

EBI RDF Platform: 6 months in

<http://www.ebi.ac.uk/rdf/>

Simon Jupp <jupp@ebi.ac.uk>

What is EMBL-EBI?

- Part of the European Molecular Biology Laboratory
- International, non-profit research institute
- Europe's hub for biological data services and research
- Based in Hinxton, Cambridge



Data resources at EMBL-EBI

Genes, genomes & variation

European Nucleotide Archive
1000 Genomes

Ensembl
Ensembl Genomes

European Genome-phenome Archive
Metagenomics portal

Gene, protein & metabolite expression

ArrayExpress

Expression Atlas

Metabolights
PRIDE

Literature & ontologies

Europe PubMed Central
Gene Ontology
Experimental Factor Ontology

Protein sequences, families & motifs

InterPro

Pfam

UniProt

Molecular structures

Protein Data Bank in Europe
Electron Microscopy Data Bank

Chemical biology

ChEMBL

ChEBI

Reactions, interactions & pathways

IntAct

Reactome

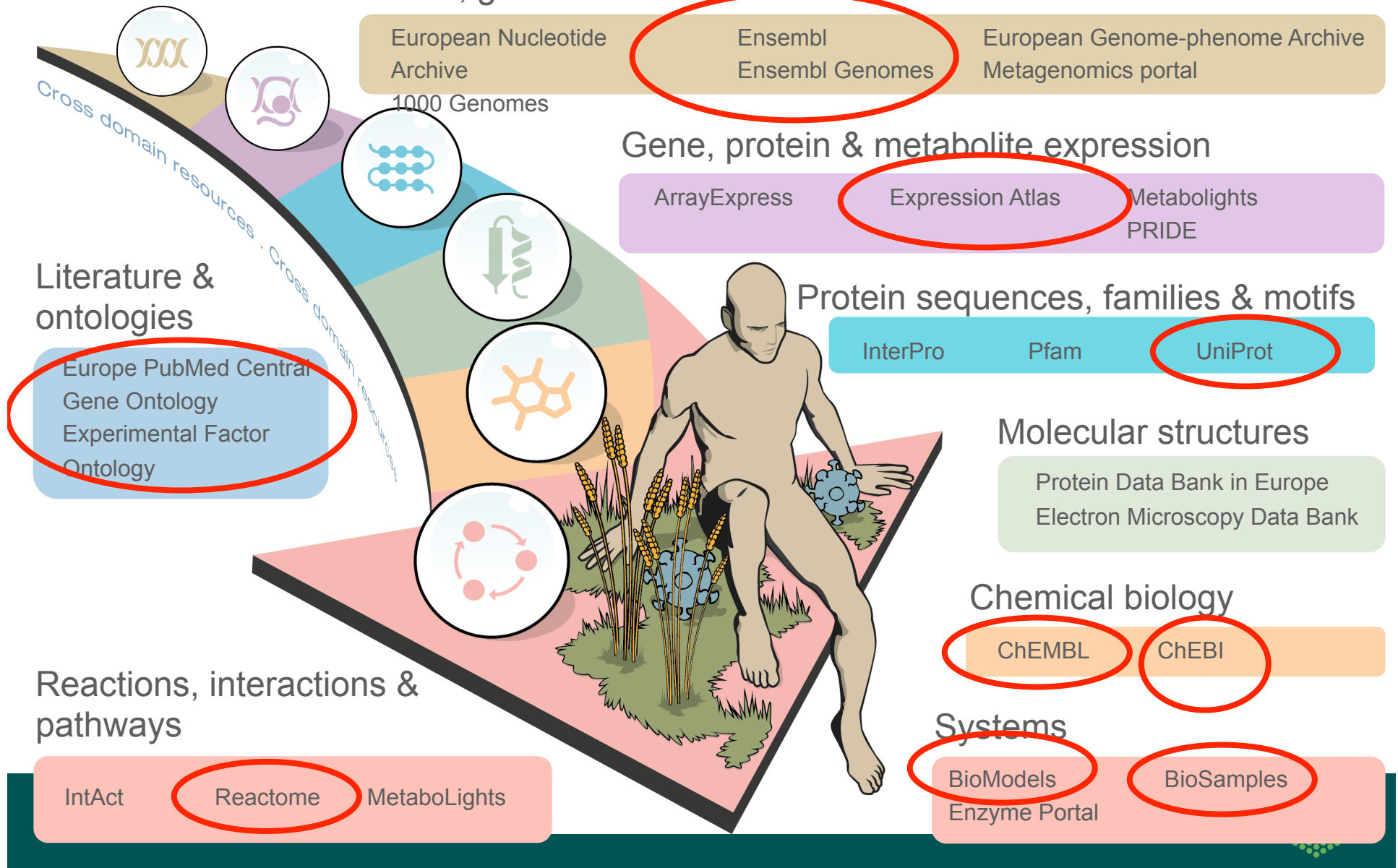
MetaboLights

Systems

BioModels

Enzyme Portal

BioSamples



Why RDF at EMBL-EBI?

- Interest for a number of years from some of our user base particularly some industry partners
- Betas and pilots starting up within individual EBI projects without coordination
- Overall feeling that technology is maturing and life science RDF community is growing
- EBI RDF has been available elsewhere but users had concerns of stability and faithfulness to source

Components for each resource

- **Stable** public SPARQL endpoint
 - Stable = production level service
- Consistent documentation with
 - Example queries
 - Schema diagram
 - Data set description (VOID)
- Bulk downloads
 - Full RDF download from FTP site
- Hosted via Linked Data Browser (lodestar – common look and feel across resources)

The screenshot displays the EMBL-EBI RDF Platform interface. The top navigation bar includes 'Services', 'Research', 'Training', and 'About us'. The main content area is titled 'Expression Atlas SPARQL Endpoint' and features a text input field for 'Enter SPARQL Query'. Below the input field, a SPARQL query is shown, including prefixes for 'rdf', 'xsd', 'owl', 'obo', 'obo:owl', 'obo:owl:IR_1', 'obo:owl:IR_2', 'obo:owl:IR_3', 'obo:owl:IR_4', 'obo:owl:IR_5', 'obo:owl:IR_6', 'obo:owl:IR_7', 'obo:owl:IR_8', 'obo:owl:IR_9', 'obo:owl:IR_10', 'obo:owl:IR_11', 'obo:owl:IR_12', 'obo:owl:IR_13', 'obo:owl:IR_14', 'obo:owl:IR_15', 'obo:owl:IR_16', 'obo:owl:IR_17', 'obo:owl:IR_18', 'obo:owl:IR_19', 'obo:owl:IR_20', 'obo:owl:IR_21', 'obo:owl:IR_22', 'obo:owl:IR_23', 'obo:owl:IR_24', 'obo:owl:IR_25', 'obo:owl:IR_26', 'obo:owl:IR_27', 'obo:owl:IR_28', 'obo:owl:IR_29', 'obo:owl:IR_30', 'obo:owl:IR_31', 'obo:owl:IR_32', 'obo:owl:IR_33', 'obo:owl:IR_34', 'obo:owl:IR_35', 'obo:owl:IR_36', 'obo:owl:IR_37', 'obo:owl:IR_38', 'obo:owl:IR_39', 'obo:owl:IR_40', 'obo:owl:IR_41', 'obo:owl:IR_42', 'obo:owl:IR_43', 'obo:owl:IR_44', 'obo:owl:IR_45', 'obo:owl:IR_46', 'obo:owl:IR_47', 'obo:owl:IR_48', 'obo:owl:IR_49', 'obo:owl:IR_50', 'obo:owl:IR_51', 'obo:owl:IR_52', 'obo:owl:IR_53', 'obo:owl:IR_54', 'obo:owl:IR_55', 'obo:owl:IR_56', 'obo:owl:IR_57', 'obo:owl:IR_58', 'obo:owl:IR_59', 'obo:owl:IR_60', 'obo:owl:IR_61', 'obo:owl:IR_62', 'obo:owl:IR_63', 'obo:owl:IR_64', 'obo:owl:IR_65', 'obo:owl:IR_66', 'obo:owl:IR_67', 'obo:owl:IR_68', 'obo:owl:IR_69', 'obo:owl:IR_70', 'obo:owl:IR_71', 'obo:owl:IR_72', 'obo:owl:IR_73', 'obo:owl:IR_74', 'obo:owl:IR_75', 'obo:owl:IR_76', 'obo:owl:IR_77', 'obo:owl:IR_78', 'obo:owl:IR_79', 'obo:owl:IR_80', 'obo:owl:IR_81', 'obo:owl:IR_82', 'obo:owl:IR_83', 'obo:owl:IR_84', 'obo:owl:IR_85', 'obo:owl:IR_86', 'obo:owl:IR_87', 'obo:owl:IR_88', 'obo:owl:IR_89', 'obo:owl:IR_90', 'obo:owl:IR_91', 'obo:owl:IR_92', 'obo:owl:IR_93', 'obo:owl:IR_94', 'obo:owl:IR_95', 'obo:owl:IR_96', 'obo:owl:IR_97', 'obo:owl:IR_98', 'obo:owl:IR_99', 'obo:owl:IR_100'. The query is executed, and the results are displayed in a table format. The table has columns for 'id', 'label', 'description', 'type', and 'value'. The first row shows 'musculus' as the label and 'Mus musculus' as the description. The second row shows 'organism' as the type and 'A material entity that is an individual living system, such as animal, plant, bacteria or virus, that is capable of replicating or reproducing, growth and maintenance in the right environment. An organism may be unicellular or made up, like humans, of many billions of cells divided into specialized tissues and organs.' as the value. The interface also includes a 'Related to' section with a 'subClassOf' link.

Numbers of triples

The screenshot shows the EMBL-EBI RDF Platform website. The main header features the EMBL-EBI logo and navigation links for Services, Research, Training, and About us. Below this is the 'RDF Platform' title and a secondary navigation bar with links for RDF Platform, Services, Documentation, FAQ, and About, along with a Feedback button.

The main content area includes a paragraph explaining the platform's goal: "The EBI RDF Platform aims to bring together the efforts of a number of EMBL-EBI resources that provide access to their data using Semantic Web technologies. It provides a unified way to query across resources using the W3C SPARQL query language. We welcome comments or questions via our feedback form."

Below the text is a section titled "Current RDF resources" which contains a table with three columns: Services, Quick links, and Example query. The table lists six services with their respective triple counts highlighted in green boxes:

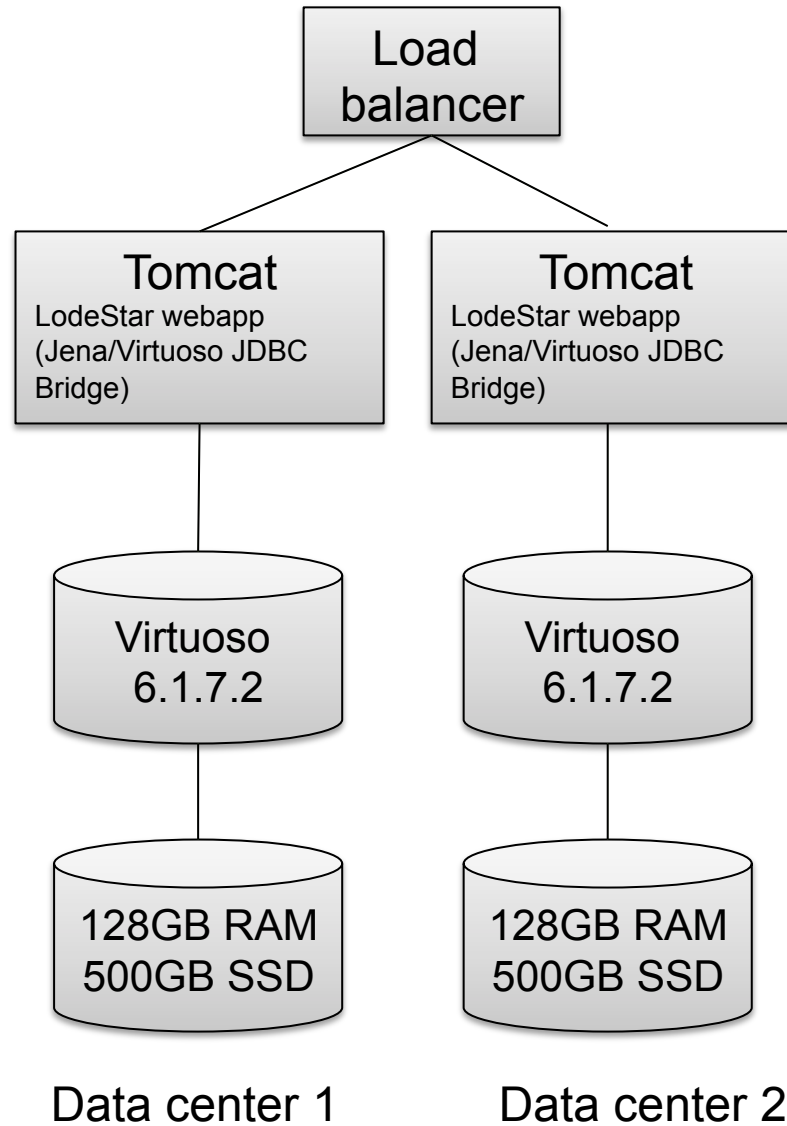
Services	Quick links	Example query
BioModels	11,385,558 triples	ions to plex
BioSamples	102,520,402 triples	f
ChEMBL	374,762,364 triples	1.
Expression Atlas	447,149,547 triples	is
Reactome	12,487,422 triples	Service description Pathways that references Inculg
UniProt	9,024,662,088 triples	Service description What are the preferred gene name

At the bottom of the table, there is a "Feedback" link.

On the right side of the page, there is a sidebar titled "RDF Platform" with a dropdown menu containing the following items:

- RDF Platform
- About the technology
- Getting started
- About the project
- EBI RDFApp Competition - win an iPad Mini!

Technical Infrastructure



1 Virtuoso instance per resource

Loading times and benchmarking

- Loading tests run 8 virtual CPU core, 128GB RAM, 3TB SAN storage on flash-based Violin array.
 - RedHat Enterprise Linux version 6, virtualised on VMWare infrastructure
- Tested Sesame, OWLIM-lite, OWLIM-SE, Virtuoso 6 and 7

Dataset	Triples	OWLIM-SE	Virtuoso 7	Virtuoso 6
Expression Atlas 13.01	339,991,787	1 hr 45 mins	24 mins	25 mins
ChEMBL 14	97,717,782	32 mins	5 mins	2 mins 23 seconds
UniProt 2013_01	6,419,569,246	57 hrs 40 mins	5 hrs 11 mins	N/A
BioModels 24	5,788,106	4 mins	20 seconds	15 seconds
Reactome 42	12,927,678	9 mins	30 seconds	21 seconds

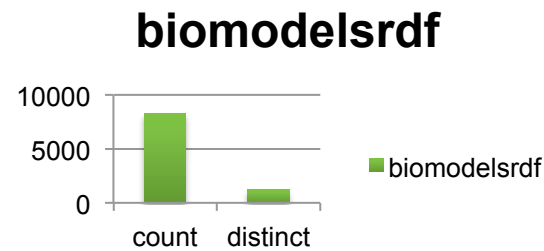
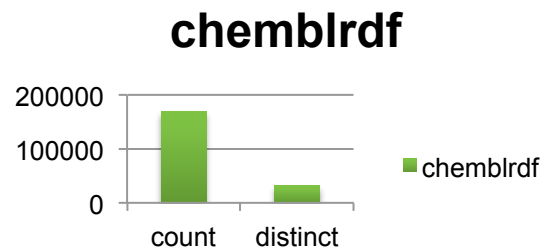
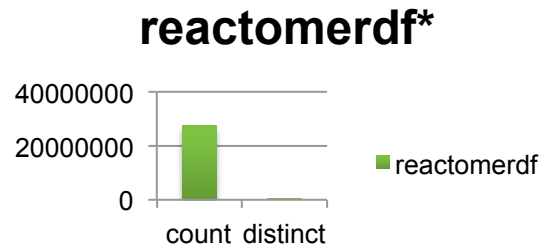
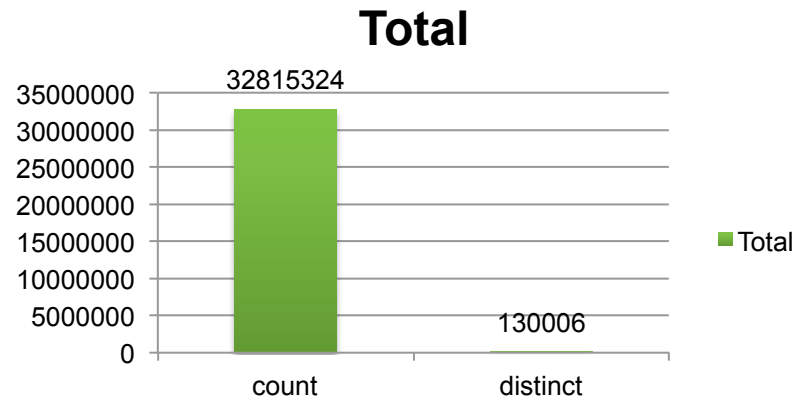
Usage analysis

- Load balancer access logs
 - No. of user
 - Number of queries
- Log every SPARQL query
 - Including execution time
- Query analysis using Jena
 - Path size
 - Used SPARQL operators
 - UNION, OPTION, GROUP, FILTER, SERVICE, SUBQUERY
 - Predicate vocabularies

Usage

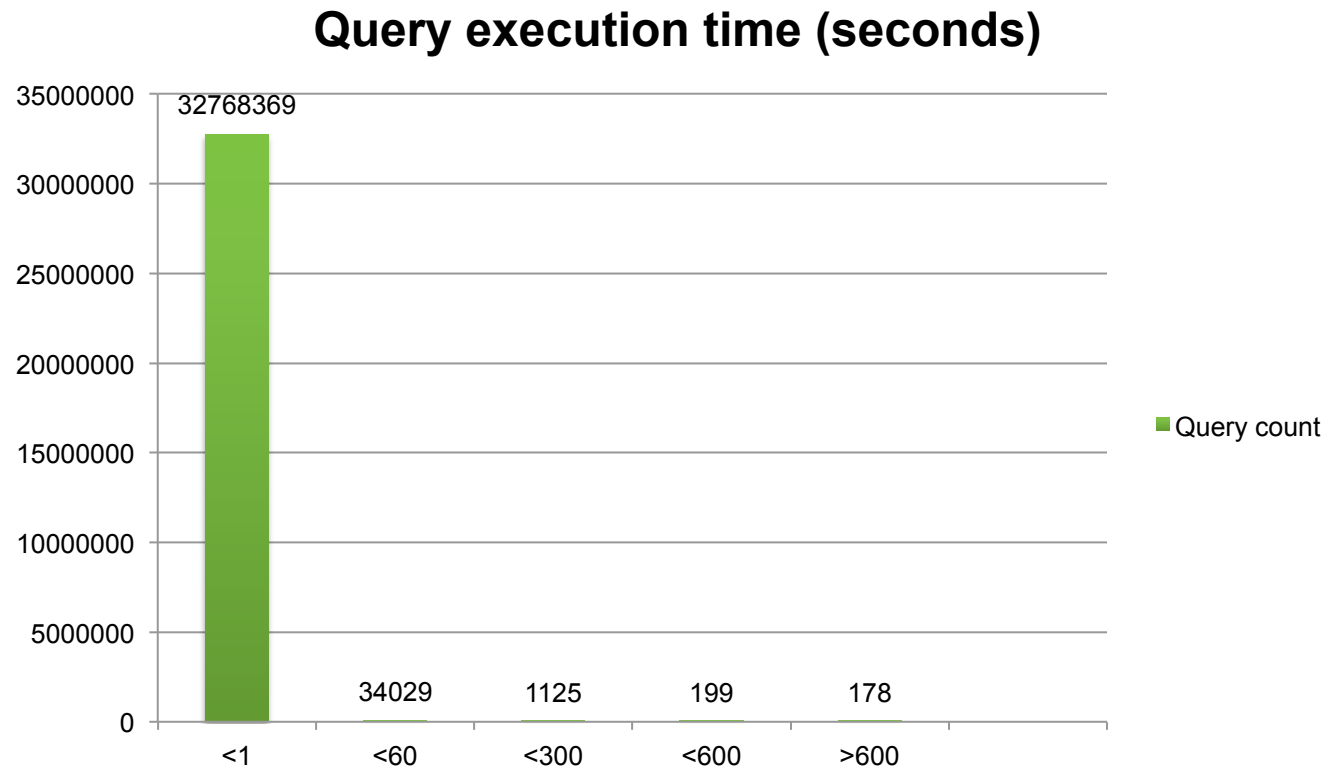
- ~500 unique visitors per month
- ~55 millions hits to site
- Most usage coming from UK
 - Internal pipelines and API calls that use the RDF platform
- Some datasets being downloaded and hosted externally behind other public APIs
 - OpenPhacts project and ChEMBL RDF
- Most resources already have good modes of access

Number of queries (6 months)



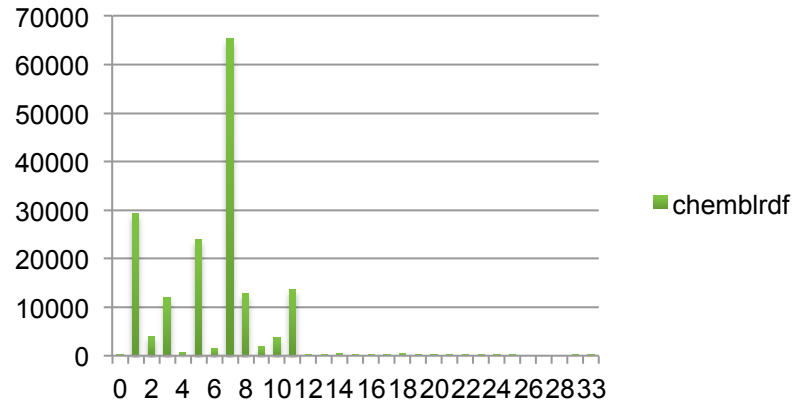
*Excludes ~90 million queries generated
By internal pipeline between Dec 1-4th 2013

Query execution time

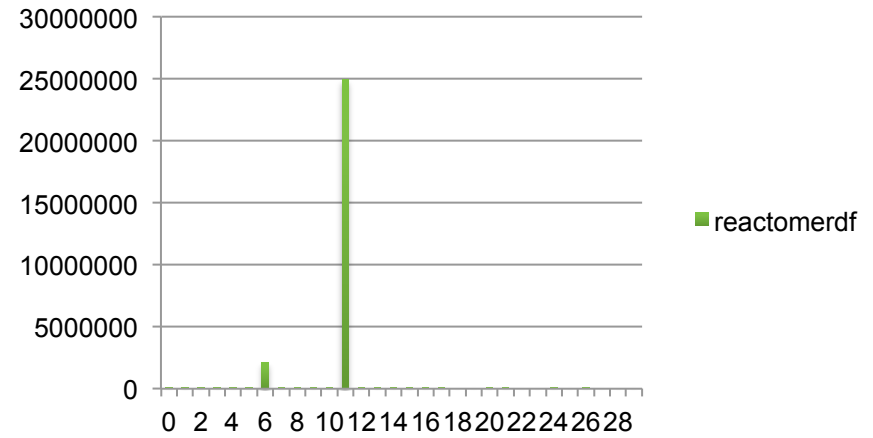


Query path size

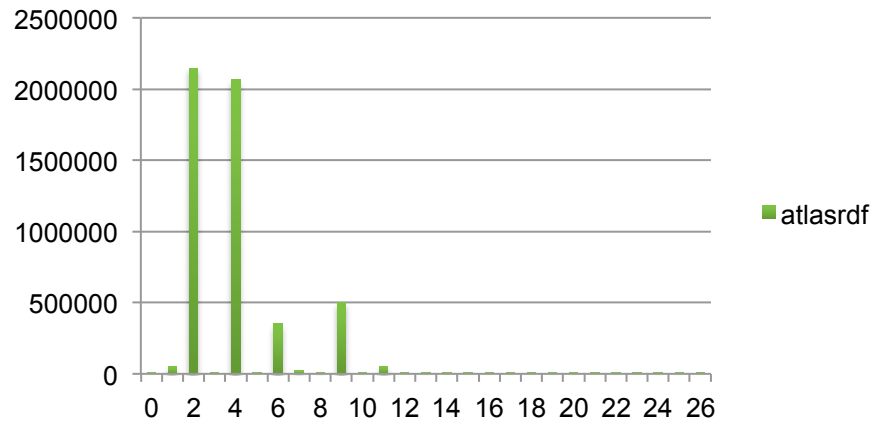
chemblrdf



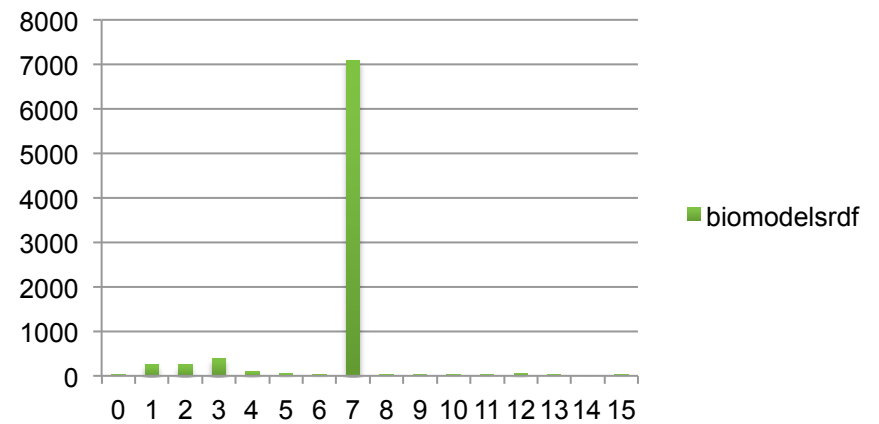
reactomerdf



atlasrdf

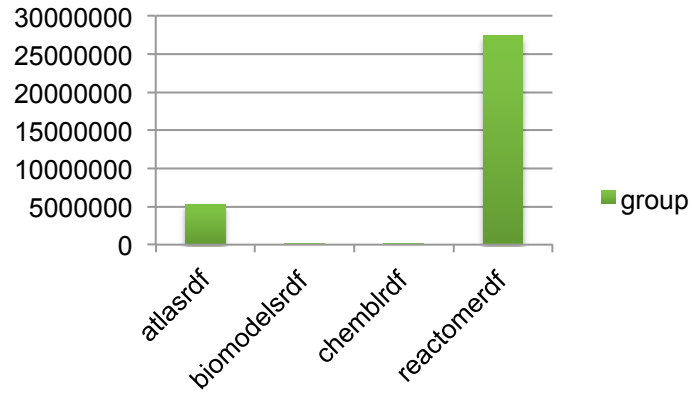


biomodelsrdf

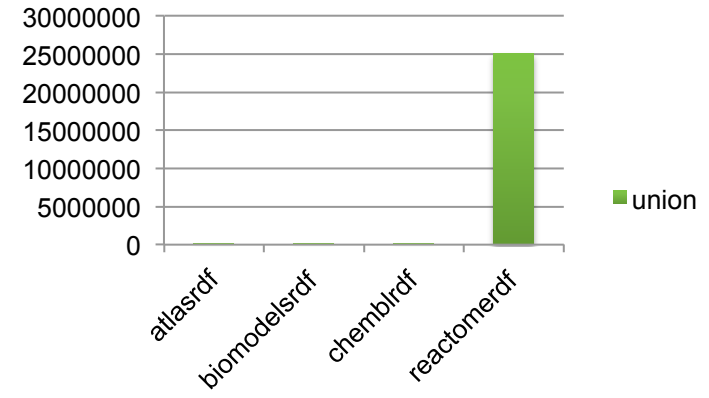


SPARQL elements

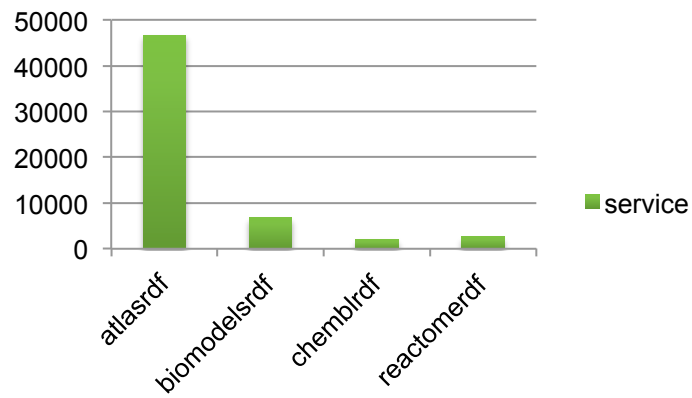
group



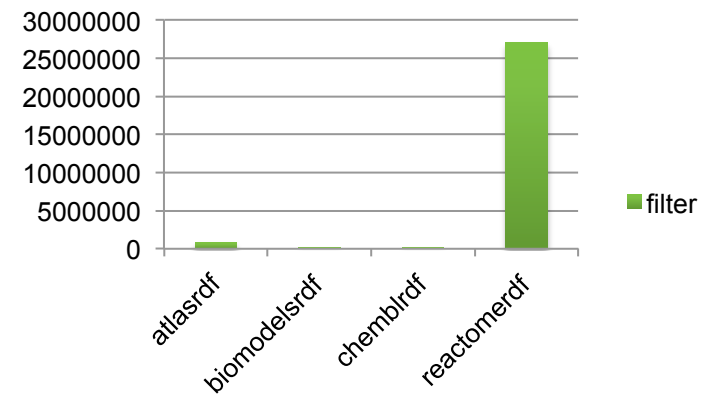
union



service

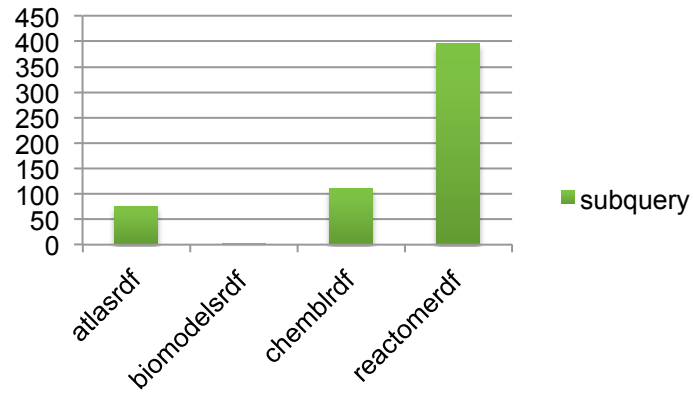


filter

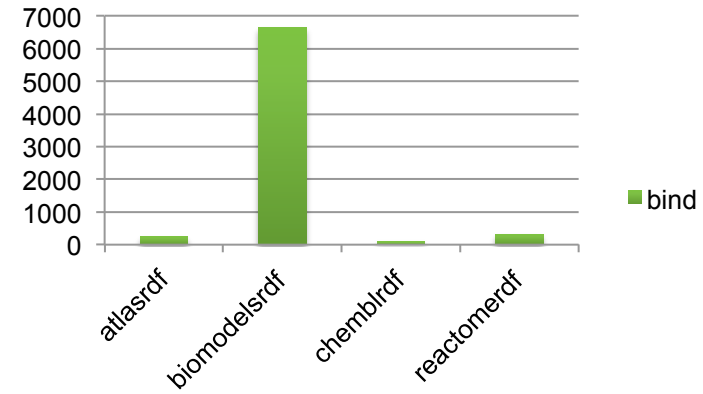


SPARQL elements cont...

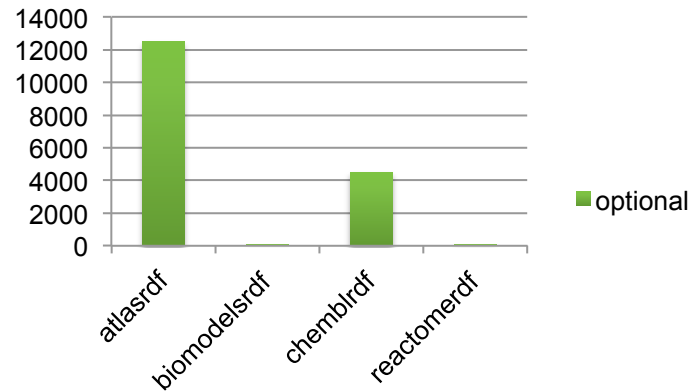
subquery



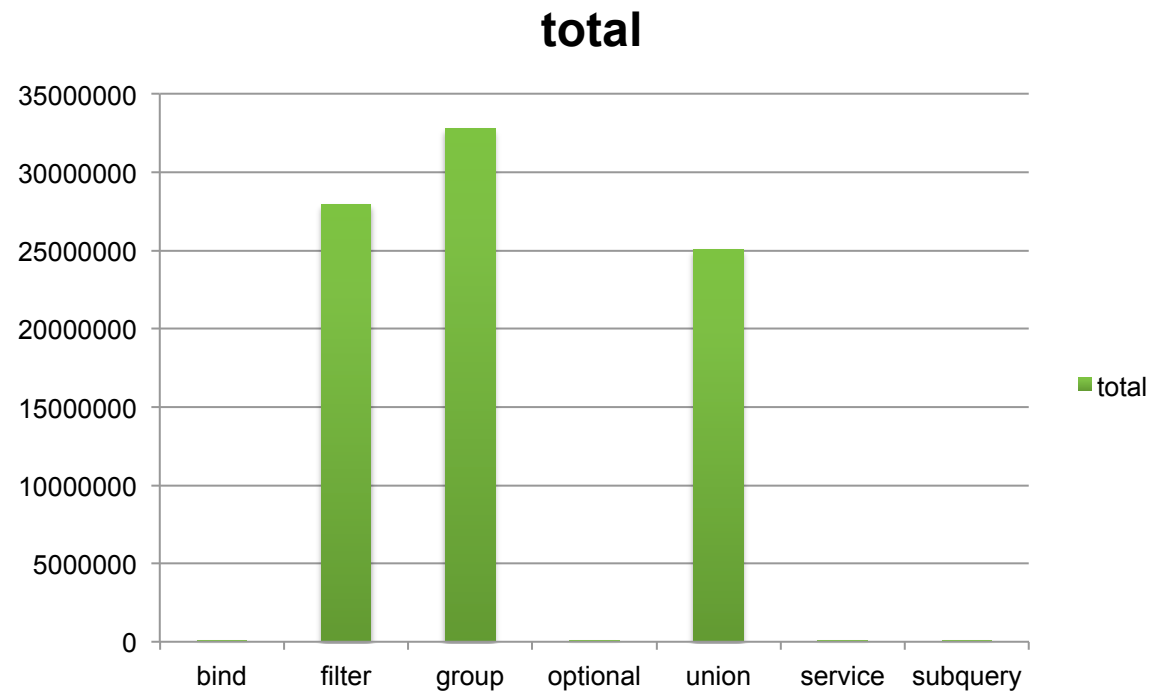
bind



optional

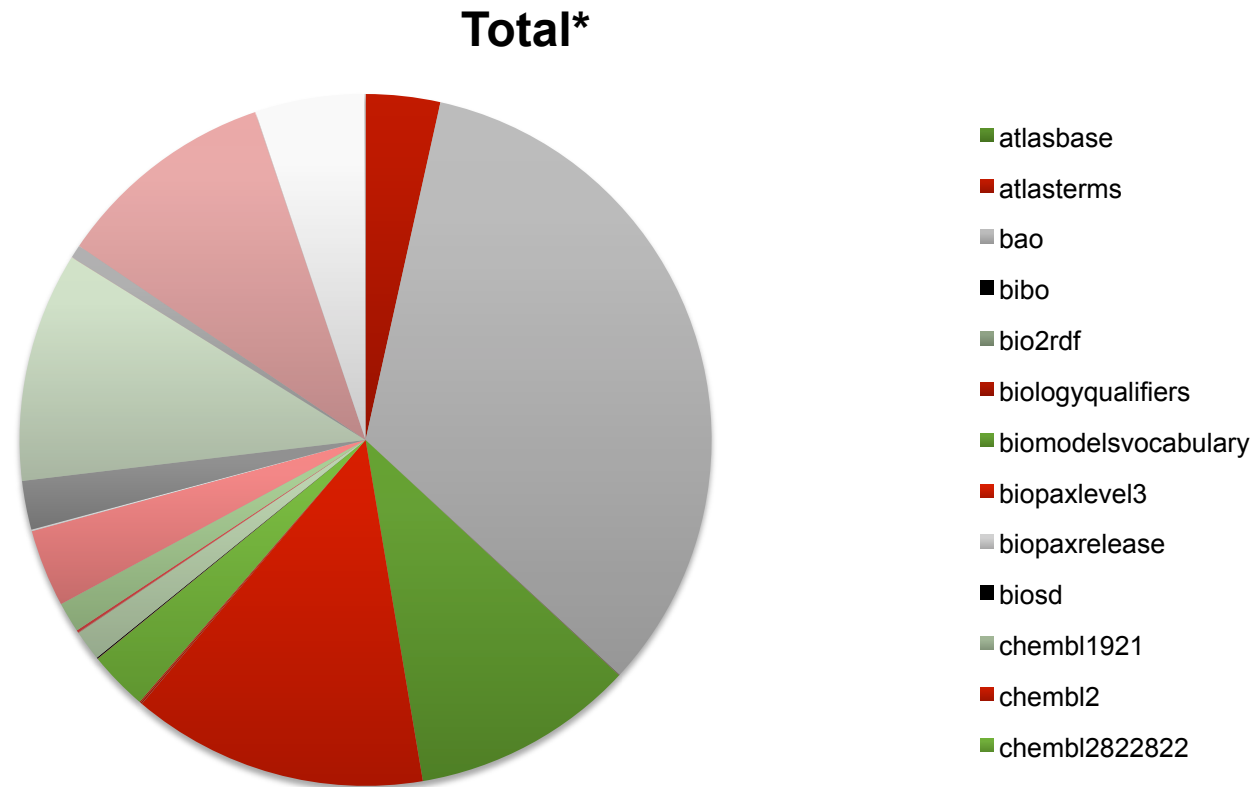


SPARQL elements



Predicate vocabularies

- 47 different vocabularies in total
- 577 relations



* Excludes rdf vocabulary

Performance problems

- Federated querying
 - One triple store vs multiple federate stores
- We use lots of ontologies (3.9 million OWL class usage)
 - Transitive queries still hard to answer
 - Get **metformin** associated **pathways** with **differentially expressed genes**, find any **proteins** that are **targets** for known **diabetes drugs**
- Some big datasets coming in 2014
 - Ensembl gene database. 78 different species
 - Many billions of triples

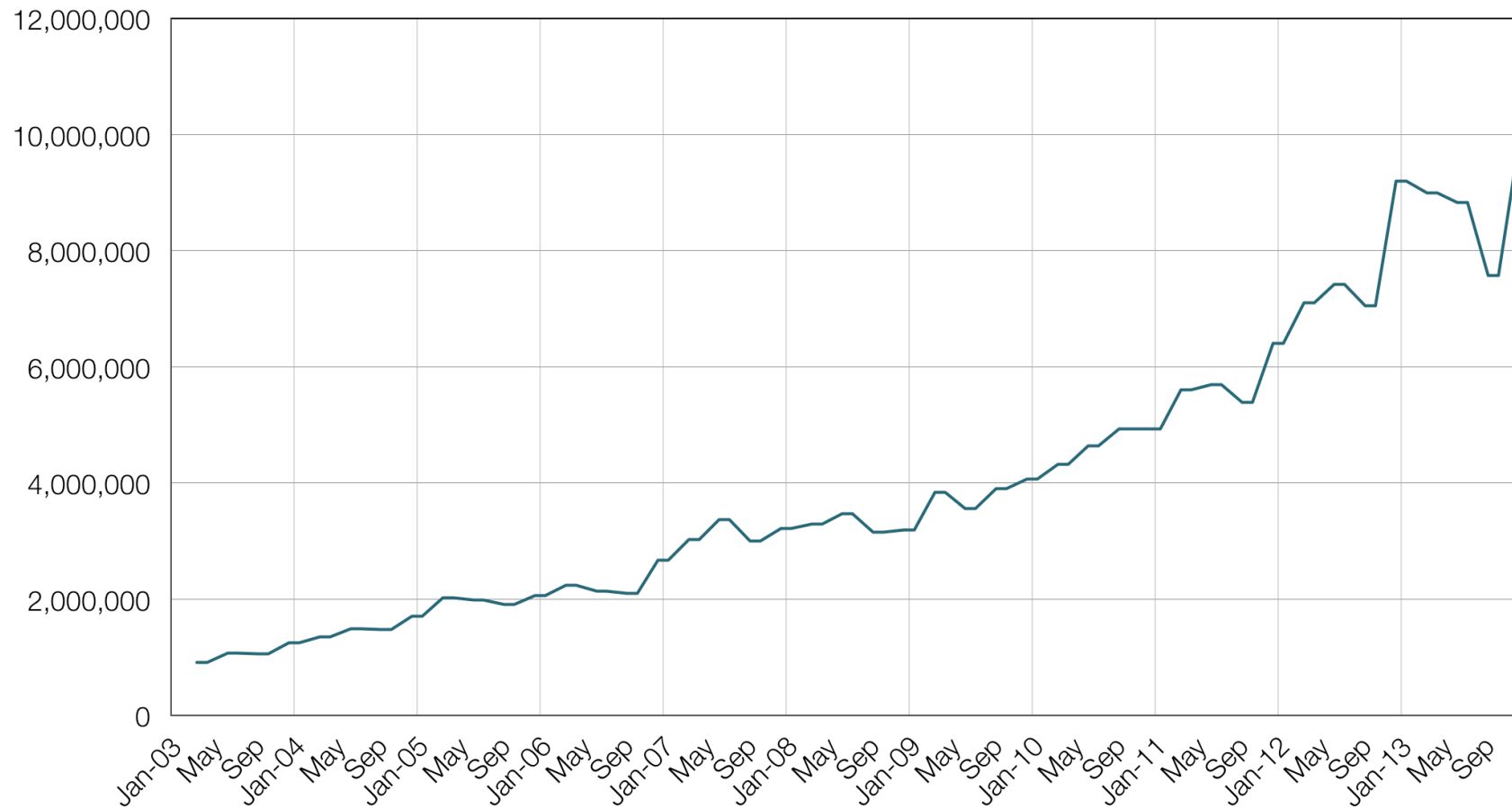
How can we help

- All datasets are publicly available form FTP site
 - Currently ~1 billion triples
 - More datasets planned in 2014 expect triples to rise
- We have complex biological queries currently difficult to ask
 - We need better inference support
 - We need federated queries to scale
- Open for collaboration on projects to help develop and push this technology

Acknowledgments and Funding

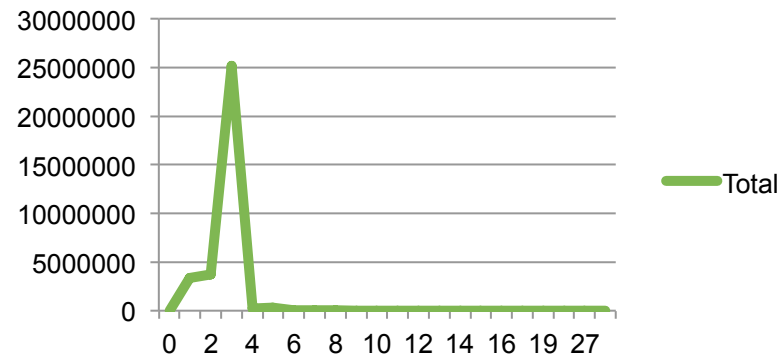
- RDF platform: James Malone, Andy Jenkinson, Mark Davies, Marco Brandizi, Sarala Wimalaratne, Leyla Garcia, Jerven Bolleman
- OntoText and OpenLink
- EMBL
- European Commission:
 - BioMedBridges [284209]
 - Diachron [601043]
- OpenPhacts
- NCBO through NIH NCBC grant U54-HG004028

Requests per day, 2003-2013

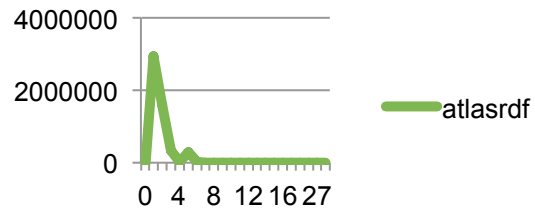


Selected variables

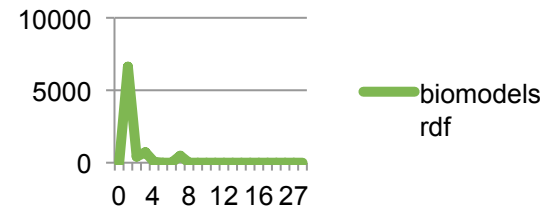
Total



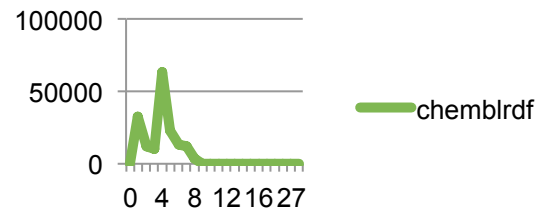
atlasrdf



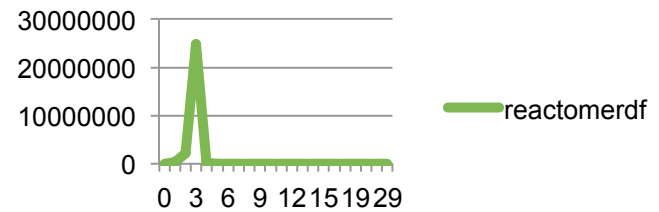
biomodelsrdf



chemblrdf

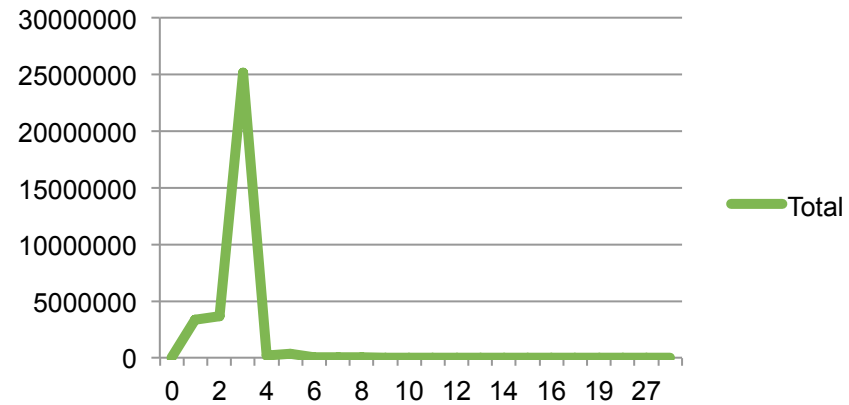


reactomerdf

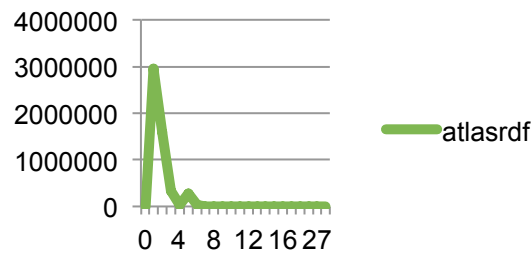


Query variables

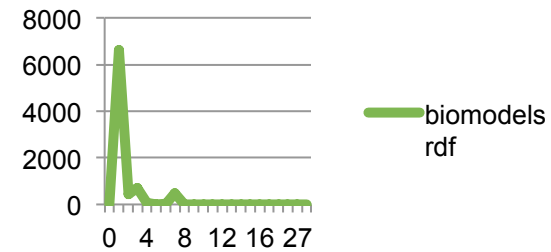
Total



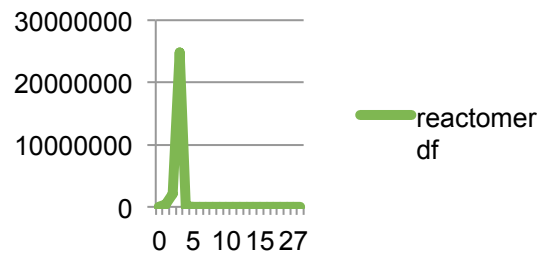
atlasrdf



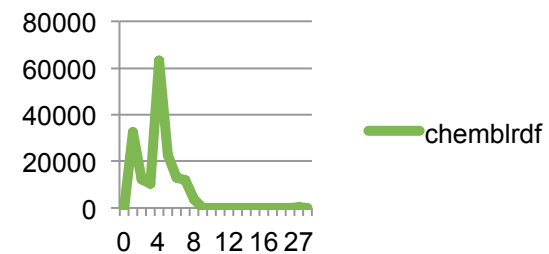
biomodelsrdf



reactomerdf



chemblrdf



RDF platform Released October 2013

EMBL-EBI Services Research Training About us




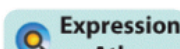




RDF Platform

RDF Platform Services Documentation FAQ About Feedback

The *EBI RDF Platform* aims to bring together the efforts of a number of EMBL-EBI resources that provide access to their data using [Semantic Web technologies](#). It provides a unified way to query across resources using the [W3C SPARQL](#) query language. We welcome **comments or questions** via our [feedback form](#).

Current RDF resources

Services	Quick links	Example query
 BioModels	<ul style="list-style-type: none"> ◦ Service description ◦ SPARQL endpoint ◦ Documentation ◦ RDF download 	All model elements with annotations to acetylcholine-gated channel complex (GO:0005892)
 BioSamples	<ul style="list-style-type: none"> ◦ Service description ◦ SPARQL endpoint ◦ Documentation ◦ RDF download 	Samples derived from known kinds of listeria organism
 ChEMBL	<ul style="list-style-type: none"> ◦ Service description ◦ SPARQL endpoint ◦ Documentation ◦ RDF download 	Find drug-like (but currently not approved) molecules which bind 7TM1 GPCRs with high affinity
 Expression Atlas	<ul style="list-style-type: none"> ◦ Service description ◦ SPARQL endpoint ◦ Documentation ◦ RDF download 	Under what experimental conditions is <u>Ensembl gene ENSG00000129991 (TNNI3)</u> expressed?
 Reactome	<ul style="list-style-type: none"> ◦ Service description ◦ SPARQL endpoint ◦ Documentation ◦ RDF download 	Pathways that references <u>Insulin (P01308)</u>
 UniProt	<ul style="list-style-type: none"> ◦ Service description ◦ SPARQL endpoint ◦ Documentation ◦ RDF download 	What are the preferred <u>gene name and disease annotations of all human UniProt entries that are known to be involved in a disease?</u>

RDF Platform

- ▾ RDF Platform
- About the technology
- Getting started
- About the project
- EBI RDFApp Competition - win an iPad Mini!



[Services](#) > [Gene Expression Atlas](#)

Gene Expression Atlas

The Functional Genomics Production Team are pleased to announce the beta publication of gene expression data from the Gene Expression Atlas as RDF Linked Open Data. The Expression Atlas Linked Dataset is an alternative API to the Gene Expression Atlas data. The purpose of this API is to enable richer queries over the data, it also supports federated queries over other linked datasets, including ChEMBL, Reactome, BioModels, BioSample, [Uniprot](#), [BioPortal](#) and [Bio2RDF](#).

The primary interface to the Expression Atlas RDF data is via the [SPARQL endpoint](#).

Gene Expression Atlas

- [SPARQL endpoint](#)
- [Atlas documentation](#)

Full VOID dataset description at <http://rdf.ebi.ac.uk/dataset/atlas/13.07>

Title	Gene Expression Atlas RDF
Description	RDF representation of all the experiments loaded into the Gene Expression Atlas database.
Version	13.07
Issued	July 19 2013
Number of triples	447149547

Contact

If you are interested in this work, please sign up to our [EFO user mailing list](#) where we will announce updates.

If you have questions or would like more information, you can email Simon Jupp [jupp \[at\] ebi.ac.uk](mailto:jupp@ebi.ac.uk) and James Malone [malone \[at\] ebi.ac.uk](mailto:malone@ebi.ac.uk) directly.

Atlas RDF Related Links

- [Experimental Factor Ontology \(EFO\)](#)
- [Gene Expression Atlas](#)
- [ArrayExpress](#)
- [Functional Genomics Production Team](#)



[Services](#) > [Atlas](#) > Atlas SPARQL endpoint

Expression Atlas SPARQL Endpoint

Enter SPARQL Query

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX owl: <http://www.w3.org/2002/07/owl#>
PREFIX dcterms: <http://purl.org/dc/terms/>
PREFIX obo: <http://purl.obolibrary.org/obo/>
PREFIX sio: <http://semanticscience.org/resource/>
PREFIX efo: <http://www.ebi.ac.uk/efo/>
PREFIX atlas: <http://rdf.ebi.ac.uk/resource/atlas/>
PREFIX atlasterms: <http://rdf.ebi.ac.uk/terms/atlas/>

SELECT DISTINCT ?experiment ?description WHERE
{?experiment
a atlasterms:Experiment ;
dcterms:description ?description ;
atlasterms:hasAssay
[atlasterms:hasSample
[atlasterms:hasSampleCharacteristic
[ atlasterms:propertyType ?propertyType ;
atlasterms:propertyValue ?propertyValue ]
]
]
}
filter regex (?description, "diabetes", "i")
```

RDFS inference?

Output: HTML

Results per page: 25

Submit Query Reset

Example Queries

- [Query 1](#)
Get experiments where the sample description contains diabetes
- [Query 2](#)
Get differentially expressed genes where factor is asthma
- [Query 3](#)
Show expression for ENSG00000129991 (TNNI3)
- [Query 4](#)
Show expression for ENSG00000129991 (TNNI3) with its GO annotations from Uniprot (Federated query to <http://beta.sparql.uniprot.org/sparql>)
- [Query 5](#)
For the genes differentially expressed in asthma, get the gene products associated to a Reactome pathway
- [Query 6](#)

RDF Platform

Services > Atlas > Atlas Linked Data

About: [Mus musculus](#)

http://purl.obolibrary.org/obo/NCBITaxon_10090

Identifier : "10090"

Type: [organism](#)

A material entity that is an individual living system, such as animal, plant, bacteria or virus, that is capable of replicating or reproducing, growth and maintenance in the right environment. An organism may be unicellular or made up, like humans, of many billions of cells divided into specialized tissues and organs.
more types...



Expression Atlas RDF

- [Dataset home](#)
- [SPARQL endpoint](#)
- [Documentation](#)

Related to

subClassOf



Atlas documentation

The schema for the atlas datasets is depicted below. It shows how resources are connected in the RDF graph. Each resource is types according to the atlas terms ontology. The stable URI for the latest version of this ontology is <http://rdf.ebi.ac.uk/terms/atlas>.

This ontology acts as our internal schema and provides a basic description of the resources in the Atlas RDF dataset. In order to facilitate integration of the Atlas RDF data with other datasets, we provide an additional ontology that maps the atlas terms ontology to several external ontology development projects. We currently map to the Semantic Science Integration Ontology (SIO), the Ontology for Biomedical Investigations (OBI), the Experimental Factor Ontology (EFO) and the ontology of bioinformatics operations, topics, types of data including identifiers, and formats (EDAM). A table of the mappings can be seen [here](http://rdf.ebi.ac.uk/terms/atlas-mapping). The full atlas terms mapping ontology can be found at <http://rdf.ebi.ac.uk/terms/atlas-mapping>. We are happy to extend our mapping to other relevant ontologies and are keen to discuss the integration of this dataset with other similar resources.

Documentation

- ▾ Atlas documentation
 - Atlas contact
 - Atlas download
 - Atlas examples
 - AtlasRDF R package
- BioModels documentation
- ChEMBL documentation
- Reactome documentation
- UniProt documentation
- Using the SPARQL Endpoints
- Example SPARQL queries
- Programmatic access
- Technical documentation

