# A PDF Storage Backend for Protégé

**Henrik Eriksson**

**Linköping University**

# Storage of the Pizza example



Project and ontology files

# How do you package an ontology?

- **Gift wrapping?**



- **Document packaging**

# Persistent storage in Protégé

- **Files**
  - Serialization
  - Protégé Frames: CLIPS-like/XML
  - Protégé OWL: XML-based

- **Databases**

Voluminous

Verbose

Slow parsing & writing

Multiple file (e.g., .pprj, .owl)

There is a storage problem here

Linköpings universitet

# Background: Semantic Documents

- **Combining documents with knowledge representation**
  - Like semantic web, but for "real" documents

- **Problem: Large amounts of information is available electronically, but it is**
  - difficult to find the right information when the search query is complex, and
  - difficult to navigate content-rich information.

- **Goal**
  - Semantic description of document content (i.e., a meta-model for documents)
  - Support for systematic authoring of complex electronic documents
  - Adding support for PDF to Protégé – a PDF tab for Protégé

Linköpings universitet

# One Document—Many Applications

On-screen viewing

Printing

Semantic document

Decision support

SAGE Diabetes Guideline

Reasoning

Metadata: Guideline logic

Semantic search

Consistency check

One format for all applications

Linköpings universitet

# Semantic Documents

- **Knowledge representation**
  - **Semantic web: OWL**
  - **Ontologies**

- **Document models**
  - **Adobe's Portable Document Format (PDF)**
  - **Extensible Metadata Platform (XMP)**
  - **MS Word, RTF (?)**

- **Functions**
  - **Semantic search based on metadata**
  - **Reasoning, inference**



Statistics documents (PDF)

XMP markup

Report publication database

Document retrieval

Semantic search

Reasoning engine

Functions

Linköpings universitet

# PDFTab: Annotation tool for Protégé



Annotation tool

Protégé

Adobe Acrobat (PDF)

# Lightweight semantic documents

- **Semantic documents are nice, but**
  - sometimes too heavy
  - advanced tools required (heavy)

- **The PDF backend provides**
  - a new save method
  - a compact storage format
  - storage using standard PDF attachments
  - file access through standard PDF tools (e.g., Acrobat)

Linköpings universitet

# PDF Attachments

- **Little known feature of PDF**

- **Just like e-mail attachments**

Linköpings universitet

# The "Secrets" of the Portable Document Format (PDF)

- **Open and documented format**

- **PDF files contain something like a file system**
  - Indexing for fast random access
  - Like the .doc format of MS Word

- **Extendible file layout**
  - Custom additions

- **Different object and streams with support for text, binary data, compression, and encryption**

**Document (PDF)**

○ ○ ○ ○ **Objects**

**Streams**

**Pages**

**Metadata**

**Index (xref)**

Linköpings universitet

# Internal PDF Structure



```
                    ┌──────────────┐
                    │   Document   │
                    └──────┬───────┘
                           │
                    ┌──────┴───────┐
                    │ Root/Catalog │
                    └──────┬───────┘
        ┌──────────┬───────┴───────┬──────────┐
   ┌────┴───┐ ┌────┴───┐     ┌─────┴────┐ ┌───┴────┐
   │ Pages  │ │Outlines│     │ Metadata │ │ Names  │
   └────┬───┘ └────────┘     └─────┬────┘ └───┬────┘
        │                          │          │
   ┌────┴─────┐             ┌──────┴──┐  ┌─────┴────────┐
   │ Contents │             │   XMP   │  │Embedded files│
   └──────────┘             └─────────┘  └──────────────┘
```

# Inserting ontologies in documents

Storage backend

Linköpings universitet

# Experimental implementation

- **New knowledge base format/project type**

# Resulting PDF document

# Scenarios

- **Generated documents**

```
                                    ┌─────────────────┐
                              ┌----->│ PDF generation  │────┐
                              ¦      └─────────────────┘    │
                              ¦                             ▼
┌──────────────┐   ┌──────────────┐   ┌──────────────┐   ┌──────────────┐   ┌──────────────┐
│   Ontology   │──>│  Testing &   │──>│  Validation  │──>│   Protégé    │──>│   Document   │
│ development  │   │   revising   │   │              │   │     Save     │   │ publication  │
└──────────────┘   └──────────────┘   └──────────────┘   └──────────────┘   └──────────────┘
```

- **Authored documents**

```
      ┌───────────────────────────────┐
      │                               ▼
┌──────────────┐   ┌──────────────┐   ┌──────────────┐
│  Authoring   │──>│   Editing    │──>│ PDF conversion│────┐
└──────────────┘   └──────────────┘   └──────────────┘    │
                                                          ▼
                                                   ┌──────────────┐   ┌──────────────┐
                                                   │   Protégé    │──>│   Document   │
                                                   │     save     │   │ publication  │
                                                   └──────────────┘   └──────────────┘
                                                          ▲
┌──────────────┐   ┌──────────────┐   ┌──────────────┐   │
│   Ontology   │──>│  Testing &   │──>│  Validation  │───┘
│ development  │   │   revising   │   │              │
└──────────────┘   └──────────────┘   └──────────────┘
```
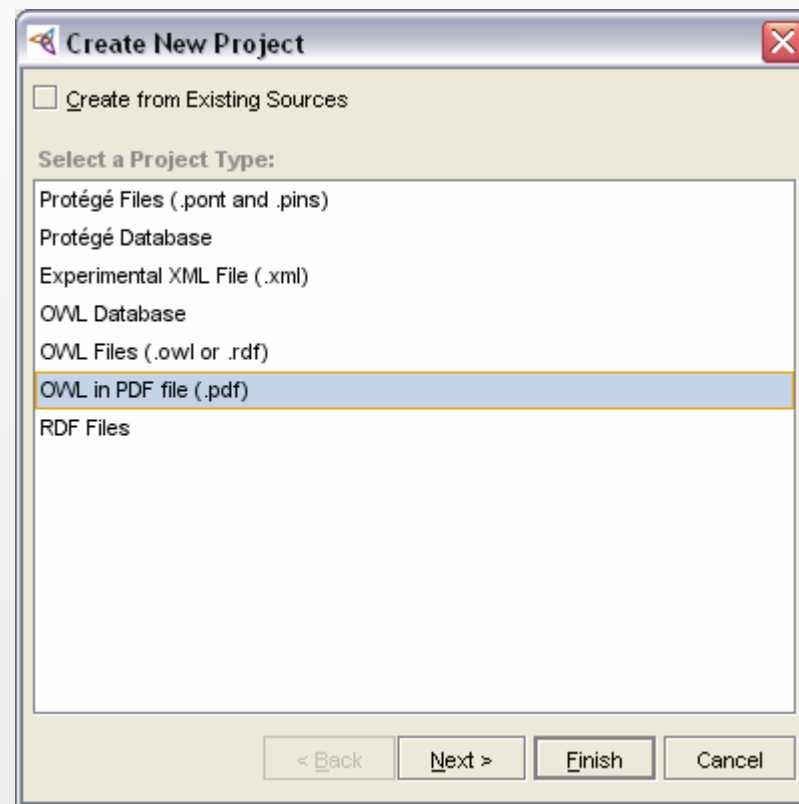
# Discussion

- **Architecture for storage (packaging) formats**
  - **Other formats possible**
  - **Examples: zip, tar, tgz, …**

- **Implementation issues**
  - **Currently "research prototype"**
  - **API changes/additions/debugging required**
    - **pdfbox, OWL plug-in, Protégé core**
  - **One PDF kb format required for each major storage type**
    - **Example: PDF-Protégé-Frames, PDF-Protégé- OWL, PDF-Protégé-RDFS**
    - **Should really be separated in a general PDF filter (more API changes required)**

Linköpings universitet

# Summary



- **Semantic documents**
    - Combine printable documents with ontologies and knowledge bases
    - Combined documentation (human-readable) and reasoning (machine-readable)
    - One document with several applications

- **PDF storage backend**
    - Lightweight semantic documents
    - Attaching ontology files to PDF documents
    - Straightforward access from Acrobat

Linköpings universitet