

EBI is an Outstation of the European Molecular Biology Laboratory.

Proteomics data sharing: why?

- Since the Enlightenment, Science is built upon the sharing of knowledge and information
“If I have looked further, it is because I stood on the shoulder of Giants”
-- Sir Isaac Newton, paraphrasing Bernard of Chartres
- Data producers are not always the best data analysts
Sharing of data allows analysts access to real data, and in turn allows better analysis tools to be developed
- Meta-analysis of data can recycle previous findings for new tasks
putting findings in the context of other findings increases their scope
- Sharing data allows independent review of the findings
When actual replication of an experiment is often impossible, a re-analysis or spot checks on the obtained data become vitally important
- Simple economics
“Information, no matter how expensive to create, can be replicated and shared at little or no cost.” -- Thomas Jefferson

Simply sharing data is not enough...

Table 1. Identities of stress-induced proteins

Spot ID	Synonym	Function
1202	SCO0525	Hypothetical p
3307	SCO2988	UDP-glucose
3509	SCO2180	Putative dihyd
6413	SCO6027	Probable acet
6419	SCO1494	3-Dehydroquin
6823	SCO5477	Putative oligo
118	SCO1340	Conserved hy
1104	SCO2368	Conserved hy
1617	SCO5373	ATP synthase
2601	SCO5373	ATP synthase
3616	SCO5371	ATP synthase
5721	SCO4814	Bifunctional p
1515	SCO2180	Putative dihyd
1616	SCO3661	Putative chap
2706	SCO3671	Heat shock p
2906	SCO5999	Aconitase
3504	SCO1936	Putative trans
5310	SCO0506	NH(3)-depend
7417	SCO5477	Putative oligo
505	SCO1998	30S ribosoma
1711	SCO1352	Xaa-pro amin
2618	SCO0681	Putative ferred
2722	SCO1998	30S ribosoma
4407	SCO5113	Oligopeptide
4509	SCO2390	Beta-ketoacyl
1803	SCO2181	2 Oxoglutarat
2113	SCO4277	Hypothetical p
3101	SCO3899	Hypothetical p
4309	SCO1081	Putative elect
4512	SCO5212	3-Phosphoshi
5514	SCO3629	Putative aden

+ SOD1

Table 1. Identification of exosomal proteins based on MALDI-TOF peptide mass fingerprinting or MS/MS-derived sequences

Band (Fig. 1)	Protein Name	Identification Method ^a	Accession Number ^b	Molecular Mass (kDa)	Matching Peptides	Sequence Coverage (%)
A	Mac-1 α -chain = CD11b					
1	Complement C3 ^c	MS/MS (7)	4093220			
1	PK-120 ^c	MS/MS (3)	Not in databases			
1	α 2-Macroglobulin ^f	MS/MS (2)	Not in databases			
2	Plastinogen ^e	MS	P06868	91	28	37
3	Alix	MS	6755002 ^d	96	26	34
3	Mac-1 β -chain = CD18	MS/MS (6)				
4		MS	P11835	85	27	38
4	hsp90- β = hsp84	MS/MS (1)				
5	Serum albumin ^e	MS	P11499	83	30	38
		MS	P02769	69	42	66
		MS/MS (3)				
B	hsc73					
B, C	MFG-E8/lactadherin					
6	Tubulin β	MS	P05218	50	20	44
7	Annexin VII = synexin	MS	Q07076	50	6	13
		MS/MS (3)				
7	Bovine coagulation factor X ^c	MS/MS (2)	P00743	54		
7	PEPF ^e	MS/MS (3)	Q95121	46		
7	Tumor susceptibility protein (tsq) 101	MS/MS (2)	3184260 ^d	44		
7	Rab GDP dissociation inhibitor (GDI) 3	MS	Q61598 ^e	51	10	21
7	Elongation factor (EF) 1- α -1	MS/MS (2)	P10126	50		
7	EIF-4A-II	MS/MS (2)	P10630	46		
8	Annexin I	MS	P10107	39	7	25
8	Reverse transcriptase/pol (murine leukemia virus)	MS/MS (1)	61790 ^d			
D	γ -Actin					
E	G protein G _{2e} subunit					
F	Annexin II					
9	Annexin V	MS	P48036	36	16	54
		MS/MS (4)				
10	Annexin IV	MS	P97429	36	20	63
		MS/MS (4)				
10	Galactin-3 = Mac-2	MS	P16110	27	11	37
10		MS/MS (6)				
11	Syntenin	MS	2197106 ^d	32	17	35
		MS/MS (6)				
G	Gag polyprotein (murine leukemia virus)					
G	MHC class II β -chain					
12	14-3-3 protein η	MS	P11576	28	21	68
		MS/MS (4)				
12	14-3-3 protein γ 6	MS	P35215	28	20	63
		MS/MS (2)				
12	14-3-3 protein γ	MS/MS (1)	3065929 ^d			
13	Apolipoprotein A-I ^c	MS	P15497	30	25	67
H	CD9					
14	Thioredoxin peroxidase II	MS	P35700	22	8	43
		MS/MS (6)				
14	Rab 11	MS/MS (1)	P46638	24		
14	κ -Casein ^e	MS/MS (2)	P02668	21		
15	Rab-7	MS	P51150	24	5	26
		MS/MS (3)				
16	Fermitin light chain ^f	MS	O46415	20	15	73
16	Rap1-B	MS	P09526	21	14	57
17	Cofilin	MS	P18760	19	10	50
18	Histone H3	MS	Z85979 ^e	15	7	45
19	Histone H2B	MS	P10853	14	13	82
19	Histone H2A	MS	P20670	14	12	67
20	Histone H4	MS	S0626 ^e	11	15	90
20	Profilin I	MS	P10924	15	11	60
21	Hemoglobin γ -chain ^f	MS	P02081	16	16	74
21	Hemoglobin α -chain ^f	MS	P01966	15	9	66

	(kDa/pI)	Accession no.	Species
	Experimental	Theoretical	
hase			
Stress			
CS	69.07/5.8	BAA07338	<i>Arabidopsis thaliana</i>
CS			
CS	62.82/6.4	<u>P42893</u>	<i>Oryza sativa</i>
CS			
CS	65.29/6	<u>Q43097</u>	<i>Lotus japonicus</i>
CS			
EtOH	62.82/6.4	<u>P42893</u>	<i>Oryza sativa</i>
EtOH			
EtOH	67.71/6.9	<u>Q00775</u>	<i>Salanum tuberosum</i>
EtOH			
EtOH	67.69/7.1	BAA77351	<i>Triticum aestivum</i>
EtOH			
HS	62.64/8.5	Q42608	<i>Brassica rapa</i>
HS			
HS	54.96/5.2	Q36881	<i>Acetabularia acetabulum</i>
HS			
HS	57.44/7	P30567	<i>Gossypium hirsutum</i>
HS			
HS	57.94/8	P37215	<i>Lycopersicon esculentum</i>
NaCl			
NaCl			
NaCl	49.59/7.1	S33520	Soybean
NaCl	43.04/6.1	P51110	<i>Lycopersicon esculentum</i>
NaCl			
NaCl	38.79/6.2	P51110	<i>Vitis vinifera</i>
NaCl			
P1	27.54/8.8	BAB03428	<i>Oryza sativa</i>
P1	27.54/8.8	BAB03428	<i>Oryza sativa</i>
P1	24.36/8.6	BAA02870	<i>Oryza sativa</i>
P1	26.58/6.4	P09886	<i>Rhus sativum</i>
P1			
P1			
P1			
	56.77/6.1	P55238	<i>Hordeum vulgare</i>

A nuance: *available* data vs. *accessible* data

When data is only made available as arbitrarily formatted PDF tables, it carries important limitations

- Source data (e.g., mass spectra) are not made available
 - No peer review validation possible (*the numbers game*)
 - Very little raw materials for testing innovative *in silico* techniques are available (*data hoarding*)
- Automated (re-)processing of the results (e.g., identifications) is impossible (*eliminating objective technique comparison*)
- Data producers do not actually feed their results and knowledge back to the community (*evading responsibility for the results*)

Accessibility requires proper infrastructure

- **Community supported, standardized data formats**

Necessary to allow efficient access to the data

- **Controlled vocabularies (CV's) and ontologies**

To provide unambiguous context and metadata to the actual data, as well as enabling powerful queries to be performed on the data

- **Minimal reporting requirements for specific data types**

Ensures the presence of certain bits of information without which interpretation is ambiguous, hampered or impossible

- **Publicly available, online repositories**

Bioinformatics grew up along side the internet, and this is reflected in the successful online data sharing mechanisms already in place in the life sciences. *The repositories should implement the standards, use the CV's and ontologies, and adhere to the minimal requirements.*

Community standards for proteomics



The Human Proteome Organisation (HUPO)
Proteomics Standards Initiative (PSI)



<http://www.psidev.info>

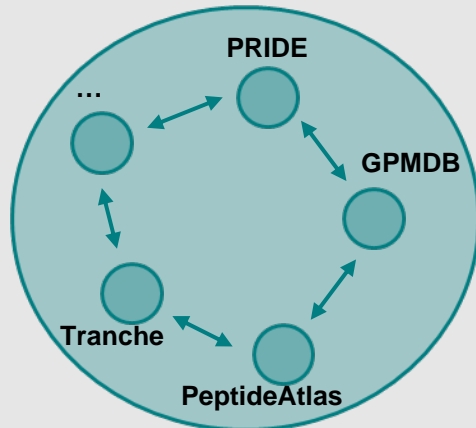
- Creates standard formats, CV's and ontologies, and minimal requirements
- Composed of several workgroups

<i>Molecular Interactions</i>	<i>(MI)</i>	<i>PSI-MI format v2.5</i>
<i>Mass Spectrometry</i>	<i>(MS)</i>	<i>mzData, mzML format</i>
<i>Sample Processing</i>	<i>(SP)</i>	<i>alpha stage</i>
<i>Gel Electrophoresis</i>	<i>(Gel)</i>	<i>GelML format</i>
<i>Proteomics Informatics</i>	<i>(PI)</i>	<i>analysisXML format</i>
<i>Protein Modifications</i>	<i>(Mod)</i>	<i>PSI-MOD ontology</i>

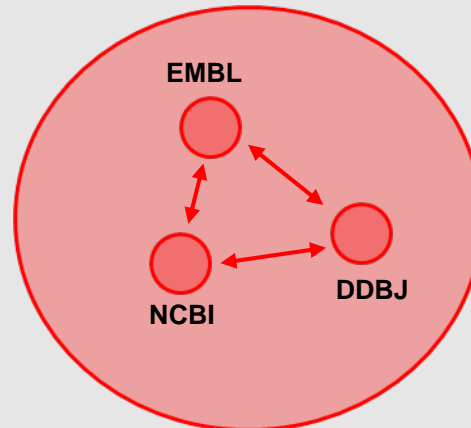
Furthermore, superstructures must be built

Often, multiple repositories will emerge more or less simultaneously in a particular field. By exchanging data, and by collaborating on data acquisition an increase in coverage as well as a more comprehensive dataset is obtained by each individual resource.

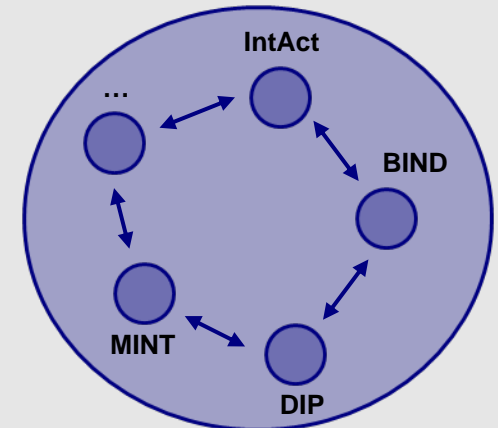
Such superstructures do require additional infrastructure, however.



mass spec
ProteomExchange



sequence databases
(INSDC)

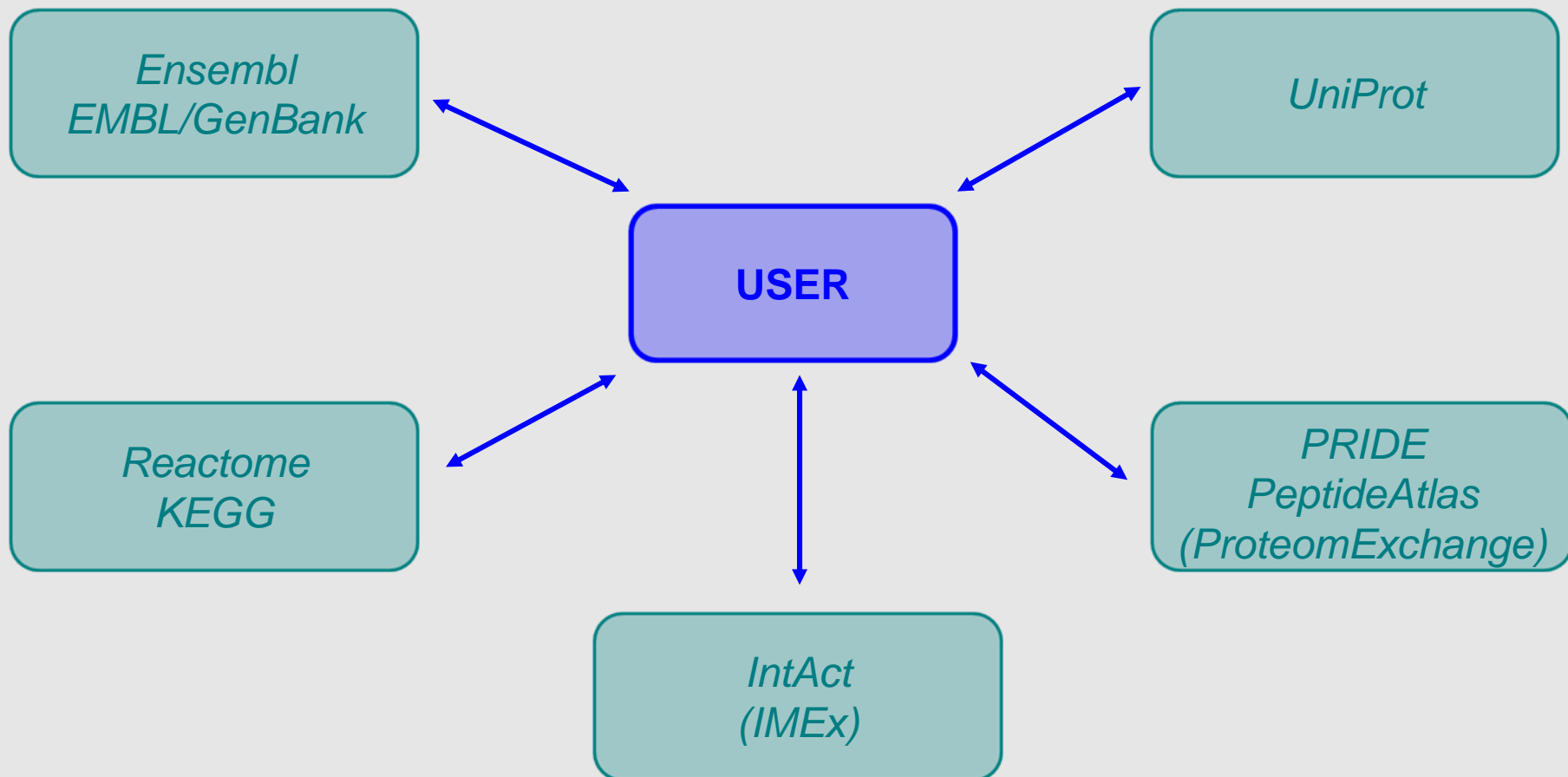


interactions
IMEx

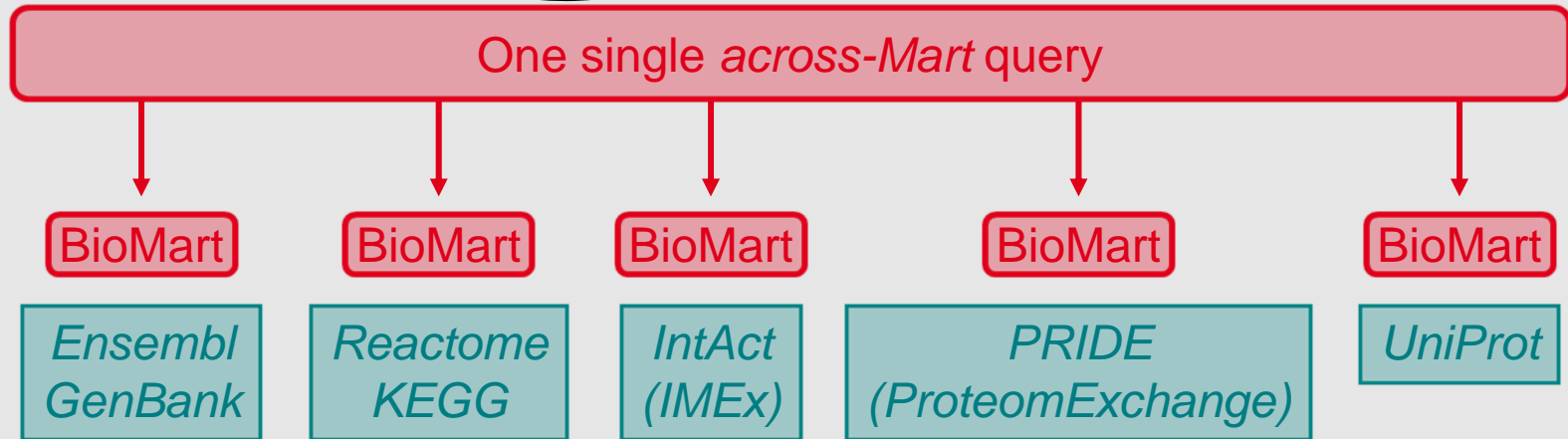
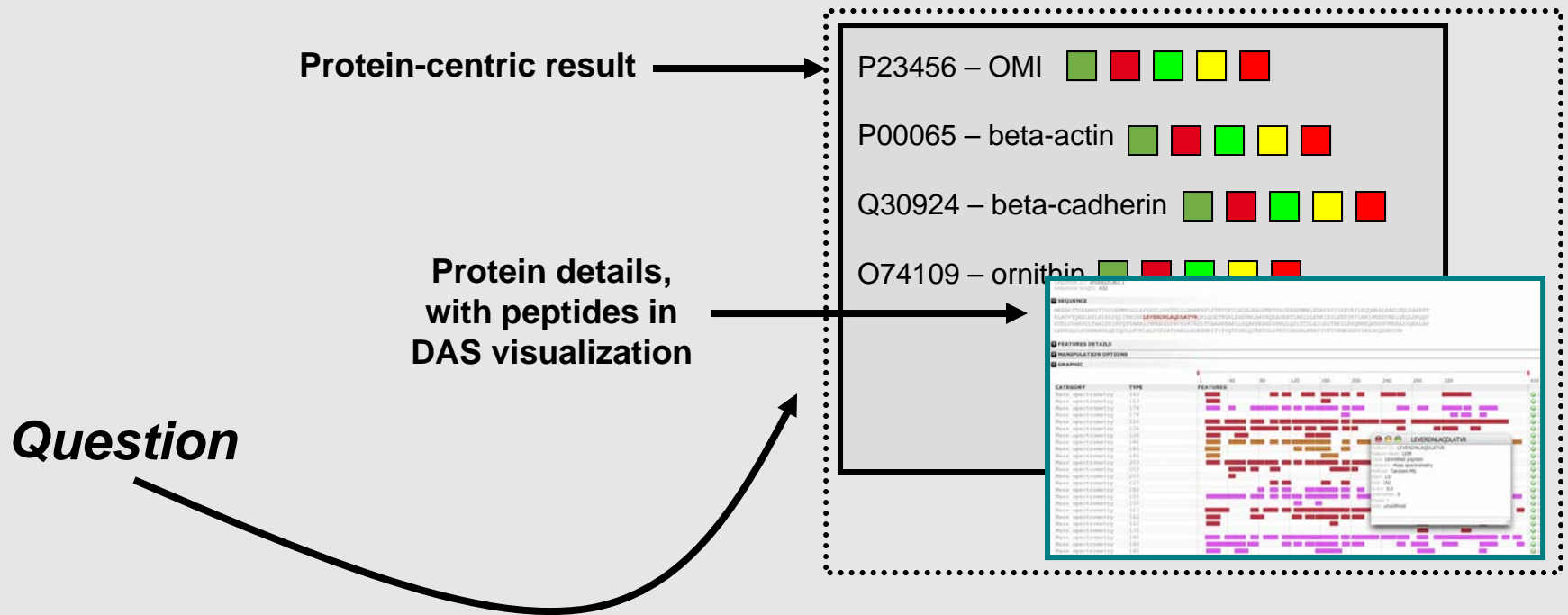
The final goal: cross-domain integration

Current situation

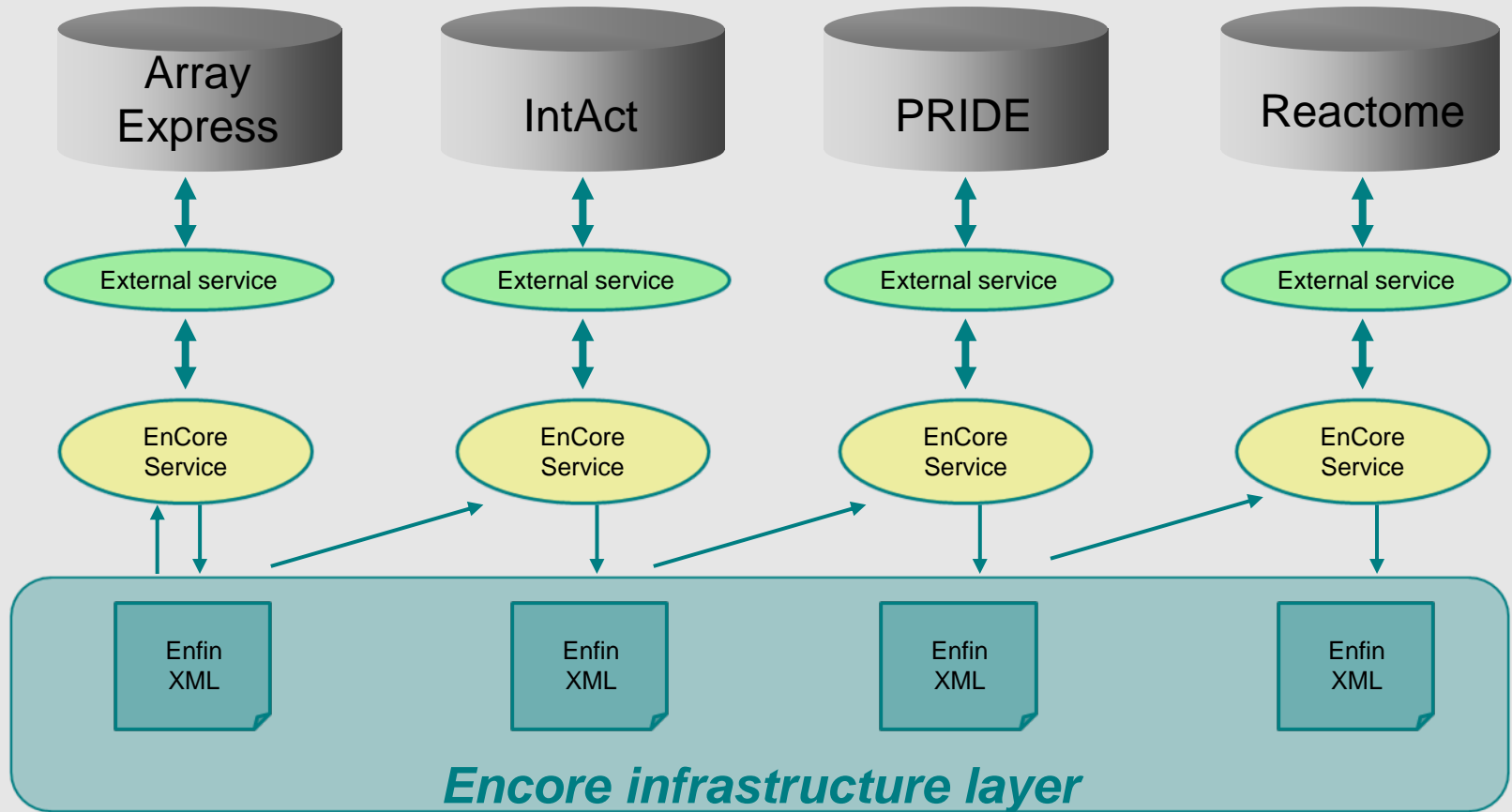
(requires user expertise in querying all the different repositories)



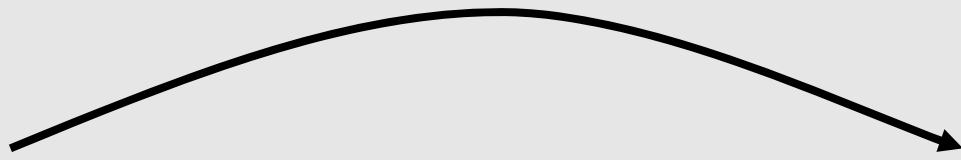
BioMart as a common interface layer



Enfin/Encore as a cross-domain aggregator



Question



**EnVision interface
for results**



How do we make this all happen?

- **Journal guidelines**

Journal guidelines heavily influence the decisions taken by authors; by first requesting and subsequently mandating data submission to established repositories, they provide an important stick.

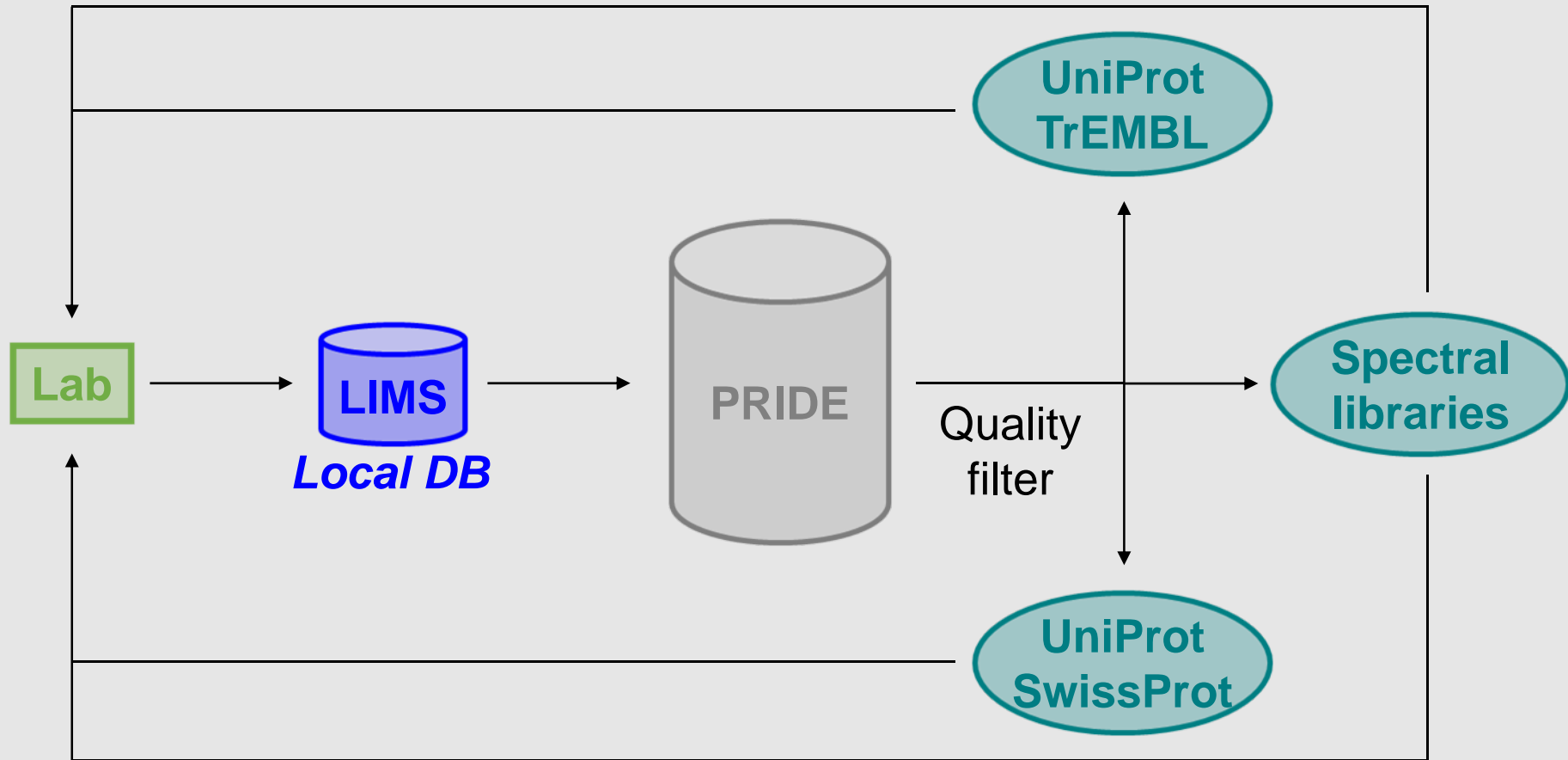
- **Funder support and guidelines**

Funders contribute both sticks and carrots. The sticks lie in the grant application guidelines; they can require a plan for data management and dissemination. The carrot is in providing specific funding for this aspect of science.

- **Data repositories**

The availability of reliable, freely available repositories is key; submission thresholds should be kept low and added value needs to be provided. Furthermore, feedback loops need to be established in order to ensure that accumulated data flows back to the user community. Repositories thus provide mostly carrots.

An example of a proteomics feedback loop



Funding example: EC FP6 ProDaC grant



- Coordination Action grant in the 6th Framework Programme of the European Commission
- Good example of a coordination between: (i) funder, (ii) data producers, (iii) repository providers, (iv) standards organization, and (v) journals
- Workpackages

WP1: standards for data representation

WP2: Standards implementation

WP3: Data integration tools

WP4: Proteomics repository adaptation

WP5: Data flow management

WP6: Proteomics data exploitation

(WP7: project management

PSI standards development

provide compatible software

enabling data submission

ensuring repository compliance

overall data gathering and submission

establishing data feedback loops

overall project management)