



openMIN7ED



CALL FOR PARTICIPATION: BioCreative V.5. BeCalm (Biomedical annotation metaserver) task (<http://www.becalm.eu>)

Background

There is an increasing demand in being able to effectively access, evaluate, compare, visualize and integrate multiple text mining systems in order to process natural language document collections. Several BioCreative tasks tried to promote the development of online text annotation servers. In particular, the BioCreative Meta-Server (BCMS) was the first distributed prototype platform requesting, retrieving and unifying biomedical textual annotations.

Despite the relevance of those previous efforts, several crucial aspects have not been sufficiently or only partially addressed (including continuous evaluation, extraction of textual content from heterogeneous sources, harmonization of multiple biomedical text annotations and visualization and comparative assessment of automatic and manual annotations).

This inspired the BeCalm task.

BioCreative V.5. BeCalm tasks

We are now inviting text mining and language processing researchers world-wide to implement text Annotation Servers and to provide their results for the following tasks:

- **CEMP (Chemical Entity Mention recognition).** This task requires the recognition of chemical named entity mentions in text. This implies to return the start and end indices corresponding to all the chemical entity mentions. The training set for this task will consist of 21,000 patent abstracts and the test set will consist of 9,000 patent abstracts. The same evaluation strategy as for the BioCreative CHEMDNER task will be used (i.e., F-score).

- **GPRO (Gene and Protein Related Object recognition).** For this task teams have to recognize mentions of gene and protein related objects (named as GPROs) in running text. The training set for this task will consist of 21,000 patent abstracts and the test set will consist of 9,000 patent abstracts. The same evaluation strategy as for the BioCreative CHEMDNER task will be used (i.e., F-score).

- **TIPS (Technical interoperability and performance of annotation servers).** This novel BioCreative task will focus the first time specifically on the technical aspects of making available and evaluating text Annotation Servers (ASs) for continuous named entity recognition. This is an open task, in the sense that it is not restricted to a particular named entity recognition type/tool. Moreover, for this task we allow that the participant Annotation Servers may be fully developed in-house or integrate/adapt third party recognition software as building block components (see list of initial existing candidate bio-NER tools [here](#)).

For this task we will reinforce a minimal set of functional specifications (metadata info) and the use of a common communication protocol for serializing and distributing text annotations. It is in line with the efforts of ELIXIR /EXCELERATE in benchmarking the ELIXIR catalog of methods.

Specifically, three levels of evaluation will be considered for this task:

- **technical** level: *stability* (i.e. ability to respond to continuous requests, to respond within stipulated time window, and to provide updated server status information), *response time* (i.e. time taken to respond to a request measured in terms of the number and contents of the

requested documents and the volume of predictions returned) and *batch processing* (i.e. ability to respond to requests with a varied number of documents) capabilities will be monitored.

- **data** level: ability to return NER annotation results as structured data (represented in one or several of the following formats XML/BioC, JSON/BioCJSON or TXT/TSV); ability to retrieve and process documents from different providers (e.g. patents, PubMed abstracts and PMC full-texts).
- **functional specification** level: metadata requirements, following functional specifications inspired by the OpenMinTeD interoperability project. *Mandatory* metadata include server name, institution/company, server administrator, programming language (main language, if using several), integration of third-party recognition software, recognised annotation types (e.g. chemical entities, genes, proteins, diseases, organisms, cellular lines and types, and mutations) supported annotation formats (e.g. XML/BioC, JSON/BioCJSON or TXT/TSV) and version control. *Important* metadata include software license, specification of third-party recognition software (if any), dedicated vs shared server, and relevant publications. Finally, *interesting* metadata include operating system, distributed processing, hardware characteristics (i.e. number of processors and RAM information).

Annotation Servers should implement a REST (Representational State Transfer) API application that listens and responds to the metaserver requests (see <http://becalm.eu/api>). To support the implementation of Annotation Servers, the organizers will host a team registration page and a system for managing user accounts and mailing list/newsletter.

In order to enable a more robust evaluation setting that facilitates direct comparison and visualization of different online and offline predictions the BeCalm (Biomedical Annotation Metaserver) platform will be used (<http://www.becalm.eu>). Participating teams do not need to send results for all of three sub-tasks. They can send results only for individual sub-tasks.

The CHEMDNER-patents task (BioCreative V - <http://www.biocreative.org>) is a community challenge on named entity recognition of chemical compounds and genes/proteins in patents (<http://www.biocreative.org/tasks/biocreative-v/track-2-chemdner/>).

The teams that participated in previous biomedical text mining tasks posed at the BioCreative competition heavily used machine learning approaches, like CRFs and SVMs. The evaluation of these systems in online annotation environments/settings, where systems have to manage periodical and varied annotation requests, is of critical importance to promote widespread real-world use of these systems. Also, it will enable the continuous evaluation of these systems, which may highlight new user needs and technology updates.

Registration and participation

Teams interested in the Biocreative V.5 BeCalm task should register at the 'Participant sign in' at <http://www.becalm.eu/>

Tentative timeline/phases

- *Registration* phase: announcement to the community, OpenMinted partners, previous BioCreative task participants and call for participation. During this phase we will provide additional infrastructure and documentation details for additional annotation server participants.

-*Offline participation* phase: release of a dataset to registered teams for which they can upload offline predictions to the BeCalm platform.

-*Online participation* phase: annotation requests will be sent periodically to the registered Annotation Servers; responses should be returned within a pre-specified time window; server status will be continuously monitored (e.g. operating, overload, or shutdown).

- *Evaluation* phase: Overall analysis of the generated results and survey responses according to the different task evaluations. Presentation of the obtained results at the evaluation workshop together with technical details on the annotation server implementations by a selected number of participating teams.

Biocreative V.5 workshop proceedings and journal special issue

In line with previous BioCreative challenges, the result of participating teams will be presented at the BioCreative V.5. evaluation workshop that will take place in Madrid (Spain) in March 2017. It will include an overview talk presenting the datasets used and results obtained by the participating teams, and a number of teams will be invited to present their systems. We plan to have also a session where teams, task organizers and domain experts will discuss the obtained results and future steps. Finally, during the poster session, all teams will be able to present their participating strategies.

Participating teams will be invited to contribute to the: Proceedings of the BioCreative V.5. Challenge Evaluation Workshop. A selected number of top performing teams will also be invited to contribute with a system description paper to a special issue of a relevant journal in the field. All successfully participating teams will be invited to contribute to a journal paper describing the BioCreative V.5 challenge.

Task Organizers

- Martin Krallinger, Spanish National Cancer Research Centre
- Anália Lourenço, University of Vigo
- Martin Pérez-Pérez, University of Vigo
- Gael Pérez-Rodríguez, University of Vigo
- Florentino Fdez-Riverola, University of Vigo
- Alfonso Valencia, Spanish National Cancer Research Centre

Advisory Board

- Sophia Ananiadou, National Centre for Text Mining (NaCTeM)
- Cecilia Arighi, University of Delaware
- Donald Comeau, National Center for Biotechnology Information (NCBI), NIH,
- Rezarta Islamaj Doğan, National Center for Biotechnology Information (NCBI), NIH
- Lynette Hirschman, MITRE Corporation
- Zhiyong Lu, National Center for Biotechnology Information (NCBI), NIH
- Johanna McEntyre, EMBL - EBI
- Fabio Rinaldi, Institute of Computational Linguistics, University of Zurich
- Angus Roberts, University of Sheffield

Relevant references

- Martin Pérez-Pérez, Gael Pérez-Rodríguez, Obdulia Rabal, Miguel Vazquez, Julen Oyarzabal, Florentino Fdez-Riverola, Alfonso Valencia, Martin Krallinger, Anália Lourenço (2016) The Markyt visualisation, prediction and benchmark platform for chemical and gene entity recognition at BioCreative/CHEMDNER challenge. *Database*, baw120
- Florian Leitner, et al. (2008) Introducing meta-services for biomedical information extraction. *Genome biology*. 9(2), 1.
- Comeau, D. C., Doğan, R. I., et al. (2013). BioC: a minimalist approach to interoperability for biomedical text processing. *Database*, bat064.
- Arighi, C. N., Lu, Z., et al. (2011). Overview of the BioCreative III workshop. *BMC bioinformatics*, 12(8), 1.
- Katayama, T., Arakawa, K., et al. (2010). The DBCLS BioHackathon: standardization and interoperability for bioinformatics web services and workflows. *Journal of biomedical semantics*, 1(1), 1.
- Belleau, F., Nolin, M. A., Tourigny, N., Rigault, P., & Morissette, J. (2008). Bio2RDF: towards a mashup to build bioinformatics knowledge systems. *Journal of biomedical informatics*, 41(5), 706-716.