

# Facilitating reproducible research:



from the perspective of a public data archive,  
ArrayExpress, at EMBL-EBI

Amy Tang

Functional Genomics Curation and  
Training Project Leader

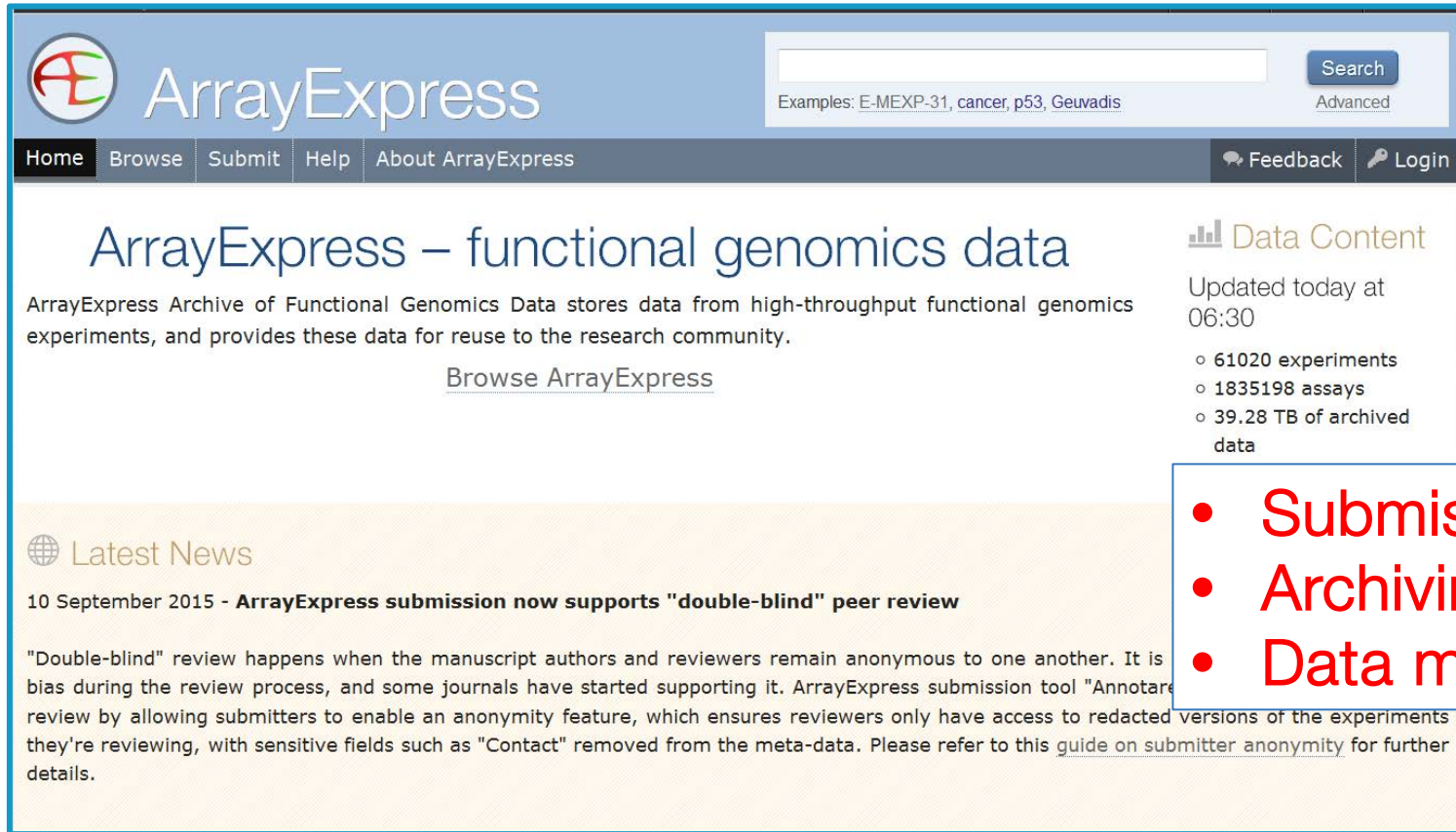
EMBL-EBI

# Talk outline

- What is ArrayExpress? 
- Facilitating reproducibility
  1. Submission timing is key 
  2. The ideal world: standards, and ontologies.
  3. When standards fail --- why we need curators
  4. When the dust settles --- meta-data updates
- Further challenges

# ArrayExpress at EMBL-EBI

[www.ebi.ac.uk/arrayexpress](http://www.ebi.ac.uk/arrayexpress) (daily release at 6am UK time)



The screenshot shows the ArrayExpress website. At the top, there is a header with the ArrayExpress logo and name, a search bar with a 'Search' button, and a navigation menu with links: Home, Browse, Submit, Help, About ArrayExpress, Feedback, and Login. Below the header, the main content area features the title 'ArrayExpress – functional genomics data' and a description: 'ArrayExpress Archive of Functional Genomics Data stores data from high-throughput functional genomics experiments, and provides these data for reuse to the research community.' A link 'Browse ArrayExpress' is provided. To the right, a 'Data Content' section shows statistics: 'Updated today at 06:30', '61020 experiments', '1835198 assays', and '39.28 TB of archived data'. Below this, a 'Latest News' section dated '10 September 2015' announces that 'ArrayExpress submission now supports "double-blind" peer review'. The news text explains that 'Double-blind' review keeps authors and reviewers anonymous to each other to prevent bias, and that ArrayExpress has implemented this by allowing submitters to enable an anonymity feature that redacts sensitive information like contact details from the meta-data. A callout box on the right side of the screenshot lists three key features: Submission, Archiving, and Data mining.

ArrayExpress – functional genomics data

ArrayExpress Archive of Functional Genomics Data stores data from high-throughput functional genomics experiments, and provides these data for reuse to the research community.

[Browse ArrayExpress](#)

**Data Content**

Updated today at 06:30

- 61020 experiments
- 1835198 assays
- 39.28 TB of archived data

**Latest News**

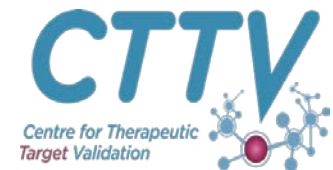
10 September 2015 - **ArrayExpress submission now supports "double-blind" peer review**

"Double-blind" review happens when the manuscript authors and reviewers remain anonymous to one another. It is designed to prevent bias during the review process, and some journals have started supporting it. ArrayExpress submission tool "Annotare" supports double-blind review by allowing submitters to enable an anonymity feature, which ensures reviewers only have access to redacted versions of the experiments they're reviewing, with sensitive fields such as "Contact" removed from the meta-data. Please refer to this [guide on submitter anonymity](#) for further details.

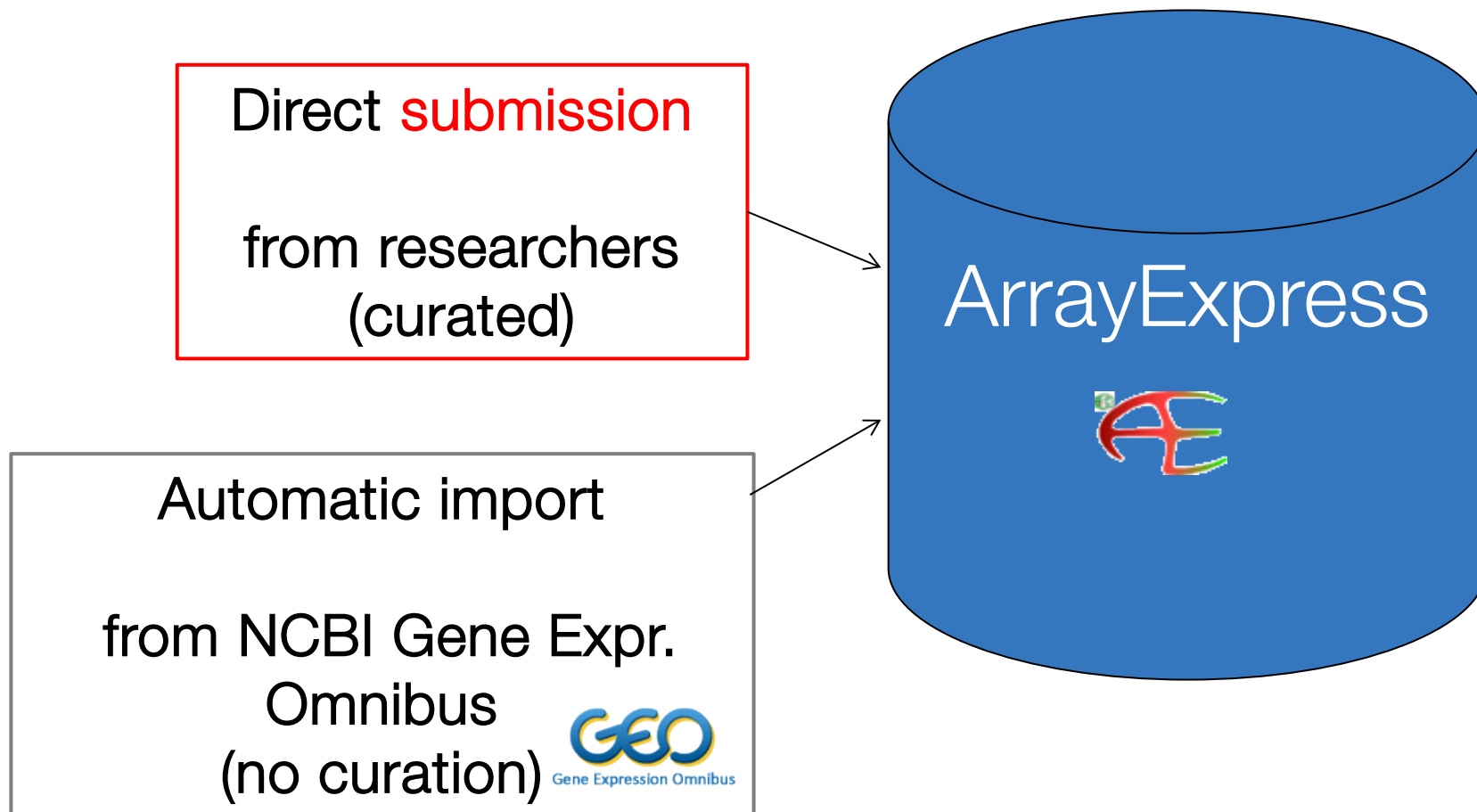
- Submission
- Archiving
- Data mining

“PubMed for data sets”

# Find us on the Wellcome Genome Campus



# ArrayExpress data sources

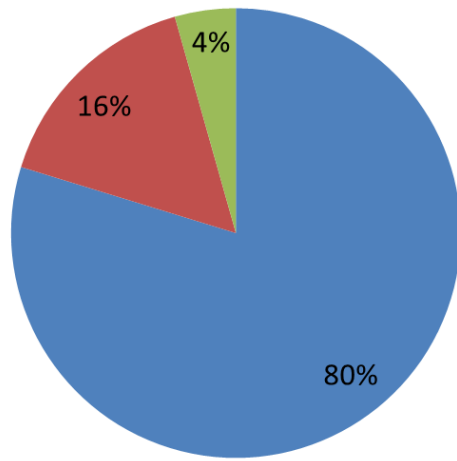




# Data volume in ArrayExpress

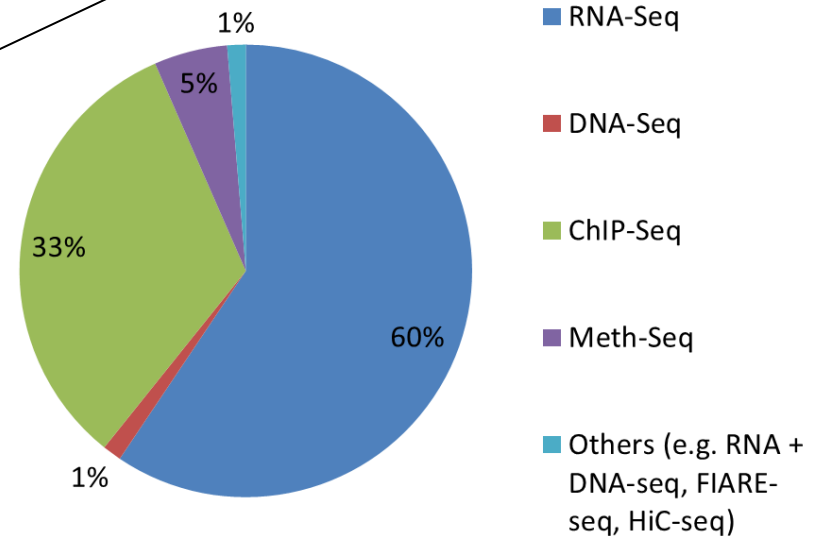
~61600 experiments, ~16% direct submissions, the rest imported

■ Microarray ■ HTS ■ others/mixed tech



Microarray vs HTS

RNA-, DNA-, ChIP-seq breakdown

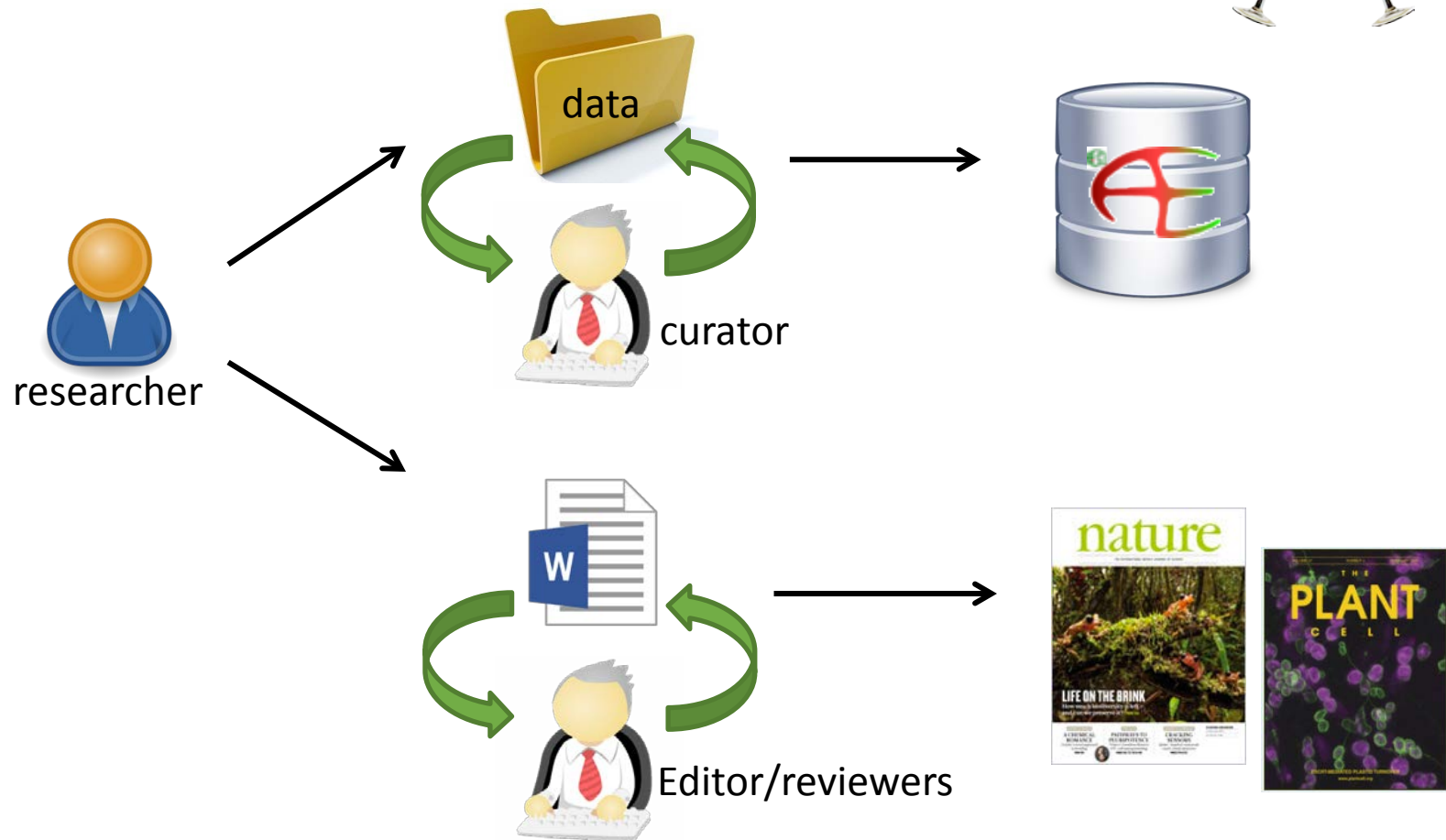


(Pie charts as of 13 Nov 2015)

# Submission timing

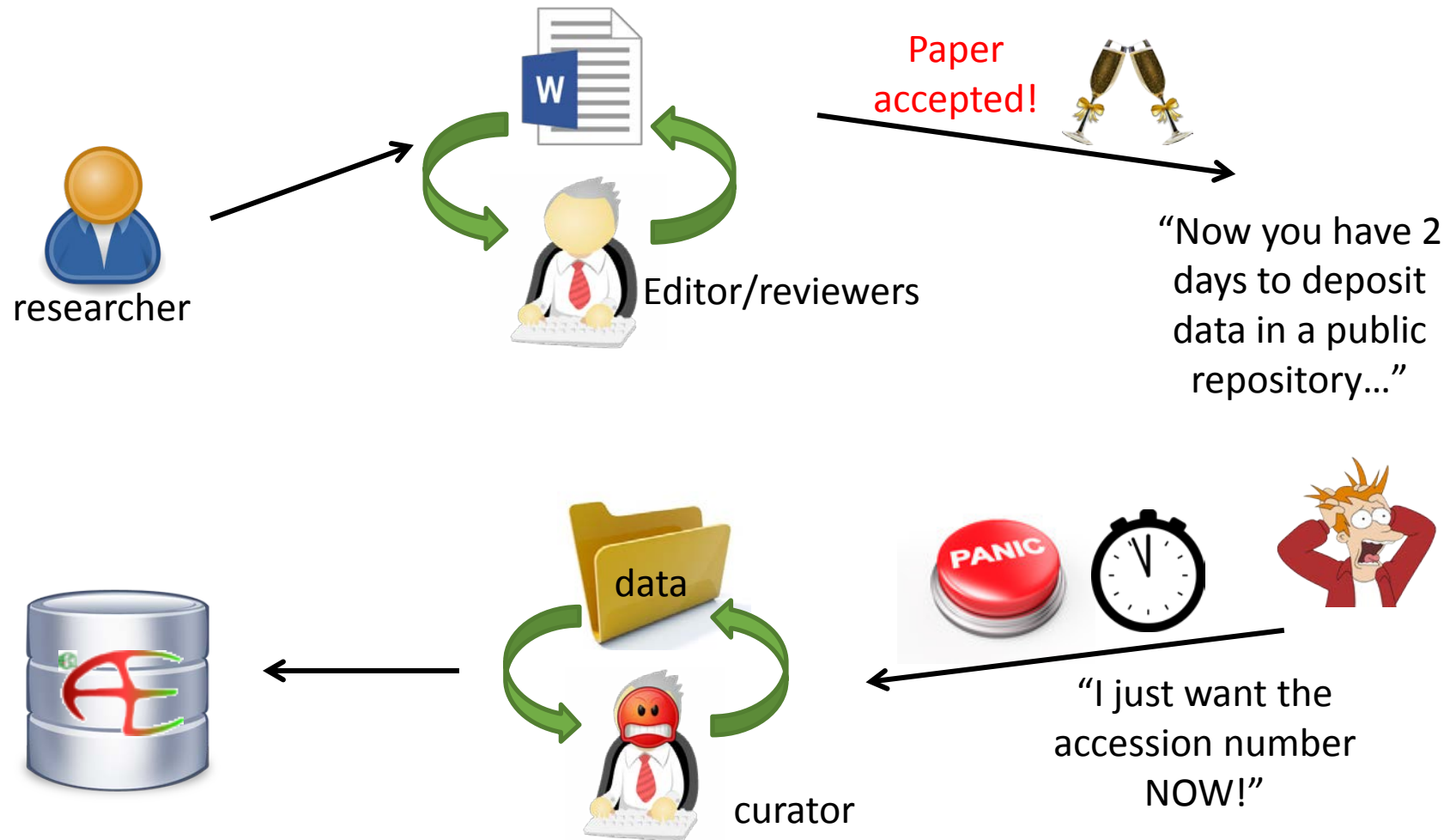
# Submit to ArrayExpress – expected timing

Submission → Curation/Review → Publication





# Submit to ArrayExpress – common scenario



# Submitting to ArrayExpress – when?

*As explained, our paper is **conditionally accepted in Science** and they want the accession numbers before the proof stage. Since the data files and metadata have now been submitted could you please assign and send us an accession number?*

*“The issue is that we need an accession number **as soon as possible.**”*

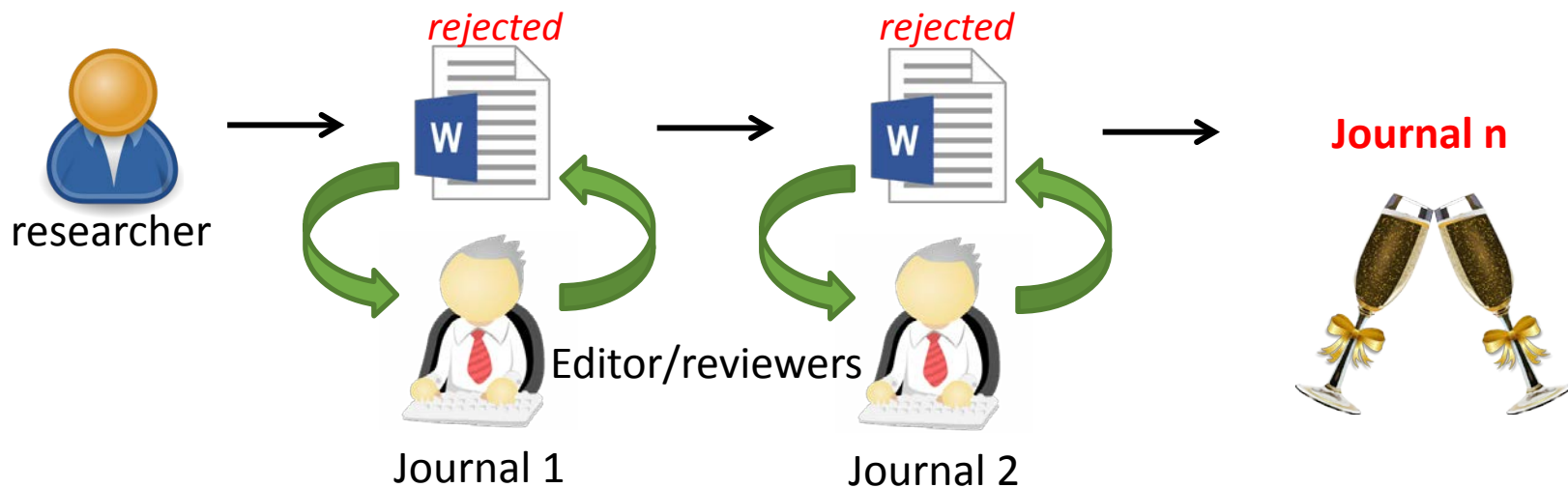
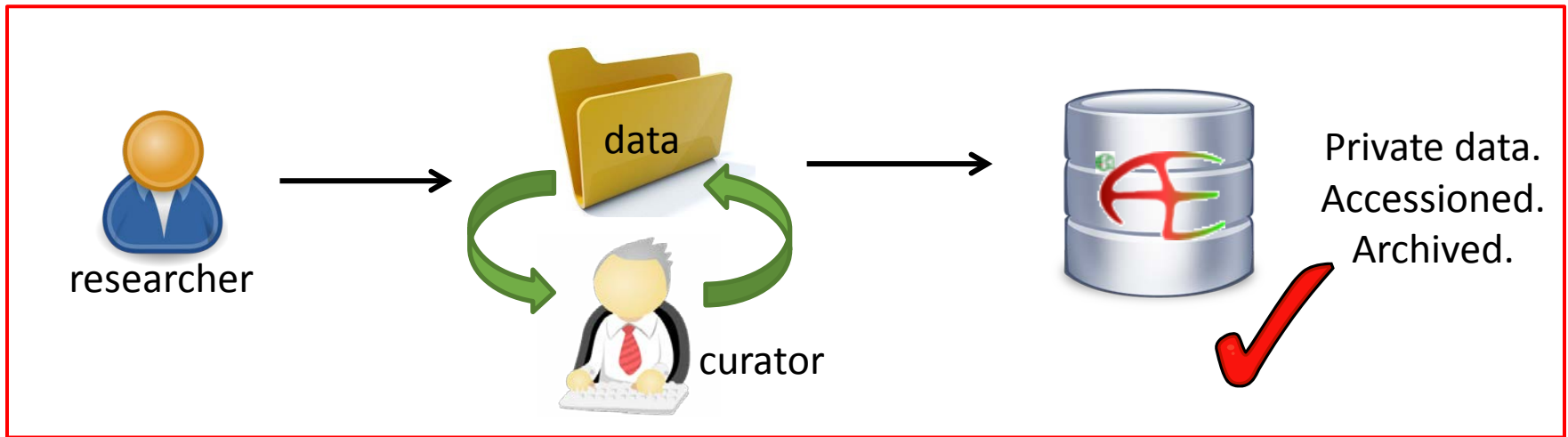
*Dear Sir/Madam,*

*I have done a MicroArray submission today entitled: “xxxxx.” I would like to request you for **providing me its accession number today, as I intend to submit the concerned manuscript by today.***

*Many thanks for your time and considerations.*



# Submit to ArrayExpress – do it early, park it



# Standards and ontologies

# Data standards – where we're at in 2001...



© 2001 Nature Publishing Group <http://genetics.nature.com>

*commentary*

## 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>

Microarray analysis has become a widely used tool for the generation of gene expression data on a genomic scale. Although many significant results have been derived from microarray studies, one limitation has been the lack of standards for presenting and exchanging such data. Here we present a proposal, the Minimum Information About a Microarray Experiment (MIAME), that describes the minimum information required to ensure that microarray data can be easily interpreted and that results derived from its analysis can be independently verified. The ultimate goal of this work is to establish a standard for recording and reporting microarray-based gene expression data, which will in turn facilitate the establishment of databases and public repositories and enable the development of data analysis tools. With respect to MIAME, we concentrate on defining the content and structure of the necessary information rather than the technical format for capturing it.

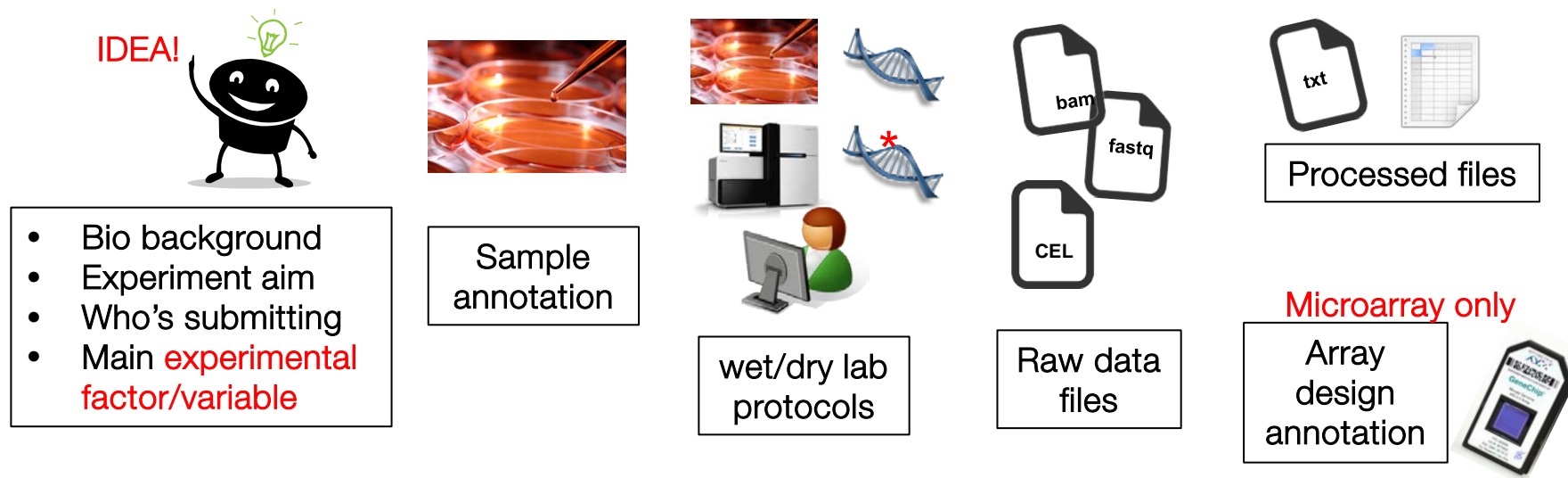
### Introduction

After genome sequencing, DNA microarray analysis<sup>1</sup> has become

cult, because at present, microarrays do not measure gene expression levels in any objective units. In fact, most measurements report

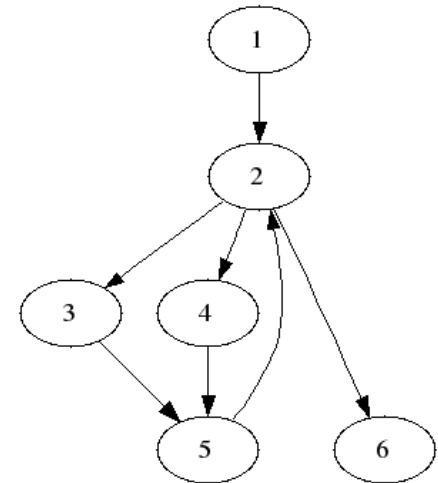
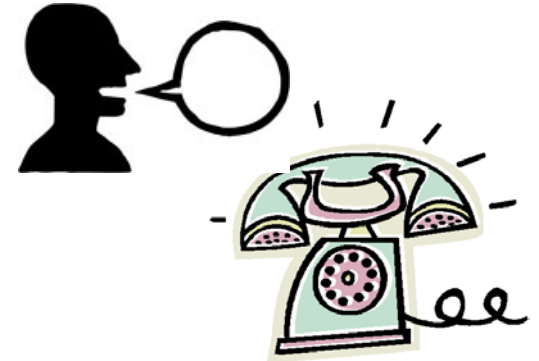
# Data standards

- **MIAME** = **M**inimal **I**nformation **A**bout a **M**icroarray **E**xperiment ([http://www.mged.org/Workgroups/MIAME/miame\\_2.0.html](http://www.mged.org/Workgroups/MIAME/miame_2.0.html))
- **MINSEQE** = **M**inimal **I**nformation about a high-throughput **N**ucleotide **SEQ**uencing **E**xperiment (<http://www.mged.org/minseqe>)
- **Aim:** capture information from every stage of an experiment





# Various ways submitters tried to comply with our standards



# Fulfil standards with submission tool “Annotare”

[http://www.ebi.ac.uk/fgpt/annotare\\_help/](http://www.ebi.ac.uk/fgpt/annotare_help/)

Webforms asks for information based on the content checklist:  
MIAME/MINSEQE

**EXPERIMENT: UNACCESSIONED** Help Feedback Validate Submit to ArrayExpress

Experiment Description Samples and Data IDF Preview SDRF Preview

General Information

Contacts

Title \* RNA-Seq CAGE (Cap Analysis of Gene Expression) analysis of human tissues in RIKEN FANTOM5 project

Description \* This experiment captures the expression data reported by the RIKEN FANTOM5 project ( <http://fantom.gsc.riken.jp/5/> ), focusing on tissue/organism part data which was deposited in the sequence read archive (SRA) under study accssion DRP001031 ( <https://www.ebi.ac.uk/ena/data/view/DRP001031> ). The samples in this experiment can also be found on a dedicated page of the FANTOM website: [http://fantom.gsc.riken.jp/5/sstar/Browse\\_samples](http://fantom.gsc.riken.jp/5/sstar/Browse_samples) . Since this is CAGE analysis, gene expression data is reported by FANTOM5 in TPMs (tags per million) for gene promoters.  
(at least 50 characters)

ArrayExpress Experiment Type \* RNA-seq of coding RNA

**Annotare 2.0**

**EXPERIMENT: E-MTAB-2407** Help Feedback

Experiment Description Samples and Data IDF Preview SDRF Preview

Create samples, add attributes and experimental variables

Create extracts and assign ENA library info

Upload and assign data files

Protocols

High-throughput sequencing

Sample Attributes and Variables... Add Sample Delete Samples Fill Down Value Import Values

<input type="checkbox"/>	Name	Organism	Disease	Material Type	Organism Part	Sex	Individual
<input type="checkbox"/>	Sample 1	Homo sapiens	kidney neoplasm	organism part	kid	male	patient 1
<input type="checkbox"/>	Sample 2	Homo sapiens	normal	organism part	kidney (UBERON_0002113)		patient 1
<input type="checkbox"/>	Sample 3	Homo sapiens	kidney neoplasm	organism part		male	patient 2
<input type="checkbox"/>	Sample 4	Homo sapiens	kidney neoplasm	organism part		male	patient 3
<input type="checkbox"/>	Sample 5	Homo sapiens	normal	organism part		male	patient 2
<input type="checkbox"/>	Sample 6	Homo sapiens	normal	organism part		male	patient 3

Automatically generates  
MAGE-TAB  
format  
spreadsheet

# Sample annotation



- **State the “obvious”**

e.g. mouse strain, dose of drug, bio. reps vs tech reps

- **Include confounding variables**

e.g. sex/gender in clinical studies


- **Avoid context-specific acronyms**

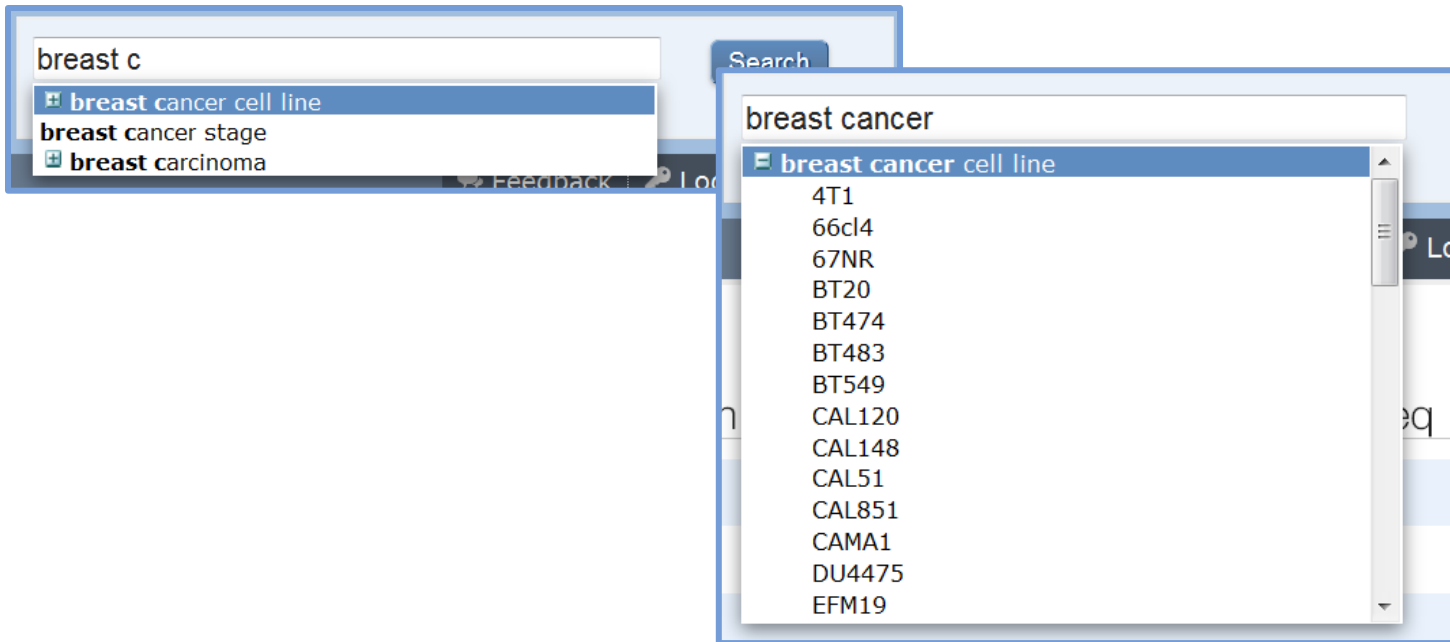
e.g. “m” =  ?  ?  ? or 

- **Consult submitters to define required meta-data for emerging technologies**

e.g. single-cell RNA-seq

# Annotate samples with ontology terms

- Ontology = controlled vocabulary with hierarchical relationships
- We mainly use Experimental Factor ontology (EFO)  Experimental Factor Ontology



- We request new terms regularly

# Annotate samples with ontology terms – why?

## Avoid acronyms and typographical errors

- Choose terms from dropdown, e.g. NCBI Taxonomy species

Organism
Escherichia coli O157
Escherichia coli O158
Escherichia coli O159
Escherichia coli O160

## Avoid ambiguity – standardise terminology

- Developmental stage (different nomenclature systems)
- Mouse phenotype (IMPC project)
- Plant ontology (various plant “treatments”)

# How hard is it to agree on how we say “female”?

*female*

*F*

*femme*

2



# How many ways can you say “female”?



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(gynocious)	remale	metafemale
f	femele	semi-engorged female	sterile female
famale	female, pooled	sexual oviparous female	normal female
femal	femalen	sterile female worker	sf
female	females	strictly female	vitellogenic replete female
female - worker	females only	tetraploid female	worker
female (alate sexual)	gynocious	thelytoky	hexaploid female
female (calf)	healthy female	female (gynocious)	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*

# How many ways can you say “male”?



37 year old male	initial phase male	male fetus	six males mixed
600 yr. old male	m	male plant	stallion
adult male	make	male, 8 weeks old	steer
bull	makle	male, castrated	sterile male
castrated male	mal e	male, pooled	strictly male
cm	male	males	tetraploide male
dioecious male	male (7-2872)	man	type i males
diploid male	male (7-3074)	men	type ii males
drone	male (m-a)	normale male	virgin male
engorged male	male (m-o)	ram	winged and wingless males
fertile male	male caucasian	rooster	young male
four males mixed	male child	s1 male sterile	
individual male	male fertile	sex: male	

male (note: this sample was originally provided as a \female\ 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 male individual)

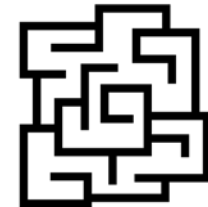
*Courtesy of N. Silvester, European Nucleotide Archive, EMBL-EBI*

# Protocols



## ➤ Include full protocols

- Quoting pre-published papers – can't verify.
- Quoting paper 1 → paper 2 → paper n : trail lost!
- “Where is the protocol in the 50-page PDF!?”
- Often too generic:



*“Standard protocol” ;*

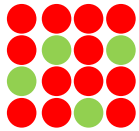
*“... carried out according to manufacturer's instructions”*

# Protocols



- Include supporting files, e.g.
  - Very detailed protocols (run over pages)
  - spike-in sequences in single-cell RNA-seq
  - gtf file for transcript quantification
  - scripts used to process the data
  - README file describing supporting files

# Raw data – rigorous check



Correct array design assignment?

## Fastq files



- Trimmed reads? (read length = 1bp!)
- Check mate pairs (if paired end)

## Bam files

- Contain “all” reads (mapped + unmapped)

# NGS raw data – stable archiving at INSDC

- Standard SRA-format fastq
- Daily data exchange/mirroring



International Nucleotide  
Sequence Database  
Collaboration (INSDC)  
running the sequence  
read archive (SRA)





# When standards aren't met...

# The gatekeepers, aka the curators



Amy Tang



Maria Keays



Melissa Burke



Anja Fullgrabe



Laura Martinez



Satu Koskinen

Curators/Bioinformaticians

# How submitters see us ....



Amy Tang



Maria Keays



Melissa Burke



Anja Fullgrabe



Laura Martinez



Satu Koskinen

Curators/Bioinformaticians

# Curator: the “gap filler”



Sample	genotype	treatment
sample 1	DT	1
sample 2	DT	1
sample 3	DT	2
sample 4	DT	2
sample 5	DS	1
sample 6	DS	1
sample 7	DS	2
sample 8	DS	2



- Genotype of what gene? What is “DT”?
- What is treatment “1” ?
- Any controls?



# Curator: the “cleaner”



Sample Attributes and Variables...					Add Sample	Delete Samples	Fill Down Value	Import Values
<input type="checkbox"/>	Name	Organism	Material Type	Description				
<input type="checkbox"/>	Sample 1	Homo sapiens	cell	islet cells isolated from 23-year-old male with type I diabetes				
<input type="checkbox"/>	Sample 2	Homo sapiens	cell	islet cells isolated from 23-year-old male with type I diabetes				
<input type="checkbox"/>	Sample 3	Homo sapiens	cell	islet cells isolated from 28-year-old female with type I diabetes				
<input type="checkbox"/>	Sample 4	Homo sapiens	cell	islet cells isolated from 24-year-old female with type I diabetes				
<input type="checkbox"/>	Sample 5	Homo sapiens	cell	islet cells isolated from 21-year-old male with type I diabetes				
<input type="checkbox"/>	Sample 6	Homo sapiens	cell	islet cells isolated from 25-year-old male with type I diabetes				

Sample Attributes and Variables...					Add Sample	Delete Samples	Fill Down Value	Import Values
<input type="checkbox"/>	Name	Material Type	Organism	Disease	Cell Type	Sex	Age (year)	
<input type="checkbox"/>	Sample 1	cell	Homo sapiens	type I diabetes mellitus	islet cell	male	23	
<input type="checkbox"/>	Sample 2	cell	Homo sapiens	type I diabetes mellitus	islet cell	male	22	
<input type="checkbox"/>	Sample 3	cell	Homo sapiens	type I diabetes mellitus	islet cell	female	28	
<input type="checkbox"/>	Sample 4	cell	Homo sapiens	type I diabetes mellitus	islet cell	female	24	
<input type="checkbox"/>	Sample 5	cell	Homo sapiens	type I diabetes mellitus	islet cell	male	21	
<input type="checkbox"/>	Sample 6	cell	Homo sapiens	type I diabetes mellitus	islet cell	male	25	

# Curator: the “investigator”



Sample Name	Sex
Patient A	male
Patient A	female
Patient B	male
Patient B	male

- Organism: *Equus caballus*
- Age: 70 (years)

- Sex: male
- Organism part: endometrium

- Lung sample from one mouse
- Sex = “mixed sex”.





# When the dust (finally) settles...

# Contact info and citation

- Contact update – people move!
- Citation update
- Release policy: we crawl literature for published ArrayExpress accessions; turn private data sets to public if the paper is out.



# Continuing challenges

- We have no power to reject data sets 
- Submitters don't treat us like journal editors/reviewers
- Deposition may be mandatory, standards aren't
- Non-obvious errors will pass without suspicion
- Secure funding to run archival databases 

# Faces behind ArrayExpress

Robert  
Petryszak



Ugis  
Sarkans



Amy Tang



Maria Keays



Melissa Burke



Anja Fullgrabe

Curators/Bio-  
informaticians



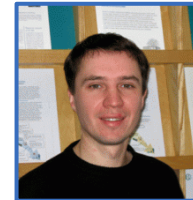
Satu Koskinen



Laura Martinez



Ahmed Ali



Nikolay Kolesnikov  
("Mr Annotare")



Catherine Snow



Mirosław Dyląg

Data management

web