

# data model experiences from various EU projects MAGETAB, FuGE, XGAP, Pheno-OM

Metadata Capture Symposium  
November 10, Utrecht

Morris Swertz and biodata modelers around the world\*

BBMRI-NL, EU-GEN2PHEN, EU-CASIMIR, EU-EURATRANS, LifeLines, EU-SYSGENET, EU-PANACEA, NBIC and other consortia

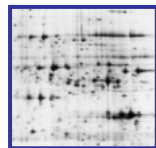
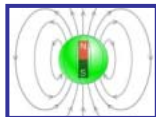
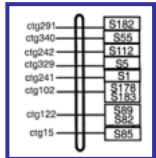
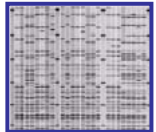
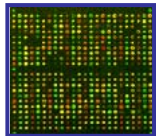
*\*Genomics Coordination Center Groningen*



# Genomics Coordination Center – UMC/University Groningen



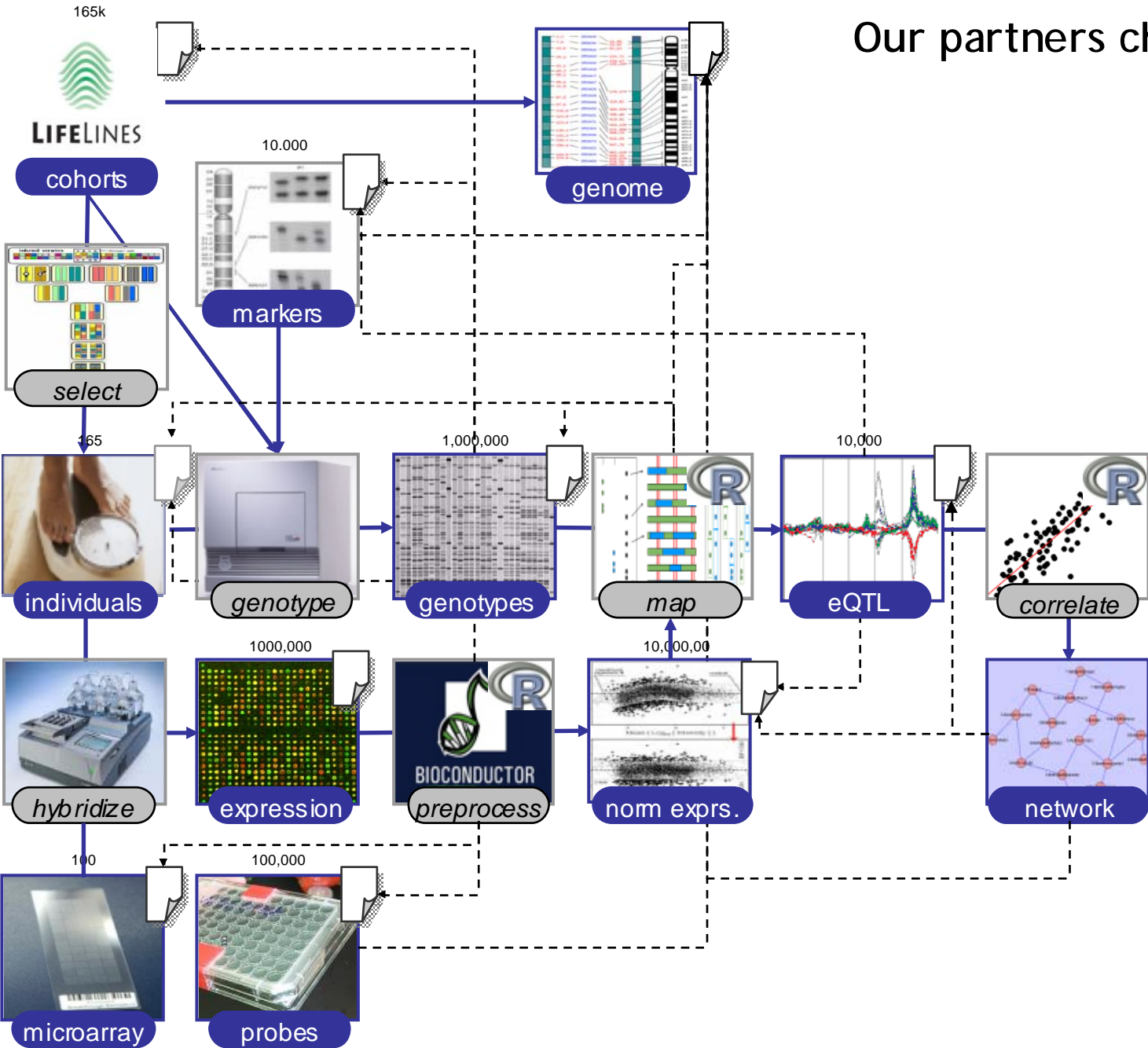
*etc.*



*etc.*



# Our partners challenges



# Some models

- MAGE-TAB for microarrays
  - Oldest model, born from experience microarray community
- Pheno-OM for phenotypes
  - Recent model, but present now as it shows core concepts
  - Born from XGAP and PaGE
- XGAP for all phenotypes/genotypes
  - Extended version of Pheno-OM and FuGE
- FuGE for functional genomics experiments
  - Exemplar for its protocol – protocolApplication structure
- Others like ISA-TAB (aka MAGE-TAB++), CHADO, HL7, endless list

=> I will show convergence of these models

# MAGE-TAB (format)

Born out of failed MAGE-OM (100+ pages spec)

Originally for microarrays

Contents:

- IDF Investigation description file
- SDRF Sample-Data-Relationship file
- Data files: raw format or matrix files
- Optional: array design file(s)

Widely used in many production systems!

<http://www.mged.org/mage-tab/>

<http://www.ebi.ac.uk/microarray/doc/help/MAGE-TAB.html>

Rayner et al (2006) BMC Bioinformatics 7, 489

# IDF

<b>Investigation Title</b>	University of Heidelberg H sapiens TK6		
<b>Experimental Design</b>	genetic_modification_design	time_series_design	
<b>Experimental Factor Name</b>	Genetic Modification	Incubation Time	
<b>Experimental Factor Type</b>	genetic_modification	time	
<b>Experimental Factor Term Source REF</b>	MGED Ontology	MGED Ontology	
<b>Person Last Name</b>	Maier	Fleckenstein	Li
<b>Person First Name</b>	Patrick	Katharina	Li
<b>Person Email</b>	patrick.maier@radonk.ma.uni-heidelberg.de		
<b>Person Phone</b>	+496213833773		
<b>Person Address</b>	Theodor-Kutzer-Ufer 1-3		
<b>Person Affiliation</b>	Department of Radiation Oncology, University of Heidelberg		
<b>Person Roles</b>	submitter; investigator	investigator	investigator
<b>Person Roles Term Source REF</b>	MGED Ontology	MGED Ontology	MGED Ontology
<b>Quality Control Type</b>	biological_replicate		
<b>Quality Control Term Source REF</b>	MGED Ontology		
<b>Replicate Type</b>	biological_replicate		
<b>Replicate Term Source REF</b>	MGED Ontology		
<b>Date of Experiment</b>	2005-02-28		
<b>Public Release Date</b>	2006-01-03		
<b>PubMed ID</b>	12345678		
<b>Publication Author List</b>	Patrick Maier; Katharina Fleckenstein; Li Li; Stephanie Laufs; Jens Zeller; Stefan Fruehauf; Carsten Herskind; Frederik Wenz		
<b>Publication Status</b>	submitted		
<b>Experiment Description</b>	Gene expression of TK6 cells transduced with an oncoretrovirus expressing MDR1 (TK6MDR1) was compared to untransduced TK6 cells and to TK6 cell transduced with an oncoretrovirus expressing the Neomycin resistance gene (TK6neo). Two biological replicates of each were generated and the expression profiles were determined using Affymetrix Human Genome U133 Plus2.0 GeneChip microarrays. Comparisons between the sample groups allow the identification of genes with expression dependent on the MDR1 overexpression.		
<b>Protocol Name</b>	GROWTHPRCL10653		TRANPRCL10656
<b>Protocol Type</b>	grow	action	bioassay_data_transformation
<b>Protocol Description</b>	TK6 cells were grown in suspension in RPMI 1640 medium supplemented with 10% fetal calf serum (Invitrogen, Karlsruhe, Germany). The cells were routinely maintained at 37 C and 5% CO2.	TK6 cells were grown in suspension in RPMI 1640 medium supplemented with 10% fetal calf serum (Invitrogen, Karlsruhe, Germany). The cells were routinely maintained at 37 C and 5% CO2.	Mixed Model Normalization with SAS Micro Array Solutions (version 1.3).
<b>Protocol Parameters</b>	media; time	Extracted Product; Affymetrix	
<b>Protocol Term Source REF</b>	MGED Ontology	MGED Ontology	
<b>SDRF File</b>	e-mexp-428_tab.txt		
<b>Term Source Name</b>	Cell Type Ontology	MGED Ontology	NCI Metathesaurus
<b>Term Source File</b>	http://obo.sourceforge.net/cgi-bin/detail.cgi?cell	http://mged.sourceforge.net/ontologies/MGEDontology.php	http://ncimeta.nci.nih.gov/indexMetaphrase.html
<b>Term Source Version</b>		1.3.0.1	

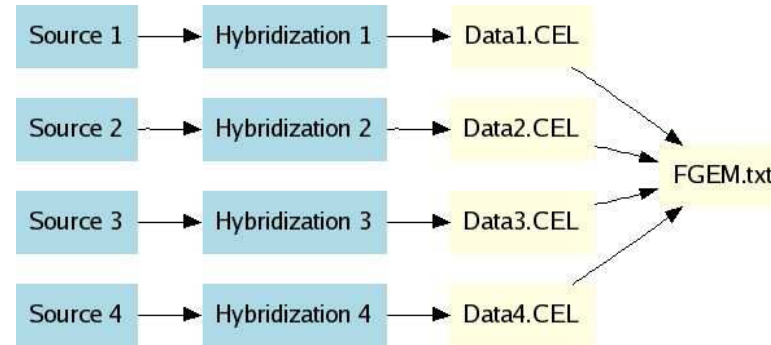
Repeating columns

Separated Values

Ontology Sources

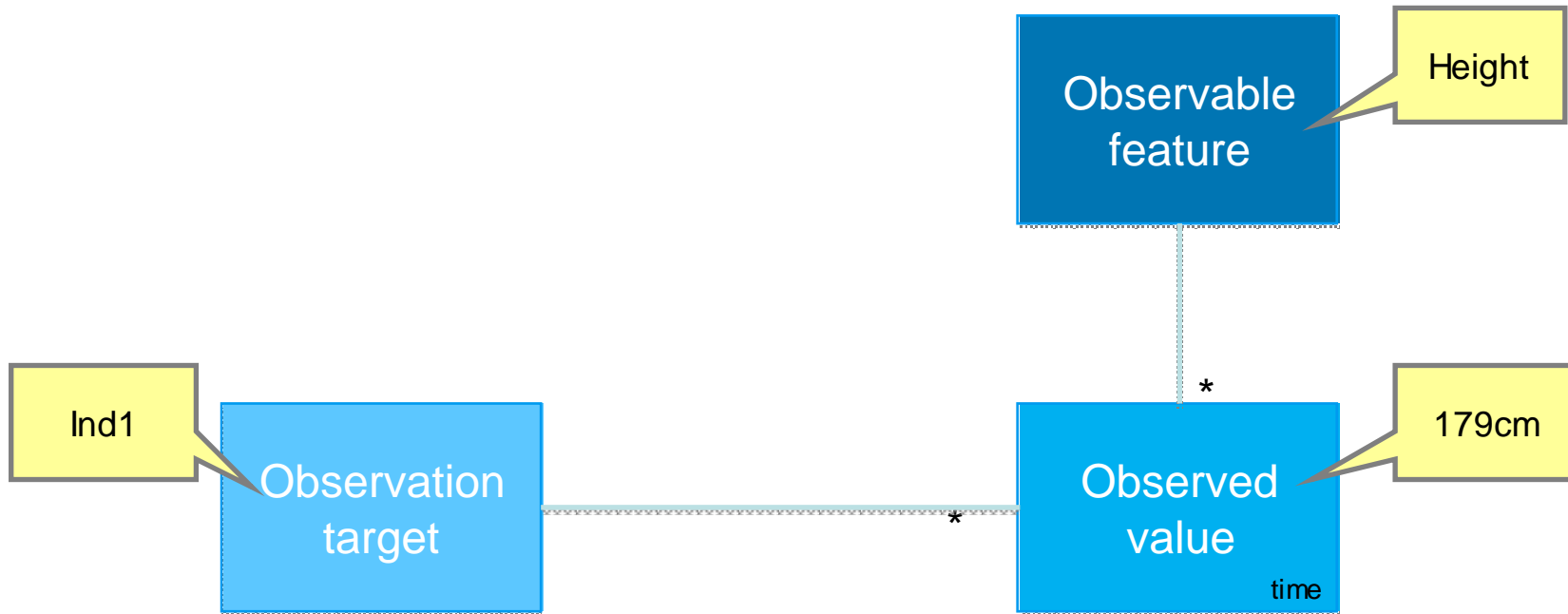


# SDRF



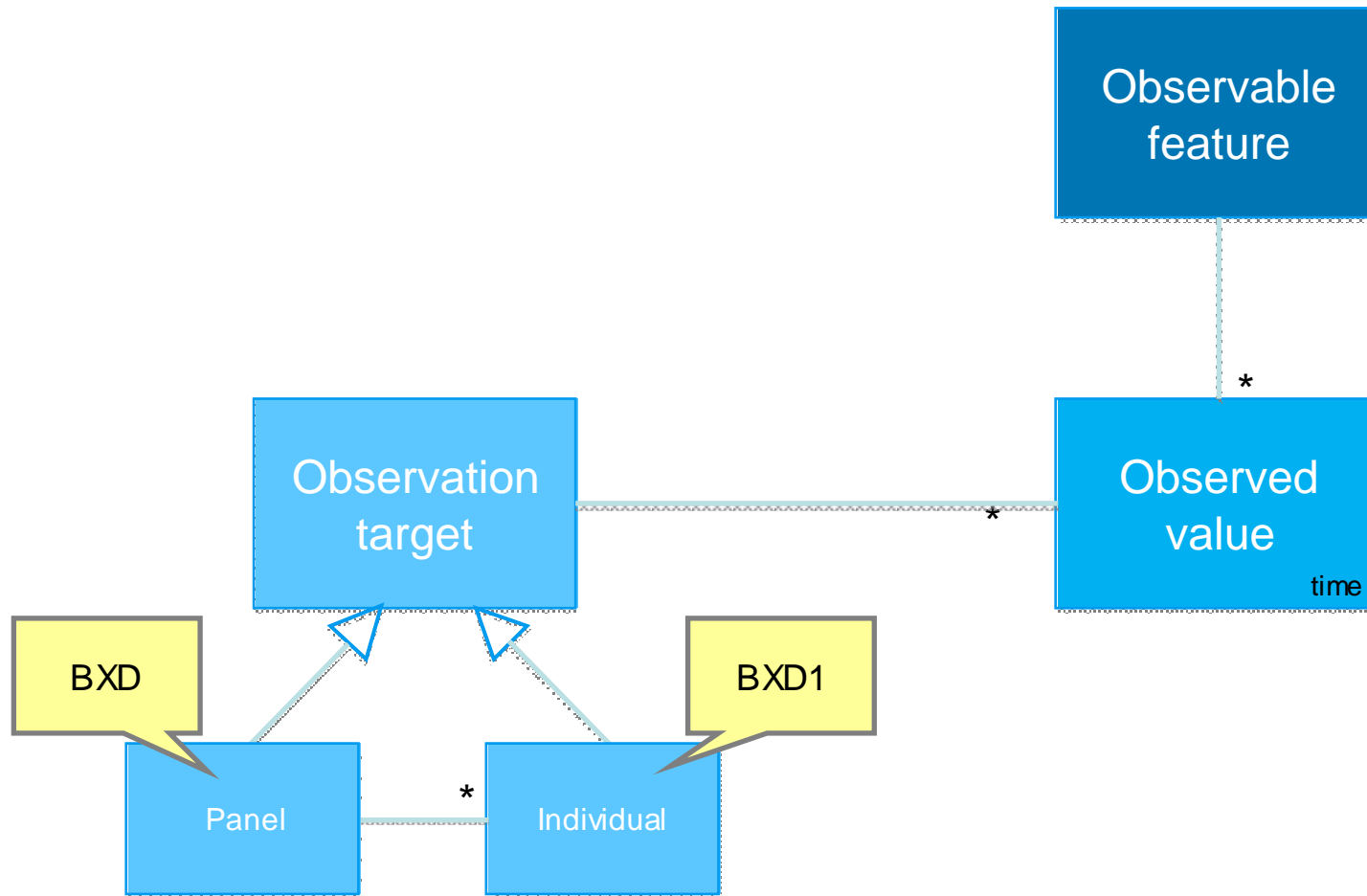
<b>Source Name</b>	<b>Protocol REF</b>	<b>Hybridization Name</b>	<b>Array Data File</b>	<b>Derived Array Data Matrix File</b>
Source 1	P-XMPL-10	Hybridization 1	Data1.CEL	FGEM.txt
Source 2	P-XMPL-10	Hybridization 2	Data2.CEL	FGEM.txt
Source 3	P-XMPL-10	Hybridization 3	Data3.CEL	FGEM.txt
Source 4	P-XMPL-10	Hybridization 4	Data4.CEL	FGEM.txt

# Pheno-OM ... for phenotypes

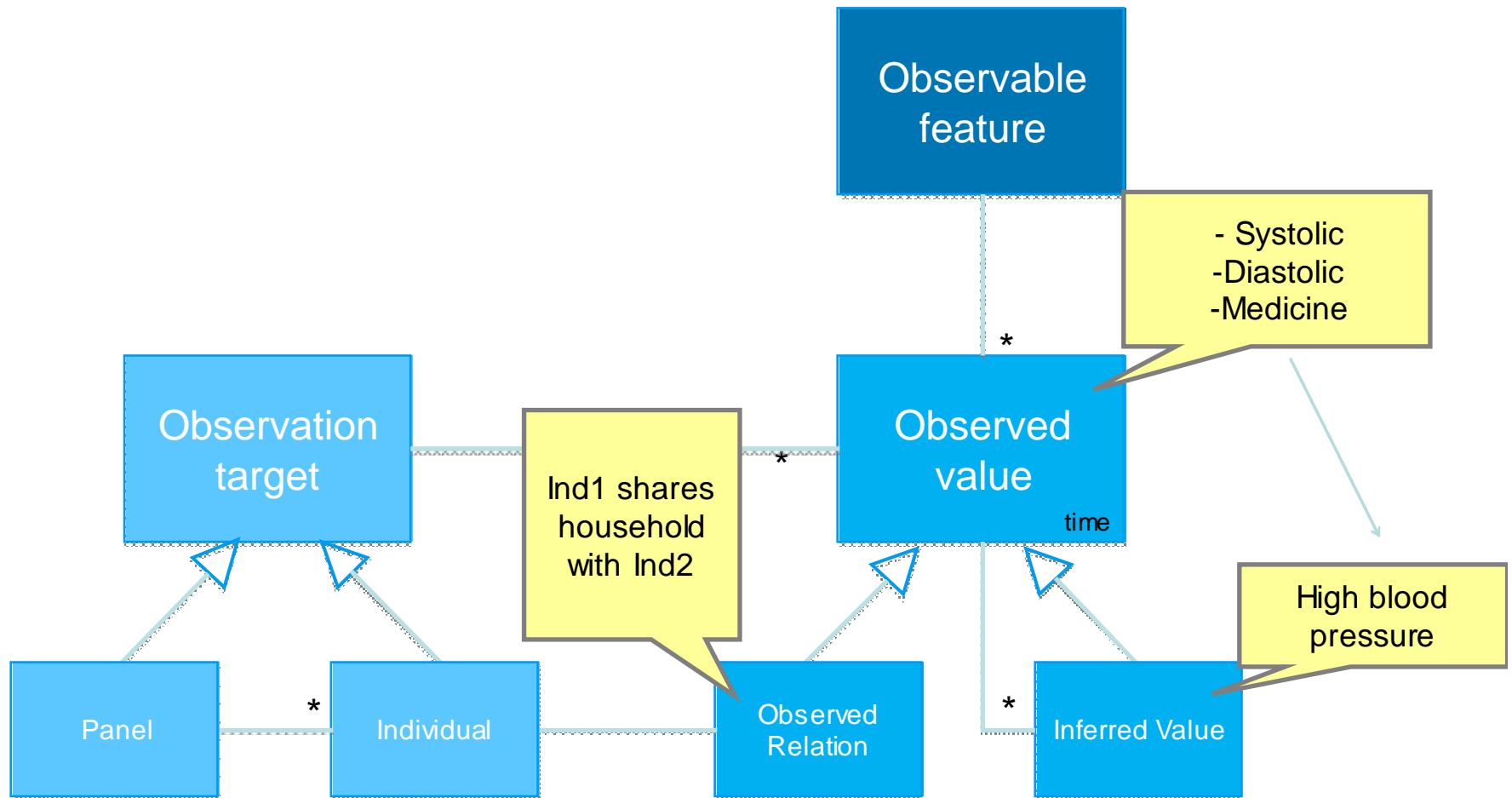




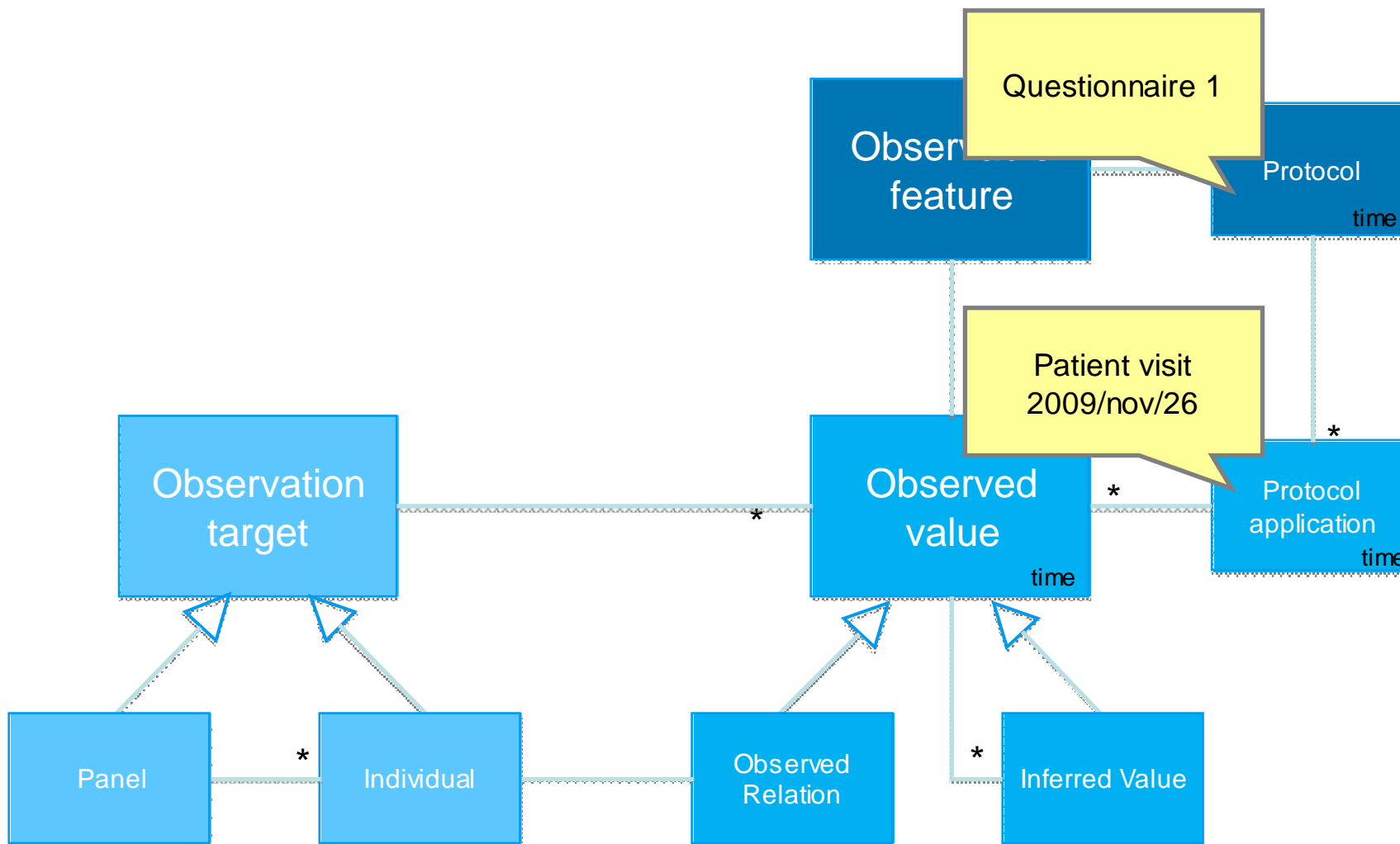
# Individual and panel level data



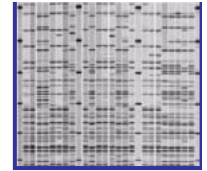
# Values, relations, inferences



# Protocols/SOPs used ...

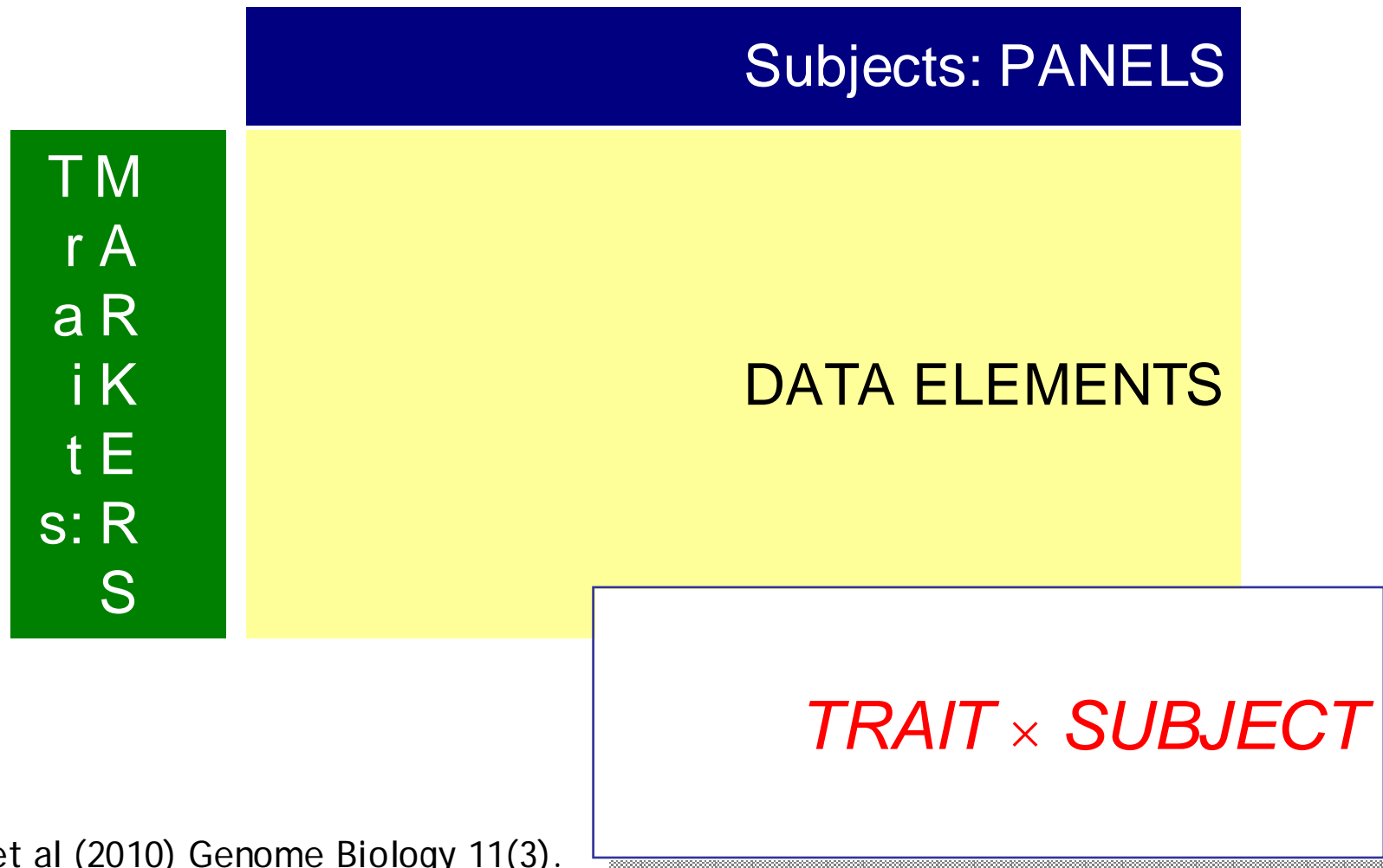


# XGAP for genotypes (superset of pheno-OM)



Data in matrices

*Genotype data*



# Minimal model

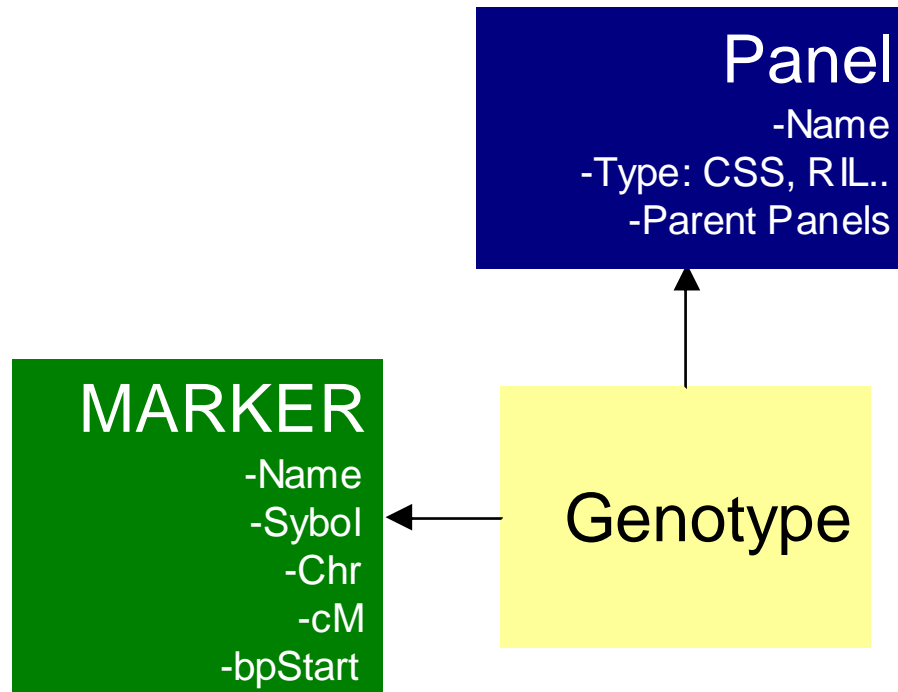
- Annotations in tables, e.g. Marker

## MARKER

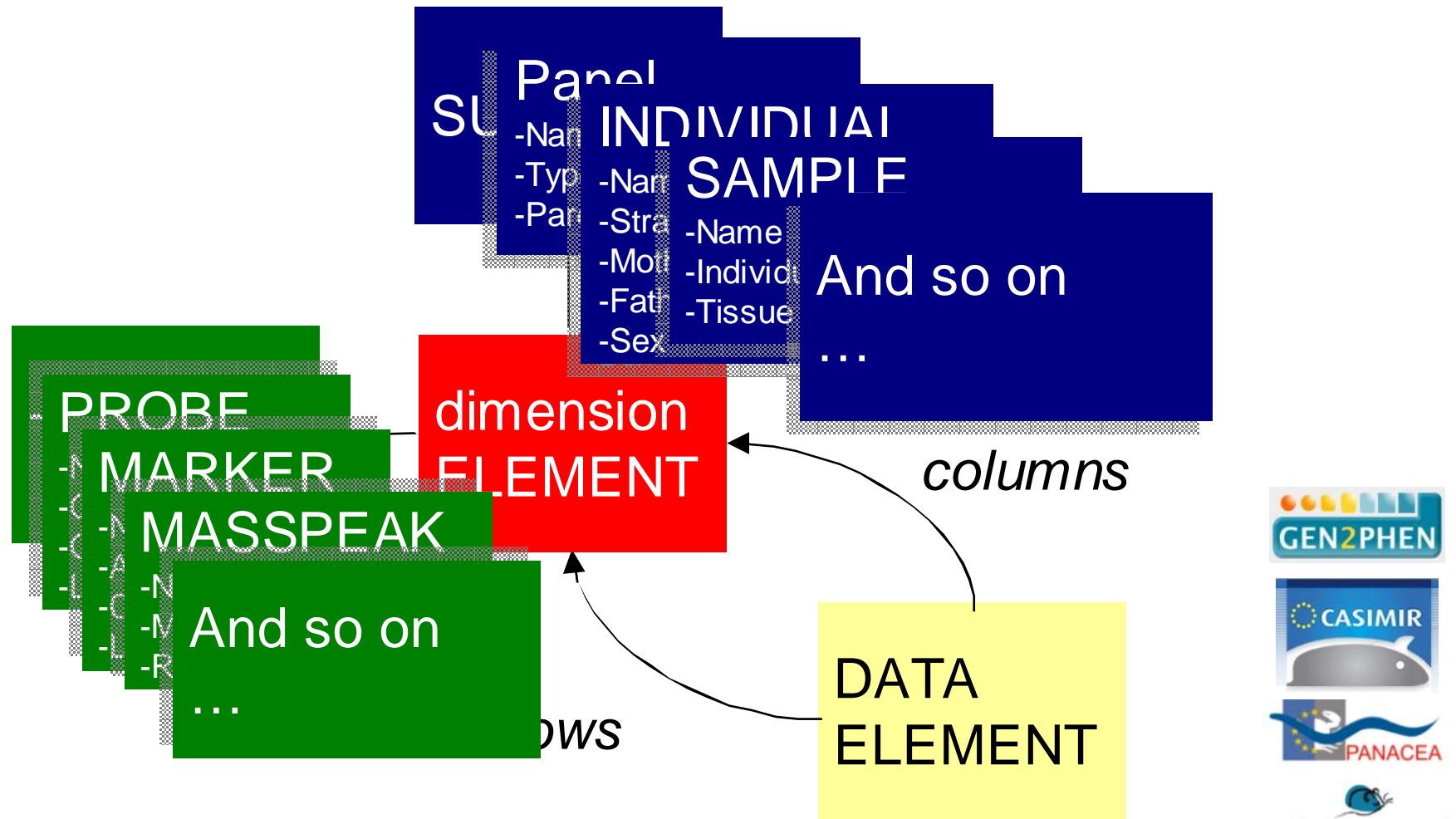
-Name  
--Symbol  
-Chr  
-cM  
-bpStart

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



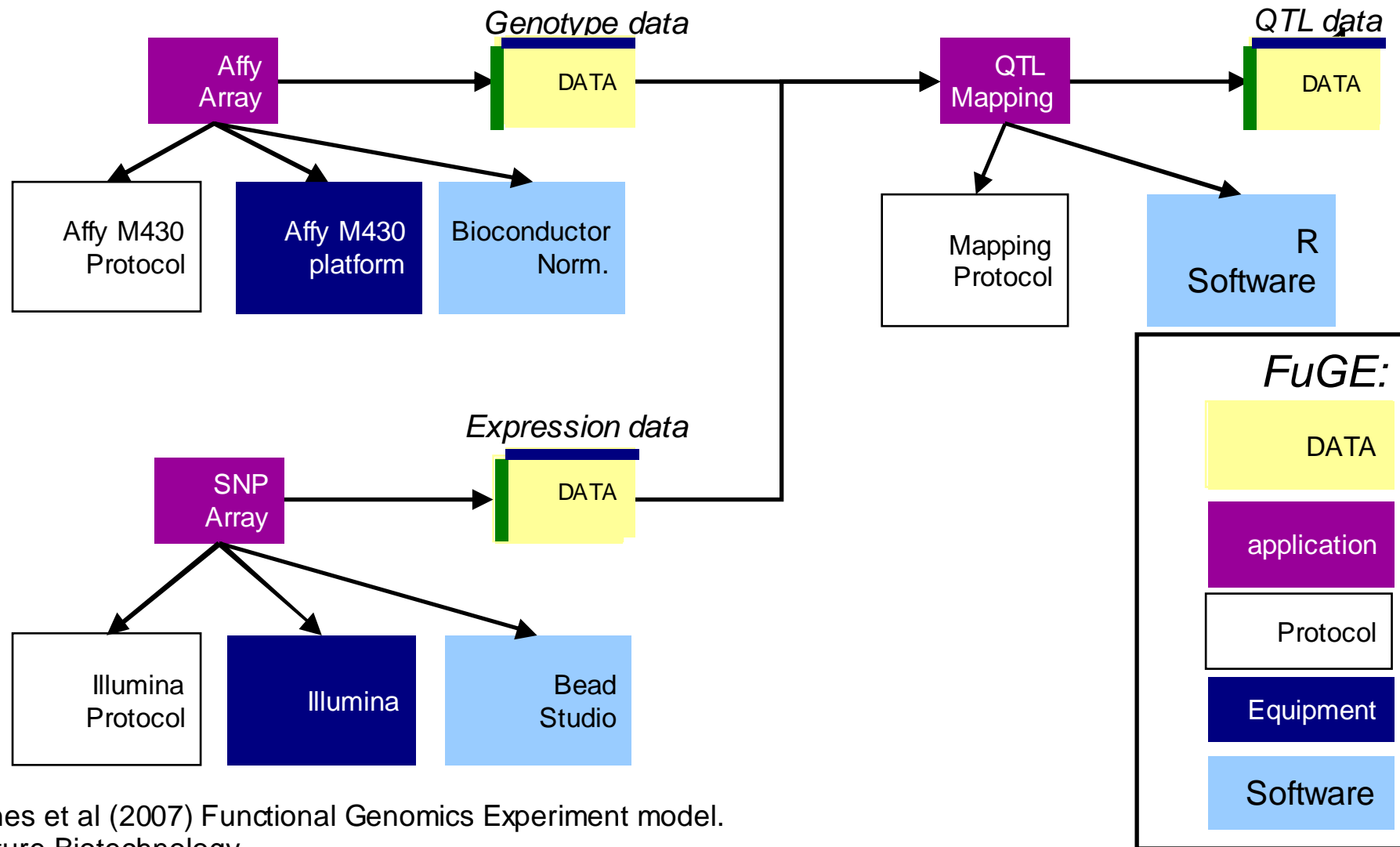
# XGAP (extension based variation mechanism)





# XGAP extends on FuGE

/



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

# ISA-TAB(generic model)

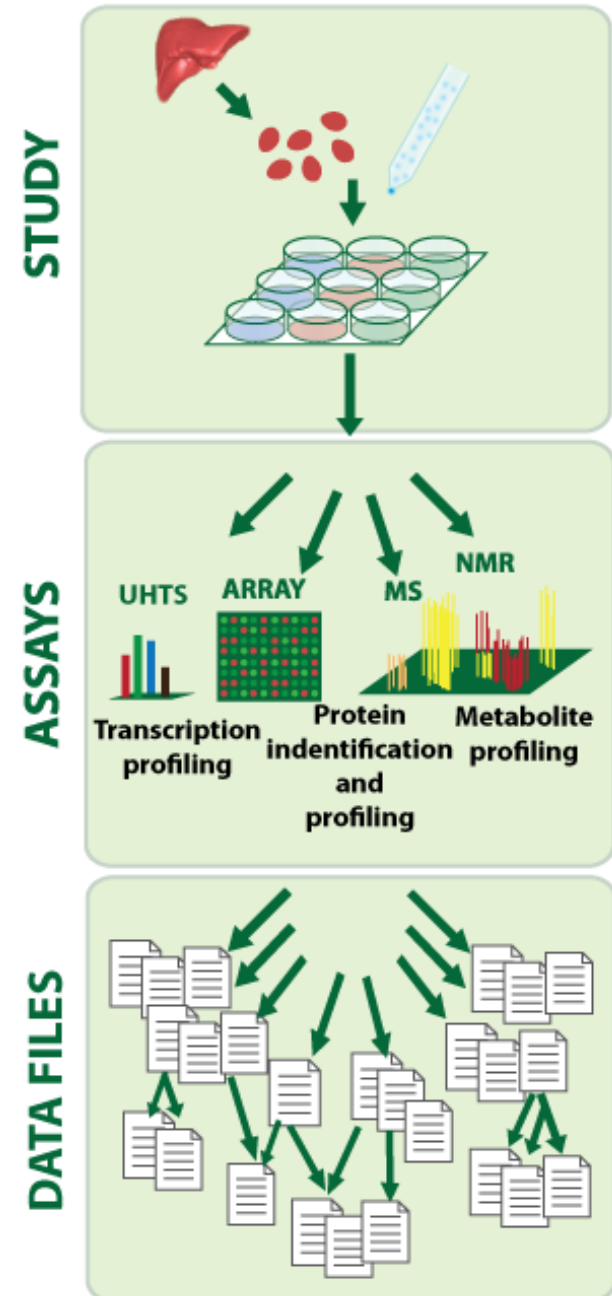
Differs from MAGE-TAB

- Nested investigations (as studies)
- To have templates assays
- More aligned to FuGE
- But some find it still difficult
  - Reflects largely the complexity of the experiments not that of the format.

ISA =

- Investigation
- Study (Investigation component)
- Assay (a component of Study)
- Data files
  - Used at the Harvard Stem Cell Discovery Lab, NERC and Sysmo-DB

## Example of experimental WORKFLOW



## Current work: merging the models

<b>Pheno-OM</b>	<b>FuGE + XGAP</b>	<b>MAGE-TAB</b>	<b>PaGE</b>
Investigation	Investigation	Investigation	Investigation
ObservableFeature	Trait	~Characteristic	Observable_feature
ObservationTarget	Subject	Source	Abstract_observation_target
Individual	Individual	NA	Individual
Panel	Panel	NA	Panel
ObservedValue	DataElement	CharacteristicValue	ObservedValue
ObservedRelationship	NA	NA	NA
InferredValue	via ProtocolApplication	NA	ObservedValue (self assoc)
OntologySource	OntologySource	OntologySource	Ontology_source
OntologyTerm	OntologyIndividual	OntologyTerm	Ontology_term
Protocol	Protocol	Protocol	Observation_method
ProtocolApplication	ProtocolApplication	ProtocolApplication	Observation_method
Code	NA	~Characteristic	NA

# Discussion for afternoon

- How to properly map this to RDF
- Map data values to RDF/OWL
  - Most models ontology enabled
- Map data structures for RDF/OWL
  - So annotate data model
- Some experiences with MOLGENIS D2RQ
  - Report of GEN2PHEN internship student