

IntelliGenWiki: An Intelligent Semantic Wiki for Life Sciences

Bahar Sateli¹, Marie-Jean Meurs^{1,2}, Greg Butler^{1,2}, Justin Powlowski^{2,3}, Adrian Tsang^{2,4}, René Witte¹✉

¹Department of Computer Science and Software Engineering; Concordia University, Montréal, Canada

²Centre for Structural and Functional Genomics; Concordia University, Montréal, Canada

³Department of Chemistry and Biochemistry; Concordia University, Montréal, Canada

⁴Department of Biology; Concordia University, Montréal, Canada

Motivation and Objectives

The rapid growth of the scholarly literature makes the management and curation of the available information a labor-intensive and time-consuming task for researchers, during which significant knowledge can be easily missed. To address this problem, efforts have been made to use Natural Language Processing (NLP) techniques as a means to (semi-)automatically improve the exhaustive analysis of the available information. In order to make these NLP techniques more end-user friendly and integrate them with knowledge management workflows, we developed IntelliGenWiki, a novel combination of a wiki system with state-of-the-art techniques from the NLP and Semantic Computing domains. Wikis are well known as an easy-to-use, collaborative platform for creating and organizing knowledge. For example, the Gene Wiki project (Huss III et al, 2010) applies community intelligence to the annotation of gene and protein functions. However, existing approaches rely on a manual analysis of the literature. With IntelliGenWiki, we aim to leverage the collaborative nature of wikis by introducing new Human-AI collaboration patterns: Our goal is to provide text mining assistants that work together with humans on literature analysis tasks, like curation or the generation of semantic metadata, which can be used in an Linked Open Data context. IntelliGenWiki is based on an open service-oriented architecture: it can be applied to different projects by deploying custom NLP analysis pipelines suitable for the specific task and domain. Here, we demonstrate the benefits of this approach within a collaborative literature curation context.

Methods

We first describe the general workflow for working with NLP assistants, followed by a description of the underlying architecture.

Workflow. IntelliGenWiki provides a standard wiki user interface. From any wiki page (Fig. 1, top), users can ask for “Semantic Assistants” from the menu (Fig. 1, left), which will result in a dynamically injected user interface from which assistants can be selected and executed (Fig. 1, bottom). The user can now select an appropriate assistant from the list and invoke it on one or multiple pages of the wiki, gathered in a so-called “collection”. This will invoke the selected NLP pipeline on the set of wiki pages. The results (e.g., detected entities) are stored in the user’s place of choice and made persistent in the wiki repository (Fig. 1, middle). Thereby, all updated pages become immediately available to all wiki users for collaborative adjustment, modification and further refinement of the results.

Architecture. Technically, IntelliGenWiki combines NLP analysis pipelines developed in the General Architecture for Text Engineering (GATE) (Cunningham et al, 2011) with MediaWiki, <http://www.mediawiki.org> (Last accessed: 26.09.2012), a widely-used wiki engine. These pipelines are published as standard web services through the Semantic Assistants framework (Witte and Gitzinger, 2008). The Wiki-NLP integration is based on a service-oriented architecture that seamlessly introduces these NLP web services into wiki systems (Sateli and Witte, 2012). This allows wiki users to benefit from text mining techniques directly within their wiki environment, without the need for switching to an external application. Additionally, we support the generation of semantic metadata from NLP analysis results. This metadata is formally represented in the wiki through the Semantic MediaWiki (SMW) extension: <http://semantic-mediawiki.org/> (Last accessed on Sept 26, 2012). This formal representation of the available wiki knowledge can be exploited by exporting it in form of RDF triples. It can also be queried directly within the wiki using SMW inline queries. For example, users could write queries to retrieve



Sysop [my talk](#) [my preferences](#) [my watchlist](#) [my contributions](#) [log out](#)

page
discussion
edit
history
delete
move
protect
watch
refresh

PMID: 20709852

[Contents \[show\]](#)

Characterization of a Cellobiohydrolase (MoCel6A) Produced by Magnaporthe oryzae [\[edit\]](#)

PMID: 20709852

Authors: Machiko Takahashi,1‡ Hideyuki Takahashi,1‡ Yuki Nakano,1 Teruko Konishi,2 Ryohei Terauchi,1 and Takumi Takeda1*

Iwate Biotechnology Research Center, Kitakami, Iwate 024-0003, Japan,1 University of the Ryukyus, Department of Bioscience and Biotechnology, Faculty of Agriculture, 1 Senbaru Nishihara, Okinawa 903-0213, Japan2

*Corresponding author. Mailing address: Iwate Biotechnology Research Center, 22-174-4, Narita, Kitakami, Iwate 024-0003, Japan. Phone: 81 (197) 68-2911. Fax: 81 (197) 68-3811. E-mail: ttakeda@ibrc.or.jp

‡M.T. and H.T. contributed equally to this work.

Received March 10, 2010; Accepted July 30, 2010.

Full Text [\[edit\]](#)

Abstract [\[edit\]](#)

Three GH-6 family cellobiohydrolases are expected in the genome of *Magnaporthe grisea* based on the complete genome sequence. Here, we demonstrate the properties, kinetics, and substrate specificities of a *Magnaporthe oryzae* GH-6 family cellobiohydrolase (MoCel6A). In addition, the effect of cellobiose on MoCel6A activity was also investigated. MoCel6A contiguously fused to a histidine tag was overexpressed in *M. oryzae* and purified by affinity chromatography. MoCel6A showed higher hydrolytic activities on phosphoric acid-swollen cellulose (PSC), β -glucan, and cellooligosaccharide derivatives than on cellulose, of which the best

These results suggest that enhancement or inhibition of hydrolytic activities by cellobiose is dependent on the reaction mixture pH.

PMID: 20709852 [\[PubMed - indexed for MEDLINE\]](#) PMID: PMC2950481 [Free PMC Article](#) [mycoMINE on PMID:20709852_Abstract \(View\)](#) [\[View\]](#)

Content	Type	Start	End	Features
cellobiohydrolase	Enzyme	89	106	<ul style="list-style-type: none"> ■ enzyme_alias: cellobiohydrolase ■ BRENDA_SystematicName: 4-beta-D-glucan cellobiohydrolase ■ BRENDA_EcNumber: 3.2.1.91 ■ abbreviation_alias: - ■ google_search: http://www.google.com/search?q=cellobiohydrolase ■ BRENDA_RecommendedName: cellulose 1,4-beta-cellobiosidase ■ SwissProt_ID: O68438 ■ BRENDA's page: http://www.brenda-enzymes.org/php/result_flat.php4?ecno=3.2.1.91

This page was last modified on 28 March 2012, at 19:15. This page has been accessed 24 times. [Privacy policy](#) [About G-nWiki](#) [Disclaimers](#) 

Available Assistants
Results Target
Global Settings
Console

Step 1. Select the service your wish to execute on your collection.
Once you add this page to your collection, you can continue browsing as your collection is saved.

Available Assistants Select a service

Runtime Parameters Select a service

mycoMINE
 IR Information Extractor
 Information Extractor
 OrganismTagger

Collection

Add
Clear

Fig. 1: The wiki interface with integrated text analysis services (bottom), showing automatically added, NLP-extracted entities (middle), together with original content (top)

literature that contains a certain type of entities, such as enzymes or organisms.

Results and Discussion

To test the effectiveness of NLP assistants in a wiki environment, we deployed an IntelliGenWiki installation within the Genozymes project: <http://www.fungalgenomics.ca> (Last accessed on Sept 20, 2012). The task we aimed to support in the project is biomedical literature curation for lignocellulose research. For this experiment, we deployed the mycoMINE NLP pipeline (Meurs et al, 2012), which automatically extracts knowledge from the literature on fungal enzymes by using semantic text mining approaches combined with ontological resources. We manually pre-filled the wiki with a corpus of 30 documents composed of PubMed abstracts and their corresponding full-text papers, selected by two expert biocurators. These biocurators provided us with their average time needed for curation without support on the same task. They performed the corpus curation through the wiki using mycoMINE to automatically extract relevant entities, and they kept track of the time spent on each document. The time for abstract selection (triage task) decreased from 1min. (without support) to 20sec. (using IntelliGenWiki), and from 37.5min (without support) to 30.6min (using IntelliGenWiki) for full paper selection (curation task), showing a productivity enhancement of 67% and 20%, re-

spectively. The results gathered from this experiment confirm the usability and the effectiveness of our approach.

The IntelliGenWiki system, including the NLP integration back-end, is available as open source software from <http://www.semanticsoftware.info/intelligenwiki>.

Acknowledgements

Funding for this work was provided by NSERC, Genome Canada and Génome Québec. Caitlin Murphy and Sherry Wu are acknowledged for their participation in the evaluation task.

References

1. Cunningham H, Maynard D, et al (2011) Text Processing with GATE (Version 6), University of Sheffield, Department of Computer Science
2. Huss III J. W., et al (2010) The Gene Wiki: Community Intelligence Applied to Human Gene Annotation, *Nucleic Acids Research* 38, p. 633–639. doi:10.1093/nar/gkp760
3. Meurs MJ, Murphy C, et al (2012) Semantic Text Mining Support for Lignocellulose Research, *BMC Medical Informatics and Decision Making* 12(Suppl 1):S5. doi:10.1186/1472-6947-12-S1-S5
4. Sateli B and Witte R (2012) Natural Language Processing for MediaWiki – The Semantic Assistants Approach, In 8th International Symposium on Wikis and Open Collaboration (WikiSym 2012). Linz, Austria.
5. Witte R and Gitzinger T (2008) Semantic Assistants – User-Centric Natural Language Processing Services for Desktop Clients, In Asian Semantic Web Conference (ASWC 2008), Springer LNCS 5367, pp.360–374. doi:10.1007/978-3-540-89704-0_25