Wright State University

# CORE Scholar

4-11-2008

# RDB2RDF: Incorporating Domain Semantics in Structured Data

Satya S. Sahoo
*Wright State University - Main Campus*

Follow this and additional works at: https://corescholar.libraries.wright.edu/knoesis

Part of the Bioinformatics Commons, Communication Technology and New Media Commons, Databases and Information Systems Commons, OS and Networks Commons, and the Science and Technology Studies Commons

kno.e.sis

COLLECTING THE DOTS | CONNECTING THE DOTS

# RDB2RDF: Incorporating Domain Semantics in Structured Data

Satya S. Sahoo

Kno.e.sis Center, Computer Science and Engineering Department,
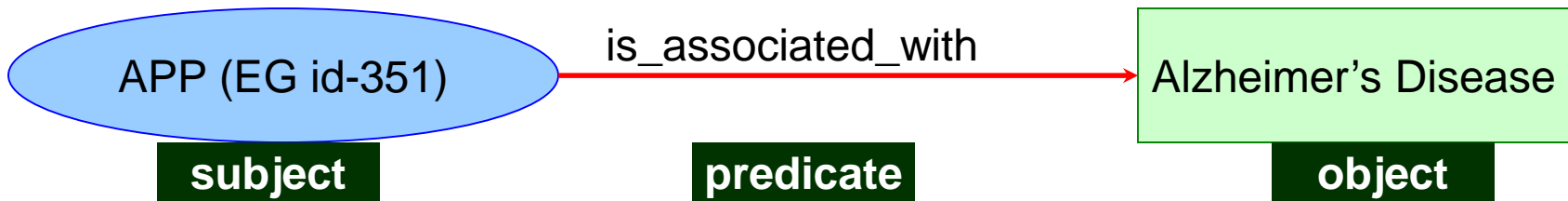
Wright State University, Dayton, OH, USA

# Acknowledgements

- Dr. Olivier Bodenreider (U.S NLM, NIH)
- Dr. Amit Sheth (Kno.e.sis Center, Wright State University)
- Dr. Joni L. Rutter (NIDA, NIH)
- Dr. Karen J. Skinner (NIDA, NIH)
- Lee Peters (U.S NLM, NIH)
- Kelly Zeng (U.S NLM, NIH)

# Outline

- RDB to RDF – Objectives
- Method I: RDB to RDF without ontology
- Application I: Genome ↔ Phenotype
- Method II: RDB to RDF with ontology
- Application II: Genome ↔ Biological Pathway integration
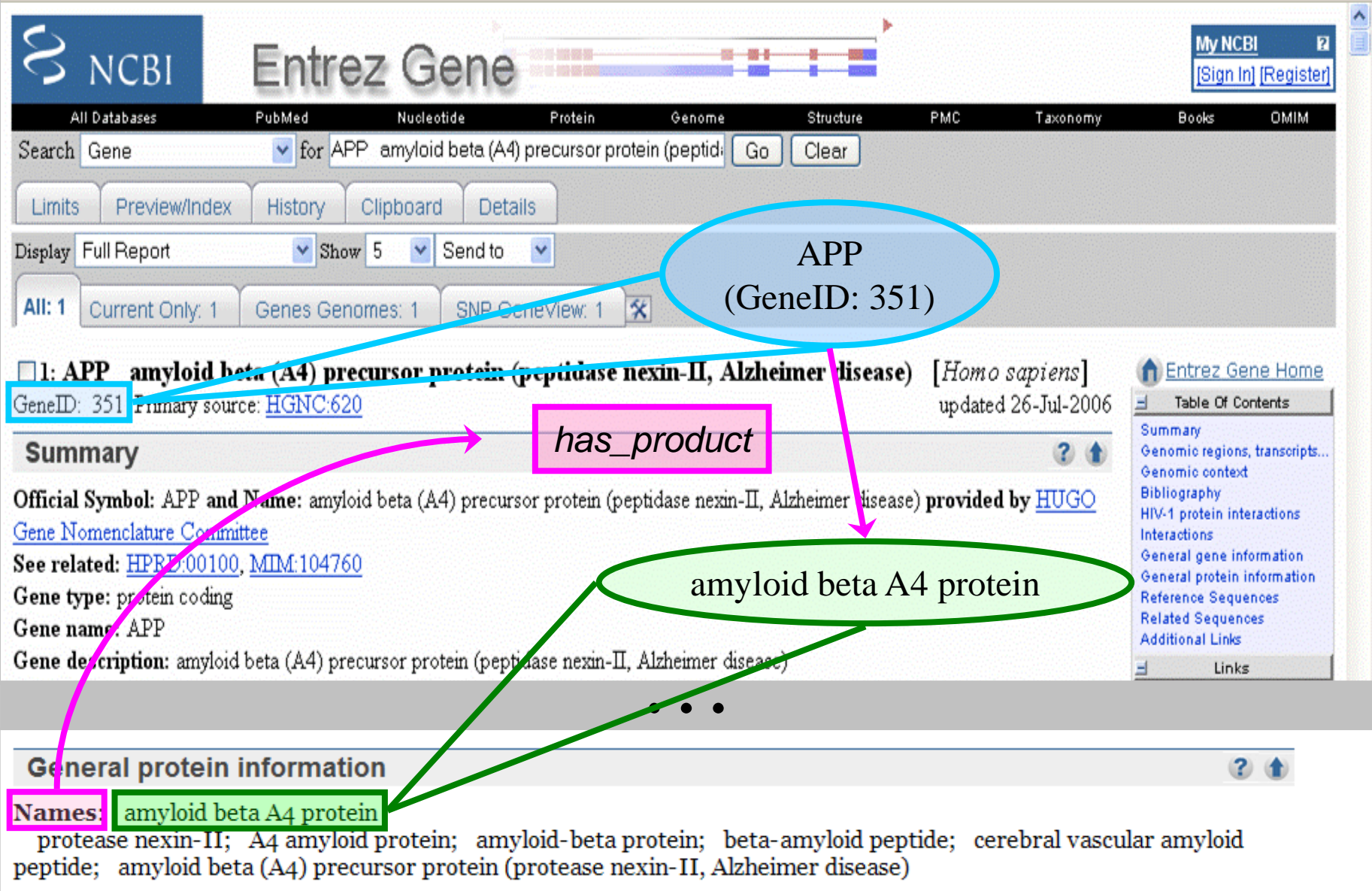- Conclusion

- RDF data model



- RDF enables modeling of logical relationship between entities

- Relations are at the heart of Semantic Web*

- RDF data - Logical Structure of the information

- Reasoning over RDF data → knowledge discovery

*Relationships at the Heart of Semantic Web: Modeling, Discovering, and Exploiting Complex Semantic Relationships,

Relationship Web: Blazing Semantic Trails between Web Resources

# Outline

- RDB to RDF – Objectives
- **Method I: RDB to RDF without ontology**
- **Application I: Genome ↔ Phenotype**
- Method II: RDB to RDF with ontology
- Application II: Genome ↔ Biological Pathway integration
- Conclusion

- NCBI Entrez Gene: gene related information from sequenced genomes and model organisms*
  - o 2 million gene records
  - o Gene information for genomic maps, sequences, homology, and protein expression
  - o Available in XML, ASN.1 and as a Webpage

*http://www.ncbi.nlm.nih.gov/sites/entrez/

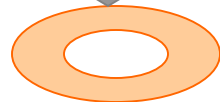# Entrez Gene Web Interface

- Mapped 106 elements tags out of 124 element tags to named relations

- 50GB XML file → 39GB RDF file (411 million RDF triples)

- Oracle 10g release 2 with part of the 10.2.03 patch

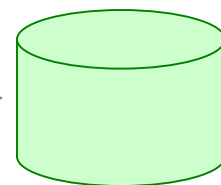- On a machine with 2 dual-core Intel Xeon 3.2GHz processor running Red Hat Enterprise Linux 4 (RHEL4)

## From *glycosyltransferase* to *congenital muscular dystrophy**

WRIGHT STATE
UNIVERSITY

# Outline

- In collaboration with National Institute on Drug Abuse (NIH)

- List of 449 human genes putatively involved with nicotine dependence (identified by Saccone et al.*)

- Understand gene functions and interactions, including their involvement in biological pathways

- List of queries:
  - *Which genes participate in a large number of pathways?*
  - *Which genes (or gene products) interact with each other?*
  - *Which genes are expressed in the brain?*

*S.F. Saccone, A.L. Hinrichs, N.L. Saccone, G.A. Chase, K. Konvicka and P.A. Madden et al., Cholinergic nicotinic receptor genes implicated in a nicotine dependence association study targeting 348 candidate genes with 3713 SNPs, Hum Mol Genet 16 (1) (2007), pp. 36–49

WRIGHT STATE
UNIVERSITY

- Method I: cannot answer query "*Which genes participate in a large number of pathways?*"

- Need to specify a particular instance of gene or pathway as starting point in RDF graph
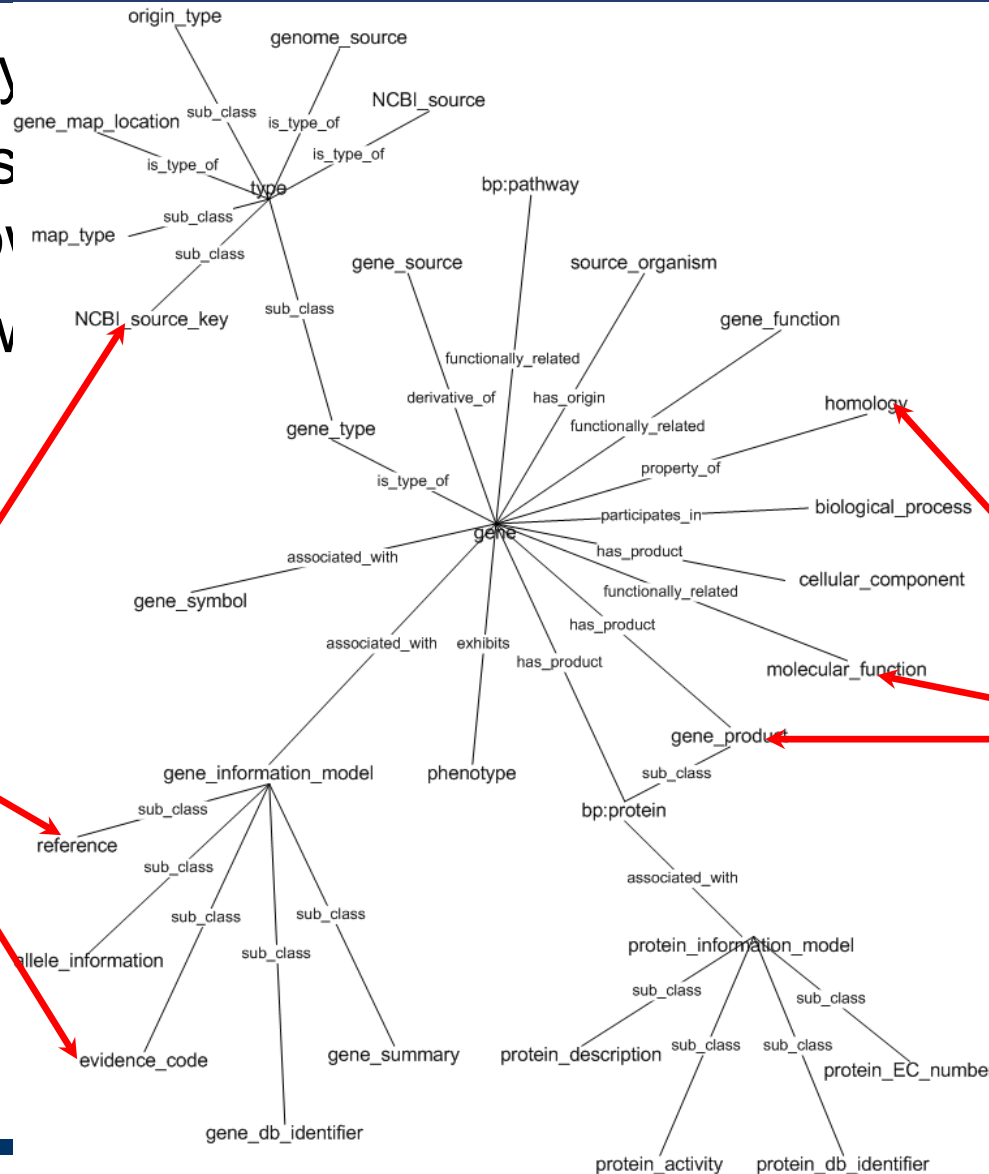
- Need to *classify* RDF instance data – Schema + Instance



SCHEMA

INSTANCE

- No ontology
- Created a s                                                    Entrez Gene – Entrez Kno
- Integrated w                                                    l pathway data)



Information model concepts

Domain concepts

# Outline

- RDB to RDF – Objectives
- Method I: RDB to RDF without ontology
- Application I: Genome ↔ Phenotype
- Method II: RDB to RDF with ontology
- Application II: Genome ↔ Biological Pathway integration
- **Conclusion**

- Application driven approach for RDB to RDF – Biomedical Knowledge Integration

- Explicit modeling of <span style="color:red">domain semantics</span> using named relations for

  o Accurate context based querying

  o Enhanced reasoning using relations based logic rules

- Use of ontology as reference knowledge model

- GRDDL compatible approach (using XSLT stylesheet) for transformation of RDB to RDF

- More information at:

http://knoesis.wright.edu/research/semsci/application_domain/sem_life_sci/bio/research/

*Thank you*