



# Converting biological information to the W3C Resource Description Framework (RDF): Experience with Entrez Gene

Presented By: Satya Sanket Sahoo

Mentor: Dr. Olivier Bodenreider

# Outline

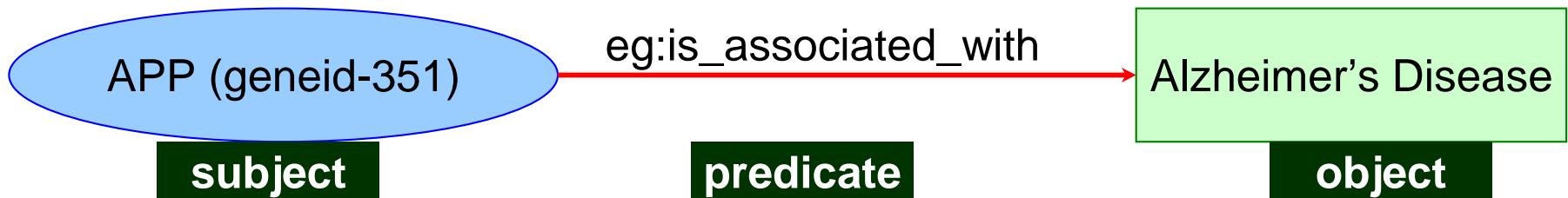
- Motivation
- RDF – Background
- Implementation technique
- Inference
- Unique identifiers
- Issues and challenges

# Motivation: knowledge management

- Concentrate on the logical structure of data
- Explicit definition of terms and relationships
- Information integration – one universe for data from diverse background
- Inference: use existing knowledge to infer implicit knowledge

# Resource Description Framework

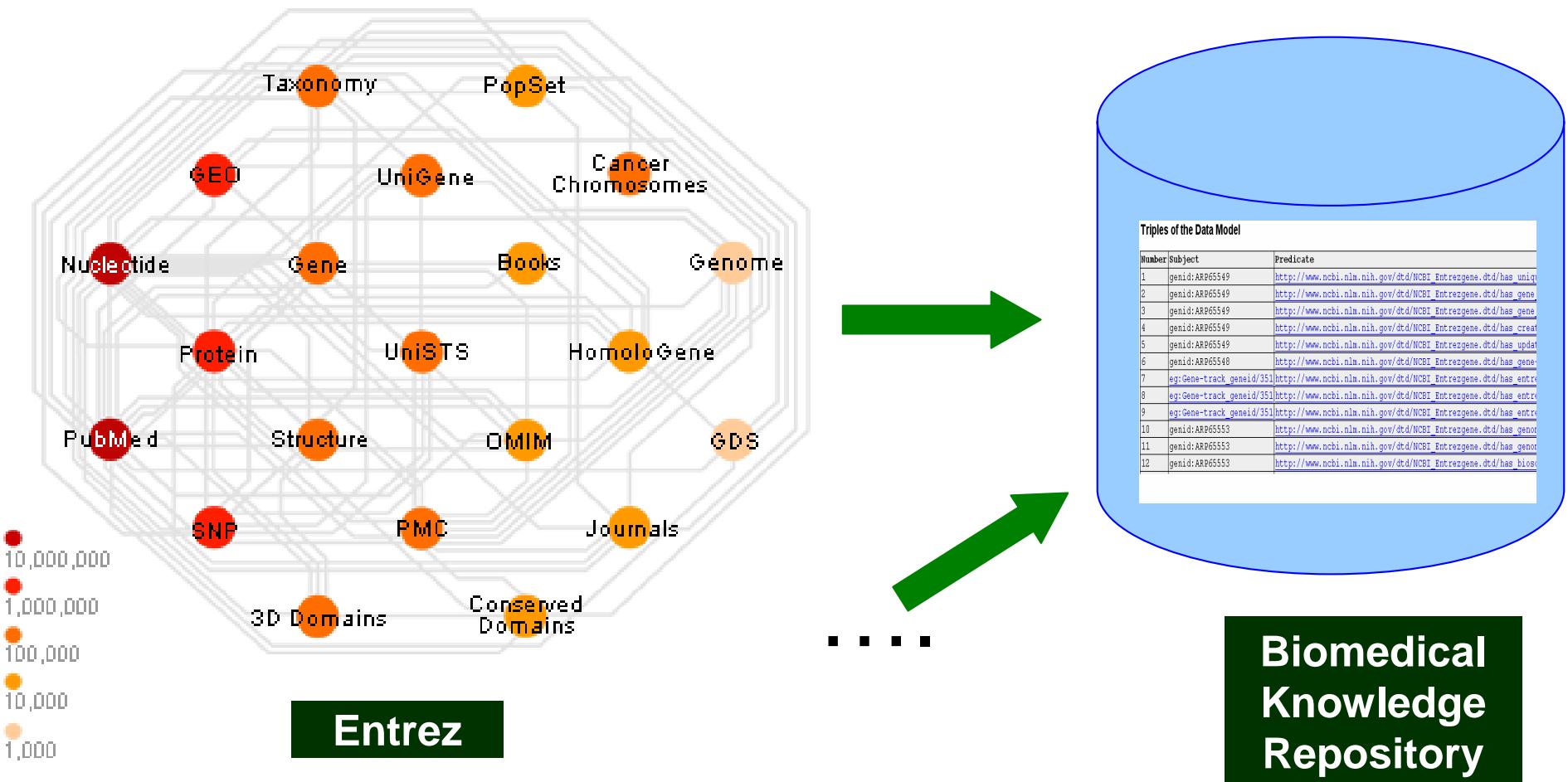
- All information represented as a ‘triple’



- Advantages include:
  - machine ‘understandable’
  - enables inference
  - represents the logical structure of the data
  - integration of data under one universe

Namespace - eg = [http://www.ncbi.nlm.nih.gov/dtd/NCBI\\_Entrezgene.dtd/](http://www.ncbi.nlm.nih.gov/dtd/NCBI_Entrezgene.dtd/)

# RDF – contd.



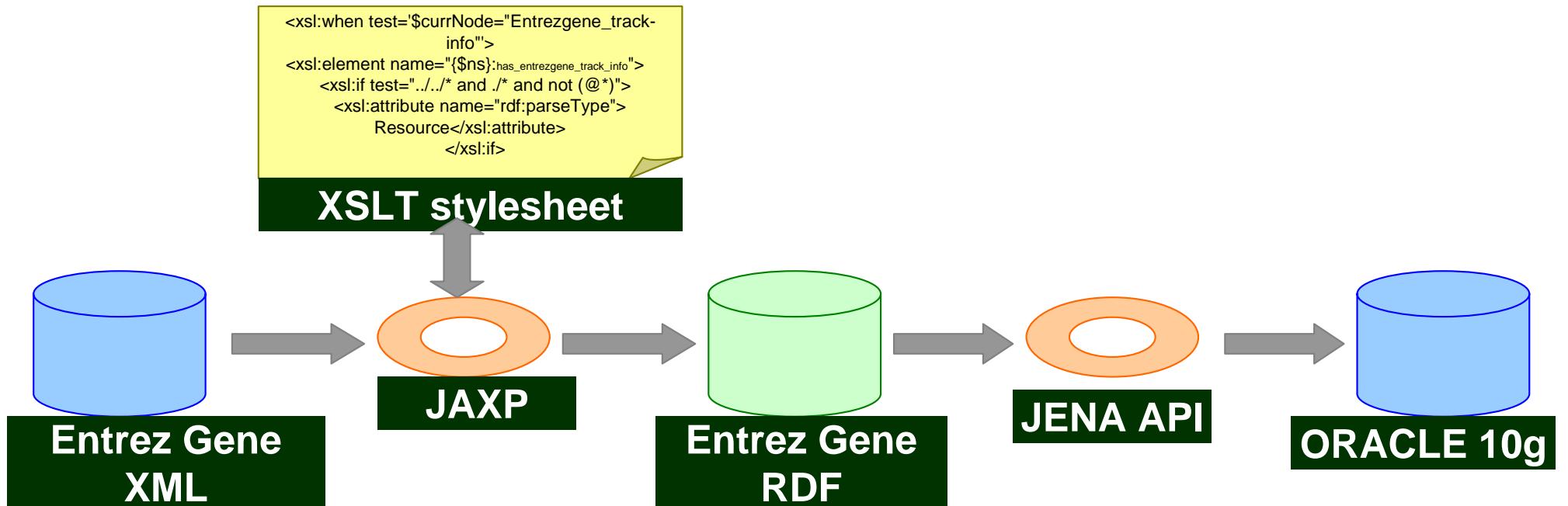
# RDF – contd.

- RDF triples can be thought as normalized assertions
- Similar to normalization of text
- But, instead of **lexical resemblance** RDF triples enable **semantic resemblance**

# Implementation: Entrez Gene XML to RDF

- Mapped element tags to more meaningful relations
- Started building an ontology of relationships
- Using XSLT stylesheet and XPath expressions converted XML to RDF
- The RDF reflects the nesting structure of terms in the Entrez gene records

# Implementation: Entrez Gene XML to RDF



- Modular - Separates application code from transformation framework
- Extensible – specific stylesheets may be used to for each of the Entrez databases
- Flexible – changes in application logic or transformation logic are separate

# Implementation

**Entrez Gene**

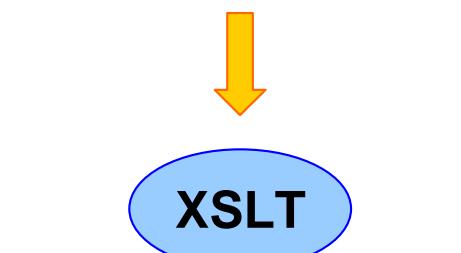
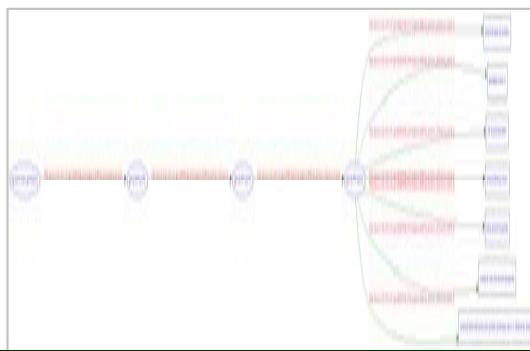


```

- <Entrezgene-Set>
- <Entrezgene_prot>
- <Prot-ref>
- <Prot-ref_name>
  <Prot-ref_name_E>amyloid beta A4 protein</Prot-ref_name_E>
  <Prot-ref_name_E>protease nexin-II</Prot-ref_name_E>
  <Prot-ref_name_E>A4 amyloid protein</Prot-ref_name_E>
  <Prot-ref_name_E>amyloid-beta protein</Prot-ref_name_E>
  <Prot-ref_name_E>beta-amyloid peptide</Prot-ref_name_E>
  <Prot-ref_name_E>cerebral vascular amyloid peptide</Prot-ref_name_E>
- <Prot-ref_name>
  amyloid beta (A4) precursor protein (protease nexin-II, Alzheimer disease)
</Prot-ref_name>
</Prot-ref>
</Entrezgene_prot>
</Entrezgene-Set>

```

**Entrez Gene XML**

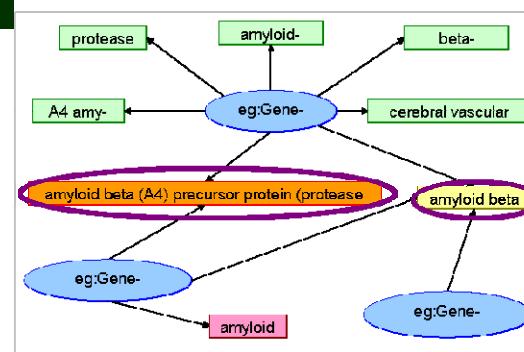


```

<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:eg="http://www.ncbi.nlm.nih.gov/dtd/NCBI_Entrezgene.dtd">
  <rdf:Description rdf:about="eg:Gene-track_geneid/351">
    <eg:has_entrezgene_protein rdf:type="Resource">
      <eg:has_protein_reference rdf:type="Resource">
        <eg:has_protein_reference_name rdf:type="Resource">
          <eg:has_protein_reference_name_E>amyloid beta A4 protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>protease nexin-II</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>A4 amyloid protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid-beta protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>beta-amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>cerebral vascular amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid beta (A4) precursor protein (protease nexin-II, Alzheimer disease)</eg:has_protein_reference_name_E>
        </eg:has_protein_reference_name>
      </eg:has_protein_reference>
    </eg:has_entrezgene_protein>
  </rdf:Description>
</rdf:RDF>

```

**Entrez Gene RDF**



# Web interface

NCBI Entrez Gene

All Databases PubMed Nucleotide Protein Genome Structure PMC Taxonomy Books OMIM

Search Gene for APP amyloid beta (A4) precursor protein (peptidase nexin-II, Alzheimer disease) [Homo sapiens] Go Clear

Limits Preview/Index History Clipboard Details

Display Full Report Show 5 Send to

All: 1 Current Only: 1 Genes Genomes: 1 SNP GeneView: 1

1. APP - amyloid beta (A4) precursor protein (peptidase nexin-II, Alzheimer disease) [Homo sapiens]  
GeneID: 351 Primary source: HGNC:620 updated 26-Jul-2006

**Summary**

Official Symbol: APP and Name: amyloid beta (A4) precursor protein (peptidase nexin-II, Alzheimer disease) provided by HUGO Gene Nomenclature Committee  
See related: HPRD:00100, MIM:104760  
Gene type: protein coding  
Gene name: APP  
Gene description: amyloid beta (A4) precursor protein (peptidase nexin-II, Alzheimer disease)

...  
**General protein information**

**Names:** amyloid beta A4 protein  
protease nexin-II, A<sub>4</sub> amyloid protein; amyloid-beta protein; beta-amyloid peptide; cerebral vascular amyloid peptide; amyloid beta (A4) precursor protein (protease nexin-II, Alzheimer disease)

Entrez Gene Home  
Table Of Contents  
Summary  
Genomic regions, transcripts...  
Genomic context  
Bibliography  
HIV-1 protein interactions  
Interactions  
General gene information  
General protein information  
Reference Sequences  
Related Sequences  
Additional Links  
Links

# Implementation

The screenshot shows the Entrez Gene page for the APP gene. The search bar at the top has "APP amyloid beta (A4) precursor protein (peptidase nexin-II, Alzheimer disease)" entered. Below the search bar, the gene summary for APP is displayed, including its official symbol (APP), name (amyloid beta (A4) precursor protein (peptidase nexin-II, Alzheimer disease)), and gene ID (351). The page also includes sections for Gene Homologs, Summary, and References.

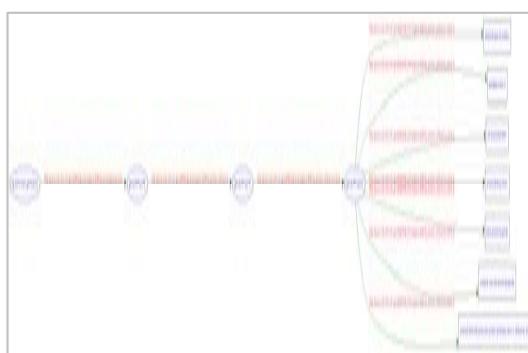
Entrez Gene



```

<-> EntrezGene-Set
  - <Entrezgene_prot>
    - <Prot-ref>
      - <Prot-ref_name>
        <Prot-ref_name_E>amyloid beta A4 protein</Prot-ref_name_E>
        <Prot-ref_name_E>protease nexin-II</Prot-ref_name_E>
        <Prot-ref_name_E>A4 amyloid protein</Prot-ref_name_E>
        <Prot-ref_name_E>amyloid-beta peptide</Prot-ref_name_E>
        <Prot-ref_name_E>beta-amyloid peptide</Prot-ref_name_E>
        <Prot-ref_name_E>cerebral vascular amyloid peptide</Prot-ref_name_E>
      - <Prot-ref_name>
        <Prot-ref_name_E>amyloid (A4) precursor protein (protease nexin-II, Alzheimer disease)</Prot-ref_name_E>
      - </Prot-ref_name>
    - </Prot-ref>
  - </Entrezgene_prot>
</EntrezGene-Set>
  
```

Entrez Gene XML



Entrez Gene RDF graph



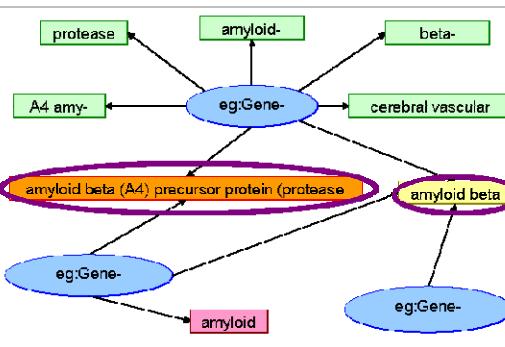
XSLT



```

<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:eg="http://www.ncbi.nlm.nih.gov/dtd/NCBI_Entrezgene.dtd">
  <rdf:Description rdf:about="eg:Gene-track_geneid/351">
    <eg:has_entrezgene_protein rdf:type="Resource">
      <eg:has_protein_reference rdf:type="Resource">
        <eg:has_protein_reference_name rdf:type="Resource">
          <eg:has_protein_reference_name_E>amyloid beta A4 protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>protease nexin-II</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>A4 amyloid protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid-beta peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>beta-amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>cerebral vascular amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid (A4) precursor protein (protease nexin-II, Alzheimer disease)</eg:has_protein_reference_name_E>
        </eg:has_protein_reference_name>
      </eg:has_protein_reference>
    </eg:has_entrezgene_protein>
  </rdf:Description>
</rdf:RDF>
  
```

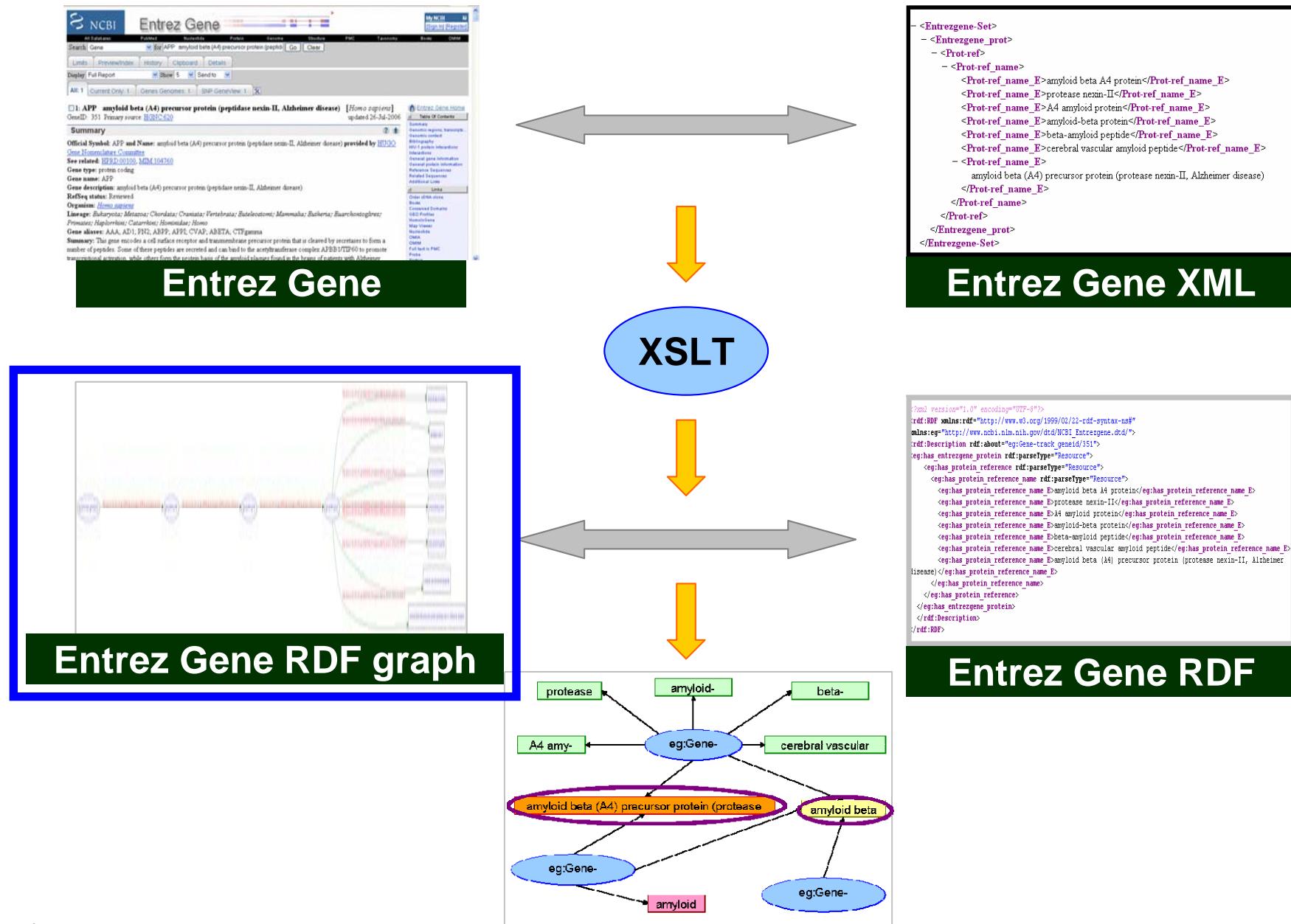
Entrez Gene RDF



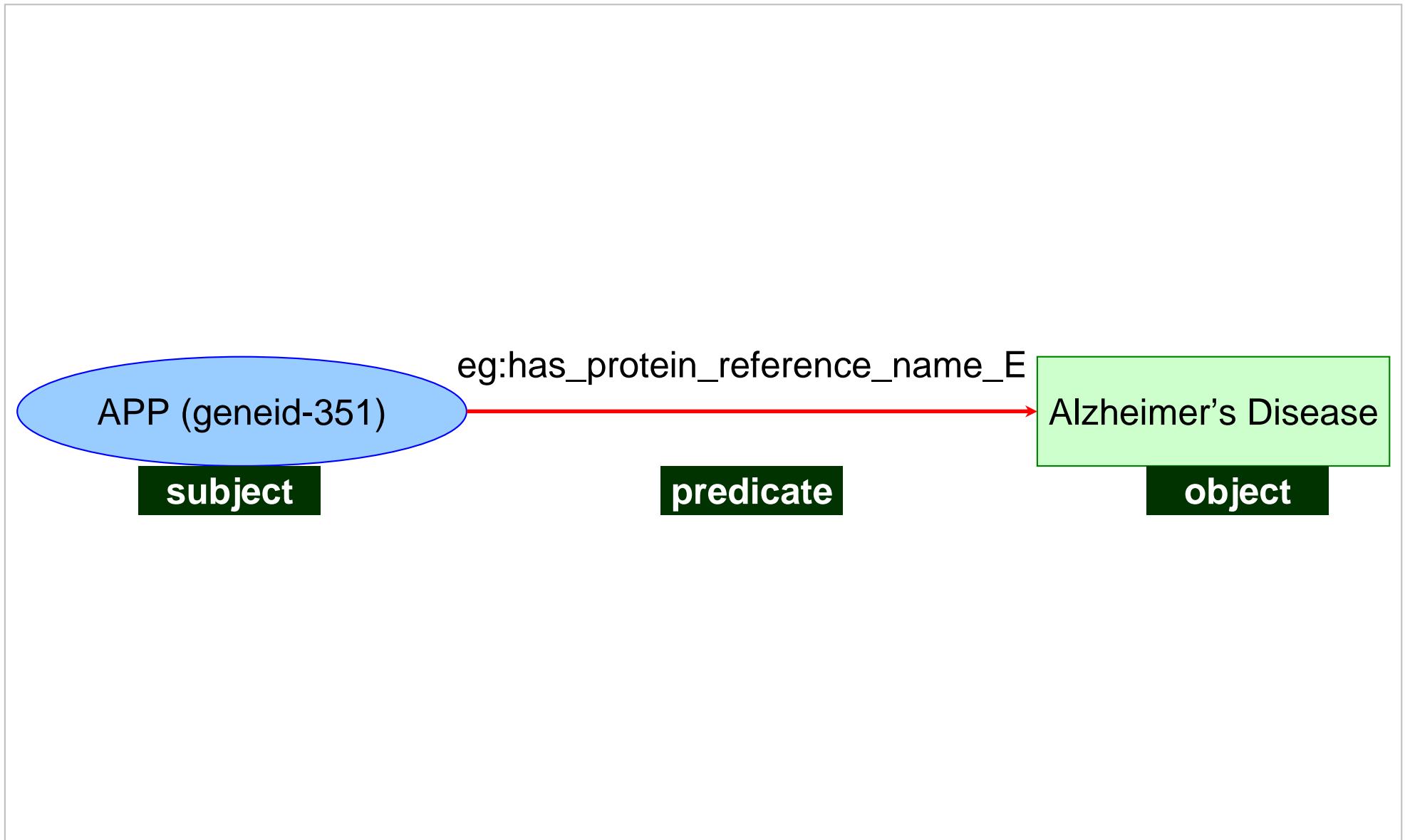
# XML

```
<?xml version="1.0"?>
<!DOCTYPE Entrezgene-Set PUBLIC "-//NLM//DTD NCBI-Entrezgene, 21st January 2005//EN" "NCBI_Entrezgene.dtd">
<Entrezgene-Set>
<Entrezgene>
<Entrezgene_track-info>
  <Gene-track>
    <Gene-track_geneid>351</Gene-track_geneid>
  </Gene-track>
</Entrezgene_track-info>
<Entrezgene_prot>
  <Prot-ref>
    <Prot-ref_name>
      <Prot-ref_name_E>amyloid beta A4 protein</Prot-ref_name_E>
      <Prot-ref_name_E>protease nexin-II</Prot-ref_name_E>
      <Prot-ref_name_E>A4 amyloid protein</Prot-ref_name_E>
      <Prot-ref_name_E>amyloid-beta protein</Prot-ref_name_E>
      <Prot-ref_name_E>beta-amyloid peptide</Prot-ref_name_E>
      <Prot-ref_name_E>cerebral vascular amyloid peptide</Prot-ref_name_E>
      <Prot-ref_name_E>amyloid beta (A4) precursor protein (protease nexin-II, Alzheimer disease)</Prot-ref_name_E>
    </Prot-ref_name>
  </Prot-ref>
</Entrezgene_prot>
</Entrezgene>
</Entrezgene-Set>
```

# Implementation



# RDF Graph



# RDF Graph



Entrez Gene RDF graph  
(W3C Validator Site - <http://www.w3.org/RDF/Validator/>)

# Implementation

The screenshot shows the NCBI Entrez Gene page for gene ID 351. The search term is "APP amyloid beta (A4) precursor protein (peptidase nexin-II, Alzheimer disease)". The page displays various sections including Summary, Gene Structure, and References.

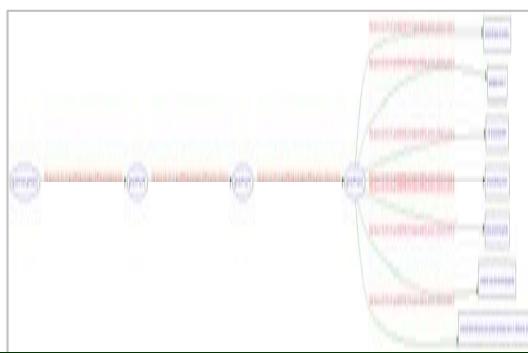
Entrez Gene

```

<- <Entrezgene-Set>
  - <Entrezgene_prot>
    - <Prot-ref>
      - <Prot-ref_name>
        <Prot-ref_name_E>amyloid beta A4 protein</Prot-ref_name_E>
        <Prot-ref_name_E>protease nexin-II</Prot-ref_name_E>
        <Prot-ref_name_E>A4 amyloid protein</Prot-ref_name_E>
        <Prot-ref_name_E>amyloid-beta protein</Prot-ref_name_E>
        <Prot-ref_name_E>beta-amyloid peptide</Prot-ref_name_E>
        <Prot-ref_name_E>cerebral vascular amyloid peptide</Prot-ref_name_E>
      - <Prot-ref_name>
        <Prot-ref_name_E>amyloid beta (A4) precursor protein (protease nexin-II, Alzheimer disease)</Prot-ref_name_E>
      </Prot-ref_name>
    </Prot-ref>
  </Entrezgene_prot>
</Entrezgene-Set>

```

Entrez Gene XML



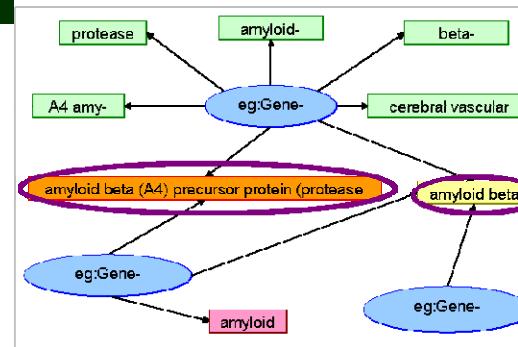
Entrez Gene RDF graph

```

<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:eg="http://www.ncbi.nlm.nih.gov/dtd/NCBI_Entrezgene.dtd">
  <rdf:Description rdf:about="eg:Gene-track_geneid/351">
    <eg:has_entrezgene_protein rdf:type="Resource">
      <eg:has_protein_reference rdf:type="Resource">
        <eg:has_protein_reference_name rdf:type="Resource">
          <eg:has_protein_reference_name_E>amyloid beta A4 protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>protease nexin-II</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>A4 amyloid protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid-beta protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>beta-amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>cerebral vascular amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid beta (A4) precursor protein (protease nexin-II, Alzheimer disease)</eg:has_protein_reference_name_E>
        </eg:has_protein_reference_name>
      </eg:has_protein_reference>
    </eg:has_entrezgene_protein>
  </rdf:Description>
</rdf:RDF>

```

Entrez Gene RDF



# RDF

```
<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:eg="http://www.ncbi.nlm.nih.gov/dtd/NCBI_Entrezgene.dtd/">
  <rdf:Description rdf:about="eg:Gene-track_geneid/351">
    <eg:has_entrezgene_protein rdf:parseType="Resource">
      <eg:has_protein_reference rdf:parseType="Resource">
        <eg:has_protein_reference_name rdf:parseType="Resource">
          <eg:has_protein_reference_name_E>amyloid beta A4 protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>protease nexin-II</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>A4 amyloid protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid-beta protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>beta-amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>cerebral vascular amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid beta (A4) precursor protein (protease nexin-II, Alzheimer
disease)</eg:has_protein_reference_name_E>
        </eg:has_protein_reference_name>
      </eg:has_protein_reference>
    </eg:has_entrezgene_protein>
  </rdf:Description>
</rdf:RDF>
```

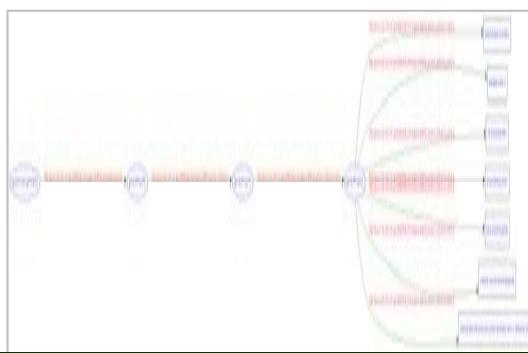
# Implementation

The screenshot shows the NCBI Entrez Gene page for gene ID 351. The search term is "APP amyloid beta (A4) precursor protein (peptidase nexin-II, Alzheimer disease)". The page displays various sections including Summary, Gene Structure, and Reference Sequences.

Entrez Gene

```
<Entrezgene-Set>
  - <Entrezgene_prot>
    - <Prot-ref>
      - <Prot-ref_name>
        <Prot-ref_name_E>amyloid beta A4 protein</Prot-ref_name_E>
        <Prot-ref_name_E>protease nexin-II</Prot-ref_name_E>
        <Prot-ref_name_E>A4 amyloid protein</Prot-ref_name_E>
        <Prot-ref_name_E>amyloid-beta protein</Prot-ref_name_E>
        <Prot-ref_name_E>beta-amyloid peptide</Prot-ref_name_E>
        <Prot-ref_name_E>cerebral vascular amyloid peptide</Prot-ref_name_E>
      - <Prot-ref_name>
        <Prot-ref_name_E>amyloid (A4) precursor protein (protease nexin-II, Alzheimer disease)</Prot-ref_name_E>
      </Prot-ref_name>
    </Prot-ref>
  </Entrezgene_prot>
</Entrezgene-Set>
```

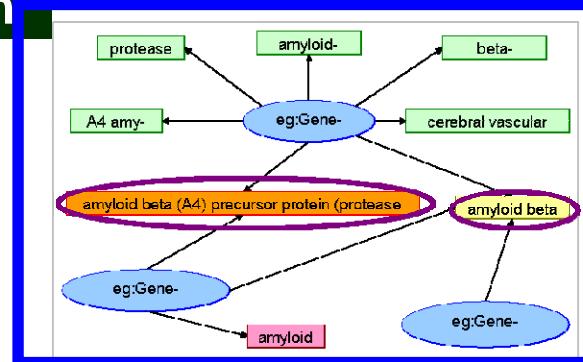
Entrez Gene XML



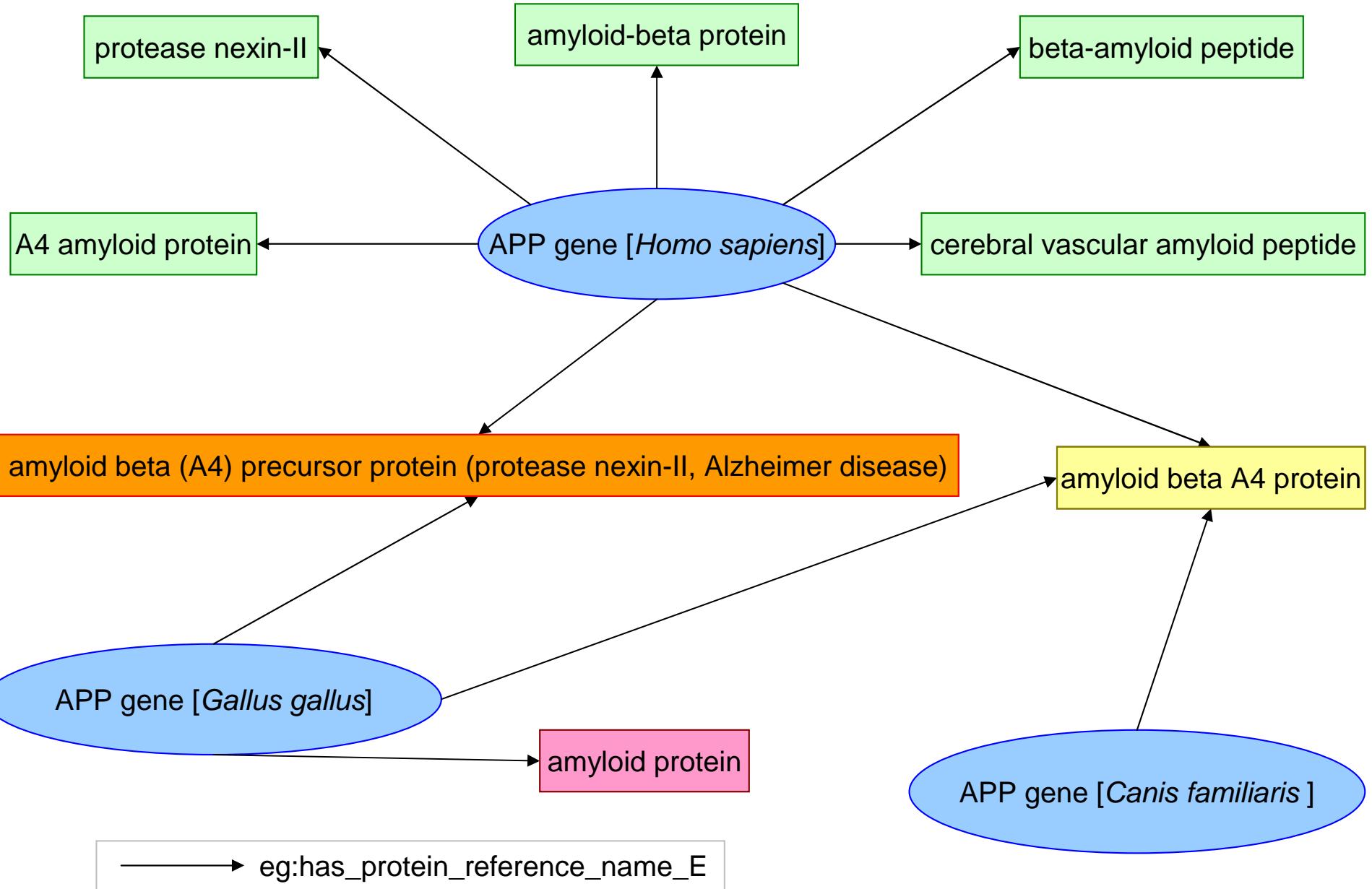
Entrez Gene RDF graph

```
<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:eg="http://www.ncbi.nlm.nih.gov/dtd/NCBI_Entrezgene.dtd">
  <rdf:Description rdf:about="eg:Gene-track_geneid/351">
    <eg:has_entrezgene_protein rdf:type="Resource">
      <eg:has_protein_reference rdf:type="Resource">
        <eg:has_protein_reference_name rdf:type="Resource">
          <eg:has_protein_reference_name_E>amyloid beta A4 protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>protease nexin-II</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>A4 amyloid protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid-beta protein</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>beta-amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>cerebral vascular amyloid peptide</eg:has_protein_reference_name_E>
          <eg:has_protein_reference_name_E>amyloid beta (A4) precursor protein (protease nexin-II, Alzheimer disease)</eg:has_protein_reference_name_E>
        </eg:has_protein_reference_name>
      </eg:has_protein_reference>
    </eg:has_entrezgene_protein>
  </rdf:Description>
</rdf:RDF>
```

Entrez Gene RDF

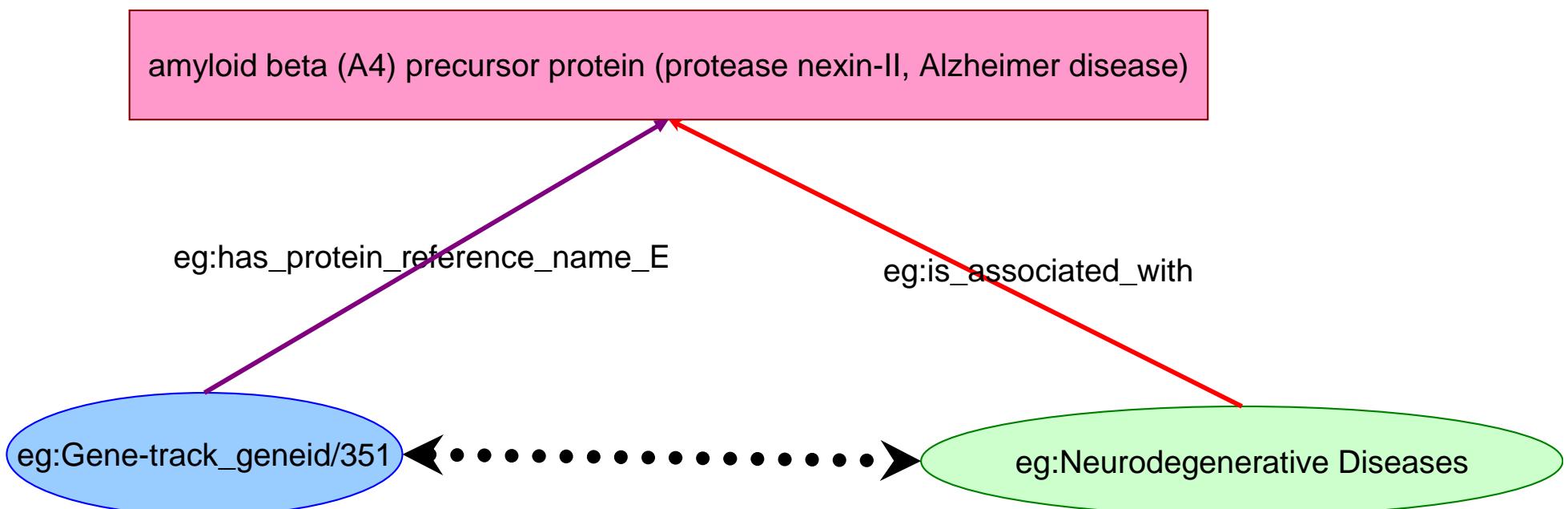


# Connecting different genes



# Inference

- Rules are objects that allow inference from RDF data [1]
- Oracle 10g allows the creation of rulebase based on RDFS (RDF Schema)



# Unique Identifier

- Identification of a resource uniquely
- Issues:
  - Can be dereferenced or not
  - Persistent or transient identifiers
- We use the Entrez Gene DTD as the namespace

[http://www.ncbi.nlm.nih.gov/dtd/NCBI\\_Entrezgene.dtd](http://www.ncbi.nlm.nih.gov/dtd/NCBI_Entrezgene.dtd)

- The possible candidates include:
  - LSID: Life Sciences Identifier
  - URI: NLM through UMLS and Entrez Gene

# Issues and Challenges

- We implemented one of the multiple approaches for transformation
- Identifier for biological entities is an issue of debate in the community
- Nesting structure, bi-directionality of relations and, circularity need to be solved
- Evolve the form of relationships used as predicate in the triples

# Special thanks to

- Kelly Zeng
- May Cheh
- Thomas C. Rindflesch
- Rob Logan
- Paul Lynch
- John Nyugen

# References

1. Maglott, D., Ostell, J., Pruitt, K.D., and Tatusova, T., "Entrez Gene: gene-centered information at NCBI", *Nucleic Acids Res.* 2005 January 1; 33(Database Issue): D54–D58.
2. Resource Description Framework (RDF), <http://www.w3.org/TR/2004/REC-rdf-primer-20040210/>
3. Rindflesch, TC, Fiszman, M., "The Interaction of Domain Knowledge and Linguistic Structure in Natural Language Processing: Interpreting Hypernymic Propositions in Biomedical Text", *Journal of Biomedical Informatics*. 2003;36(6):462-77.
4. XML Schema Language Transformation (XSLT), <http://www.w3.org/TR/xslt>
5. Alexander, N., Ravada S., "RDF Object Type and Reification in Oracle"— Technical White Paper ([http://download-east.oracle.com/otndocs/tech/semantic\\_web/pdf/rdf\\_reification.pdf](http://download-east.oracle.com/otndocs/tech/semantic_web/pdf/rdf_reification.pdf))
6. Online Mendelian Inheritance in Man, OMIM (TM). McKusick-Nathans Institute for Genetic Medicine, Johns Hopkins University (Baltimore, MD) and National Center for Biotechnology Information, National Library of Medicine (<http://www.ncbi.nlm.nih.gov/omim/>)
7. BioRDF subgroup: [http://esw.w3.org/topic/HCLSIG\\_BioRDF\\_Subgroup](http://esw.w3.org/topic/HCLSIG_BioRDF_Subgroup)
8. McBride, B. 2002. Jena: A Semantic Web Toolkit. *IEEE Internet Computing* 6, 6 (Nov. 2002), 55-59.
9. XPath: <http://www.w3.org/TR/xpath>
10. Life Sciences Identifier (LSID) project: <http://lsid.sourceforge.net/>