# Document Management using Protégé

Henrik Eriksson

Linköping University

# Approach: Semantic Documents

- Combine documents with knowledge representation
  - Like semantic web, but for "real" documents
- Semantic Documents
  - Printable electronic documents
  - Knowledge representation: Ontologies, workflows, and rules
  - An integrated format that keeps textual and computer-based guidelines together
  - Based on wide-spread document formats
- Currently supported format: PDF

# Adding Additional Information to the PDF Structure

- Ontologies inside PDF documents
- OWL-based metadata

Document (PDF)

Pages

XMP

Knowledge base (OWL)

Document

Root/Catalog

Pages | Outlines | Metadata | OWLMetadata

Contents

XMP

OWL

Added OWL statements

# PDFTab: Annotation Tool for Protégé



Annotation tool

Protégé

Adobe Acrobat (PDF)

# Tool Architecture

# Corresponding Ontology

# Document Mark Up

# Annotation Process

Tool selection

Text marking

**Document annotation** is a common technique to relate...

Document loading

Document saving

Annotation individual

Text annotation in document

Annotation ontology

# Document-centric Annotation Framework

# Ontology Structure

- **Linking documents and ontologies**
- **Standard ontology structure**
  - Annotation ontology
    - The annotation types
  - Document ontology
    - The key document parts
  - Domain ontology
    - The "regular" ontology

# Supporting multiple documents

- Architecture with multiple ontologies and ontology modules

# Case Study: Document Repository in Protégé

- Document data set
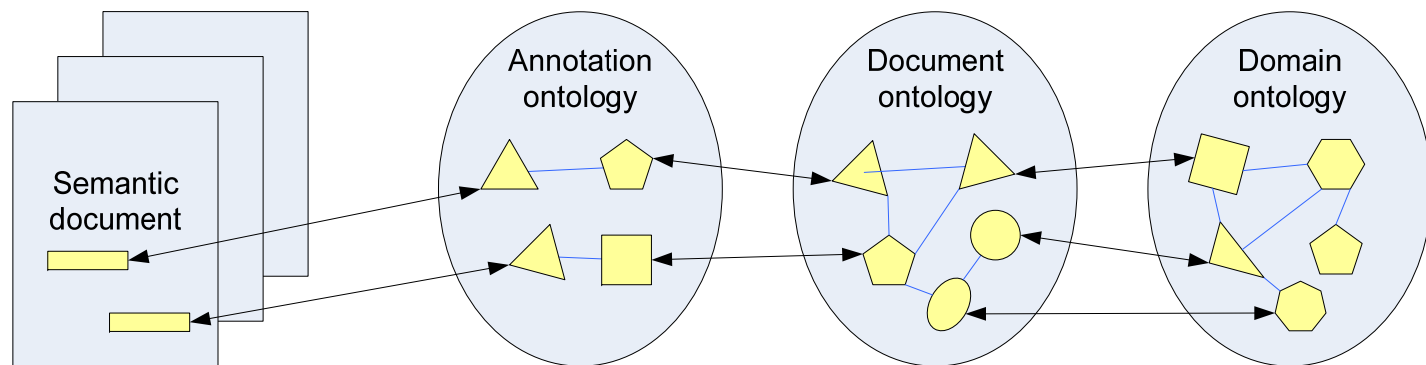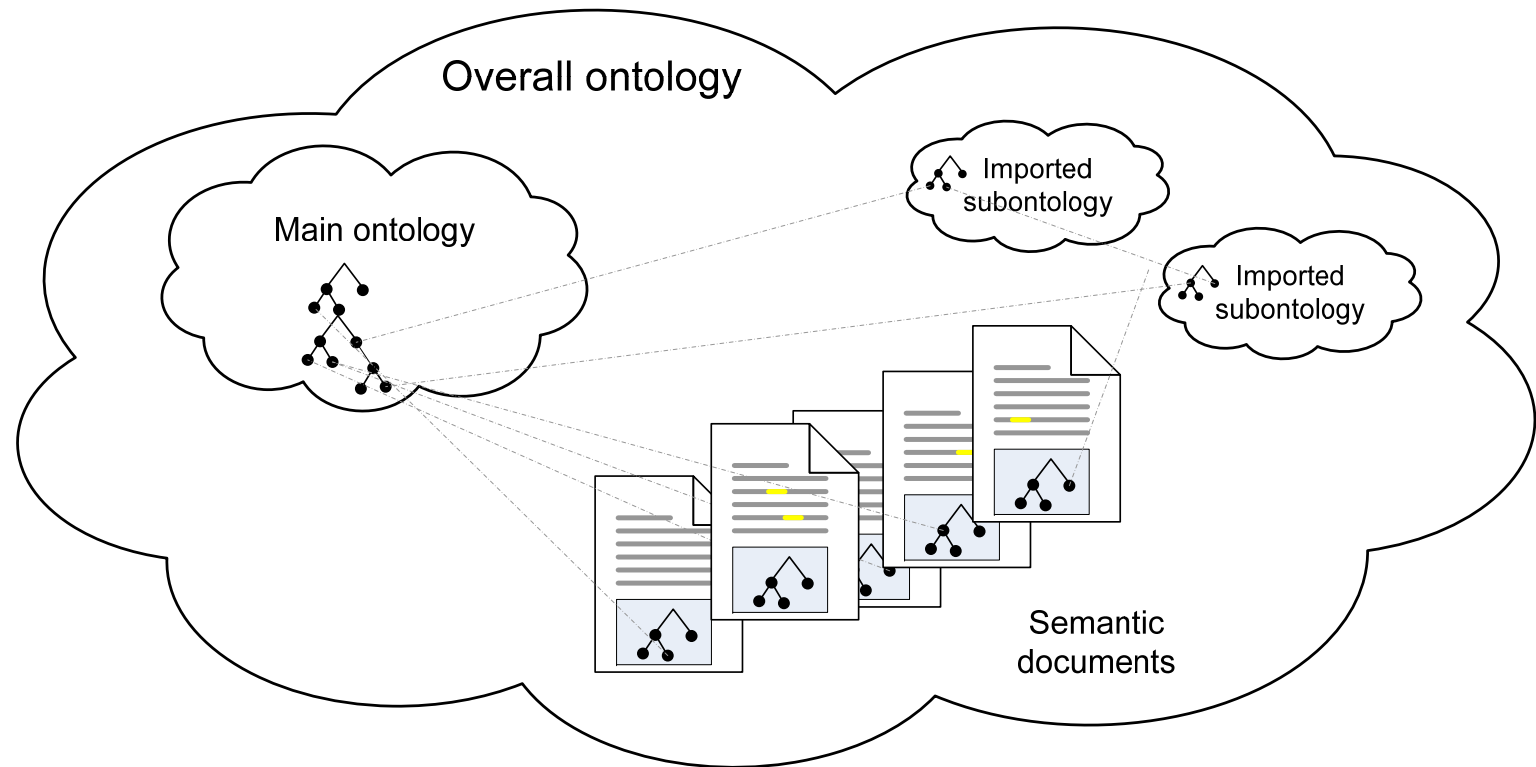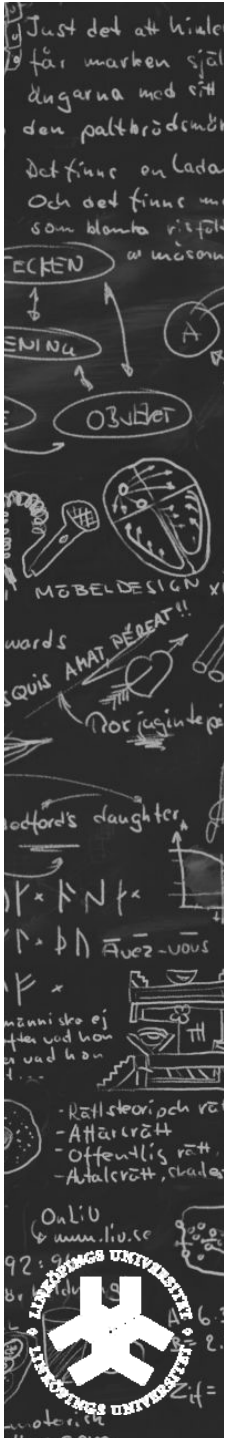  - All statistics reports (PDF) published by Statistics Sweden in 2006
  - Five volumes of Statistical Yearbook (2002–2006)
- Method
  - Document acquisition
  - Ontology development
  - Automated annotation (through annotator program)
- Number of automatically-annotated documents: 302
- Total number of annotations for these documents: 17,470

# Statistics Reports Loaded in Protégé

# Discussion

- **Scalability issues**
  - Beyond hundreds of documents
  - Too many ontologies for the current Protégé implementation
  - How can we scale to thousands or millions of documents
- **Vision: Repository storage backend**
  - Possibly backend based on a document-repository database (e.g., Dspace)
  - Normal document services and semantic services

# Summary

- Semantic Documents

- Protégé — a platform for document management

- Ontologies as model document repositories

- Furthermore, ontologies can act as document repositories

- However, large document sets will require a custom-tailored database backend