# Semantic Text Mining for Lignocellulose Research

Marie-Jean Meurs, Caitlin Murphy, Ingo Morgenstern,
Nona Naderi, Greg Butler, Justin Powlowski,
Adrian Tsang and René Witte

Department of Computer Science and Software Engineering
Centre for Structural and Functional Genomics
Department of Biology and Department of Chemistry and Biochemistry
Concordia University, Montréal (QC), Canada

DTMBIO 2011                                    Oct. 24th, Glasgow

# mycoMINE:
# a semantic infrastructure supporting biofuel research

## Automated curation of available knowledge on fungal enzymes

- curation of the existing literature
- development of ontological NLP pipelines
- integration through Web-based interfaces

## Goals

- spending less time to mine the literature for facts
- being provided with richer and semantically linked information

# Biofuel Process
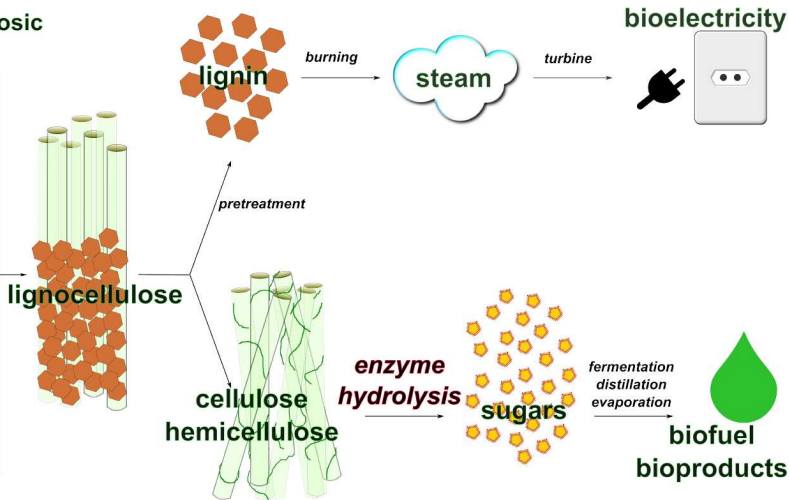
# Biofuel

## Fuels produced from biomass (lignocellulose)

- sustainable liquid fuels with low environmental impact
- promising alternative to fossil based fuels

## From cellulose to biofuel

- extraction of sugars requires to break down cellulose
- degradation of cellulose uses specific molecules called enzymes

$\Rightarrow$ discovering efficient enzymes is a key challenge

# Knowledge sources

- ever growing documentation on fungal enzymes
- many knowledge sources

## Example: the PubMed bibliographic database

- 19 million citations from over 21,000 life science journals
- linked to other databases like:
    - *Entrez Genome* provides access to genomic sequences
    - *BRENDA, The Comprehensive Enzyme Information System*, the main collection of enzyme functional data available to the scientific community

# Knowledge sources

## Querying PubMed

$\rightarrow$ collecting an often long list of relevant papers

## Analyzing the collection

$\rightarrow$ reading abstracts and sometimes full papers

## A time consuming task, difficult to handle

$\rightarrow$ significant knowledge can be missed

$\Rightarrow$ Natural Language Processing (NLP) and Semantic Web approaches
   are increasingly adopted in biomedical research.

# User Groups

Objectives ▷ Curation of fungal genes

- → guiding research and experiments
- → filling the mycoCLAP database
  - ○ DB of fungal genes encoding biochemically characterized lignocellulose-active proteins
  - ○ http://cubique.fungalgenomics.ca/mycoCLAP/

Users of our system:

- *curators*: manual curation of fungal genes
- *biology researchers*: decision about the experiments to conduct
- *experimenters*: execution of the experiments

# Semantic Entities

| Semantic Entity | Level | Definition | Example |
|---|---|---|---|
| ActivityAssayConditions | S | conditions at which the activity assay is carried out | disodium hydrogen phosphate, citric acid, pH 4.0, 37°C |
| Assay | W | name of the experimental assay | Dinitrosalicylic Acid Method (Somogyi-Nelson) |
| Enzyme | W | enzyme name | alpha-galactosidase |
| Gene | W | gene name | mel36F |
| Glycosylation | S | enzymatic process attaching glycans to organic molecules | N-glycosylation |
| Host | W | organism used to produce the recombinant protein | Escherichia coli |
| KineticAssayConditions | S | buffer, pH, temp. for the kinetic parameters determination | 0.1M (disodium hydrogen phosphate, citric acid), pH 4.0, 37°C |
| Organism | W | organism name | Gibberella sp. |
| pH | S | pH mentions | The enzyme retained greater than 90% of its original activity between pH 2.0 and 7.0 at room temperature for 3h. |
| ProductAnalysis | S | products formed from the enzyme reaction and identification method | HPLC, glucose, galactose |
| SpecificActivity | S | specific activity of the enzyme on the substrate | 11.9U/mg |
| Strain | W | strain name | F75 |
| Substrate | W | substrate name | stachyose |
| SubstrateSpecificity | S | substrate specificity mentions | The Endogluccanase from Pyrococcus furiosus had highest activity on cellopentaose. |
| Temperature | S | temperature mentions | The enzyme stability at different pH values was measured by the residual activity after the enzyme was incubated at 25°C for 3h. |

Semantic annotation types defined by the curators, applicable level (sentence, $S$ or word(s), $W$), definitions and examples

# Semantic Resources
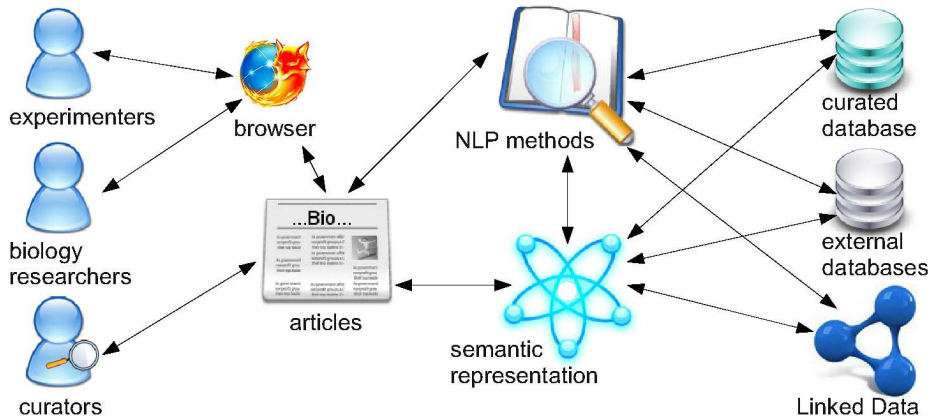
– organisms –

- the NCBI Taxonomy database
    - o http://www.ncbi.nlm.nih.gov/Taxonomy

– enzymes –

- BRENDA
    - o http://www.brenda-enzymes.org
- the UniProtKB/SwissProt database
    - o http://www.uniprot.org

# Semantic Resources

## References to the original sources integrated into the curated data:



The N-terminal amino acid sequence of also identified. The aspergillopepsin I s expressed enzyme had extra Gly-Ser of aspergillopepsin I signal-encoding seq nucleotide sequence [ 19 ]. The numbe sequence described previously [ 19 ].

The molar absorption coefficient ( e 28

Contents of the secondary structure ( a respectively.

| Enzyme | | |
|---|---|---|
| BRENDA's page0 | ▼ | http://www.brenda-enzymes.org/php/result_flat.php4?ecno=3.4.23.18 |
| BRENDA_ECNumber | ▼ | 3.4.23.18 |
| BRENDA_RecName | ▼ | Aspergillopepsin I |
| SwissProtID | ▼ | P41748 Q12567 |
| alias | ▼ | aspergillopepsin |

$\rightarrow$ facilitates semantic connections through Linked Data techniques
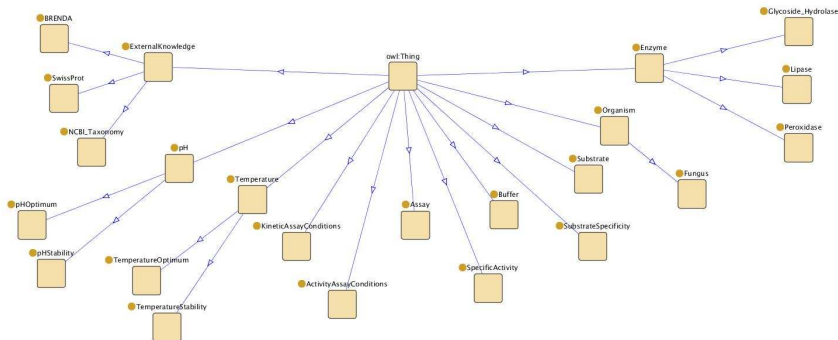
# System Architecture

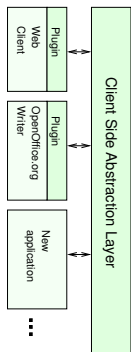# mycoMINE Ontology

# Semantic Assistants

NLP services are provided by the Semantic Assistants architecture.

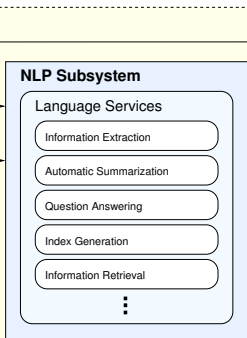# Text Mining Pipelines based on GATE

The *General Architecture for Text Engineering* (GATE):

- mature framework, more than 10 years of development
- development team at University of Sheffield, UK (gate.ac.uk)

Preprocessing steps:

- tokenization
- sentence splitting
- part-of-speech tagging

Custom pipelines:

- extract the semantic entities
- populate the OWL ontology using the OwlExporter component

# Ex1. The GATE pipeline for Organism Recognition

## Modules for organism entity detection:

- based on pattern matching to the NCBI reference taxonomy
- provide scientific names and NCBI Taxonomy Identifiers

## Modules for extraction of strain mentions:

- specific text tokenization $+$ machine learning (svm) based approach

## Resources

- external resources automatically translated for reuse in the system
- $\rightarrow$ ability to update the installation when the Taxonomy DB changes

# Ex2. The GATE pipeline for Enzyme Recognition

1. Tokens with the *-ase* enzyme suffix:

   - enzyme-specific text tokenization
   - grammar rules [JAPE language]

2. Modules for enzyme entity detection:

   - rely on automatically extracted knowledge from BRENDA
   - based on pattern matching
   - provide enzyme *EC number*, *Recommended Name*, *Systematic Name*
     and *URL* on the BRENDA website.

# Pipeline in GATE Developer

# System output — Literature Curation: results displayed in Firefox

**Semantic Assistants Sidebar**

🔀 Purification and properties of an extra...

- 🔵 ▽ Enzyme
  - ▷ beta-xylosidase (4)
  - hydrolase
  - XlnD
- 🟢 ▽ Substrate
  - oat spelt xylan
- 🔴 ▷ p H
- 🟠 ▽ Gene
  - ▷ xylA (3)
- 🟣 ▽ Organism
  - ▷ Aspergillus japonicus (2)
  - ▷ Aspergillus niger (1)
  - ▷ Pichia pastoris (1)
  - ▷ Saccharomyces cerevisiae...

## Purification and properties of an extracellular beta-xylosidase from Aspergillus japonicus and sequence analysis of the encoding gene.

Wakiyama M, Yoshihara K, Hayashi S, Ohta K.

Department of Applied Chemistry, Faculty of Engineering, University of Miyazaki, 1-1 Gakuen Kibanadai Nishi, Miyazaki 889-2192, Japan.

### Abstract

An extracellular protein exhibiting beta-xylosidase activity was purified from the culture filtrate of a filamentous fungus, Aspergillus japonicus strain MU-2, grown on oat spelt xylan. The purified enzyme was a monomeric glycoprotein with an apparent M(r) of 113.2 kDa as estimated by SDS-PAGE. beta-Xylosidase activity was optimal at pH 4.0 and 70 degrees C. The enzyme also showed beta-glucosidase and alpha-l-arabinofuranosidase activities. The genomic DNA and cDNA encoding this protein were cloned and sequenced. Southern blot analysis indicated that the beta-xylosidase gene (xylA) was present as a single copy in the genome. An open reading frame, consisting of 2412 bp, was not interrupted by introns, and it encoded a presumed signal peptide of 17 amino acids and a mature protein of 787 amino acids. The deduced amino acid sequence of the xylA gene product showed a high degree of identity (69%) to the primary structure of the Aspergillus niger beta-xylosidase XlnD that belongs to the glycoside hydrolase family 3. Moreover, the xylA gene was functionally expressed in the yeast Pichia pastoris.

PMID: **19000618** [PubMed - indexed for MEDLINE]    **Free full text**

## System output — Manual Annotation: pre-annotation in Teamware

# The Manual Annotation Process (1)

The annotation team:

- four biology researchers
- the researcher in charge of the curation task
- inter-annotator agreement over 80%

Table: Inter-Annotator Agreement and characteristics of the annotation tasks

| Task | T1 | T2 | T3 | T4 | T5 | T6 | T7 | Corpus |
|---|---|---|---|---|---|---|---|---|
| #papers | 1 | 3 | 4 | 4 | 4 | 5 | 5 | 26 |
| #entities | 321 | 1158 | 1673 | 2072 | 2248 | 2288 | 2133 | 11893 |
| IAA (%) | 80 | 81 | 81 | 82 | 87 | 84 | 88 | 85 |

# The Manual Annotation Process (2)

The gold standard corpus:

- 26 full text articles [21 freely accessible]
- manually annotated using GATE Teamware
- enzyme categories:
    - o Glycoside Hydrolase = 69%
    - o Lipase = 12%
    - o Peroxidase = 19%
- papers published between 1996 and 2011
- 23 articles available on PubMed http://www.ncbi.nlm.nih.gov/pubmed/
- 3 articles available on ScienceDirect http://www.sciencedirect.com/
- adjudication task achieved on 11 papers

Introduction     Project Context and System Architecture     Text Mining Pipelines     **Evaluation**     Conclusion

○    ○○○    ○○●
○○    ○    ○○    ○
○○    ○○○    ○○

# The Manual Annotation Process (3)

### The GATE Teamware:

# Evaluation

- Correctness evaluated in terms of precision, recall and F-measure
- Reference = 11 adjudicated papers

## Results on the four most common entities:

|  | Strict (overlaps discarded) | | | Lenient (overlaps included) | | |
|---|---|---|---|---|---|---|
|  | **Recall** | **Precision** | **F-m** | **Recall** | **Precision** | **F-m** |
| **Enzyme** | 0.79 | 0.64 | 0.71 | 0.91 | 0.75 | 0.82 |
| **Organism** | 0.87 | 0.86 | 0.87 | 0.91 | 0.91 | 0.91 |
| **pH** | 0.79 | 0.81 | 0.80 | 0.96 | 0.99 | 0.98 |
| **Temperature** | 0.70 | 0.66 | 0.68 | 0.93 | 0.88 | 0.91 |

# Conclusions

### Contributions:

- mycoMINE
- text mining pipelines combined with ontological resources

### Results:

- state-of-the-art results
- available gold standard corpus and system

$\rightarrow$ future work:
  - o user-system interaction for data validation
  - o quality assessment of the curated data
  - o impact on the Genozymes research workflow