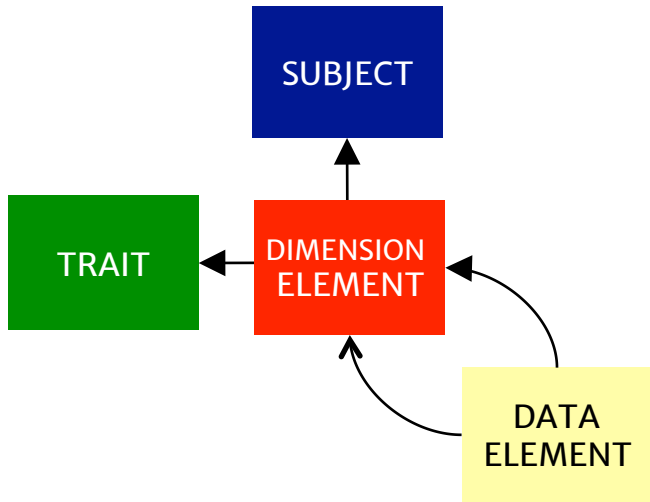
umcg



university of
 groningen

genomics coordination
 center

Using the XGAP model



1. model

```

<field name="Feature" type="xref" xref_entity="ObservationElement"
  xref_field="id" xref_label="name"
  description="References the ObservableFeature that this observation was ma
<description>FIXME: change to ObservationTarget?</description>
<field name="Target" type="xref" xref_entity="ObservationElement"
  xref_field="id" xref_label="name"
  description="References the ObservationTarget that this feature was made o
</entity>
<entity name="ObservedValue" implements="Observation">
  <description>
    Generic storage of values, relationships and optional ontology
    mapping of the value/relation. Values can be atomic observations,
    e.g., length (feature) of individual 1 (target) = 179cm (value).
    Values can also be relationship values, e.g., extract (feature) of
    sample 1 (target) = derived sample (relation).
  <br />
  Discussion: how to model sample pooling in this model?
  <br />
  More Discussion: do we want to have type specific subclasses? No,
  because you can solve this by casting during querying?
</description>
<field name="ontologyReference" type="xref" xref_entity="OntologyTerm"
  nillable="true"
  description="(Optional) Reference to the
  ontology definition or 'code' for this value (recommended for non-numeric
  values such as codes)" />
  
```

2. generate


4. import

	DPZLIL	PLUOL	LOLY	ILM4	DP10IL
RIL1	A	A	A	B	B
RIL2	B	A	B	B	A
RIL3	B	B	A	B	A
RIL4	A	B	B	A	B
RIL5	A	R	A	A	R

3. program viewer plugin

Marker 1-5 of 121

	pkP1050	pkP1101	pkP1103	pkP1052	egPE107
WSU1	0.1124274	0.1386495	0.04974491	1.208681	0.3715276
WSU2	0.2477207	0.2091639	0.6656859	0.004070799	0.5588038
WSU3	0.1799319	0.2640062	0.007156544	0.04396198	0.3165853
WSU4	0.663039	0.6303374	0.1686028	0.7328358	0.1075895
WSU5	0.6529894	0.7207575	0.2907452	0.2535538	0.1900313
WSU6	0.05398721	0.02133219	0.3999896	0.9596403	3.238696
WSU7	0.6048707	0.6055992	0.2280684	0.152925	0.3723309
WSU8	0.07370543	0.0719364	0.45418	1.144719	1.111111
WSU9	1.136192	1.220644	1.409933	1.126819	0.411111
WSU10	0.3662646	0.1562056	0.2816977	0.2554982	0.111111




Outcome: working application

Strains (panels) Chromosomes Markers Genes Measurements DerivedTraits Probes Samples

Probes

File Edit View 1 - 10 of 68,452

Search:

Micro-array probes blasted against WS220 and linked to the genes, used for gene expression phenotypes and eQTLs.

Click on to plot an item, and on to return to this list.

id	name	description	Investigation ontology	Reference	Alternative identifiers	label	Chromosome	cM
10505	WSU1	IA / SpotReport / last_match_NA / on_cumulative_start_0	Public					
10505	WSU2	25A1.8 / cea2.c.00914 / last_match_60 / on_cumulative_start_10184580	Public			clec-87		
10505	WSU3	21F3.6 / cea2.c.02677 / last_match_60 / on_cumulative_start_4912043	Public			F21F3.6		
10505	WSU4	25H2.9 / cea2.c.02801 / last_match_60 / on_cumulative_start_10567120	Public			pas-5		
10505	WSU5	56H1.4 / cea2.c.04344 / last_match_60 / on_cumulative_start_5741975	Public					
10505	WSU6	06O01.1 / cea2.c.04508 / last_match_60 / on_cumulative_start_7015970	Public			pdi-3		
10505	WSU7	20F10.2 / cea2.c.06048 / last_match_60 / on_cumulative_start_10300315	Public			T20F10.2		

Stable!

Strains (panels) Chromosomes Markers Genes Measurements DerivedTraits Probes Samples

Markers

File Edit View 1 - 10 of 1,579

Genetic markers used in one or more of the populations stored in WormQTL. Click on to view one item, and on to return to this list.

id	name	description	Investigation ontology	Reference	Alternative identifiers	label	Chromosome	cM
10347	pkP1050	CR_non_cumulative_bp_pos_169018	Public					-18
10347	pkP1101	CR_non_cumulative_bp_pos_992189	Public					-17
10347	pkP1103	CR_non_cumulative_bp_pos_1881116	Public					-11
10347	pkP1052	CR_non_cumulative_bp_pos_2818974	Public					-6
10347	egPE107	CR_non_cumulative_bp_pos_3502476	Public					-3
10347	egPF101	CR_non_cumulative_bp_pos_4338254	Public					-1
10347	pkP1054	CR_non_cumulative_bp_pos_4845516	Public					

Stable!

Marker 1-5 of 121

Probe 1-10 of 23232

Stepsize 5

Width 5

Height 10

Change size

	pkP1050	pkP1101	pkP1103	pkP1052	egPE107
WSU1	-0.1892	-0.1892	0.231	-0.8379	-0.9186
WSU2	-0.012	-0.012	0.1026	-0.2283	-0.4022
WSU3	0.0637	0.0637	0.2153	-0.1182	-0.1068
WSU4	0.0136	0.0136	0.1208	-0.1924	-0.1909
WSU5	0.054	0.054	0.1649	-0.1768	-0.1621
WSU6	0.0873	0.0873	0.1452	-0.0902	-0.0322
WSU7	-0.0529	-0.0529	-0.0248	0.0354	0.0405
WSU8	0.0629	0.0629	0.1488	-0.06	-0.1015
WSU9	0.0421	0.0421	0.0923	-0.254	-0.0614
WSU10	-0.0099	-0.0099	0.056	-0.0829	-0.0386

Dynamic!

Results

Observ-OM model

So far..

- MOLGENIS toolbox
 - Good at rapidly generating databases for **stable** data
 - Values comparable within columns
- XGAP model
 - Good at **flexible homogeneous** data sets
 - All values in data set are semantically equivalent
- *But how about..*
 - Capture **flexible heterogenous** data
 - Detailed meta-data (value type, unit, ontologies..)
 - Data lineage by protocol & application thereof

“Molecules”

“Phenotypes”

Repeat

Can we use 'entity-attribute-value' modeling?

Individual	Gender	Age	Height	Body Weight	Diastolic_Blood_Pressure	Systolic_Blood_Pressure	Cervical_Carcinoma
Individual_1	female	56	180	65	90	140	yes
Individual_2	female	45	178	75	87	130	no
Individual_3	male	65	168	100	78	125	no
Individual_4	male	35	178	45	100	150	no
Individual_5	male	34	190	55	68	134	yes
Individual_6	female	24	185	67	76	134	yes
Individual_7	female	20	179	80	102	145	yes

← The columns should be independent (more flexible)

↓ 1. model

Values NOT all comparable: Columns are NOT of the same type

```
<entity name="Locus" abstract="true">
  <description> position. Typical examples of such traits are genes,
    probes and markers. Common structure for entities that have a
    genomic</description>
  <field name="Chromosome" label="Chromosome" type="xref"
    xref_entity="Chromosome" xref_field="id" xref_label="name" nillable="false"
    description="Reference to the chromosome this
    position belongs to." />
  <field name="cM" label="cMPosition" type="decimal" nillable="true"
    description="genetic map position in centi morgan (cM)." />
  <field name="bpStart" label="Start (5')" type="long" nillable="true"
    description="numeric basepair position (5') on the chromosome" />
  <field name="bpEnd" label="End" type="long" nillable="true"
    description="numeric basepair position (3') on the chromosome" />
  <field name="Seq" type="text" nillable="true"
    description="The FASTA text representation of the sequence." />
  <field name="Symbol" type="varchar" nillable="true"
    description="todo" />
</entity>
<entity name="Chromosome" extends="ObservableFeature">
  <field name="orderNr" type="int" />
  <field name="isAutosomal" type="bool" description="Is 'yes' when number of chr
  <field name="bpLength" type="int" nillable="true" description="Lenght of the c
  <field name="Species" label="Species" type="xref" xref_entity="Species"
    xref_field="id" xref_label="name" nillable="true"
    description="Reference to the species this
    chromosome belongs to." />
</entity>
```

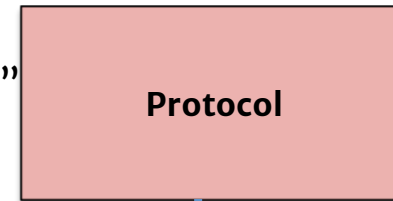


Features become columns

- Feature \sim columns, protocols \sim tables

Protocol

Measuring “Gender”, “Age”, “Height”, “Blood pressure”
in some medically standardized way (ISO-54532)



Feature

“Gender”
Categr.
M/F/O

Feature

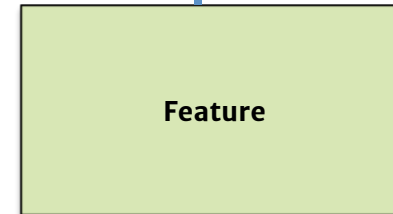
“Age”
Years
Integer

Feature

“Height”
Centimeter
Decimal

Feature

“Blood pressure”
mmHg
Decimal

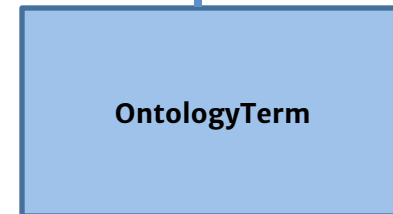


Ontology

MeSH term
“Gender Identity”

Ontology

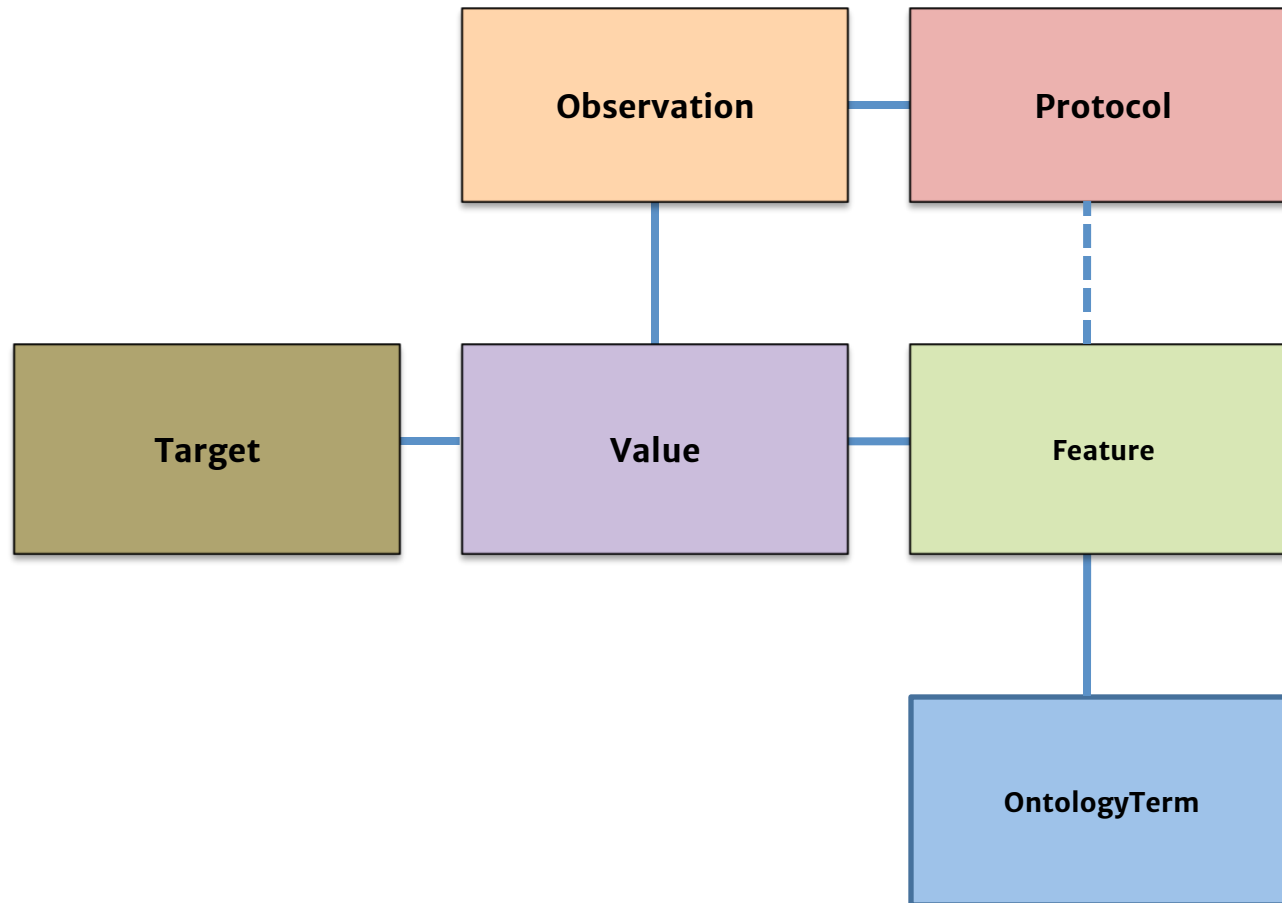
MeSH term
“Blood pressure”



Adamusiak *et al.* (2012) *Human Mutation* 33(5):867-73

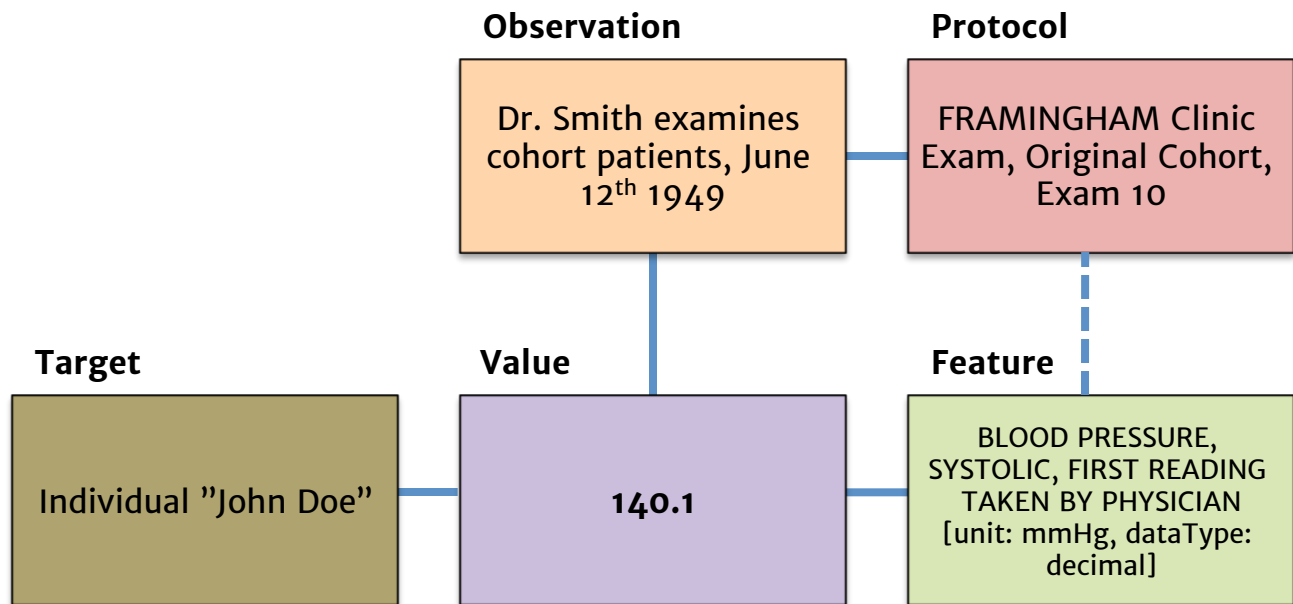
Observations become table rows

- Targets are patients, samples, groups, etc



Adamusiak *et al.* (2012) *Human Mutation* 33(5):867-73

Data example



	Negative	Positive	Doubtful	Unknown	
31	0	1	2	9	SUGAR IN URINE FC10
32	0	1	2	9	ALBUMIN IN URINE FC11
BLOOD PRESSURE (Left arm, mm Hg):					
33-38	Systolic FC12		Diastolic FC13		NURSE
44	FC14		FC15		PHYSICIAN (First reading)
45-50	FC16		FC17		PHYSICIAN (Second reading)
LUNG FUNCTION:					
51-52	FC188		TOTAL VITAL CAPACITY (Deciliters)		
53-55	FC189		FIRST SECOND VOLUME (Centiliters)		
GLUCOSE CHALLENGE:					

Ontology

MeSH: Blood Pressure

"PRESSURE of the BLOOD on the ARTERIES and other BLOOD VESSELS."

Adamusiak et al. (2012) *Human Mutation* 33(5):867-73

Individual	Gender	Age	Height	Body Weight	Diastolic_Blood_Pressure	Systolic_Blood_Pressure	Cervical_Carcinoma	Breast_Carcinoma
Individual_1	female	56	180	65	90	140	yes	no
Individual_2	female	45	178	75	87	130	no	yes
Individual_3	male	65	168	100	78	125	no	yes
Individual_4	male	35	178	45	100	150	no	yes
Individual_5	male	34	190	55	68	134	yes	yes
Individual_6	female	24	185	67	76	134	yes	no
Individual_7	female	20	179	80	102	145	yes	no
Individual_8	female	34	175	56	76	134	yes	yes
Individual_9	female	45	181	67	90	144	yes	yes
Individual_10	male	34	160	45	86	132	yes	yes
Individual_11	male	35	194	75	70	123	yes	yes
Individual_12	female	54	182	55	69	120	yes	yes
Individual_13	male	33	170	66	65	121	yes	yes
Individual_14	female	24	180	77	66	122	yes	yes
Individual_15	female	27	170	56	64	125	yes	yes
Individual_16	female	36	167	54	63	127	yes	yes
Individual_17	male	35	186	76	70	131	yes	yes
Individual_18	male	43	190	88	59	119	yes	yes
Individual_19	male	44	168	65	60	120	yes	yes
Individual_20	female	36	175	54	63	121	yes	yes

Ontologies

MeSH: Blood Pressure

"PRESSURE of the BLOOD on the ARTERIES and other BLOOD VESSELS."

Protocol

Gender ◉ Age ◉ Height ◉ Body Weight ◉ Diastolic_Blood_Pressure ◉ Systolic_Blood_Pressure ◉ Cervical_Carcinoma ◉ Breast_Carcinoma ◉

Observation

Features

Individual ◉	Gender ◉	Age ◉	Height ◉	Body Weight ◉	Diastolic_Blood_Pressure ◉	Systolic_Blood_Pressure ◉	Cervical_Carcinoma ◉	Breast_Carcinoma ◉
Individual_1	female	56	180	65	90	140	yes	no
Individual_2	female	45	178	75	87	130	no	yes
Individual_3	male	65	168	100	78	125	no	yes
Individual_4	male	35	178	45	100	150	no	yes
Individual_5	male	34	190	55	68	134	yes	yes
Individual_6	female	24	185	67	76	134	yes	no
Individual_7	female	20	179	80	102	145	yes	no
Individual_8	female	34	175	56	76	134	yes	yes
Individual_9	female	45	181	67	90	144	yes	yes
Individual_10	male	34	160	45	86	132	yes	yes
Individual_11	male	35	194	75	70	123	yes	yes
Individual_12	female	54	182	55	69	120	yes	yes
Individual_13	male	33	170	66	65	121	yes	yes
Individual_14	female	24	180	77	66	122	yes	yes
Individual_15	female	27	170	56	64	125	yes	yes
Individual_16	female	36	167	54	63	127	yes	yes
Individual_17	male	35	186	76	70	131	yes	yes
Individual_18	male	43	190	88	59	119	yes	yes
Individual_19	male	44	168	65	60	120	yes	yes
Individual_20	female	36	175	54	63	121	yes	yes

Targets

Values

Building G2P applications

MOLGENIS software

Swertz *et al*,
BMC Bioinf. (2010)

<http://www.molgenis.org>

XGAP model

Swertz *et al*,
Genome Biology (2010)

<http://www.xgap.org>

Observ-OM model

Adamusiak *et al*,
Human Mutation (2012)

<http://www.observe-om.org>

EB Registry

Van den Akker *et al*,
Human Mutation (2011)

<http://www.deb-central.org>


dystrophic eb

xQTL workbench

Arends & van der Velde *et al*,
Bioinformatics (2012)

<http://www.xqtl.org>

XQTL WORKBENCH

AnimalDB

Track and trace of animal life
events in research laboratories

<http://www.animaldb.org>



WormQTL

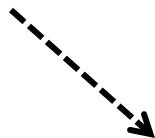
- Panacea project, *C. elegans* data
 - ~300 million measurements
- Snoek, van der Velde, Arends & Li *et al*,
Nucl. Acids Res. (2013)

<http://www.wormqtl.org>

CropQTL

Learning From Nature project,
arabidopsis thaliana data

- 1400 plants
- SNP genotypes (~70 million values)
- Classical traits, e.g. flowering time


..more



Showcase EB Registry

EB Registry: Dystr. EB mutation database



International dystrophic eb Patient Registry

International registry of patients with dystrophic epidermolysis bullosa and database of associated COL7A1 mutations

[Search](#) [Submit data](#) [Contact](#) [References](#) [Background](#) [News](#) [Login](#)

Search

Welcome to the international registry of dystrophic epidermolysis bullosa (DEB) patients and associated COL7A1 mutations.

The International Dystrophic Epidermolysis Bullosa Patient Registry contains anonymised data on both published and unpublished DEB patients, as well as their associated COL7A1 mutations and genotypes, and clinical and molecular phenotypes.

The database currently contains 590 DEB patients, of which 71 unpublished, and 395 COL7A1 mutations. Search or browse below.

Search registry

Search by typing any search term in the search field, like cDNA (e.g. "3G>T") or protein (e.g. "Arg525Ter") notations of mutations, mode of inheritance (e.g. "dominant") or specific phenotypes (e.g. "severe generalized"). Search results are shown at bottom of page.

Enter search term: [Advanced Search](#)

Show mutations

Show patients

Browse the COL7A1 gene

Click anywhere on this schematic representation of the COL7A1 gene to graphically browse the gene. With every click you will zoom in deeper on the COL7A1 gene. Mutated nucleotides are depicted in red. If the cursor is placed over the mutated nucleotide(s), the corresponding mutation is shown.

News

Patient and Mutation Update
14 Australasian patients added to Registry.

Jul 20, 2011

[More](#)

Article about DEB Registry online in Human Mutation.

Go to Pubmed to [view article](#)

Jun 16, 2011

[More](#)

Feature Update

Select type of search result tables.

Jun 9, 2011

[More](#)

Name change: International Dystrophic Epidermolysis Bullosa Patient Registry
Important update.

May 25, 2011

[More](#)

Patient and Mutation Update

22 Tunisian RDEB patients entered to database.

May 24, 2011

[More](#)

EB Registry: Dystr. EB mutation database

Phenotypic details for patient 'P10'

target

Characteristics

Age	20
Gender	m
Ethnicity	unknown
Deceased	yes
Cause of death	
MMP1 allele 1	
MMP1 allele 2	

value

feature

Protocol Application

Cutaneous

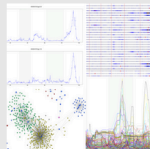
Blistering	yes
Location	generalized
Hands	unknown
Feet	unknown
Arms	unknown
Legs	unknown
Proximal body flexures	unknown
Trunk	unknown
Mucous membranes	yes
Skin atrophy	yes
Milia	unknown
Nail dystrophy	yes

Showcase WormQTL



WormQTL – Public archive and analysis web portal for natural variation data in *Caenorhabditis* spp.

WormQTL is an online scalable system for QTL exploration to service the worm community. WormQTL provides many publicly available datasets and welcomes submissions from other worm researchers.



[Find QTLs](#)



[Genome browser](#)

id	chr	start	end	score	trait
NS01	1	11892	12131	0.8379	0.5186
NS02	1	10102	10300	0.2208	0.4902
NS03	1	10637	10753	0.1182	0.1988
NS04	1	10100	10200	0.1024	0.1989
NS05	1	10541	10640	0.1768	0.1812
NS06	1	10870	10970	0.0902	0.0822
NS07	1	10528	10628	0.0248	0.0354
NS08	1	10620	10720	0.06	0.1010
NS09	1	10421	10521	0.234	0.0654
NS10	1	10300	10400	0.0824	0.0824

[Browse data](#)



[Help](#)

What can you do?

- I want to search (e)QTLs for my trait or gene
 1. Go to [Find QTLs](#)
 2. Type the name or identifier of your trait or gene and press *Search*
 3. Put any relevant hits in the shopping cart
 4. Click *Plot cart* now and explore the results
- I want to know which genes have a QTL on my favourite position
 1. Go to [Genome browser](#)
 2. Add tracks from experiments of interest
 3. Navigate to your favourite location (tip: use *open in new window*)
 4. Collect significant probe identifiers from that region
 5. Use the identifiers to do a search with [Find QTLs](#)

Phenotypes	Type of array	Sample size	Parental strains	Reference	Pubmed link	Growing temperature	Stage	Food	Medium	Dataset IDs
Gene expression	Washington State University	2x40 RILs	CB4856; N2	Li et al. 2006; Mapping determinants of gene expression plasticity by genetical genomics in <i>C. elegans</i> .	17196041	16oC and 24oC	(72h at 16 and 40h at 24); L4	OP50	NGM Plate	37 , 38
Gene expression	Affymatrix tiling array	60 RILs	CB4856; N2	Li et al. 2010; Global genetic robustness of the alternative splicing machinery in <i>Caenorhabditis elegans</i> .	20610403	24oC	(40h) L4	OP50	NGM Plate	n/a
Gene expression	Washington State University	36x3 RILs	CB4856; N2	Vinuela & Snoek et al. 2010; Genome-wide gene expression regulation as a function of genotype and age in <i>C. elegans</i> .	20488933	24oC	(40h, 96h and 214h) L4, Adult, Old	OP50	NGM Plate	3 , 5 , 6 , 7 , 8 , 9 , 10 , 11 , 12 , 13 , 14 , 15 , 16 , 17 , 18 , 19 , 20 , 21
Gene expression	Agilent 4x44k microarrays	208 RIALs	CB4856; N2	Rockman et al. 2010; Selection at linked sites shapes heritable phenotypic variation in <i>C. elegans</i> .	20947766	20oC	YA	OP50	NGM Plate	22 , 34 , 35 , 36
Feeding curves RNAi exposure	n/a	56 RILs * 12 RNAi	CB4856; N2	Elvin & Snoek et al. 2011; A fitness assay for comparing RNAi effects across multiple <i>C. elegans</i> genotypes.	22004469	20oC	Multi-generational	n/a	Liquid S-medium	24 , 32 , 33
Life-history traits	n/a	80 RILs	CB4856; N2	Gutteling et al. 2007; Mapping phenotypic plasticity and genotype-environment interactions affecting life-history traits in <i>Caenorhabditis elegans</i> .	16955112	12oC and 24oC	Egg, L4, YA	OP50	NGM Plate	25 , 26 , 27
Lifespan and pharyngeal-pumping	n/a	90 NILs	CB4856; N2	Doroszuk et al. 2009; A genome-wide library of CB4856/N2 introgression lines of <i>Caenorhabditis elegans</i> .	19542186	20oC	All; synchronised	OP50	NGM Plate	4 , 23 , 28 , 29 , 30 , 31
Lifespan, Recovery and reproduction after heat-shock	n/a	58 RILs	CB4856; N2	Rodriguez et al. 2012; Genetic variation for stress-response hormesis in <i>C. elegans</i> lifespan.	22613270	20oC and 35oC heat-shock	L4 and Adult	OP50	NGM Plate	39 , 40
Gene expression	Washington State University	CB4856 and N2	CB4856; N2	Vinuela & Snoek et al. 2012; Aging Uncouples Heritability and Expression-QTL in <i>Caenorhabditis elegans</i> .	22670229	24oC	(40h, 96h and 214h) L4, Adult, Old	OP50	NGM Plate	41 , 42 , 43

Find QTLs

- All data (175,366)
- measurement (43)
- Panel (500)
- Gene (47,360)
- Transcript (55,782)
- Chromosome (8)
- Probe (68,452)
- Sample (1,630)
- DerivedTrait (12)

seam cell

(for [ontology](#) or anatomy terms, will show the probes and related terms for that gene.

View

(9)

Found n

Your results were limited to the first 100. Please be more specific.


Probe [AGIUSA14764 / clc-2](#) reports for [WBGene00000523 - WormBase](#)
C01C10.1 / C01C10.1 / wb|C01C10.1 / non_cumu_bp_start_743339 [...more](#)

Probe [AGIUSA16119 / ceh-1](#) reports for [WBGene00000428 - WormBase](#)
F16H11.4 / F16H11.4 / wb|F16H11.4 / non_cumu_bp_start_465353 [...more](#)


Probe [AGIUSA19594 / acn-1](#) reports for [WBGene00000039 - WormBase](#)
peptidase [C42D8.5.2] / C42D8.5.2 / C42D8.5 / wb|C42D8.5.2|w [...more](#)

Probe [AGIUSA41433 / gsp-1](#) reports for [WBGene00001747 - WormBase](#)
serine/threonine protein phosphatase [F29F11.6.1] / F29F11.6 [...more](#)

Probe [AGIUSA5476 / cul-2](#) reports for [WBGene00000837 - WormBase](#)


 **Ontological terms**


- GO:0016021-integral to membrane
- GO:0005198-structural molecule activity
- GO:0005923-tight junction

 **Ontologies**

- WBbt:0005733-hypodermis
- WBbt:0005753-seam cell

 **Ontologies**

 **Ontologies**

Home Find QTLs Genome browser Browse data Help provide feedback: 

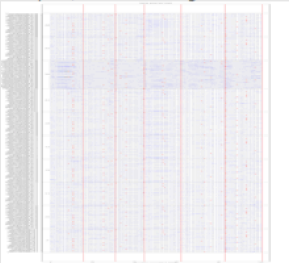
Find QTLs

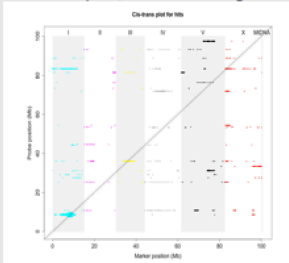
All data (175,366) daf

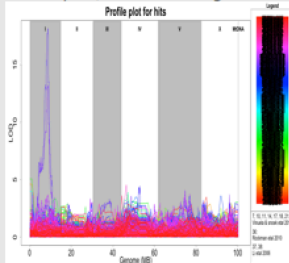
(for example: *ctl*, *daf*, *pqp-7*, *gst-27*, *Y65B4BR*, *K02B12*, *WBGene00021562*, *WBGene00006727*, *acetylcholine*, *luciferase* ...)
Gene hits, for example on [Geno Ontology](#) or anatomy terms, will show the probes and related terms for that gene.

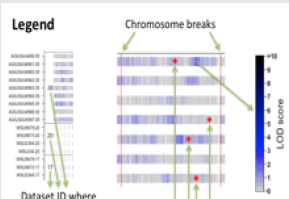
Results for my selected hits:

(get a [permanent link to these results](#))

Heatplot, click to enlarge: 

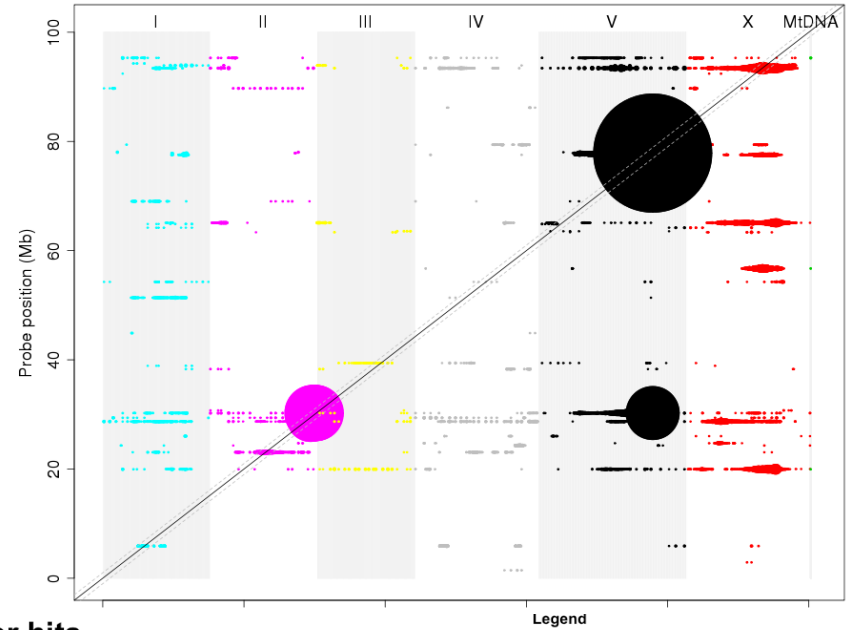
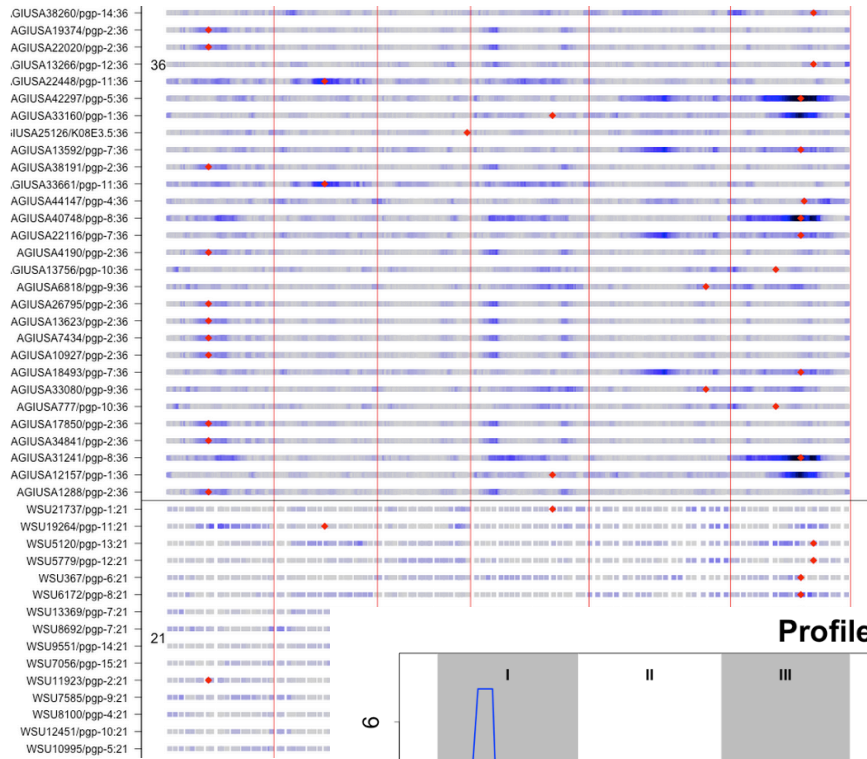
Cis-trans plot, click to enlarge: 

Profile plot, click to enlarge: 

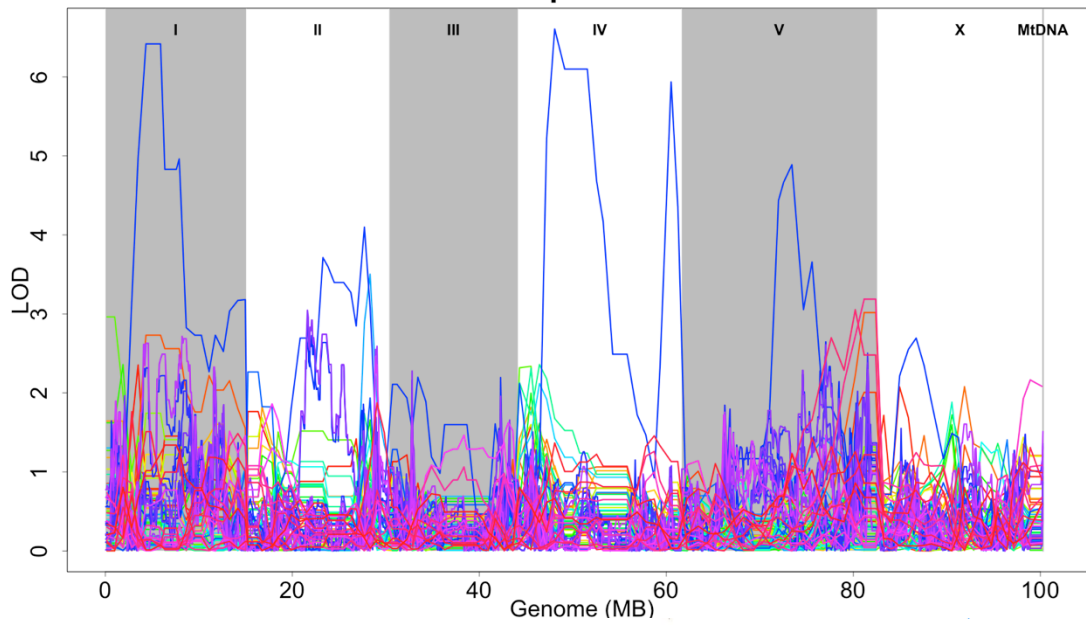
Legend, click to enlarge: 

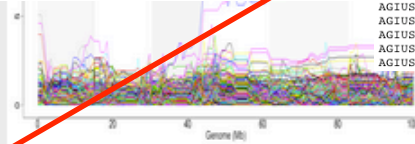
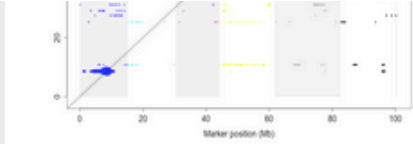
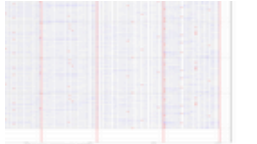
More downloads:

- Get the [Cytoscape network](#) for this plot. ([how-to import](#))
- Get the [Cytoscape nodes](#) for this plot. ([how-to import](#))
- Note: includes **significant results only**. (LOD > 3.5)
- Save both files. Import network (has LOD scores), then node attributes (*chrom*, *bploc*, *dataset*). [Example visualization](#)
- Get the generated [source data](#) for these plots.
- Get the generated [multiplot plot R script](#).
- Get the generated [cistrans R plot script](#).
- Get the generated [profile R plot script](#).



Profile plot for hits

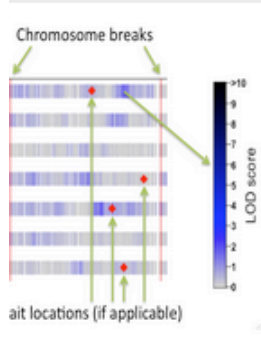




```

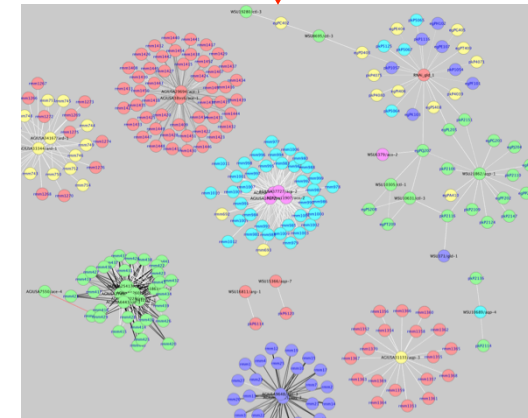
WSU6355/daf-3   QTL   pkP1101 4.6026
WSU6355/daf-3   QTL   pkP1103 3.5561
WSU6355/daf-3   QTL   pkP1050 5.0597
WSU6355/daf-3   QTL   pkP1101 5.0597
WSU7870/daf-16  QTL   egPD403 3.9109
WSU6355/daf-3   QTL   pkP1050 3.9762
WSU6355/daf-3   QTL   pkP1101 3.9762
WSU3844/daf-3   QTL   egPK601 3.6122
WSU7870/daf-16  QTL   pkP5071 3.5478
AGIUSA25467/daf-36 QTL   rmm743 3.893693484
AGIUSA25467/daf-36 QTL   rmm744 4.302316576
AGIUSA25467/daf-36 QTL   rmm745 4.302316576
AGIUSA25467/daf-36 QTL   rmm746 4.302316576
AGIUSA25467/daf-36 QTL   rmm747 4.302316576
AGIUSA25467/daf-36 QTL   rmm748 4.302316576
AGIUSA25467/daf-36 QTL   rmm749 4.302316576
AGIUSA25467/daf-36 QTL   rmm750 3.824765495
    
```

Click to enlarge:



More downloads:

- Get the [Cytoscape network](#) for this plot. ([how-to import](#))
- Get the [Cytoscape nodes](#) for this plot. ([how-to import](#))
- Note: includes **significant results only**. (LOD > 3.5)
- Save both files. Import network (has LOD scores), then node attributes (chrom, bploc, dataset). [Example visualization](#)
- Get the generated [source data](#) for these plots.
- Get the generated [multiplet plot R script](#).
- Get the generated [cistrans R plot script](#).
- Get the generated [profile R plot script](#).



[/ daf-1 \[explore deeper\]](#) - protein kinase [F29C4.1b] / F29C4.1b / F...
[/ daf-3 \[explore deeper\]](#) - F25E2.5b.3 / F25E2.5 / wb|F25E2.5b.3|wb|...
[af-11 \[explore deeper\]](#) - R0240.3 / cea2 p.107079 / blast match 60

```

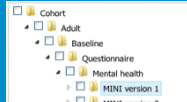
WSU6355/daf-3   trait   X   83377539   age12_int_qtl
pkP1050 marker   I   169018   age12_int_qtl
pkP1101 marker   I   992189   age12_int_qtl
pkP1103 marker   I   1881116  age12_int_qtl
WSU7870/daf-16  trait   I   10763660  age23_int_qtl
egPD403 marker  IV  47198842  age23_int_qtl
WSU3844/daf-3   trait   X   83371790  age3_qtl
egPK601 marker  X   91884289  age3_qtl
pkP5071 marker  V   76788020  age3_qtl
AGIUSA25467/daf-36 trait   V   71855469  rock_qtl
rmm743 marker   IV  52297746  rock_qtl
rmm744 marker   IV  52364173  rock_qtl
rmm745 marker   IV  52422112  rock_qtl
rmm746 marker   IV  52482040  rock_qtl
rmm747 marker   IV  52569872  rock_qtl
rmm748 marker   IV  52584390  rock_qtl
rmm749 marker   IV  52626947  rock_qtl
rmm750 marker   IV  52705021  rock_qtl
rmm751 marker   IV  52722889  rock_qtl
rmm752 marker   IV  52809758  rock_qtl
    
```

Current work

OmicsConnect running on Observ-OMX

Catalogue

Find data item and sample collections



Data

Filter individual data sets and download to Excel & SPSS

	rs11050	rs11051	rs11052	rs11053	rs11054
WSU1	-0.1892	-0.1892	0.2131	-0.8379	-0.9180
WSU2	0.0027	0.0027	0.0061	0.0298	0.0028
WSU3	0.0637	0.0637	0.2153	-0.1182	-0.1048
WSU4	0.0316	0.0316	0.1208	-0.1924	-0.1909
WSU5	0.0514	0.0514	0.1649	-0.1768	-0.1621
WSU6	0.0637	0.0637	0.0502	0.0969	0.0928
WSU7	-0.0529	-0.0529	-0.0248	0.0354	0.0405
WSU8	0.0637	0.0637	0.1649	0.0969	0.0928

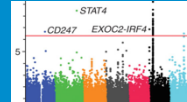
Compute

Run analysis workflows on big data compute infrastructure



GWAS Central

Explore summary level GWAS data

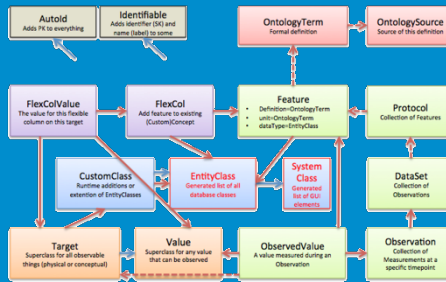


Protocol

CRFs, Questionnaires, Lab protocols, and assays



Core model



xQTL

Multi-omics association & visualization tools



NGS

Next-Generation Sequencing



XGAP

Multi-omics genotypes and phenotypes



Share

Friends, Groups and Permission management



Mutation

Explore genetic mutations and pathogenicity effects



Organization

Institutes, Departments, People, Locations & Containers



File

File storage and drivers for images and data

5-12	5P	83	08	47	14	10	00
6-1E	73	D9	02	6	GD	36	11
4-20	A3	01	76	a	PP	0	1
E-92	D8	AB	F4	b	B3	-	1
8-9B	DC	99	62	w	30	1	1
5-67	99	18	EC	o	11	92	0
8-5F	F4	CA	CB	4	94	1	1
8-91	66	65	E6	-	1	9	1
8-5F	F0	82	EC	R	6	1	1

Acknowledgements

Martijn Dijkstra
 Despoina Antonakaki
 Tomasz Adamuziak
 Rob Hastings
 Sirisha Gollapudi
 Gudmundur Thorisson
 Chao Pang
 Myles Byrne
 David van Enckvort
 Linda Mook
 Pieter Neerincx
 Ger Strikwerda
 Danny Arends
 Roan Kanninga
 Jan Bot
 George Byelas
 Yang Li
 Basten Snoek
 Noortje Festen
 Konrad Zych

And more ...

Lude Franke
 Juha Muilu
 Anthony Brookes
 Helen Parkinson
 Vincent Ferreti
 Gert-Jan van Ommen
 Jan Jurjen Uitterdijk
 Ritsert C. Jansen
 Jan Kammenga
 Cisca Wijmenga
 Paul de Bakker
 Irene Nooren
 Rob Hooft
 Salome Scholtens
 Hans Hillege
 Ronald Stolk
 Morris Swertz

And more ...

NBIC/BioAssist consortium (bioinfo)
 BBMRI-NL catalogue group(Hs)
 CTMM/TraIT consortium (Hs)
 EU-GEN2PHEN consortium (Hs)
 EU-PANACEA consortium (Ce)
 EU-BioSHARE consortitum (Hs)
 EU-CASIMIR consortium (Mm)
 EU-BioMedBridges cosortium (all)
 NL Brassica Nutr. consortium (At)
 Learning from Nature (At)
 LifeLines (Hs)
 TIFN (Hs)
 BigGrid (info)
 Target + CIT (info)

And more...



LIFELINES



TIFOOD
NUTRITION

B B M R I • N L

BiG Grid
the dutch e-science grid



Wrap-up

Summary

- Complex variation, geno-to-pheno
- Exploiting the data requires structure
- Best-of flexible and stable parts
- Support homo- and heterogeneous data

Read more

- MOLGENIS: <http://www.molgenis.org>
- MOLGENIS Compute: <http://www.molgenis.org/wiki/ComputeStart>
- xQTL: <http://www.xqtl.org>
- Adamusiak *et al* (2011) *BMC Bioinformatics*
- Akker *et al* (2011) *Human Mutation*
- Arends *et al* (2010) *Bioinformatics* 26: 2990-2992
- Brandsma *et al*, *Norsk Epidemiologi* 2012
- Snoeks *et al* (2013) *Nucleic Acids Res*
- Swertz *et al* (2010) *Genome Biology* 9;11(3): R27.
- Smedley *et al* (2008) *Briefings in bioinformatics* 9(6):532-44.
- Swertz & Jansen (2007) *Nature Reviews Genetics* 8, 235-243

Thank you!
Questions?

k.j.van.der.velde@umcg.nl

molgenis
Your database at the push of a button