# UniProt

# 100 billion+ triple KG

Jerven Bolleman

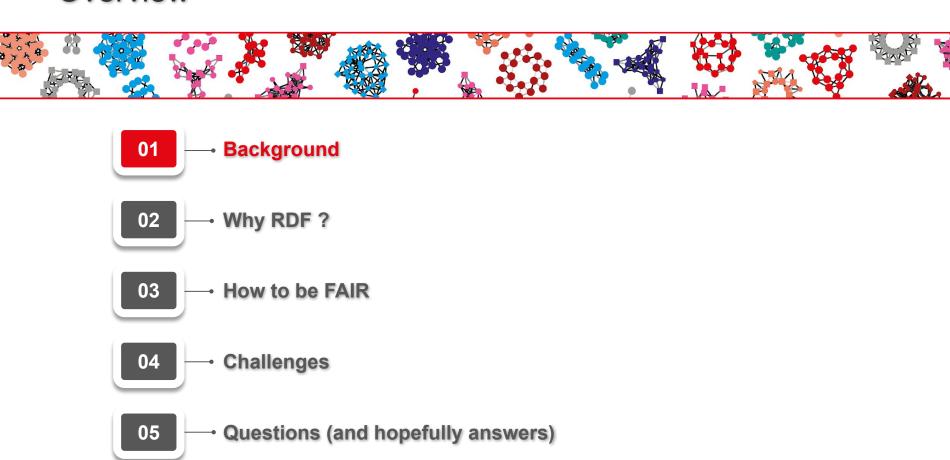Principal Software Engineer – Swiss-Prot group

SIB Swiss Institute of Bioinformatics
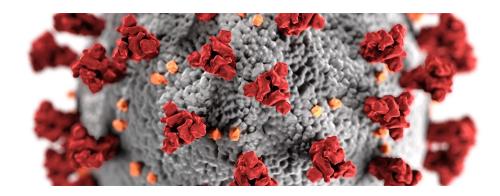
PIR

EMBL-EBI

www.sib.swiss

# Overview

# Background

Proteins and a lot of what we know about them

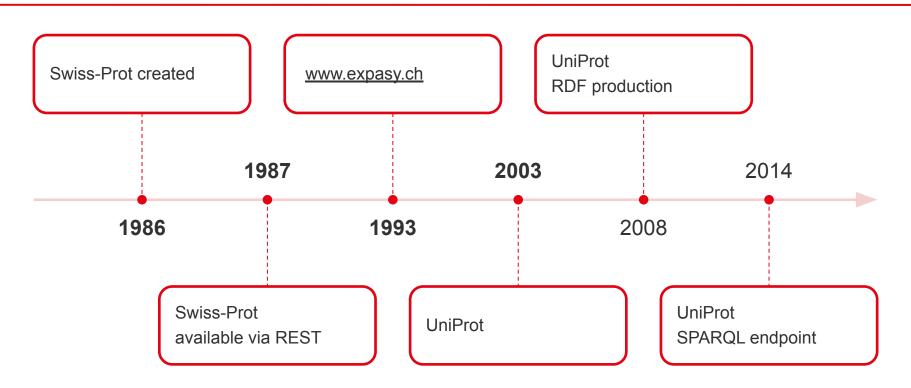Most of the parts of what makes us what we are.

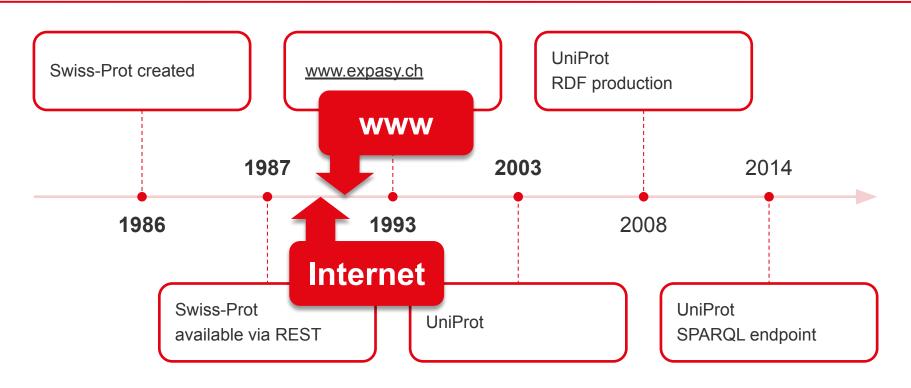Kind of important ...

# UniProt Knowledge?

- UniProtKB/Swiss-Prot (500,000)
  - Expert curated (AI accelerated)
  - Scientific literature
- UniProtKB/TrEMBL (50 billion)
  - Automated from DNA databases
  - AI annotated (learned on Swiss-Prot)
- UniRef (7 billion)
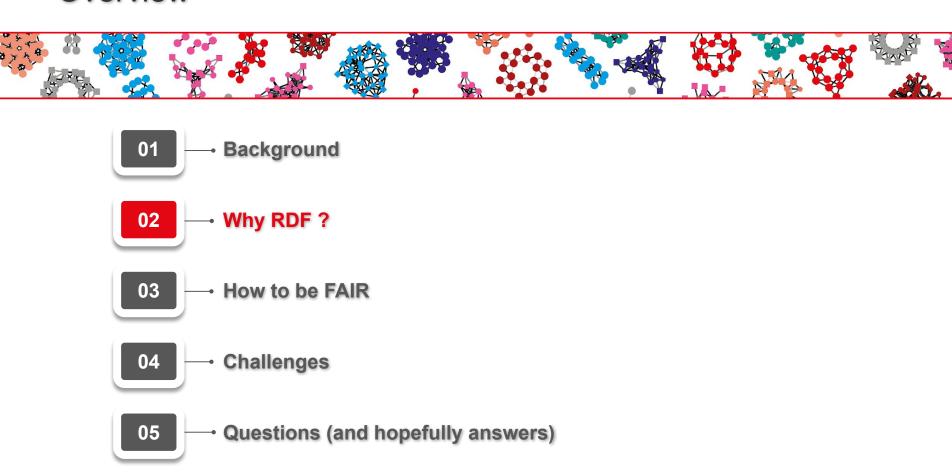- UniParc (50 billion)
  - Archive

# Background/History

Swiss-Prot created

www.expasy.ch

UniProt
RDF production

**1987**

**2003**

2014

**1986**

**1993**

2008

Swiss-Prot
available via REST

UniProt

UniProt
SPARQL endpoint

# Background/History



Swiss-Prot created

www.expasy.ch

**WWW**

UniProt
RDF production

**1987**

**2003**

2014

**1986**

**1993**

2008

**Internet**

Swiss-Prot
available via REST

UniProt

UniProt
SPARQL endpoint

# Overview

# What is 9993 ?

**Coupling of the Penicillium Duponti Acid Protease to Ethylene-Maleic Acid (1 : 1) Linear Copolymer. Preparation and Properties of the Water-Soluble Derivative**

PubMed

*DGCR2 DiGeorge syndrome critical region gene 2 [ Homo sapiens (human) ]*

NCBI Gene

I am 9993, would you like a flower

NCBI Taxonomy

# Is why we use IRIs

Coupling of the Penicillium
Duponti Acid Protease to
Ethylene-Maleic Acid (1 : 1)
Linear Copolymer.
Preparation and Properties of
the Water-Soluble Derivative

https://pubmed.ncbi.nlm.nih.gov/9993/

PubMed

*DGCR2 DiGeorge syndrome critical region*
*gene 2 [ Homo sapiens (human) ]*

http://identifiers.org/ncbigene/9993

NCBI Gene

http://identifiers.org/taxonomy/9993

I am 9993, would you like a flower

NCBI Taxonomy

## RDF not just another standard
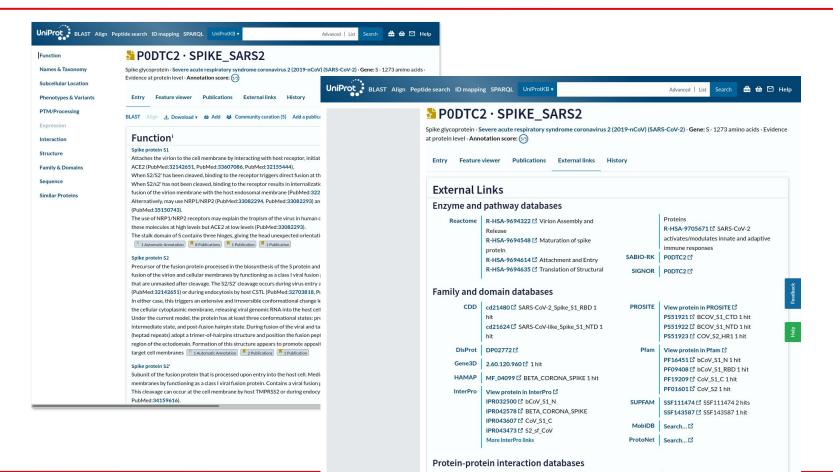
UniProt is not exhaustive
- Central point in the life sciences
- Linking to other resources

Trees have just one root (XML, JSON, Protocol Buffer etc)

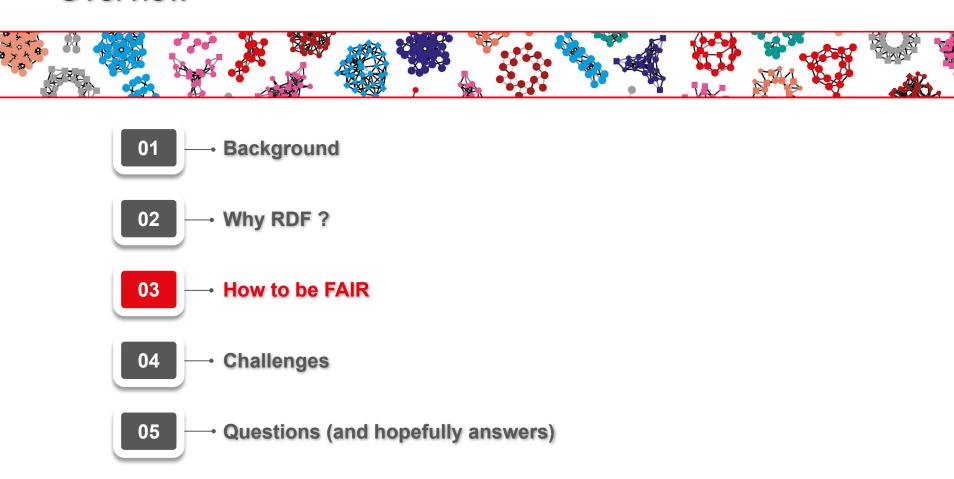Ecosystems are bigger than a single organism
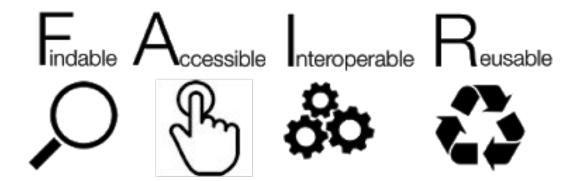
Can't dictate vendors on our users
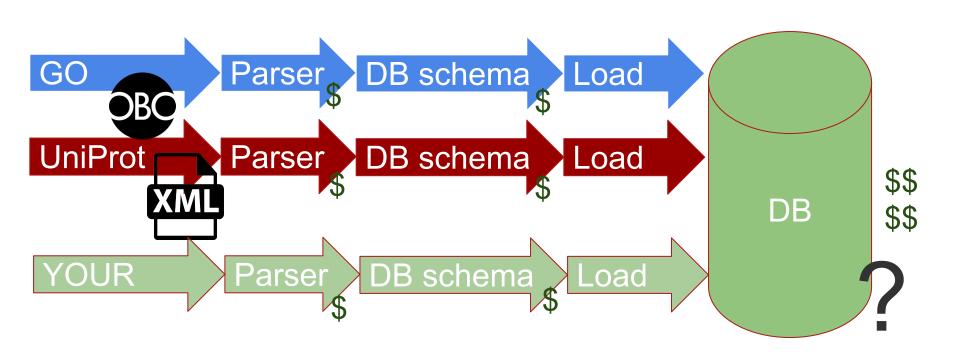
# Links

# wwPDB "another" Protein database

# Overview

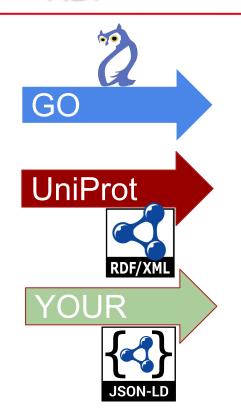# The FAIRest format of them all



F<sub>indable</sub>  A<sub>ccessible</sub>  I<sub>nteroperable</sub>  R<sub>eusable</sub>
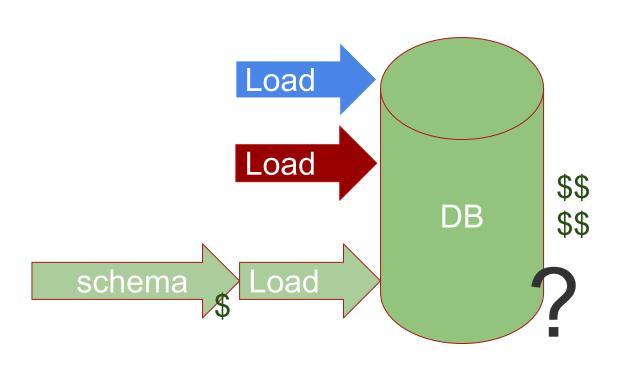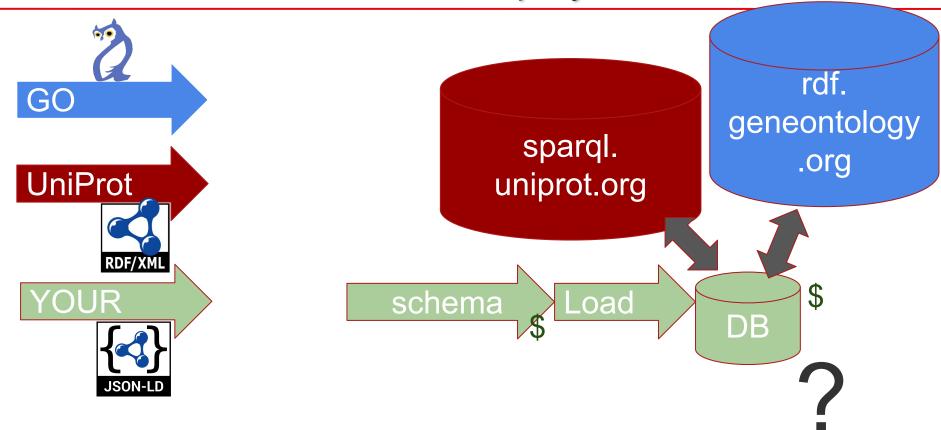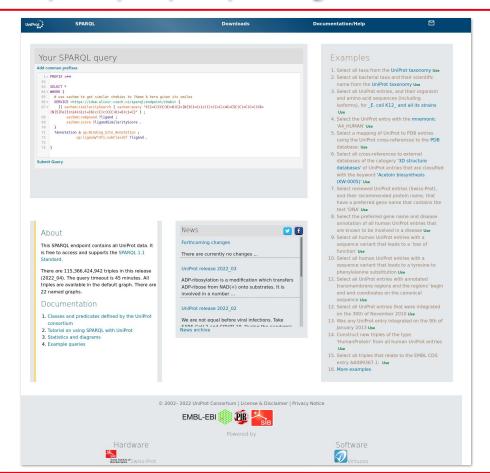
# Parse, parse, parse

# RDF

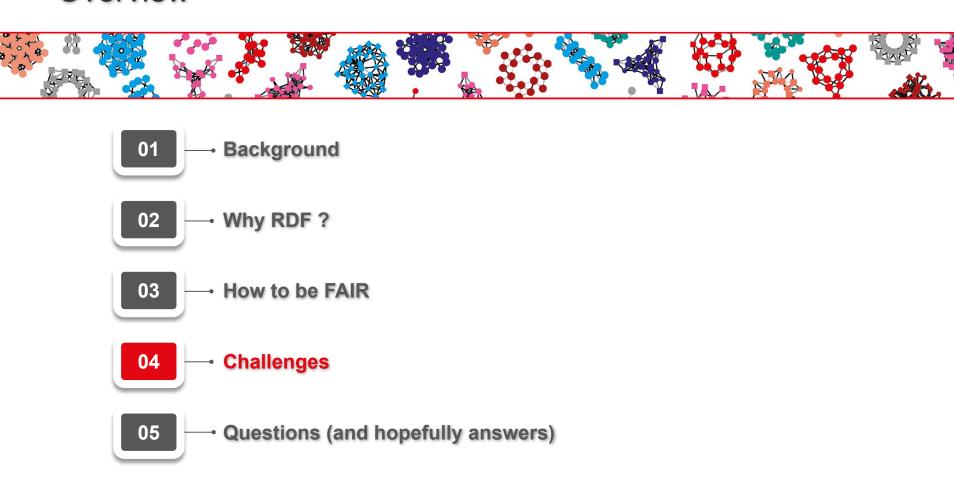# RDF + SPARQL 1.1 with build in query federation

# https://sparql.uniprot.org



Query
Examples
Documentation
Tutorials
Limits
Partners
News

# Hardware & Software

**Virtuoso 7.2+ Open Source**

**2*64 core EPYC
1 Tb Ram
24 Tb SSD NVMe**

**RDF4j  Custom Tomcat**

**Virtuoso 7.2+ Open Source**

**4*16 AMD Barcelona
256 GB Ram
8 Tb SSD Sata**

**RDF4j  Custom Tomcat**

**Virtuoso 7.2+ Open Source**

**4*16 AMD Barcelona
256 GB Ram
8 Tb SSD Sata**

# Overview

# Complexity of schema

- Entity Relation Schema Diagram
  - Large
- Biology has many concepts
  - Simplicity is misleading
  - Compare search to analysis

# Service description

- Drives documentation
- Used to rewrite hard queries
  - e.g. How many IRIs?
  - Taking care of variable and prefix use
- Note funders and database links

# Query speed

- Aggregates
- Distinct
- Order by

## Caching

- ETag: W/"2022_03"
- SERVICE
  - ETag code parses a query
  - Checks last modified of remote endpoint
  - Combined header

# Large literals

- DNA as strings
  - Multi-gigabase
  - Max int as an index is a problem
- Proteins
  - 70,000+ (tintin)

# Ecosystem

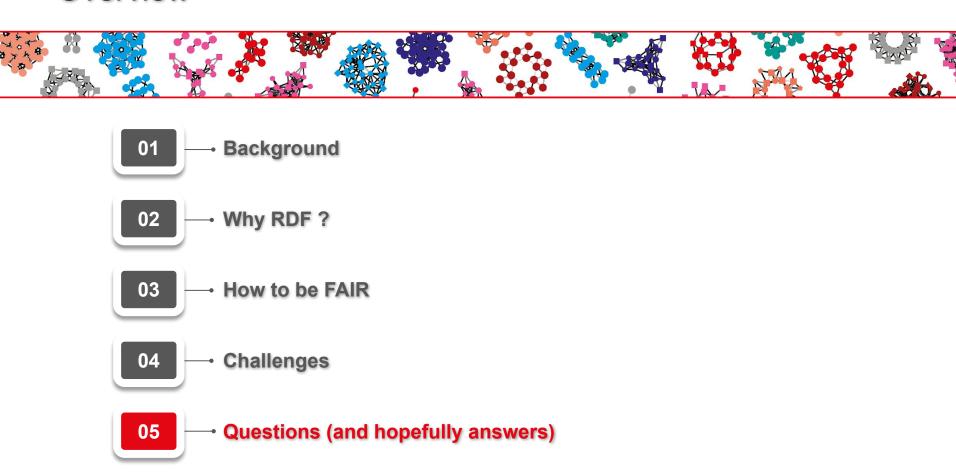- Out of date tools/reports
- Push instead of pull
- Schema.org

# Fast?

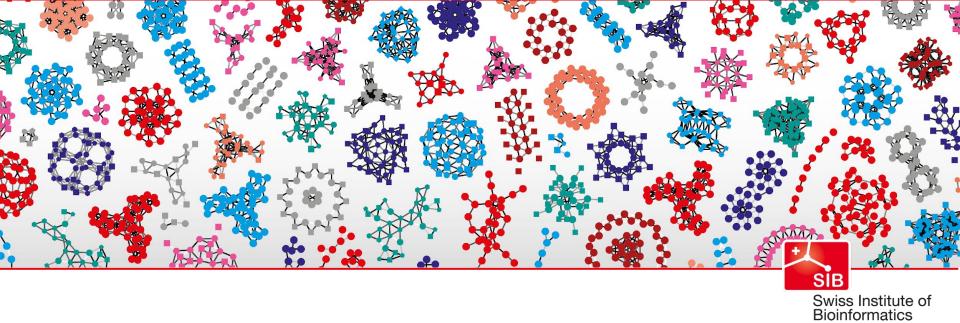- ● ZSTD
  - ○ Facebook says its fast
  - ○ Still takes a day to decompress 4TB

- ● Linked Data Fragments
  - ○ SELECT ?x WHERE {?x a ?c }

- ● IO & Network

# Overview

**UniProt: is a team effort!**