# Life Science Database Integration Using Linked Data

## Susumu Goto

Database Center for Life Science (DBCLS)
Joint Support-Center for Data Science Research (DS)
Research Organization of Information and Systems (ROIS)

International Life Science Integration Workshop
2018 / 3 / 6 @ Nakano Sunplaza, Tokyo, Japan

大学共同利用機関法人　情報・システム研究機構
データサイエンス共同利用基盤施設
Joint Support-Center for Data Science Research (DS)

大学共同利用機関法人
情報・システム研究機構
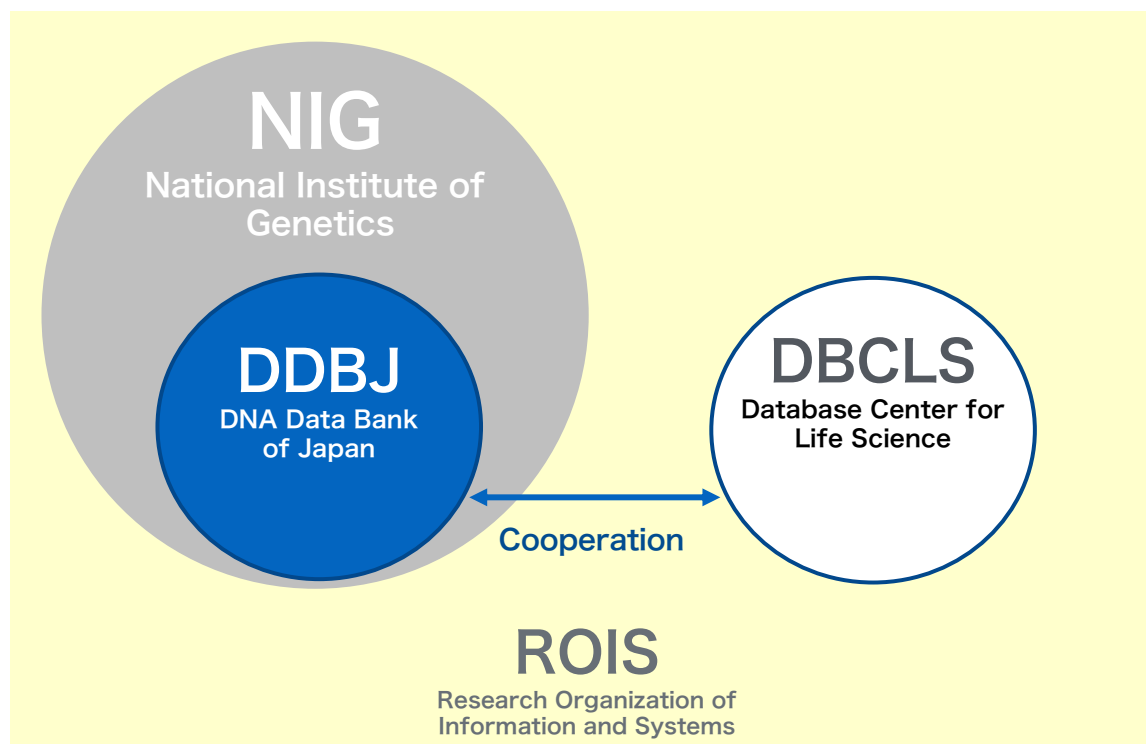Research Organization of Information and Systems

# Database Center for Life Science

- 2008-
  - Database integration based on web application
- 2011-
  - Funded by JST National Bioscience Database Center for the database integration with the FAIR principle

- Integbio DB Catalog
- LSDB Cross Search
- Life Science DB Archive
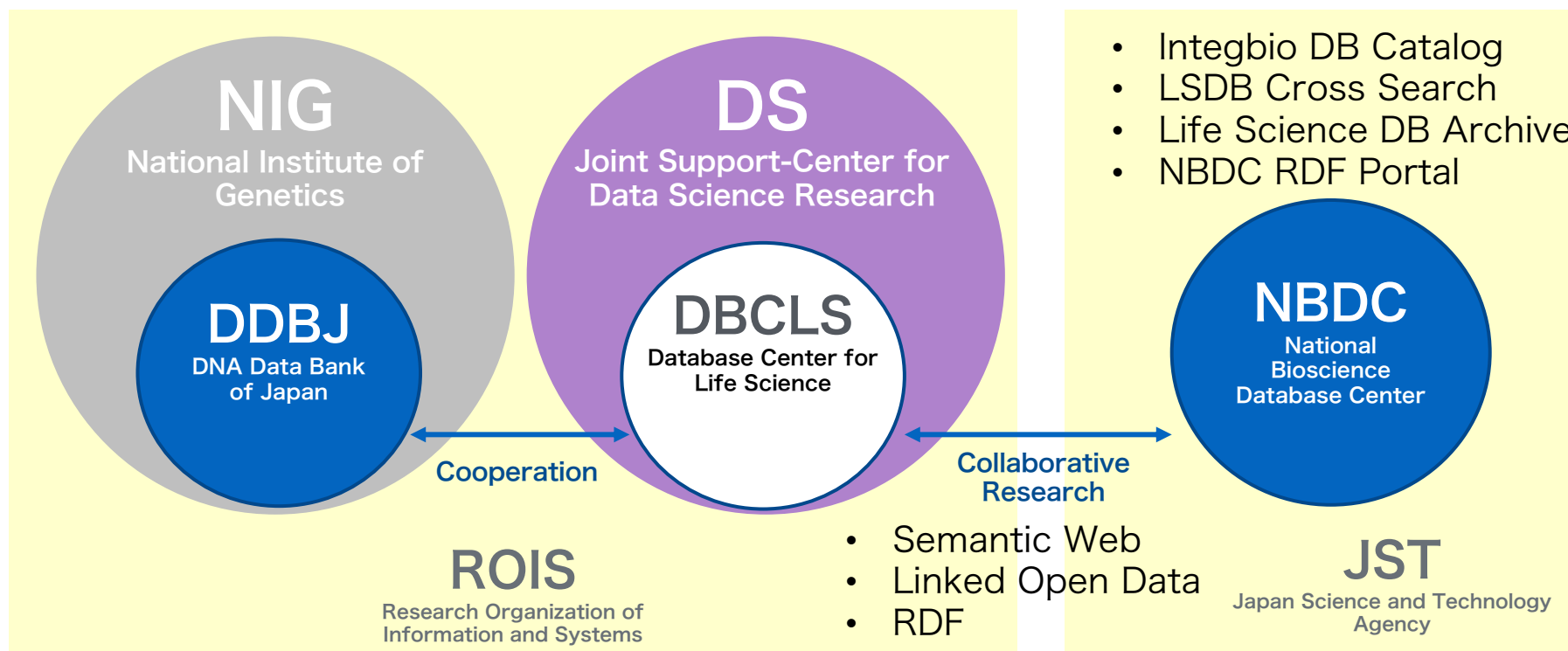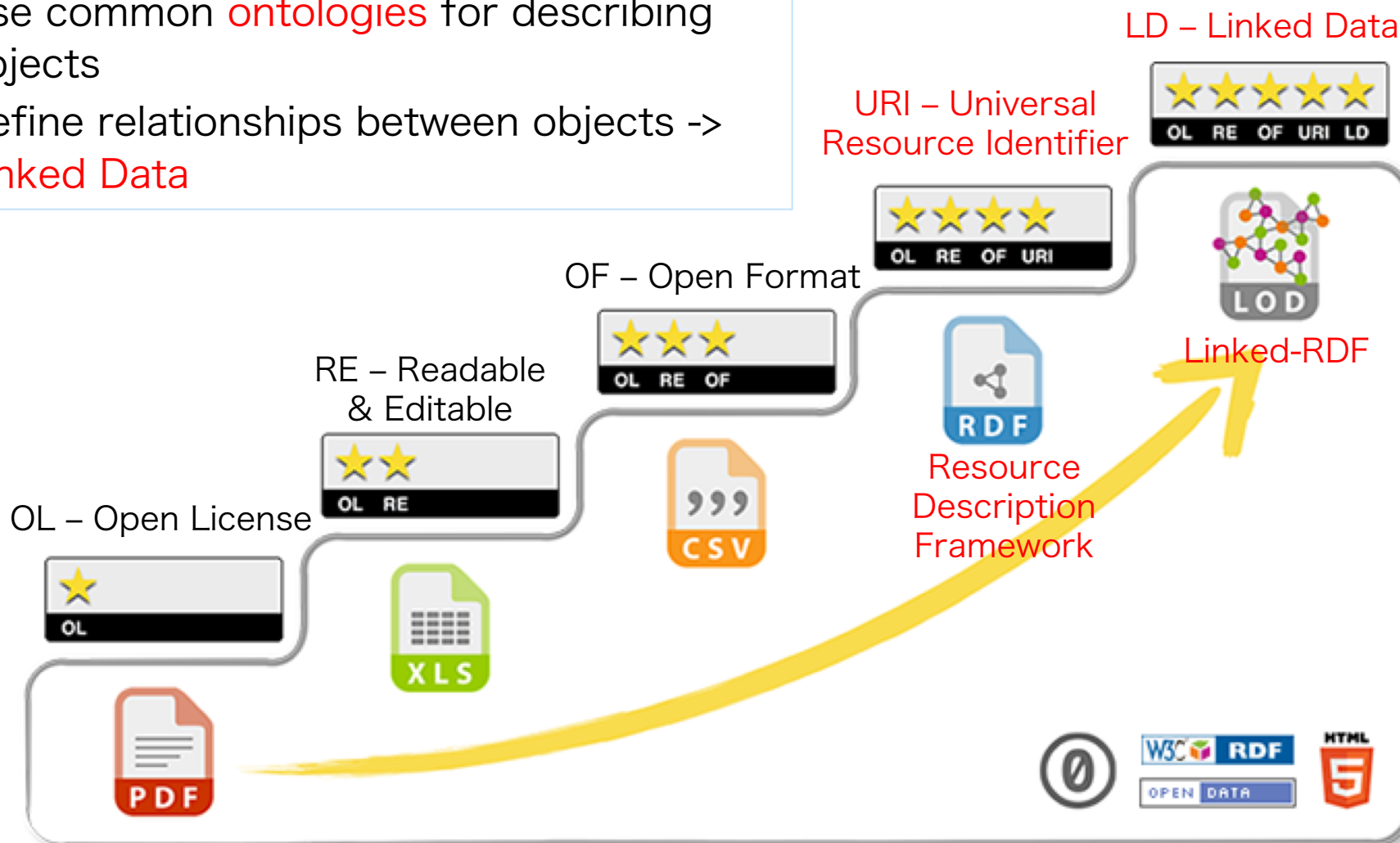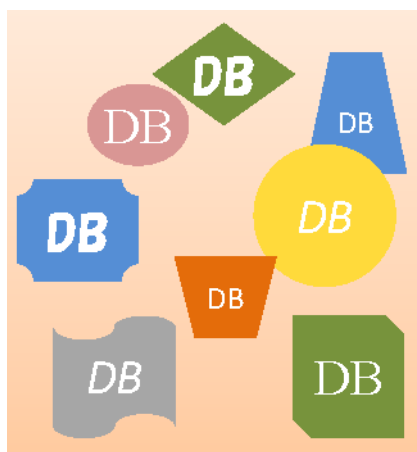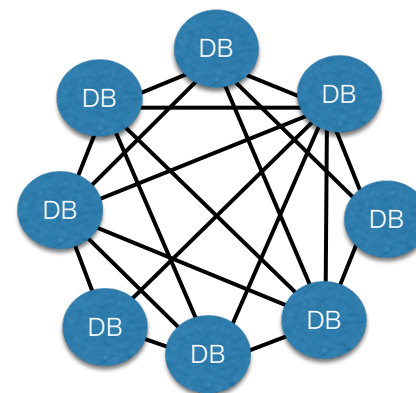- Technology development

FAIR
- Findable
- Accessible
- Interoperable
- Reusable



NIG
National Institute of Genetics

DDBJ
DNA Data Bank of Japan

DBCLS
Database Center for Life Science

Cooperation

ROIS
Research Organization of Information and Systems

# Database Center for Life Science

- 2008-
  - Database integration based on web application
- 2011-
  - Funded by JST National Bioscience Database Center for the database integration with the FAIR principle

**NIG**
National Institute of Genetics

**DDBJ**
DNA Data Bank of Japan

**DS**
Joint Support-Center for Data Science Research

**DBCLS**
Database Center for Life Science

- Integbio DB Catalog
- LSDB Cross Search
- Life Science DB Archive
- NBDC RDF Portal

**NBDC**
National Bioscience Database Center

Cooperation

Collaborative Research

**ROIS**
Research Organization of Information and Systems

- Semantic Web
- Linked Open Data
- RDF

**JST**
Japan Science and Technology Agency

3

# 5 ★ Linked Open Data

Tim Berners-Lee

- To give a unique ID to every object -> URI
- Use common ontologies for describing objects
- Define relationships between objects -> Linked Data

LD – Linked Data

URI – Universal Resource Identifier

OF – Open Format

RE – Readable & Editable

OL – Open License

Linked-RDF

Resource Description Framework

# Database Integration @ DBCLS

Ontology

RDF

Highly heterogeneous databases
using their own terms and formats

Databases integration for seamless
access and knowledge mining

- RDF: Resource Description Framework
- Triples consisting of Subject, Predicate and Object
  - Subject: ID (URI) for an object
  - Predicate: Attribute (URI) defined by an ontology
  - Object: ID (URI) or value (literal) for another object

# Database Integration @ DBCLS

# NBDC RDF Portal

- Portal site for RDF data from research groups in Japan
- 20 data sets including nine from NBDC funded databases comprising 45 billion triples (as of Nov. 2017)
- Microbial genomes, protein 3D structures, glycan structures, …
- RDF file download, SPARQL endpoints, Statistics, Metadata, …

Network of Databases

## Two important topics

- RDFyzing database guideline
  - http://wiki.lifesciencedb.jp/mw/BH14.14/RDFizingDatabaseGuidelineEnglishDraft0.1
- BioHackathons and SPARQLthons

# SPARQLthon

- Two days hackathon held every month from 2012 October.

- Theme: Life science database integration by semantic web technologies.

- >60 times in total and 1,328 (138 unique) participants from 45 institutes (15 universities, 13 research institutes, 17 private companies).

- From 2014, researchers from integrated database project funded by NBDC have attended and collaborated for creating RDF data and ontologies.

# Biohackathon

- International hackathon hosted by DBCLS/NBDC once a year in Japan from 2008

- Discuss and develop up-to-date technologies and systems for database integration and its applications

- One week intense development by international collaboration

- Summary papers have been published

- FAIR principle paper acknowledges biohackathon

# Currently Available RDF Data

| Type | RDF Data Set | Type | RDF Data Set |
|---|---|---|---|
| Gene | DDBJ | Ortholog | MBGD, PGDBj Orthology |
| Genome | Ensembl | Protein interaction | IntAct, Instruct, HINT |
| Metagenome | MicrobeDB.jp | Pathway | REACTOME, WikiPathway |
| Epigenome | KERO, ChIP-Atlas, iMETHYL | Systems biology | BioModels, SSBD |
| Genome variation | Linked ICGC, ClinVar, ExAC | Bioassay | ChEMBL, PubChem |
| Protein | UniProt | Disease | PAConto, GGDonto, DisGeNet, ClinVar, MedGen |
| Protein structure | wwPDB, BMRB, FAMSBASE | Dictionary | MeSH, Allie, LSD |
| Glycan | GlyTouCan, GlycoEpitope, WURCS | Transcriptome | ExpressionAtlas, RefEx, KERO, Open TG-GATEs |
| Chemical compound | PubChem, Nikkaji | Proteome | neXtProt, The Human Protein Atlas, jPOSTdb |
| Meta data | Quanto, integbio DB catalog, Colil, First Authors | Metabolome | MassBank, metabolonote |
| Sample | BioSamples, JCM | Ontology | BioProtal, OLS |

Red: RDF Portal、 Blue: On-going

# Tools for RDFyzing Data

## TogoDB

Converting table
data to RDB / RDF



http://togodb.org/

## D2RQ Mapper

Converting RDB to
RDF



http://d2rq.dbcls.jp/

# Database Integration @ DBCLS

# Middleware: Accessing SPARQL EPs



- **TogoStanza**: generic web framework for reusable web components
- **SPARQLList**: API for accessing SPARQL endpoints
- **SPARQL support, SPARQL builder**: web interface to support building SPARQL queries
- **YummyData**: listing and monitoring SPARQL endpoints

# YummyData: Information for SPARQL endpoint

- YummyData for endpoint information

- YummyViewer for visualization of class relationships

# Database Integration @ DBCLS

# Application: TogoGenome

- Genome database based on semantic web technology.

- Unique: implemented only by RDF data stores.

- >10,000 species including 360 eukaryotes.

- > 1 billion triples

- Genes and genomes, environmental and growth conditions, links to other DBs



Variation data

# Application: Easy access to omics data

# Application: Natural language Q&A



http://lodqa.org

# Summary

- Database integration via semantic web technology

  - RDF, Linked Open Data

  - RDF Portal and converting tools

- Tools to utilize integrated database

  - http://dbcls.jp/services

- Community for the development and utilization

  - Biohackathon

  - SPARQLthon

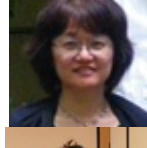  - Lecture series, TogoTV for lecture videos

# Acknowledgements

NBDC National Bioscience Database Center

JST Japan Science and Technology Agency