# The CALBC RDF Triple store: retrieval over large literature content

Samuel Croset, Christoph Grabmüller, Chen Li, Silverstras Kavaliauskas, Dietrich Rebholz-Schuhmann

croset@ebi.ac.uk

SWAT4LS

EMBL-EBI

10th December 2010, Berlin

EBI is an Outstation of the European Molecular Biology Laboratory.

# Outline

- Motivation

- Integrating multiple resources
    - CALBC Corpus
    - LexEBI
    - Public databases

- Querying the Triple Store

EMBL-EBI

# Outline

- **Motivation**

- Integrating multiple resources
    - CALBC Corpus
    - LexEBI
    - Public databases

- Querying the Triple Store

EMBL-EBI

# Why representing scientific literature in RDF?

- <u>Scientific literature:</u>
  - Primary data resource reporting novel scientific findings

- <u>Text-mining:</u>
  - Biological entities recognition
  - Population of biomedical databases through curators

- <u>RDF representation:</u>
  - Standardization of the content extracted
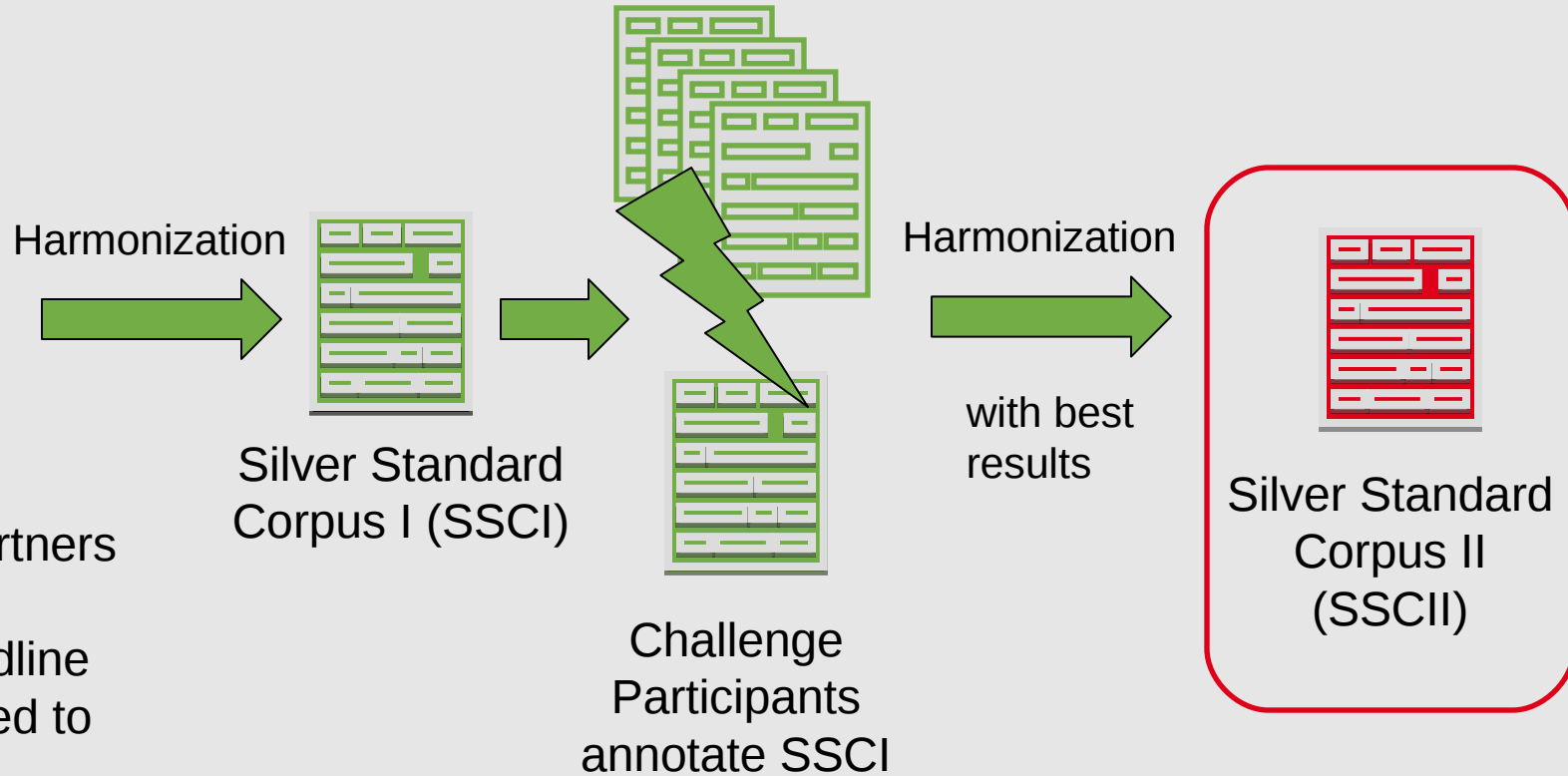  - **Exploitation of the literature in the Semantic Web**

EMBL-EBI

# Outline

- Motivation

- Integrating multiple resources
  - CALBC Corpus
  - LexEBI
  - Public databases
- Querying the Triple Store

EMBL-EBI

# CALBC Corpus

- Collaborative Annotation of a Large Biomedical Corpus

Harmonization

Harmonization

with best results

Silver Standard Corpus I (SSCI)

Silver Standard Corpus II (SSCII)

- 4 Project partners

- 150'000 Medline abstract related to Immunology annotated

Challenge Participants annotate SSCI

EMBL-EBI

# CALBC Corpus

- <u>Advantages of the CALBC Corpus:</u>

- Large-scale corpus

- 4 semantic types: Gene-Protein, Diseases, Chemicals and Species

- Generated  in a purely automatic way

-  Highly reproducible

- http://www.calbc.eu/

EMBL-EBI

# CALBC in RDF

**http://www.ebi.ac.uk/Rebholz/core/calbc/sentenceid#10605**

calbc:hasSentence

**http://www.ncbi.nlm.nih.gov/pubmed/44292**

dc:date

dc:creator

dc:identifier

1980-06-16

calbc:isIn

<urn:issn:0004-5772>

Seshadri, M S

Varkey, K

dc:title

http://www.ebi.ac.uk/Rebholz/core/corpus_calbc

"Hepatitis B surface antigen (HBsAg) positive polyarteritis nodosa. A report of two cases and review of literature"

EMBL-EBI

# CALBC in RDF

**http://www.ncbi.nlm.nih.gov/pubmed/44292**

calbc:isPartOf

**http://www.ebi.ac.uk/Rebholz/core/calbc/sentenceid#10605**

calbc:hasAnnotation

A

calbc:hasStartPosition

calbc:hasEndPosition

calbc:isEntityType

calbc:hasLabel

35

46

CHED

**"prostaglandins"**

@prefix calbc: http://www.ebi.ac.uk/Rebholz/core/calbc/
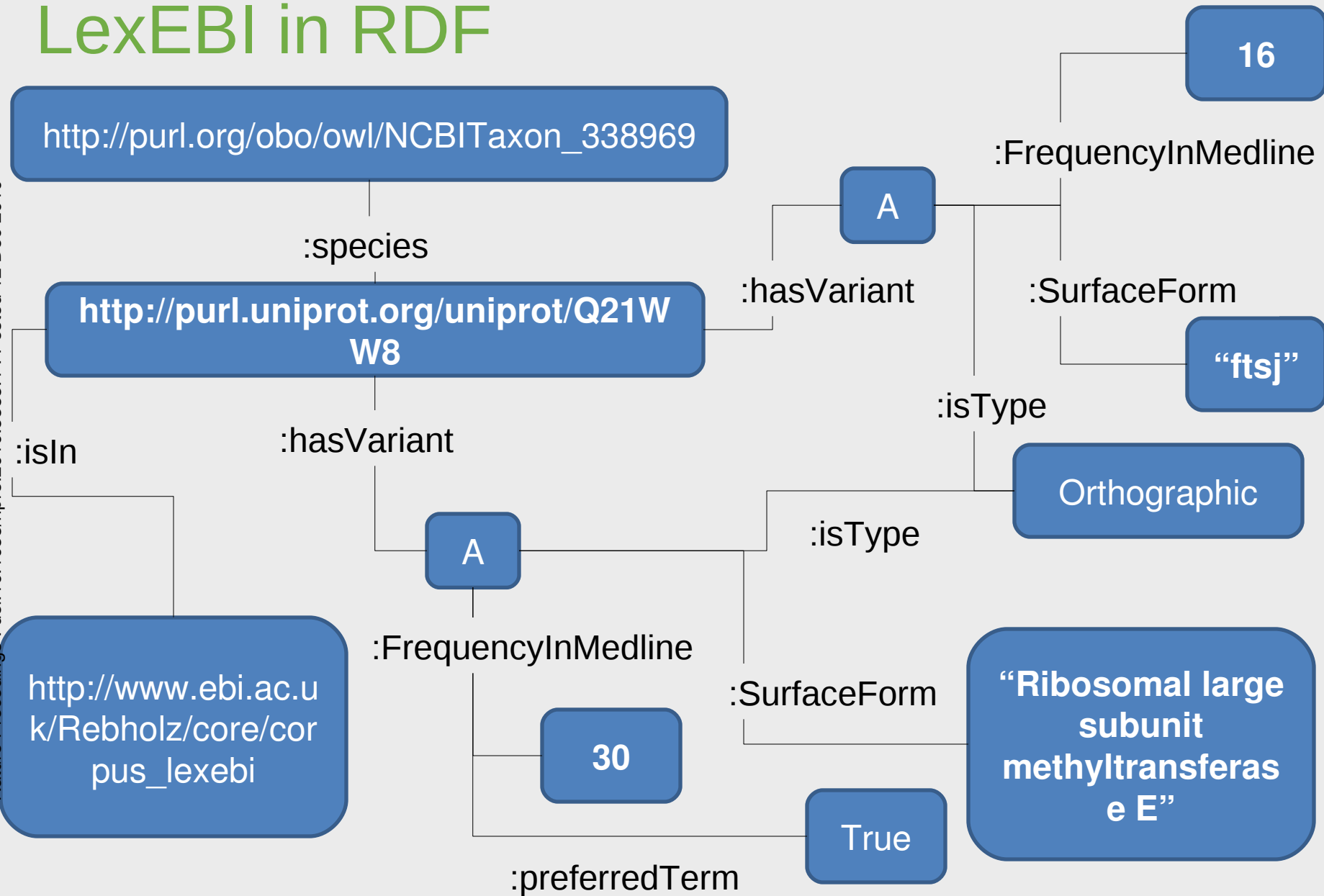
EMBL-EBI

# Outline

- **Motivation**

- **Integrating multiple resources**
  - CALBC Corpus
  - LexEBI
  - Public databases

- Querying the Triple Store

EMBL-EBI

# LexEBI

- BioThesaurus: Complete term repository for the biomedical domain
- LexEBI → XML

- <u>Features:</u>

  - Frequency count for the occurrence of the term in British National Corpus (BNC) or in MEDLINE → **Disambiguation**
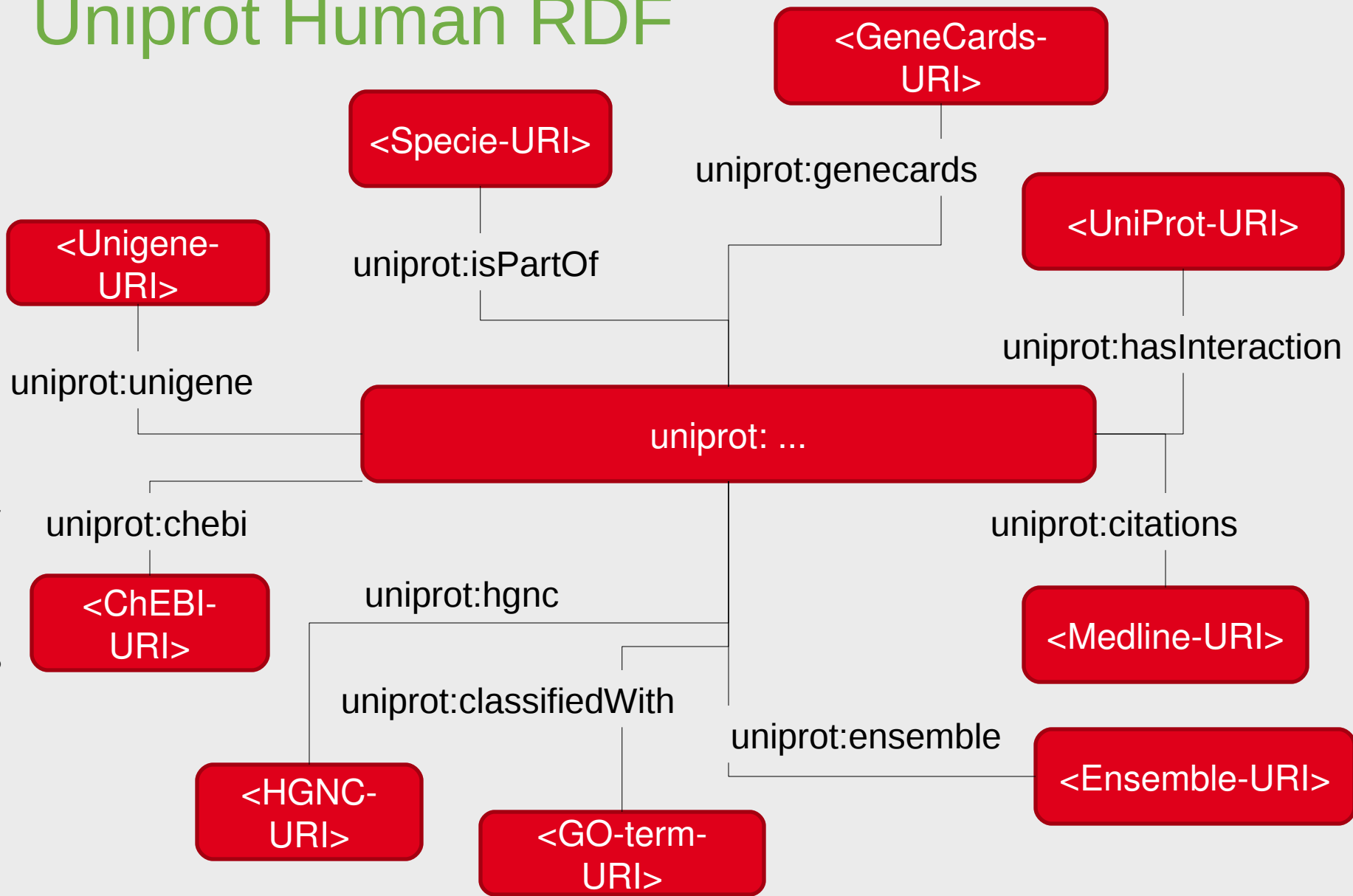  - Mapping to original resource (URI) → **Normalization**

EMBL-EBI

# LexEBI in RDF

http://purl.org/obo/owl/NCBITaxon_338969

**16**

:species

:FrequencyInMedline

**http://purl.uniprot.org/uniprot/Q21WW8**

A

:hasVariant

:SurfaceForm

:isIn

:hasVariant

:isType

**"ftsj"**

Orthographic

A

:isType

:FrequencyInMedline

http://www.ebi.ac.uk/Rebholz/core/corpus_lexebi

:SurfaceForm

**"Ribosomal large subunit methyltransferase E"**

**30**

True

:preferredTerm

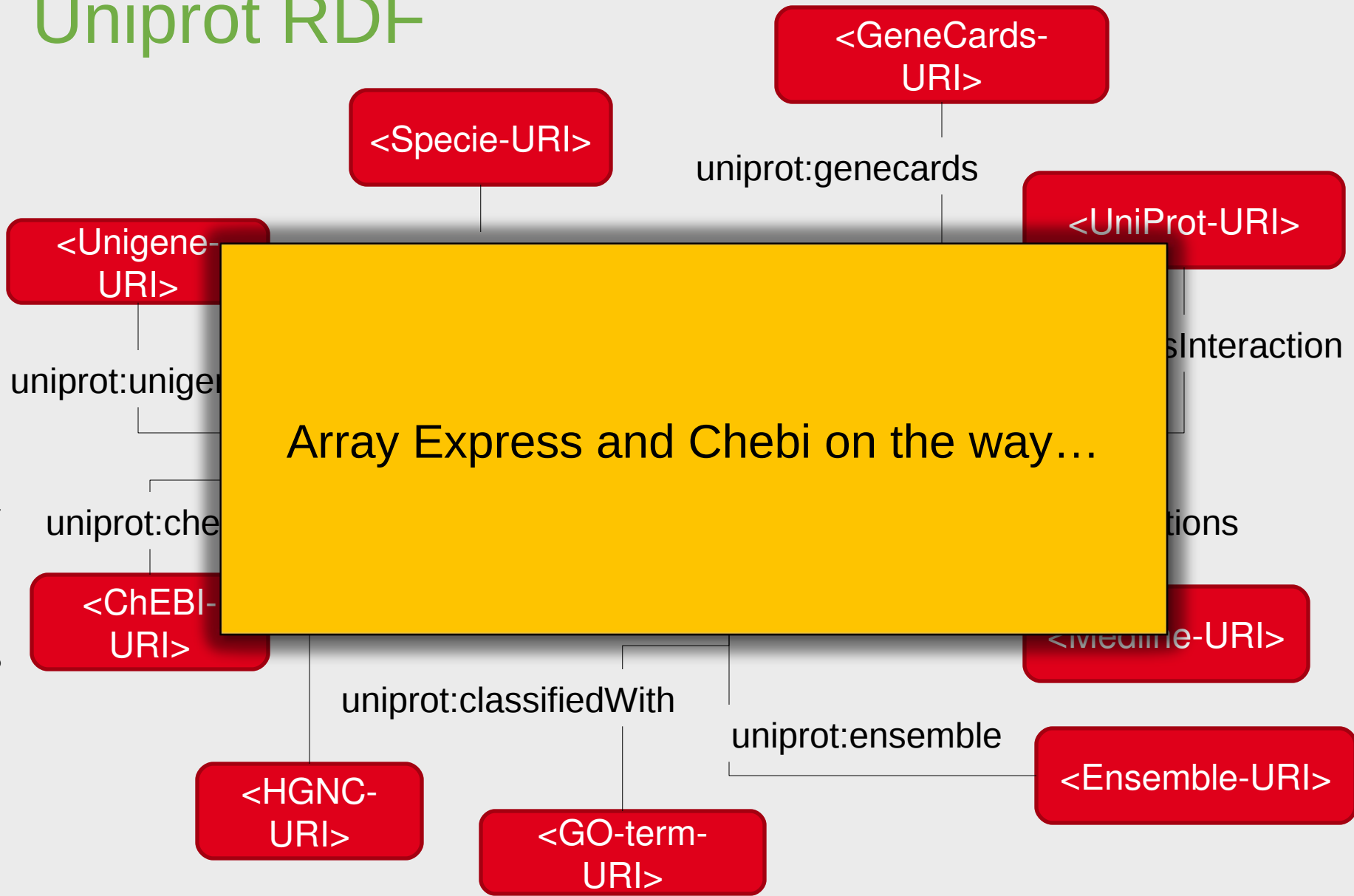@prefix : <http://www.ebi.ac.uk/Rebholz/core/lexebi#>

EMBL-EBI

# Outline

- Motivation

- Integrating multiple resources
  - CALBC Corpus
  - LexEBI
  - Public databases

- Querying the Triple Store

EMBL-EBI

# Uniprot Human RDF

<GeneCards-URI>

<Specie-URI>

uniprot:genecards

<UniProt-URI>

<Unigene-URI>

uniprot:isPartOf

uniprot:hasInteraction

uniprot:unigene

uniprot: ...

uniprot:chebi

uniprot:citations

<ChEBI-URI>

uniprot:hgnc

<Medline-URI>

uniprot:classifiedWith

uniprot:ensemble

<Ensemble-URI>

<HGNC-URI>

<GO-term-URI>

@prefix uniprot: http://purl.uniprot.org/uniprot/

EMBL-EBI

# Uniprot RDF

Array Express and Chebi on the way…

@prefix uniprot: http://purl.uniprot.org/uniprot/

EMBL-EBI

# Outline

- Motivation

- Integrating multiple resources

  - CALBC Corpus

  - LexEBI
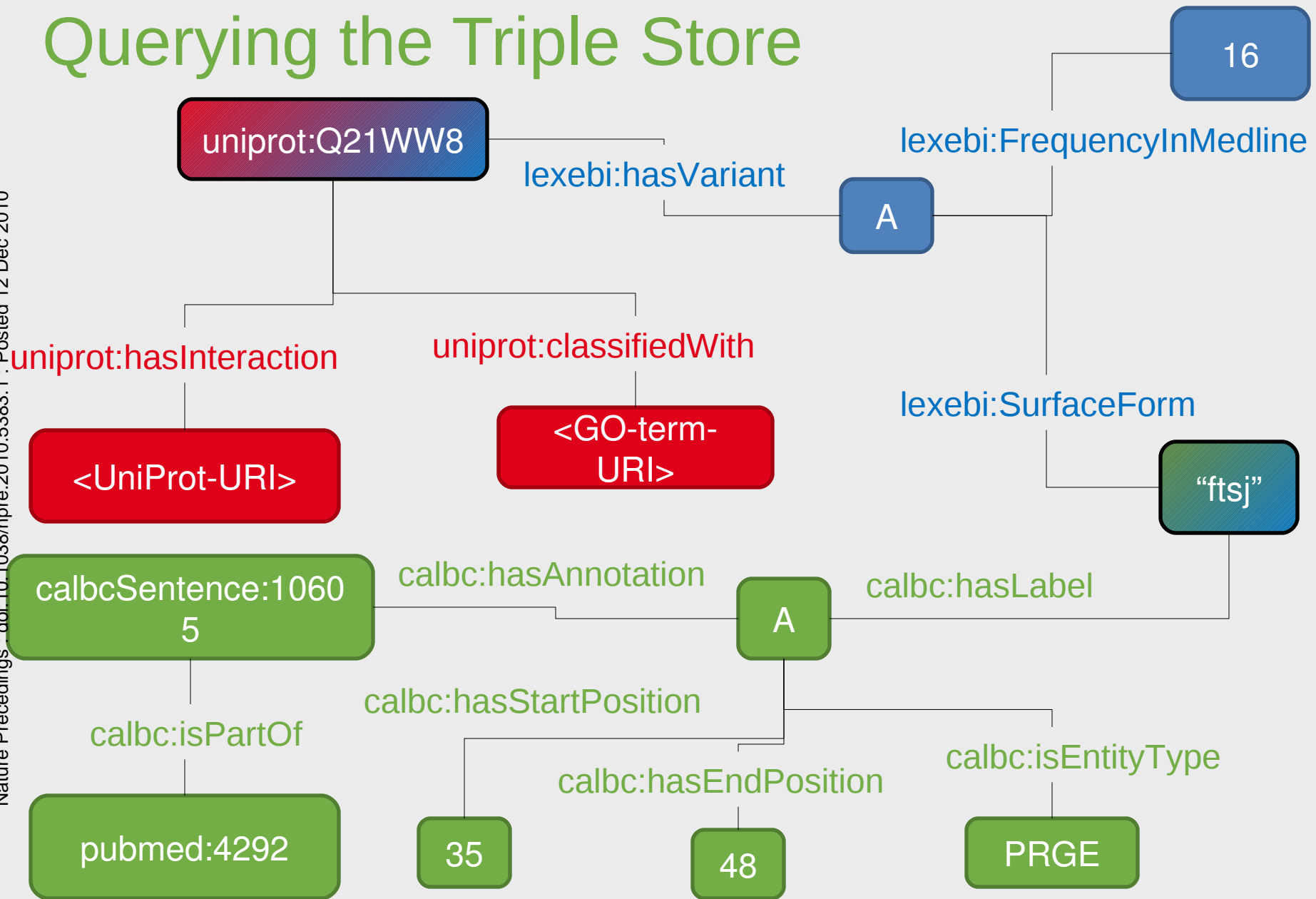
  - Public databases

- Querying the Triple Store

EMBL-EBI

# Querying the Triple Store

uniprot:Q21WW8

lexebi:hasVariant

16

lexebi:FrequencyInMedline

A

uniprot:hasInteraction

uniprot:classifiedWith

lexebi:SurfaceForm

&lt;UniProt-URI&gt;

&lt;GO-term-URI&gt;

"ftsj"

calbcSentence:10605

calbc:hasAnnotation

calbc:hasLabel

A

calbc:isPartOf

calbc:hasStartPosition

calbc:hasEndPosition

calbc:isEntityType

pubmed:4292

35

48

PRGE

EMBL-EBI

# Use cases

- Normalization of CALBC named entities
- Disambiguation of CALBC named entities
- Term collocation at the sentence level → e.g. Evidence for Gene – Disease association

- Checking consistency of bioinformatics resources from literature

EMBL-EBI

# Thank you for your attention

EMBL-EBI