

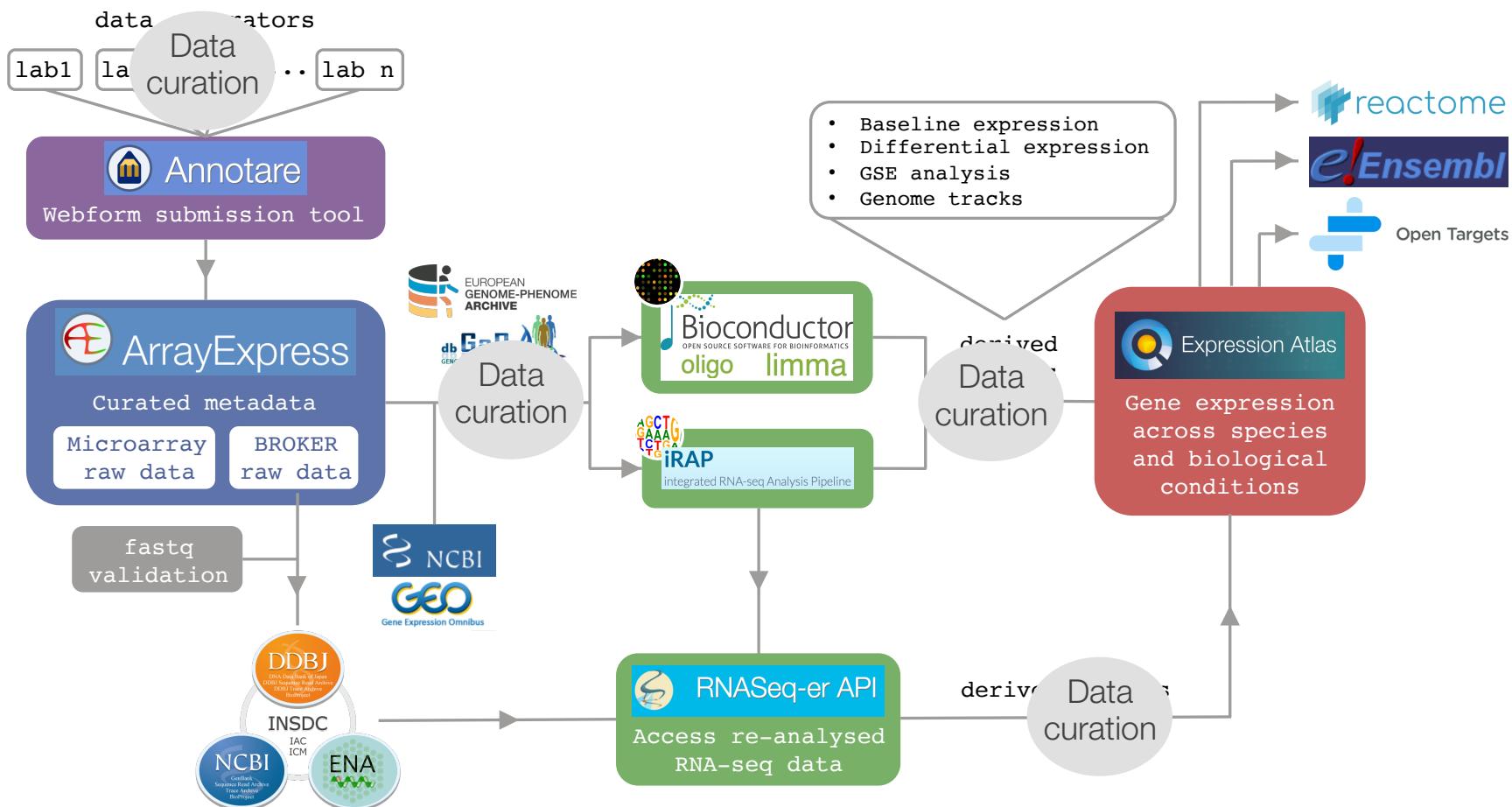
# EMBL-EBI Bioinformatics resources for exploring functional genomics data

[https://www.ebi.ac.uk/~lauhuema/  
workshop/DKFZ](https://www.ebi.ac.uk/~lauhuema/workshop/DKFZ)

Laura Huerta, PhD  
Senior Scientific Curator  
[lauhuema@ebi.ac.uk](mailto:lauhuema@ebi.ac.uk)  
9 November 2017



# Functional genomics resources at EMBL-EBI



# EMBL-EBI Bioinformatics resources for exploring functional genomics data

Data standards, curation  
and ontologies

Laura Huerta, PhD  
Senior Scientific Curator  
[lauhuema@ebi.ac.uk](mailto:lauhuema@ebi.ac.uk)  
9 November 2017



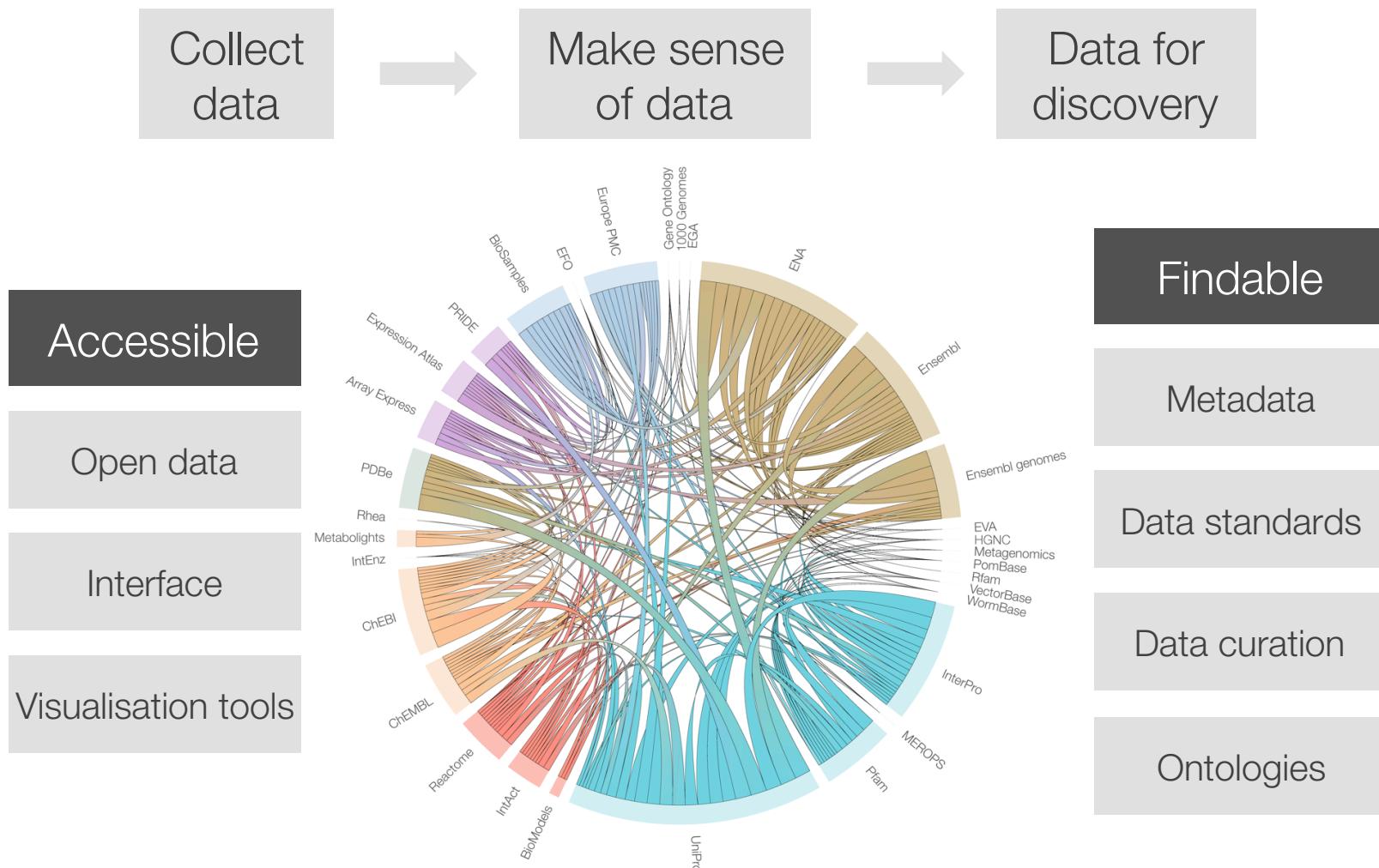
# Wunderkammer - cabinet of curiosities



Collect and preserve knowledge

[www.ebi.ac.uk/about/news/feature-story/cabinet-of-curiosities](http://www.ebi.ac.uk/about/news/feature-story/cabinet-of-curiosities)

# EMBL-EBI databases – A smart cabinet of curiosities



# EMBL-EBI databases – A smart cabinet of curiosities

- Open, online access to experimental datasets
- Annotation of datasets with adequate metadata
- Use of controlled terms in metadata annotations -> ontologies
- Mechanisms to search for metadata to find relevant experimental results

Open access

Metadata

Ontologies

Discovery

# Data standards in functional genomics

*Nature Genetics* **29**, 365 - 371 (2001)  
doi:10.1038/ng1201-365

## **Minimum information about a microarray experiment (MIAME) —toward standards for microarray data**

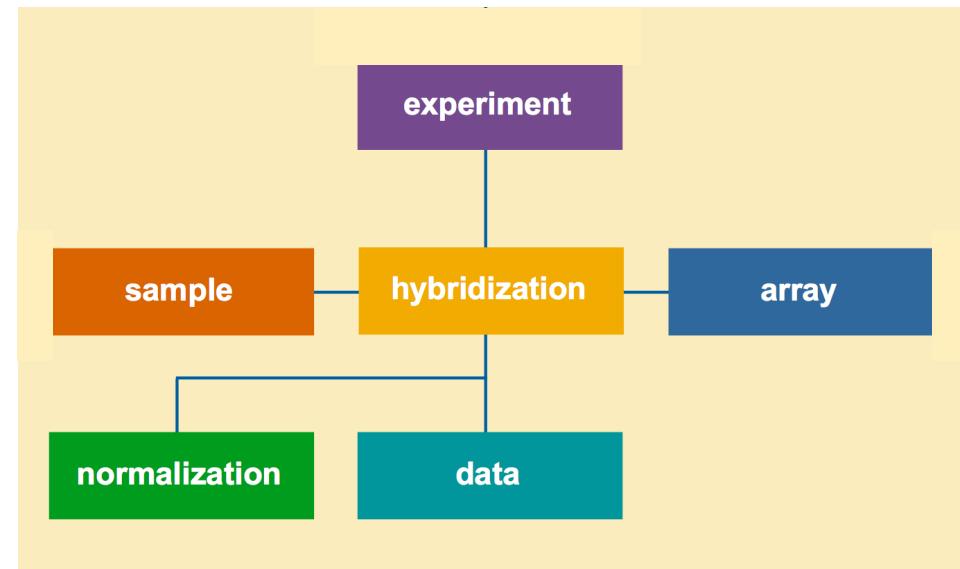
Alvis Brazma<sup>1</sup>, Pascal Hingamp<sup>2</sup>, John Quackenbush<sup>3</sup>, Gavin Sherlock<sup>4</sup>, Paul Spellman<sup>5</sup>,  
Chris Stoeckert<sup>6</sup>, John Aach<sup>7</sup>, Wilhelm Ansorge<sup>8</sup>, Catherine A. Ball<sup>4</sup>, Helen C. Causton<sup>9</sup>,  
Terry Gaasterland<sup>10</sup>, Patrick Glenisson<sup>11</sup>, Frank C.P. Holstege<sup>12</sup>, Irene F. Kim<sup>4</sup>, Victor  
Markowitz<sup>13</sup>, John C. Matese<sup>4</sup>, Helen Parkinson<sup>1</sup>, Alan Robinson<sup>1</sup>, Ugis Sarkans<sup>1</sup>, Steffen  
Schulze-Kremer<sup>14</sup>, Jason Stewart<sup>15</sup>, Ronald Taylor<sup>16</sup>, Jaak Vilo<sup>1</sup> & Martin Vingron<sup>17</sup>

*“raw data is not enough to interpret the results and to  
verify the conclusions based on microarray data analysis”*

# Data standards in functional genomics

*“MIAME describes the data and metadata that authors must provide to support conclusions drawn from a microarray investigation, so that the data obtained in the investigation can be interpreted unambiguously and the investigation can be reproduced.”*

Reproducible  
Raw data  
Metadata



# Data standards and reproducibility

The six most critical elements contributing towards MIAME are:



1. The raw data for each hybridisation (e.g., CEL or GPR files)
2. The final processed (normalised) data for the set of hybridisations in the experiment (study) (e.g., the gene expression data matrix used to draw the conclusions from the study)
3. The essential sample annotation including experimental factors and their values (e.g., compound and dose in a dose response experiment)
4. The experimental design including sample data relationships (e.g., which raw data file relates to which sample, which hybridisations are technical, which are biological replicates)
5. Sufficient annotation of the array (e.g., gene identifiers, genomic coordinates, probe oligonucleotide sequences or reference commercial array catalog number)
6. The essential laboratory and data processing protocols (e.g., what normalisation method has been used to obtain the final processed data)

<http://fged.org/projects/miame>

# Data standards and reproducibility



- 1. The description of the biological system, samples, and the experimental variables being studied:**
  - “compound” and “dose” in dose-response experiments or “antibody” in ChIP-Seq experiments, the organism, tissue, and the treatment(s) applied.
- 2. The sequence read data for each assay:**
  - read sequences and base-level quality scores for each assay; FASTQ format is recommended, with a description of the scale used for quality scores.
- 3. The ‘final’ processed (or summary) data for the set of assays in the study:**
  - the data on which the conclusions in the related publication are based, and descriptions of the data format.
- 4. General information about the experiment and sample-data relationships:**
  - a summary of the experiment and its goals, contact information, any associated publication, and a table specifying sample-data relationships.
- 5. Essential experimental and data processing protocols:**
  - how the nucleic acid samples were isolated, purified and processed prior to sequencing, a summary of the instrumentation used, library preparation strategy, labelling and amplification methodologies, alignment algorithms and data filtering plus data processing & analysis protocols.

<http://fged.org/projects/minseqe>

# Data standards and reproducibility

Capture enough information to interpret  
and reproduce the experiment

## Data standards

Experiment design

Protocols

Sample annotation

Raw data

Relationship sample-files

# Advantages of using data standards

Authoring more complete and standardised metadata...

- Will aid dataset:
  - Discovery
  - Exploration
  - Integration
  - Secondary use
- Will aid communication of scientific results
- Will aid knowledge management within research organizations
- Will make your data more FAIR

Findable

Accessible

Interoperable

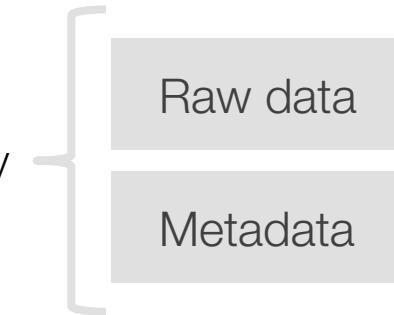
Re-usable

# Data curation

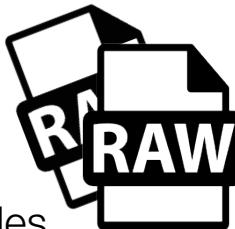
## Biocuration



Translation and integration of information relevant to biology into a database or resource



### Raw data



- unprocessed data files
- **Microarray:** files from the scanner (e.g. Affymetrix CEL files, Agilent feature extraction *txt* files, Illumina *idat* files)
- **Sequencing:** raw sequence read files (e.g. FASTQ files)

### Metadata

- Experiment description
- Experiment title
- Sample annotation
- Protocols
- Publication details (if any)
- Author contact details

# Data curation – data and files

PeerJ

✓ PEER-REVIEWED

XA21-specific induction of stress-related genes following *Xanthomonas* infection of detached rice leaves

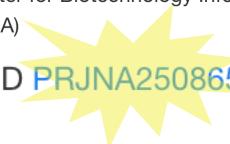


## Data Availability

The following information was supplied regarding data availability:

The National Center for Biotechnology Information Sequence Read Archive (SRA)

BioProject ID [PRJNA250865](#).



## Metadata

7 samples x 3 biological replicates

### RNA sequencing sample treatment summary

Table summarizes the experimental setup including the genotypes, time of treatment, and type of treatment used for samples used in RNA sequencing. There were three replicates for each sample for a total of 21 sequenced samples.

	A	B	C
1	Genotype	Time (hours)	Treatment
2	Kitaake	0	None
3	EFR::XA21::GFP	0	None
4	EFR::XA21::GFP	0.5	500 nM elf18
5	EFR::XA21::GFP	1	500 nM elf18
6	EFR::XA21::GFP	3	500 nM elf18
7	EFR::XA21::GFP	6	500 nM elf18
8	EFR::XA21::GFP	12	500 nM elf18

# Data curation – data and files

## Raw data



- Genotype?
- Treatment?
- Time?

[SRX873376: Other Sequencing of Japanese rice](#)

1 ILLUMINA (Illumina HiSeq 2000) run: 29.9M spots, 9G bases, 5.3Gb downloads

**Submitted by:** DOE JOINT GENOME INSTITUTE (JGI)

**Study:** Oryza sativa Japonica strain:EFR-XA21 | cultivar:Kitaake Transcriptome or Gene expression

[PRJNA250865](#) • [SRP054056](#) • [All experiments](#) • [All runs](#)

[show Abstract](#)

**Sample:** Oryza sativa cv. Kitaake EFR-XA21

[SAMN03003383](#) • [SRS843490](#) • [All experiments](#) • [All runs](#)

**Organism:** [Oryza sativa Japonica Group](#)

### Library:

*Name:* NTBA

*Instrument:* Illumina HiSeq 2000

*Strategy:* RNA-Seq

*Source:* TRANSCRIPTOMIC

*Selection:* RT-PCR

*Layout:* PAIRED

### Spot descriptor:



**Runs:** 1 run, 29.9M spots, 9G bases, [5.3Gb](#)

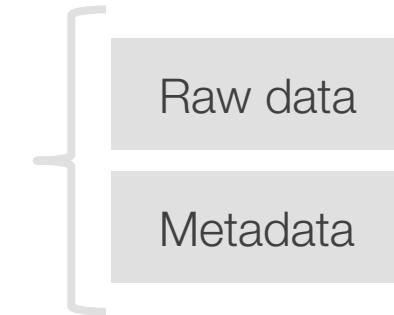
Run	# of Spots	# of Bases	Size	Published
<a href="#">SRR1799213</a>	29,893,272	9G	5.3Gb	2015-02-20

# Data curation – data and files

PeerJ

✓ PEER-REVIEWED

XA21-specific induction of stress-related genes following *Xanthomonas* infection of detached rice leaves



Raw data

Metadata

Run
SRR1799194
SRR1799197
SRR1799193
SRR1799210
SRR1799203
SRR1799205
SRR1799199
SRR1799204
SRR1799212
SRR1799208
SRR1799202
SRR1799198
SRR1799213
SRR1799201
SRR1799196
SRR1799209
SRR1799206
SRR1799195
SRR1799207
SRR1799200
SRR1799211



	FactorValue[genotype]	FactorValue[compound]	FactorValue[dose]	Unit[concentration unit]	FactorValue[time]	Unit[time unit]
SRR1799194	wild type	none		0 nanomolar		0 hour
SRR1799197	wild type	none		0 nanomolar		0 hour
SRR1799193	wild type	none		0 nanomolar		0 hour
SRR1799210	EFR:XA21:GFP	none		0 nanomolar		0 hour
SRR1799203	EFR:XA21:GFP	none		0 nanomolar		0 hour
SRR1799205	EFR:XA21:GFP	none		0 nanomolar		0 hour
SRR1799199	EFR:XA21:GFP	elf18		500 nanomolar		0.5 hour
SRR1799204	EFR:XA21:GFP	elf18		500 nanomolar		0.5 hour
SRR1799212	EFR:XA21:GFP	elf18		500 nanomolar		0.5 hour
SRR1799208	EFR:XA21:GFP	elf18		500 nanomolar		1 hour
SRR1799202	EFR:XA21:GFP	elf18		500 nanomolar		1 hour
SRR1799198	EFR:XA21:GFP	elf18		500 nanomolar		1 hour
SRR1799213	EFR:XA21:GFP	elf18		500 nanomolar		3 hour
SRR1799201	EFR:XA21:GFP	elf18		500 nanomolar		3 hour
SRR1799196	EFR:XA21:GFP	elf18		500 nanomolar		3 hour
SRR1799209	EFR:XA21:GFP	elf18		500 nanomolar		6 hour
SRR1799206	EFR:XA21:GFP	elf18		500 nanomolar		6 hour
SRR1799195	EFR:XA21:GFP	elf18		500 nanomolar		6 hour
SRR1799207	EFR:XA21:GFP	elf18		500 nanomolar		12 hour
SRR1799200	EFR:XA21:GFP	elf18		500 nanomolar		12 hour
SRR1799211	EFR:XA21:GFP	elf18		500 nanomolar		12 hour

# Data curation – correction of errors

Sample	Treatment
Culture A_control	Drug X
Culture B_treated	Drug X
Culture C_control	control
Culture D_treated	Drug X

organism > Equus caballus  
age > 79 years

sex > male  
organism part > endometrium

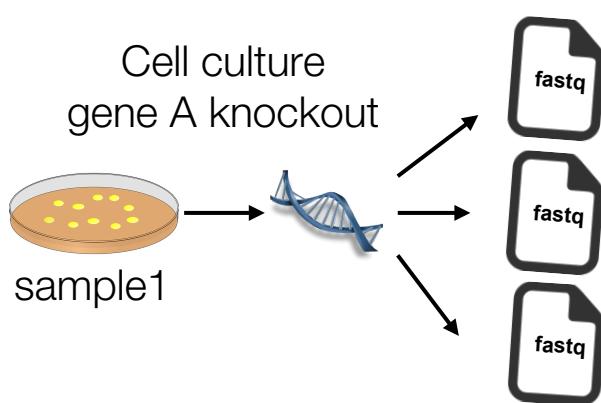
We compared knockouts  
against wild-type littermates

Sample	Genotype
Sample 1	Gene X -/-
Sample 2	Gene X -/-
Sample 3	Gene X -/-

Sample	Organism part
Sample D	heart
Sample D	lung
Sample D	liver

# Data curation – correction of errors

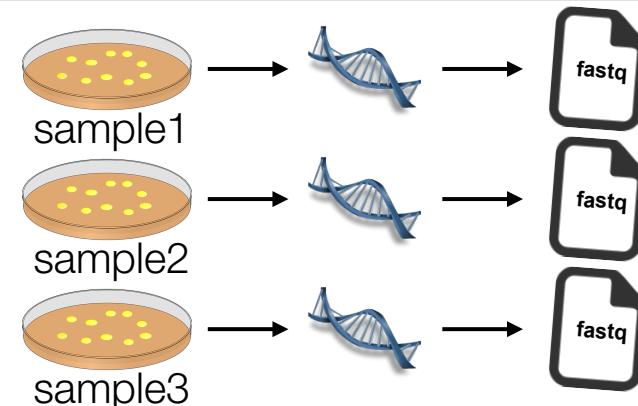
## Technical replicates



Sample	Seq. run	File
Sample 1	run_assay_1	1.fq.gz
Sample 2	run_assay_2	2.fq.gz
Sample 3	run_assay_3	3.fq.gz



## Biological replicates



Sample	Seq. run	File
Sample 1	run_assay_1	1.fq.gz
Sample 1	run_assay_2	2.fq.gz
Sample 1	run_assay_3	3.fq.gz

# Data curation – correction of errors

Sample annotation

## Hands-on activity

Let's become a biocurator,  
at least for 15 min...



In pairs

# Data curation – correction of errors

Sample annotation

<https://www.ncbi.nlm.nih.gov/geo>



GEO  
Gene Expression Omnibus

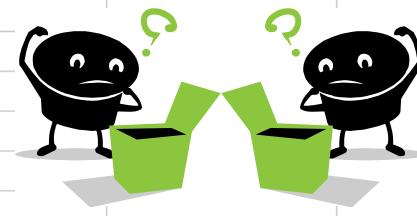
Search

- ✓ What is the experiment about?
- ✓ In which organism is the experiment done?
- ✓ What is the biological material used in the experiment?
- ✓ How many samples are authors testing?
- ✓ Is it a microarray or a RNA-sequencing experiment?
- ✓ What is the main variable that authors are studying?
- ✓ What are the two groups of samples compared?
- ✓ How many biological replicates are there in each group?
- ✓ Are there any technical replicates?

# Data curation – correction of errors

## Sample annotation

	Sample characteristics		Experimental variables		
Assay Name	Characteristics [organism]	Characteristics [cell type]	FactorValue [compound]	FactorValue [dose]	Unit [concentration unit]
GSM1087598	Homo sapiens	coronary artery endothelial cell			
GSM1087599	Homo sapiens	coronary artery endothelial cell			
GSM1087600	Homo sapiens	coronary artery endothelial cell			
GSM1087601	Homo sapiens	coronary artery endothelial cell			
GSM1087602	Homo sapiens	coronary artery endothelial cell			
GSM1087603	Homo sapiens	coronary artery endothelial cell			



GSM1087598 HCAEC\_SsnB\_rep1  
GSM1087599 HCAEC\_SsnB\_rep2  
GSM1087600 HCAEC\_SsnB\_rep3

{ Source name  
Organism  
Characteristics

HCAECs exposed to SsnB  
**Homo sapiens**  
cell type: Primary Human Coronary Artery Endothelial Cells  
treatment: control

GSM1087601 HCAEC\_DMSO\_rep1  
GSM1087602 HCAEC\_DMSO\_rep2  
GSM1087603 HCAEC\_DMSO\_rep3

{ Source name  
Organism  
Characteristics

HCAECs exposed to Vehicle Control  
**Homo sapiens**  
cell type: Primary Human Coronary Artery Endothelial Cells  
treatment: Sparstololin B

# Data curation – correction of errors

## Sample annotation

*Hi, Laura*

*The sample labels for each file are correct, e.g., these are all results for SsnB-treated cells.*

*GSM1087598 HCAEC\_SsnB\_rep1*

*GSM1087599 HCAEC\_SsnB\_rep2*

*GSM1087600 HCAEC\_SsnB\_rep3*

*The information listed under 'Characteristics' appears to be inverted; that is, 'treatment: control' are actually the SsnB-treated cells and 'treatment: Sparstolonin B' are actually the DMSO vehicle controls.*

# Annotation with ontology terms

## Ontology



A systematic way to name and organise entities, establishing relationships between the entities

Controlled vocabulary

Hierarchy (relationship)



The screenshot shows the EFO homepage with a dark blue header. On the left is the EFO logo and the text "Experimental Factor Ontology". On the right is a search bar with placeholder text "Search EFO" and a magnifying glass icon. Below the search bar are examples: "cancer", "HeLa", and "Li-Fraumeni syndrome". A navigation bar at the bottom includes links for "Home", "Browse EFO", "Submit Term", "EBI RDF Platform", and "About".

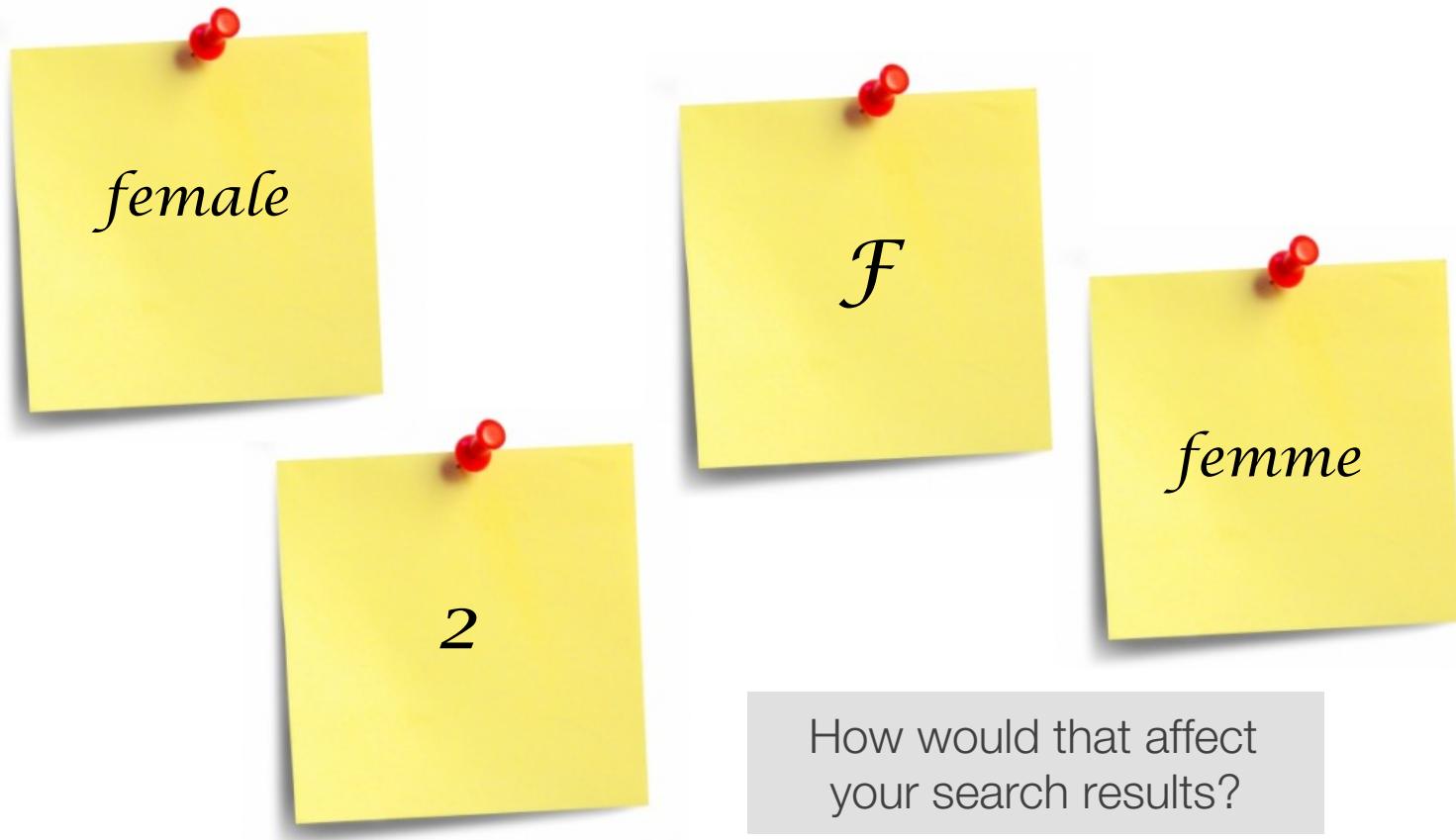
## Representing experimental variables with EFO

The **Experimental Factor Ontology** (EFO) provides a systematic description of many experimental variables available in EBI databases, and for external projects such as the NHGRI GWAS catalog. It combines parts of several biological ontologies, such as UBERON anatomy, ChEBI chemical compounds, and Cell Ontology. The scope of EFO is to support the annotation, analysis and visualization of data handled by many groups at the EBI and as the core ontology for Open Targets. EFO is developed by the EMBL-EBI Samples, Phenotypes and Ontologies Team (SPOT). We also add terms for external users when requested. If you are new to ontologies, there is a short introduction on the subject available and a blog post by James Malone on what ontologies are for.

[www.ebi.ac.uk/efo](http://www.ebi.ac.uk/efo)

# Why we annotate samples using ontologies?

How many ways can you say “female”?



How would that affect  
your search results?

# Why we annotate samples using ontologies?

18-day pregnant females	female (lactating)	individual female	worker caste (female)
2 yr old female	female (pregnant)	lgb*cc females	sex: female
400 yr. old female	female (outbred)	mare	female, other
adult female	female parent	female (worker)	female child
asexual female	female plant	monosex female	femal
castrate female	female with eggs	ovigerous female	3 female
cf.female	female worker	oviparous sexual females	female (phenotype)
cystocarpic female	female, 6-8 weeks old	worker bee	female mice
dikaryon	female, virgin	female enriched	female, spayed
dioecious female	female, worker	pseudohermaphroditic female	femlale
diploid female	female(gynoecious)	remale	metafemale
f	femele	semi-engorged female	sterile female
famale	female, pooled	sexual oviparous female	normal female
femail	femalen	sterile female worker	sf
female	females	strictly female	vitellogenetic replete female
female - worker	females only	tetraploid female	worker
female (alate sexual)	gynoecious	thelytoky	hexaploid female
female (calf)	healthy female	female (gynoecious)	female (f-o)
hen	probably female (based on morphology)		

female (note: this sample was originally provided as a \"male\" sample to us and therefore labeled this way in the brawand et al. paper and original geo submission; however, detailed data analyses carried out in the meantime clearly show that this sample stems from a female individual)",

*Courtesy of N. Silvester, European Nucleotide Archive, EMBL-EBI*

# Annotation with ontology terms

Sample description:  
“pancreatic adenocarcinoma  
cell line AsPC-1 treated  
with quercetin for 12 hours”

OLS > Experimental Factor Ontology EFO > EFO:1000044

## pancreatic adenocarcinoma

http://www.ebi.ac.uk/efo/EFO\_1000044

An adenocarcinoma which arises from the exocrine pancreas. Ductal adenocarcinoma and its variants are the most common types of pancreatic adenocarcinoma

Synonyms: adenocarcinoma of pancreas, pancreas adenocarcinoma

Tree view Term history

- experimental factor
  - material property
    - disposition
    - disease
      - endocrine system disease
        - endocrine neoplasm
          - pancreatic neoplasm
            - pancreatic carcinoma
              - pancreatic adenocarcinoma

Graph view Reset tree Show all siblings

Term info

DOID definition citation DOID:4074

NCI Thesaurus definition citation NCIT:C8294

SNOMEDCT definition citation SNOMEDCT:700423003

UMLS definition citation UMLS:C0281361

term editor Laura Huerta

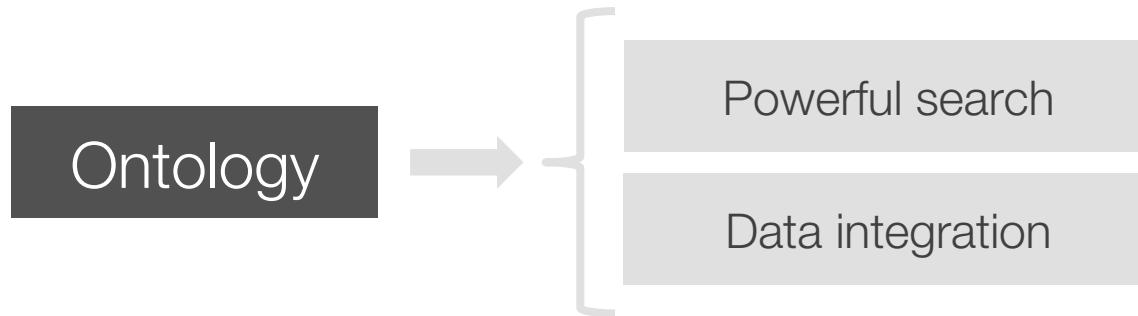
Term relations

Subclass of:

  - pancreatic carcinoma
  - has\_disease\_location some (pancreas or part\_of some pancreas)

cell line	disease	compound	time (hour)
AsPC-1	pancreatic adenocarcinoma	quercetin	12

# Ontology annotation



Smart search

Efficient search via ontology-driven query expansion

Biological conditions

lung car|

Ex lung carcinoid tumor

+ lung carcinoma

This section shows a screenshot of the EFO search interface. It features a search bar with the input "lung car|". Below the search bar, a dropdown menu lists "lung carcinoid tumor" and "lung carcinoma". The "lung carcinoma" option is highlighted with a blue border.

Complex queries

lung carcinoma

- bronchogenic carcinoma
- large cell lung carcinoma
  - Lung Sarcomatoid Carcinoma
  - Lymphoepithelioma-Like Lung Carcinoma
- lung adenocarcinoma
- non-small cell lung carcinoma
  - pulmonary mucoepidermoid carcinoma
  - pulmonary neuroendocrine tumor
  - small cell lung carcinoma
  - squamous cell lung carcinoma

This section shows a hierarchical search results tree for the query "lung carcinoma". The root node is "lung carcinoma", which branches into "bronchogenic carcinoma", "large cell lung carcinoma" (which further branches into "Lung Sarcomatoid Carcinoma" and "Lymphoepithelioma-Like Lung Carcinoma"), "lung adenocarcinoma", "non-small cell lung carcinoma" (which further branches into "pulmonary mucoepidermoid carcinoma", "pulmonary neuroendocrine tumor", "small cell lung carcinoma", and "squamous cell lung carcinoma").

# Ontology annotation – adding value to your data

<https://www.ncbi.nlm.nih.gov/geo>



Series GSE81089      Query DataSets for GSE81089

Status Public on Jun 01, 2016  
Title Next Generation Sequencing (RNAseq) of non-small cell lung cancer  
Organism *Homo sapiens*  
Experiment type Expression profiling by high throughput sequencing  
Summary Cancer testis antigens (CTAs) are of clinical interest as biomarkers and present valuable targets for immunotherapy. To comprehensively characterize the CTA landscape of non-small cell lung cancer (NSCLC), we compared RNAseq data of 199 NSCLC tissues to the normal transcriptome of 142 samples from 32 different normal organs. Of 232 CTAs currently annotated in the CTdatabase, 96 were confirmed in NSCLC. applied stringent criteria on CTAs, of which 55 genes v analysis revealed that CTA e expression is common. Immune expression of selected genes regulatory mechanism of CTA Cancer Genome Atlas. The pr was not confirmed, neither in analysis of 1117 NSCLC cases

Overall design Fresh frozen tumor tissue surgically treated 2006-2011 Sweden and 19 paired norm the regional lung cancer regis Several of the new CTAs are ;  
  
Sample characteristics values pTNM: decided by Hans Bruni Stage according to pTNM: 1= Histology diagnosis spring 20 1=squamous cell cancer 2=Ad Surgery date: the date when Age: age when surgery was p Vital date: day of death or lat Dead: 0=no 1=yes Smoking history : 1=current WHO performance status: Per

Please note that the L608T\_ data files) are associated with

Profiling cancer testis antigens in non-small-cell lung cancer.  
(PMID:27699219 PMCID:PMC5033889)

Abstract Citations BioEntities Related Articles External Links

Djureinovic D<sup>1</sup>, Hallström BM<sup>2</sup>, Horie M<sup>3</sup>, Mattsson JS<sup>1</sup>, La Fleur L<sup>1</sup>, Fagerberg L<sup>2</sup> , Brunnström H<sup>4</sup>, Lindskog C<sup>1</sup> , Madjar K<sup>5</sup>, Rahnenführer J<sup>5</sup> , Ekman S<sup>6</sup>, Ståhle E<sup>7</sup> , Koyi H<sup>8</sup>, Brandén E<sup>8</sup>, Edlund K<sup>9</sup>, Hengstler JG<sup>9</sup>, Lambe M<sup>10</sup>, Saito A<sup>3</sup>, Botling J<sup>1</sup>, Pontén F<sup>1</sup> , Uhlén M<sup>2</sup>, Micke P<sup>1</sup>

Affiliations   
[JCI Insight](#) [07 Jul 2016, 1(10):e86837]

Type: research-article, Journal Article

machines (Illumina) using the standard Illumina RNAseq protocol with a read length of 2 × 100 bases. The raw data has been uploaded together with clinical information on GEO, with the accession number

GSE81089 (<http://www.ncbi.nlm.nih.gov/geo/>).

Contributor(s) Djureinovic D, Hallström BM, Horie M, Mattsson JS, La Fleur L, Fagerberg L, Brunnström H, Lindskog C, Madjar K, Rahnenführer J, Ekman S, Ståhle E, Koyi H, Brandén E, Edlund K, Hengstler JG, Lambe M, Saito A, Botling J, Pontén F, Uhlén M, Micke P

Citation missing Has this study been published? Please [login](#) to update or notify GEO.

# Ontology annotation – adding value to your data



<https://www.ncbi.nlm.nih.gov/geo>

→ GSE81089

Sample GSM2142642	
Status	Public on Jun 01, 2016
Title	matched sample_L511N
Sample type	SRA
Source name	Human non-malignant tissue
Organism	<a href="#">Homo sapiens</a>
Characteristics	tumor (t) or normal (n): L511N

non-small cell  
lung carcinoma

Sample GSM2142481	
Status	Public on Jun 01, 2016
Title	L511T
Sample type	SRA
Source name	Human NSCLC tissue
Organism	<a href="#">Homo sapiens</a>
Characteristics	tumor (t) or normal (n): L511T stage tnM: 5 histology: 2 surgery date: 2007-04-05 age: 56 gender: female vital date: 2013-04-28 dead: 0 smoking: 3 ps who: 0

non-small cell  
lung carcinoma



Disease?

# Ontology annotation – adding value to your data

The screenshot shows the OLS interface. At the top, it says "Ontology Lookup Service". Below that, the URL is http://www.ebi.ac.uk/efo/EFO\_0003060. The main title is "non-small cell lung carcinoma". A detailed description follows: "A heterogeneous aggregate of at least three distinct histological types of lung cancer, including SQUAMOUS CELL CARCINOMA. They are dealt with collectively because of their shared treatment strategy." Under "Synonyms", there is a long list including "Lung Carcinomas, Non-Small-Cell, NSCLC, Carcinoma, Non-Small-Cell Lung, Lung Carcinoma, Non-Small-Cell, Non-Small-Cell Lung Carcinoma, Carcinomas, Non-Small-Cell Lung, Non-Small Cell Lung Cancer, Non-Small Cell Lung Carcinoma, Non-Small-Cell Lung Carcinomas, Non Small Cell Lung Carcinoma, NONSMALL CELL LUNG CARCINOMA, Carcinoma, Non-Small Cell Lung, Non-small cell lung cancer (disorder), Carcinoma, Non Small Cell Lung". On the left, there is a tree view of the ontology structure under "experimental factor". The node "non-small cell lung carcinoma" is highlighted in blue. On the right, there is a "Term info" panel with various definitions and citations from DOID, MSH, NCI Thesaurus, OMIM, SNOMEDCT, and gwas trait.

We use OLS to annotate samples using ontology terms

[www.ebi.ac.uk/ols/index](http://www.ebi.ac.uk/ols/index)

# Ontology annotation – adding value to your data



<https://www.ncbi.nlm.nih.gov/geo>

→

Sample GSM2142642	
Status	Public on Jun 01, 2016
Title	matched sample_L511N
Sample type	SRA
Source name	Human non-malignant tissue
Organism	<a href="#">Homo sapiens</a>
Characteristics	tumor (t) or normal (n): L511N lung

Sample GSM2142481	
Status	Public on Jun 01, 2016
Title	L511T
Sample type	SRA
Source name	Human NSCLC tissue
Organism	<a href="#">Homo sapiens</a>
Characteristics	tumor (t) or normal (n): L511T stage tnM: 5 histology: 2 surgery date: 2007-04-05 age: 56 gender: female vital date: 2013-04-28 dead: 0 smoking: 3 ps who: 0



Tissue?

# Ontology annotation – adding value to your data



<https://www.ncbi.nlm.nih.gov/geo>

→ GSE81089

Sample GSM2142642	
Status	Public on Jun 01, 2016
Title	matched sample_L511N
Sample type	SRA
Source name	Human non-malignant tissue
Organism	Homo sapiens
Characteristics	tumor (t) or normal (n): L511N

non-tumor

Sample GSM2142481	
Status	Public on Jun 01, 2016
Title	L511T
Sample type	SRA
Source name	Human NSCLC tissue
Organism	Homo sapiens
Characteristics	tumor (t) or normal (n): L511T tumor stage tnM: 5 histology: 2 surgery date: 2007-04-05 age: 56 gender: female vital date: 2013-04-28 dead: 0 smoking: 3 ps who: 0



Sampling site?

# Ontology annotation – adding value to your data



<https://www.ncbi.nlm.nih.gov/geo>

→ GSE81089

Sample GSM2142642	
Status	Public on Jun 01, 2016
Title	matched sample_L511N
Sample type	SRA
Source name	Human non-malignant tissue
Organism	<a href="#">Homo sapiens</a>
Characteristics	tumor (t) or normal (n): L511N

individual 511

Sample GSM2142481	
Status	Public on Jun 01, 2016
Title	L511T
Sample type	SRA
Source name	Human NSCLC tissue
Organism	<a href="#">Homo sapiens</a>
Characteristics	tumor (t) or normal (n): L511T stage tnM: 5 histology: 2 surgery date: 2007-04-05 age: 56 gender: female vital date: 2013-04-28 dead: 0 smoking: 3 ps who: 0

individual 511



Same patient?

# Data curation – adding value to your data

E-GEO-81089 - RNA-seq of 199 non-small cell lung carcinoma patients and 19 paired normal lung tissues

Status	Released on 1 June 2016, last updated on 22 August 2017
Organism	Homo sapiens
Samples (218)	<a href="#">Click for detailed sample information and links to data</a>
Protocols (3)	<a href="#">Click for detailed protocol information</a>
Description	<p>Cancer testis antigens (CTAs) are of clinical interest as biomarkers and present valuable targets for immunotherapy. To comprehensively characterize the CTA landscape of non-small cell lung cancer (NSCLC), we compared RNAseq data of 199 NSCLC tissues to the normal transcriptome of 142 samples from 32 different normal organs. Of 232 CTAs currently annotated in the CTdatabase, 96 were confirmed in NSCLC. To obtain an unbiased CTA profile of NSCLC, we applied stringent criteria on our RNAseq data set and defined 90 genes as CTAs, of which 55 genes were not annotated in the CTdatabase. Cluster analysis revealed that CTA expression is histology-dependent and concurrent expression is common. Immunohistochemistry confirmed tissue specific protein expression of selected genes. Furthermore, methylation was identified as a regulatory mechanism of CTA expression based on independent data from the Cancer Genome Atlas. The proposed prognostic impact of CTAs in lung cancer, was not confirmed, neither in our RNAseq-cohort nor in an independent meta-analysis of 1117 NSCLC cases. Fresh frozen tumor tissue from 199 patients diagnosed with NSCLC and surgically treated 2006-2010 at the Uppsala University Hospital, Uppsala, Sweden and 19 paired normal lung tissues. Clinical data were retrieved from the regional lung cancer registry. Several of the new CTAs are poorly characterized Sample characteristics values represent; pTNM: decided by Hans Brunnstrom, pathologist in Lund Spring 2013 Stage according to pTNM: 1=1a 2=1b 3=2a 4=2b 5=3a 6=3b 7=IV Histology diagnosis spring 2013 HB: 1=squamous cell cancer 2=AC unspecified 3=Large cell/ NOS Surgery date: the date when sample arrived at Patologen UAS Age: age when surgery was performed Vital date: day of death or latest contact Dead: 0=no 1=yes Smoking history : 1=current 2=ex &gt;year 3=never WHO performance status: Performance status 0-4 Please note that the L608T_2122, L771T_1 data columns (in the processed data files) are associated with L608T and L771T samples, respectively.</p>
Experiment type	RNA-seq of coding RNA
Contacts	<a href="mailto:bjorn.hallstrom@scilifelab.se">✉ Bjorn M Hallstrom &lt;bjorn.hallstrom@scilifelab.se&gt;</a> , Akira Saito, Cecilia Lindskog, Dijana Djureinovic, Elisabeth Stahle, Eva Branden, Fredrik Ponton, Hans Brunnstrom, Hirsh Koyi, Jan G Hengstler, Johan Botling, Johanna S Mattsson, Jorg Rahnenfuhrer, Karolina Edlund, Katrin Madjar, Linn Fagerberg, Linnea La Fleur, Masafumi Horie, Mathias Uhlen, Mats Lambe, Patrick Micke, Simon Ekman
Citation	<a href="#">Profiling cancer testis antigens in non-small-cell lung cancer</a> . Djureinovic D, Hallstrom BM, Horie M, Mattsson JS, La Fleur L, Fagerberg L, Brunnstrom H, Lindskog C, Madjar K, Rahnenfuhrer J, Ekman S, Stahle E, Koyi H, Branden E, Edlund K, Hengstler JG, Lambe M, Saito A, Botling J, Ponton F, Uhlen M, Micke P. <i>Europe PMC</i> 27699219
MINSEQE	* * * - *
Exp. design	*
Protocols	*
Variables	*
Processed	-
Seq. reads	*
Files	Investigation description Sample and data relationship Additional data (1)  <a href="#">Click to browse all available files</a>
	<a href="#">E-GEO-81089.idf.txt</a> <a href="#">E-GEO-81089.sdrf.txt</a> <a href="#">E-GEO-81089.additional.1.zip</a>
Links	<a href="#">ENA - SRP074349, GEO - GSE81089</a> <a href="#">Send E-GEO-81089 data to </a>

MAGE-TAB format

# Data curation – adding value to your data

E-GEO-81089 - RNA-seq of 199 non-small cell lung carcinoma patients and 19 paired normal lung tissues

Sample Attributes										Links to Data	
Source Name	organism	individual	sex	age	organism part	disease	sampling site	clinical information	ENA	FASTQ	
GSM2142480 1	Homo sapiens	L504	female	65 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142481 1	Homo sapiens	L511	female	56 (year)	lung	non-small cell lung carcinoma	tumor tissue	non-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142481 1	Homo sapiens	L511	female	56 (year)	lung	non-small cell lung carcinoma	tumor tissue	non-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142642 1	Homo sapiens	L511	female	56 (year)	lung	non-small cell lung carcinoma	non-malignant tissue	non-smoker	<a href="#">n</a>	<a href="#">d</a>	
GSM2142642 1	Homo sapiens	L511	female	56 (year)	lung	non-small cell lung carcinoma	non-malignant tissue	non-smoker	<a href="#">n</a>	<a href="#">d</a>	
GSM2142482 1	Homo sapiens	L529	female	79 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142482 1	Homo sapiens	L529	female	79 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142483 1	Homo sapiens	L530	female	76 (year)	lung	non-small cell lung carcinoma	tumor tissue	smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142483 1	Homo sapiens	L530	female	76 (year)	lung	non-small cell lung carcinoma	tumor tissue	smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142484 1	Homo sapiens	L531	male	74 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142484 1	Homo sapiens	L531	male	74 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142485 1	Homo sapiens	L532	female	69 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142485 1	Homo sapiens	L532	female	69 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142643 1	Homo sapiens	L532	female	69 (year)	lung	non-small cell lung carcinoma	non-malignant tissue	non-smoker	<a href="#">n</a>	<a href="#">d</a>	
GSM2142643 1	Homo sapiens	L532	female	69 (year)	lung	non-small cell lung carcinoma	non-malignant tissue	non-smoker	<a href="#">n</a>	<a href="#">d</a>	
GSM2142486 1	Homo sapiens	L534	male	66 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142486 1	Homo sapiens	L534	male	66 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142487 1	Homo sapiens	L535	male	65 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142487 1	Homo sapiens	L535	male	65 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142488 1	Homo sapiens	L538	male	56 (year)	lung	non-small cell lung carcinoma	tumor tissue	smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142488 1	Homo sapiens	L538	male	56 (year)	lung	non-small cell lung carcinoma	tumor tissue	smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142489 1	Homo sapiens	L539	male	78 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142489 1	Homo sapiens	L539	male	78 (year)	lung	non-small cell lung carcinoma	tumor tissue	ex-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142490 1	Homo sapiens	L541	female	70 (year)	lung	non-small cell lung carcinoma	tumor tissue	non-smoker	<a href="#">t</a>	<a href="#">d</a>	
GSM2142490 1	Homo sapiens	L541	female	70 (year)	lung	non-small cell lung carcinoma	tumor tissue	non-smoker	<a href="#">t</a>	<a href="#">d</a>	

# EMBL-EBI databases – A smart cabinet of curiosities

- Open, online access to experimental datasets
- Annotation of datasets with adequate metadata -> curation
- Use of controlled terms in metadata annotations -> ontologies
- Mechanisms to search for metadata to find relevant experimental results

Open access

Metadata

Ontologies

Discovery

# EMBL-EBI Bioinformatics resources for exploring functional genomics data

# Data standards, curation and ontologies

# Laura Huerta, PhD

## Senior Scientific Curator

[lauhuema@ebi.ac.uk](mailto:lauhuema@ebi.ac.uk)

9 November 2017