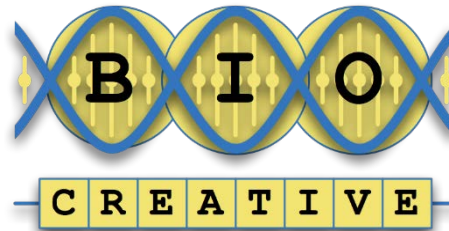# BioCreative Challenges
## Information Retrieval and Text Mining for Biology
## SIB

June 4, 2015.

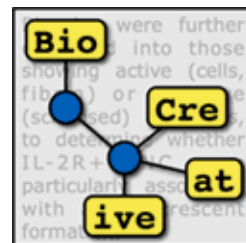**Cecilia N. Arighi, PhD**

Research Associate Professor

Protein Information Resource

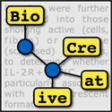CBCB, University of Delaware

arighi@dbi.udel.edu

http://www.biocreative.org

- Overview of BioCreative Effort

- Evolution of Tasks

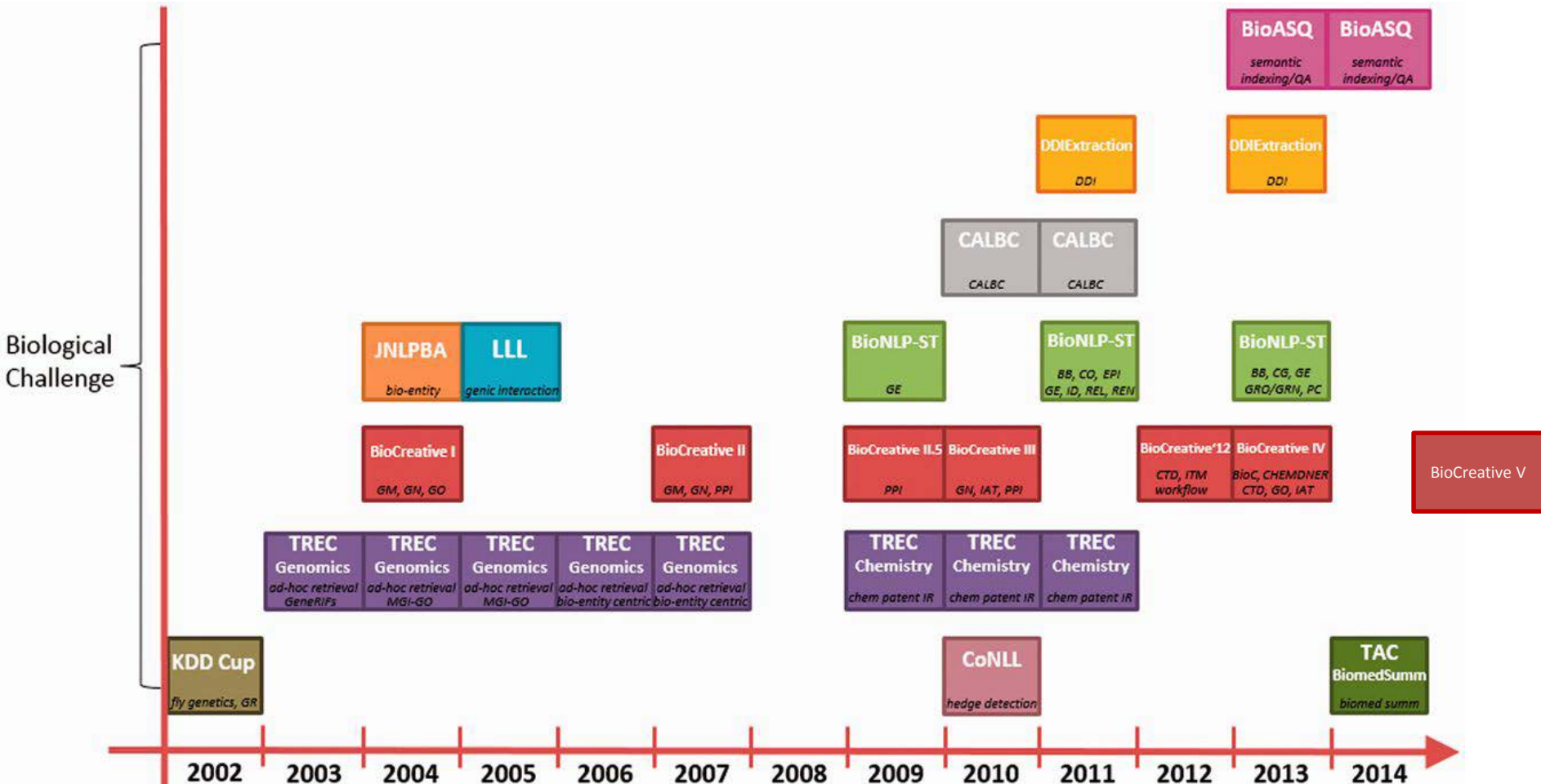- Impact of BioCreative

- BioCreative User Interactive Task

- Concluding Remarks

Community-wide effort for evaluating text mining systems applied to the biomedical domain
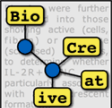
Collaborative and interdisciplinary effort

**Adapted from Chung-Chi Huang, and Zhiyong Lu Brief
Bioinform 2015;bib.bbv024**

Briefings in
**Bioinformatics**

# What is special about BioCreative?

Intends to:

- o Attract researchers from both natural language processing and biomedical domain

- o Address problems of importance to the biology and bioinformatics community (focus on biocuration)

- o Create legacy training and test data suites that could be used for development and benchmarking of future applications

- o Allow the assessment of the state-of-the-art on real biological tasks

# • Design of tracks based on user needs

**Generic biocuration workflow**

**GO curation:** identification of articles with curatable GO information

**PPI curation:** selection of relevant PPI articles from PubMed

**CTD:** identification of gene, chemical, disease and link to CTD vocabulary

**PPI curation:** Extracting interactant pair, PPI method

**ChemDNER:** Identification of chemical entities

**Source**

**Triage**

**Full curation**

RELEVANT ARTICLE

**Entity Detection**

human

AKT1

Akt1

Taxonomy

SPECIES/ TAXON

UniProt

PROTEIN DBs

Entrez Gene

GENE DBs

Experimental evidence

AKT1

BCL10

**Relation/ Evidence**

anti bait co-IP

interacts_with

*Courtesy of Lu and Hirschman*

# Gene normalization (linking a gene mention to database identifier)

BCI Abstracts

BCII Abstracts

BCIII Full-length

Species are known

Closer to real scenario

# Chemical recognition
(identifying compound names in text)

BCIV Abstracts

BCV Patents

Noisy data

Some problems faced by shared-task challenges:
  o Many different formats
  o Many new projects start over
  o An atmosphere of competition

## The Needs?
  o Common format
  o Simple-to-learn software to access the format
  o Sufficient resources to motivate users

## A Solution
  o A convenient format to share text documents and annotations
  o A library to promote interoperability of data and tools

# • **Fosters interoperability**

BioC, a BioCreative interoperability initiative, is a simple extensible XML language format to share text data and annotations

Goals:
o  simplicity
o  interoperability
o  broad use and reuse

| BioC Implementations | BioC Tools | BioC Corpora |
|---|---|---|
| C++<br>Java<br>SWIG-Python<br>SWIG-Perl<br>Python<br>Ruby<br>Go | Natural Language Tools:<br>  sentence segmenting<br>  tokenizing<br>  part-of-speech tagging<br>  lemmatization<br>  dependency parsing<br><br>NER<br>  diseases<br>  mutations<br>  chemicals<br>  species<br>  genes/proteins<br><br>Manual annotation<br>Sentence simplification | PMC-BioC<br>Disease NER<br>BioNLP Shared task<br>Abbrev. definition<br>WBI repository<br>SRL<br>iSimp<br>Metabolites |

Syntax is XML defined by a Document Type Definition (DTD)
Key file describes content of XML file

http://bioc.sourceforge.net/

PMID: 22187158

Tat mostly activated the MIP-1alpha expression in a p65-dependent manner.

Gene Name

```
<annotation id ="G0">
  <infon key="type">Gene_name</infon>
  <location offset="0" length="3" />

  <text> Tat</text>
</annotation>

<annotation id ="G1">
  <infon key="type">Gene_name</infon>
  <location offset="25" length="10" />

  <text>MIP-1alpha</text>
</annotation>

<annotation id ="G2">
  <infon key="type">Gene_name</infon>
  <location offset="52" length="3" />

  <text>p65</text>
</annotation>
```
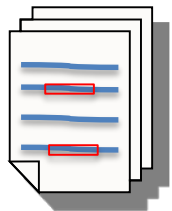
http://www.ncbi.nlm.nih.gov/CBBresearch/Dogan/BioC_ISMB2014.pd.

BioNLP evaluations have focused on isolated tasks, they have emphasized 'off-line' accuracy measures
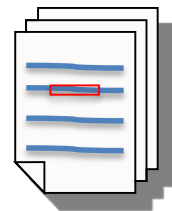
But in the biocuration world….

Documents automatically annotated, or retrieved by TM

Reviewed by curator

Curated document

Information stored in DB



database

**Interactive Task (IAT):** Evaluation of text mining systems by potential users and report on performance and usability
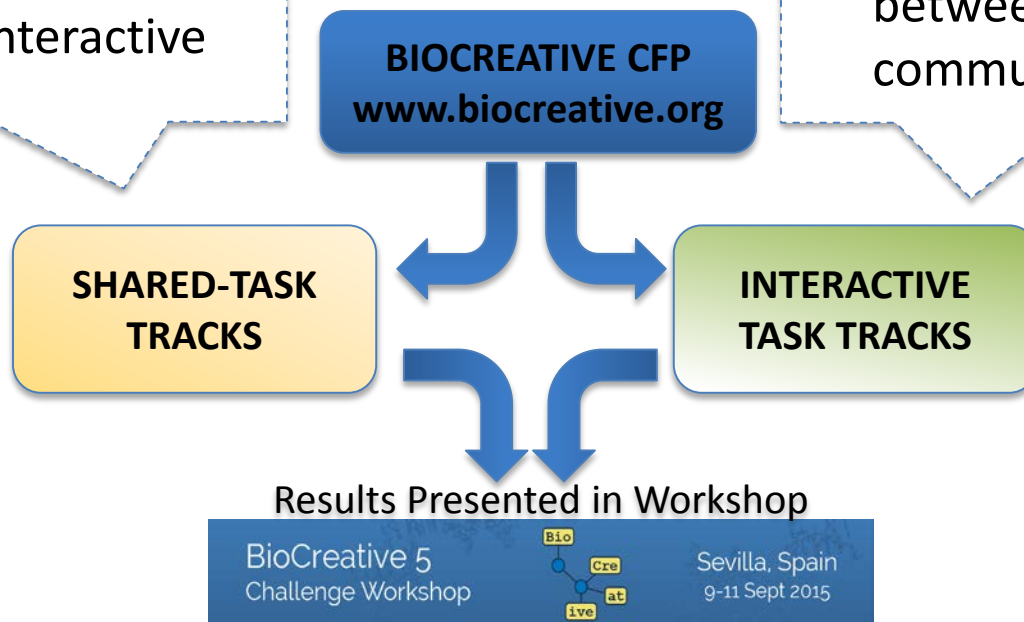
# Competitive

- Task relevant to biomedical domain
- Drive state-of-the-art TM tools development
- Provide building modules for systems in interactive task

# Non-competitive

- System development for literature curation tasks
- Tested by users
- Foster interaction between bioNLP and user communities

**BIOCREATIVE CFP**
**www.biocreative.org**

**SHARED-TASK TRACKS**

**INTERACTIVE TASK TRACKS**

Results Presented in Workshop

BioCreative 5
Challenge Workshop

Sevilla, Spain
9-11 Sept 2015

| | | Tasks | Publications |
|---|---|---|---|
| **BC I** | Spain, 2004 | • Gene Mention<br>• Gene Normalization<br>• GO | BMC Bioinformatics 2005, 6 (Suppl 1) |
| **BC II** | Spain, 2007 | • Gene Mention<br>• Gene Normalization<br>• Protein-protein Interaction | Genome Biology 2008, 9 (Suppl 2) |
| **BC II.5** | Spain, 2009 | Protein-Protein Interaction:<br>• Interactor Normalization<br>• Interaction Pair<br>• Article Categorization | IEEE Transactions in Computational Biology and Bioinformatics 2010 |
| **BC III** | USA, 2010 | • Gene Normalization<br>• User Interactive Task<br>• Protein-protein Interaction | BMC Bioinformatics 2011 |
| **BC 2012** | USA, 2012 | • CTD<br>• Biocuration Workflow<br>• User Interactive Task | Database Virtual Issue 2012 |
| **BC IV** | USA, 2013 | • Interoperability<br>• ChemDNER<br>• CTD<br>• GO<br>• User Interactive Task | Database Virtual Issue 2014 Chemical Informatics |

http://www.biocreative.org/

# Track 1- Collaborative BioCurator Assistant Task (BioC)

**Goal:**
Build a complete system to assist BioGrid curation

**Task Organizers:**
Sun Kim, Rezarta Islamaj Doğan, Donald C. Comeau, W. John Wilbur (NCBI), Andrew Chatr-aryamontri (BioGrid)
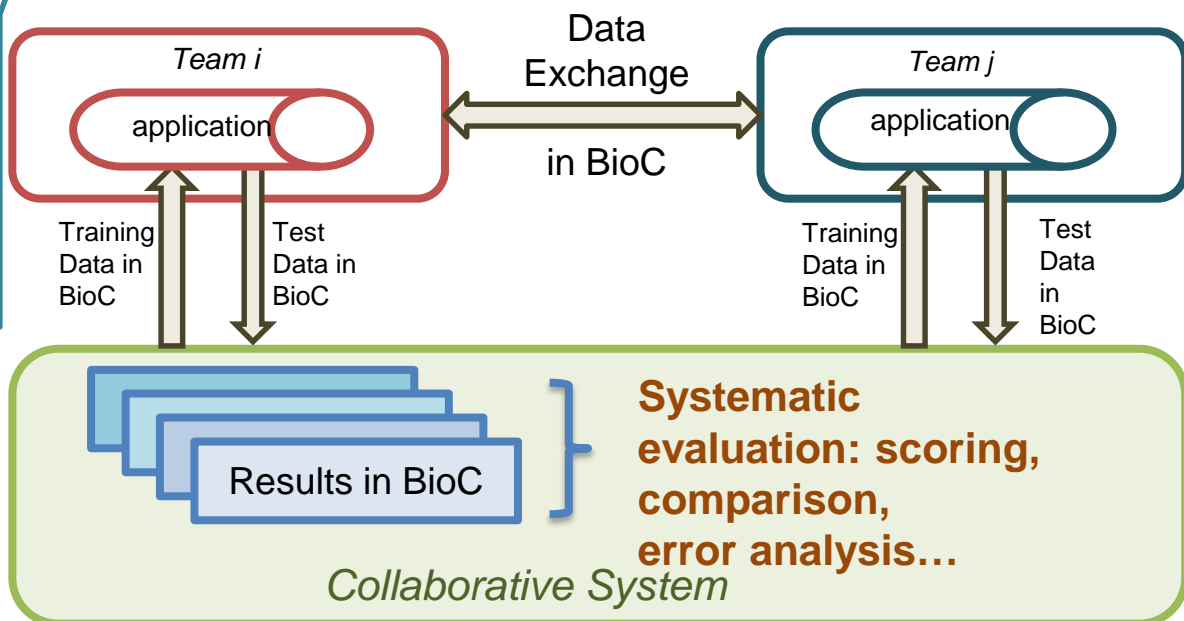
**Subtasks:**

- Protein, organism
- Physical, genetic interaction
- Experimental method
- Visualization tool

Corpora:
- PubMed abstracts
- PubMed Central full text articles

*Team i*
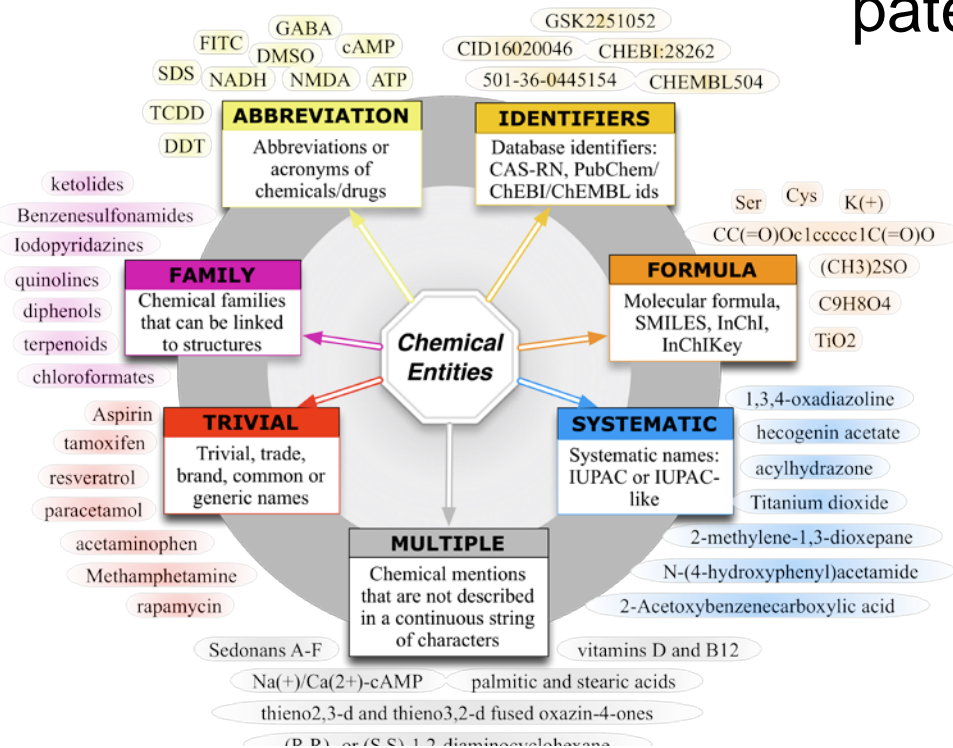application

Data Exchange
in BioC

*Team j*
application

Training Data in BioC

Test Data in BioC

Training Data in BioC

Test Data in BioC

Results in BioC

**Systematic evaluation: scoring, comparison, error analysis…**

*Collaborative System*

**Goal:** Automatic extraction of chemical and biological data from patents



*Task Organizers:*

Martin Krallinger & Alfonso Valencia (CNIO), Florian Leitner (UPM), Obdulia Rabal & Julen Oyarzabal    (CIMA)

Builds on successful task on abstracts 3,000 abstracts test set, 91% inter-annotator agreement. F-score 87.39%

Corpora: 30,000 manually annotated medicinal chemistry patent abstracts

*Subtasks:*
- **CEMP subtask** (chemical entity mention in patents)
- **CPD subtask** (chemical passage detection): the detection patent abstracts mentioning  chemicals (text classification/triage)
- **GPRO subtask** (gene and protein related object task): gene/protein mentions

# Track 3- Chemical-disease relation (CDR)

**Goal:** Advance the field in relation extraction from biomedical literature

*Task Organizers:*
Zhiyong Lu (NCBI),
Thomas Wiegers (CTD)

*Subtasks:*

- Disease Named Entity Recognition
- Chemical-induced disease relation extraction

## What is BEL?

Biological Expression Language
Computable knowledge representation

*Task Organizers:*
Fabio Rinaldi (UZurich)
Juliane Fluck ( Fraunhofer )



**Corpus**: 50 biological networks, 180,000 relationships

## *Subtasks:*

- Generation of BEL statements given the evidence
- Find evidence for a given BEL statement

http://www.openbel.org/content/bel-lang-language-structure

# Track 5- User Interactive Task

***Task Organizers:***
Cecilia Arighi, Qinghua Wang (PIR, UDel) and Lynette Hirschman (MITRE)

Evaluation of text mining tools by users

It is a *demonstration interactive* task

Need to involve users

User Advisory Group (UAG)

A diverse sample of end users with multiple text mining needs

- Help to develop end user requirements for interactive text mining tools
- Serve as users for the interactive task
- Assist in corpora annotation for biocreative tasks
- Help in recruiting biocurators

# UAG BioCreative V

## Chairs: Cecilia Arighi and Zhiyong Lu

Andrew Chatr-aryamontri

Raul Rodriguez-Esteban

Stan Laulederkind

Sherri Matis-Mitchell

Johanna McEntyre

Peter McQuilton

Evangelos Pafilis

Sandra Orchard

Sangya Pundir

Mary Schaeffer

Kimberly Van Auken

**Workshops at scientific meetings**
**Webinars (BioC)**
**Publications**

**BioC corpora**
**GN corpora**
**PPI corpora**
**CTD/CDR corpora**
**ChEMDNER corpora**

**Education & Outreach**

**V**

**I**

**Shared test collections**

**BioNLP Challenge Evaluations**

**II**

**New algorithm development & Improved results**

**Push into real-world applications**

**Interactive Task**
**BioC collaborative task**

**IV**

**Form & strengthen research communities**

**III**

**Modified from Chung-Chi Huang, and Zhiyong Lu Brief Bioinform 2015;bib.bbv024**

**Briefings in Bioinformatics**

BioCreative has promoted the development for state-of-the-art solutions

- Gene Mention: AIIA-GMT
- Gene Normalization: GNAT, GenNorm, ProMiner
- PPI triage: PIE
- GO categorization: GoCat

A variety of methods have been applied:

- Markov models
- Machine learning
- Rule-based
- Naïve Bayes classifiers
- Support Vector Machine

**Threshold Average Precision (TAP-k)**

In Gene normalization BioCreative III

Derivative of Mean average precision (MAP) with a threshold determined by the first k errors in the ranked list.

TAP-$k$ is able to measure ranking, reflect the user tolerance of prediction errors (false positives), as well as make use of confidence scores.

**Hierarchical Precision, Recall and F scores**

In Gene Ontology task BioCreative IV

Given the hierarchical nature of GO, considers common parent terms in computer-predicted and human-annotated GO terms

Successful BioNLP-user interactions through BioCreative

**PubTator**
Used by NLM for indexing. Currently being used by UniProt curators

**TagTog**
Gene indexing in Flybase (Cejuela et al., PMID:24715220)

**ODIN**
In PharmGKB workflow (Rinaldi et al., PMID: 22529178) now being tested on RegulonDB for BioCreative V

**RLIMS-P**
Phosphogrid curation (Torii et al., PMID:25122463)

# Impact/contributions

|  | # articles in ePMC |
|---|---|
| BioCreative Editorials | 6 |
| BioCreative mentioned in title or abstract | 148 |
| BioCreative is found in reference section | 389 |

| Top 10 MeSH terms in articles referencing BioCreative | Frequency |
|---|---|
| Humans | 119 |
| Natural Language Processing | 98 |
| Algorithms | 81 |
| Databases (Factual, genetic, protein) | 79 |
| Software | 78 |
| Computational Biology/methods | 69 |
| Vocabulary, Controlled | 61 |
| Artificial Intelligence | 60 |
| Information Storage and Retrieval/methods | 51 |
| Data Mining/methods | 50 |

# Gene Normalization in BioCreative

Convention for protein naming are different in different organisms. Differ in number of synonyms, in complexity of names,

| | | BC I | | BC II | BC III |
|---|---|---|---|---|---|
| |  |  |  |  |  |
| F-Score | 0.92 | 0.82 | 0.79 | 0.81 | 0.50 |
| IAA % | 87 | 91 | 69 | 91 | - |
| Average synonyms per identifier | 1.86 | 2.94 | 2.48 | 5.5 | |
| Average synonym length in words | 1 | 1.47 | 2.77 | 2.17 | |

http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2559987/

http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3269937/

Overall statistics of the annotated corpus grouped by data sets

| Data set | Articles | Genes (unique) | GO terms (unique) | Evidence text passages w.r.t. GO\|Gene\|Unique |
|---|---|---|---|---|
| Training set | 100 | 316 | 611 | 2440\|2478\|1858 |
| Development set | 50 | 171 | 367 | 1302\|1238\|964 |
| Test set | 50 | 194 | 378 | 1763\|1677\|1253 |
| Total | 200 | 681 | 1356 | 5505\|5393\|4075 |

| Genes | GO terms | Exact match | | | Hierarchical match | | |
|---|---|---|---|---|---|---|---|
| | | P | R | $F_1$ | hP | hR | $hF_1$ |
| 172 | 860 | **0.117** | 0.157 | **0.134** | 0.322 | 0.356 | **0.338** |
| 172 | 1720 | 0.092 | 0.245 | **0.134** | 0.247 | 0.513 | 0.334 |
| 172 | 3440 | 0.057 | **0.306** | 0.096 | 0.178 | **0.647** | 0.280 |
| 50 | 2639 | 0.018 | 0.075 | 0.029 | 0.064 | 0.190 | 0.096 |
| 46 | 1747 | 0.024 | 0.065 | 0.035 | 0.087 | 0.158 | 0.112 |
| 23 | 37 | 0.108 | 0.006 | 0.012 | **0.415** | 0.020 | 0.039 |

IAA GO term selection
47% strict
62% hierarchical

Database: The Journal of Biological Databases and Curation

o   The Tasks have evolved to resemble more the real scenario

o   Improvements have been achieved in many tasks

o   Combination of methods usually improves the performance

o   Although results are not of sufficient quality to use as an entirely automated process, output from these tools can provide a head start for curators

# Can text mining tools help in Biocuration

Our idea is to expose text mining systems to biocurators so they can provide feedback on the system and become adopters in the future

intuitive    efficient    usable

satisfying

**Partial**

Follow pre-defined tasks aimed at testing system usability

↓

Complete user survey

↓

Collate outputs

**Full**

Training via demo, examples, help document, annotation guidelines, and output format

Practice with examples, report bugs

Is biocurator familiar with system and annotation ?

No

Yes

Dataset selected by domain expert (or coordinator)

1/2

1/2

Non TM-assisted Annotation

TM-assisted Annotation

Complete user survey

Collect outputs and calculate metrics

Based on Questionnaire for User Interface Satisfaction (QUIS)

**Five main categories**:
1. Overall reaction
2. System's ability to help complete tasks
3. Design of application
4. Learning to use the application
5. Usability

**Goal:** Try to find correlation of response to questions in survey with overall system satisfaction to learn  what aspects are important to users

http://ir.cis.udel.edu/biocreative/survey.html

http://ir.cis.udel.edu/biocreative/survey2.html

**1-Match between system and the real world (of Biocuration)**

- The system should speak the users' language rather than system-oriented terms
- Systems should follow standards of its user community
- Sentence vs. Document level annotation

**2-Testing the Systems NOT the Users**

- Participants not being tested. But in the context of this activity we need to distinguish the participants into curation novice vs. expert because it has an impact on the performance

**3-Documentation: Annotation guidelines and tutorials**

- Provide detail annotation guidelines for the task
- Provide tutorial for system training with hands-on examples

## 4-System performance and functionalities in interface

- System performance is a key aspect to biocurators, but coupling results with functionalities that assist in easily correcting or finding additional information is very important for an interactive system

## 5-System Output

- To be useful for curation annotated results should be exported in standard formats that can be further utilized in the curation workflow. Tab-delimited and BioC formats were requested.
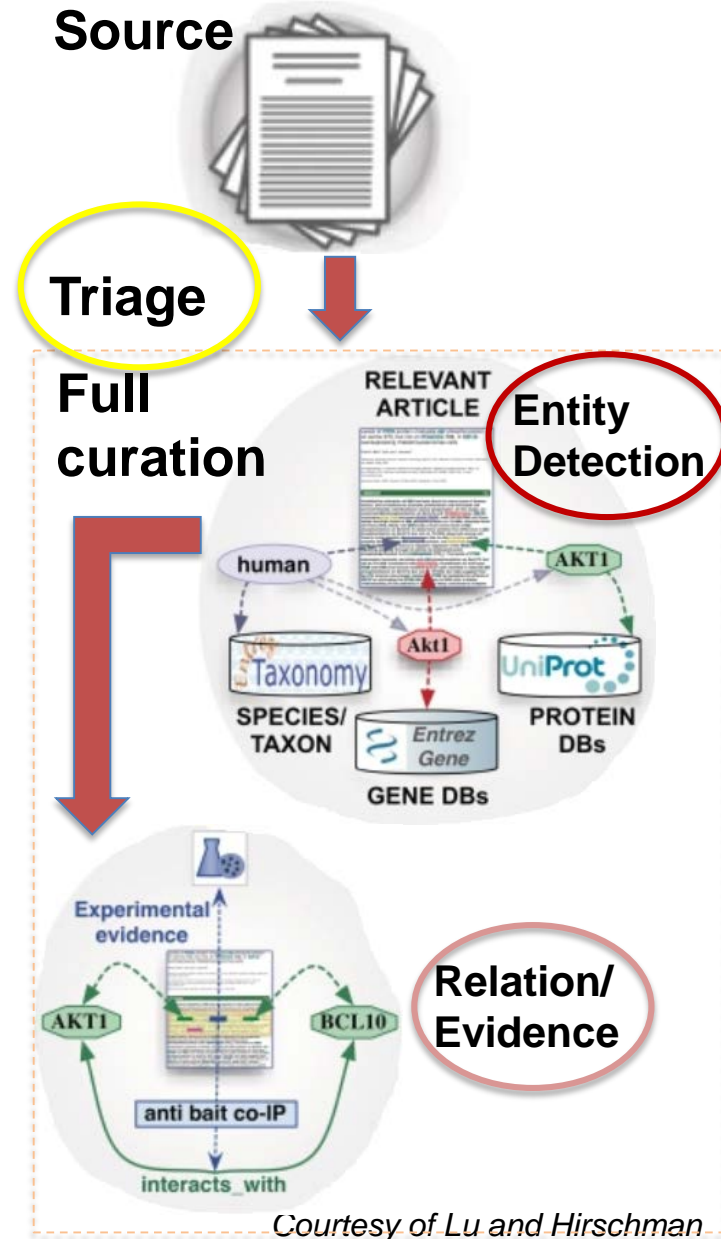
## Metrics:

Time on task (objective)
Preference via survey (subjective)

| System | Description of the tool |
|---|---|
| **Cell Finder** | Annotation of gene, expression relation and cell type in text snippets from a set of articles |
| **Ontogene** | Detection of Gene/Chemical/Diseases and their interactions |
| **MarkerRIF** | Retrieval of articles about biomarkers, and extraction of disease and biomarker (gene) with normalization |
| **SciKnowMine** | Triage based on pre-trained categories of interest in full length articles |
| **BioQRator** | Retrieval based on relevance on protein-protein interaction information and annotation of protein pair |
| **RLIMS-P** | Triage on protein phosphorylation. Annotation of kinase, substrate and site with normalization. |
| **Egas** | Identification and extraction of protein-protein intearaction events described over PubMed abstracts related to neuropathological disorders |
| **tagtog** | Annotation of gene names within full-text documents especially machine-predicted documents |
| **Argo** | Annotation of metabolic process-related named entities, namely chemical entities and genes or gene products |



Source

Triage

Full curation

Entity Detection

RELEVANT ARTICLE

human — AKT1 — Akt1

Taxonomy — UniProt

SPECIES/ TAXON — Entrez Gene — PROTEIN DBs

GENE DBs

Relation/ Evidence

Experimental evidence

AKT1 — BCL10

anti bait co-IP

interacts_with

*Courtesy of Lu and Hirschman*

**Cecilia Arighi**

## Recruitment of Biocurators

Call for participation via International Society for Biocuration (ISB) mailing list, and the ISB meeting and BioCreative websites
Personal invitation

## What's in it for Biocurators?

- Exposure to state-of-the-art text mining systems
- Direct communication and interaction with developers
- Contribution to tools that meet the needs of biocurators
- Adoption of text mining tool
- Potential publication in peer reviewed journal
- Focus on a set of articles that will add to their curation effort

A)

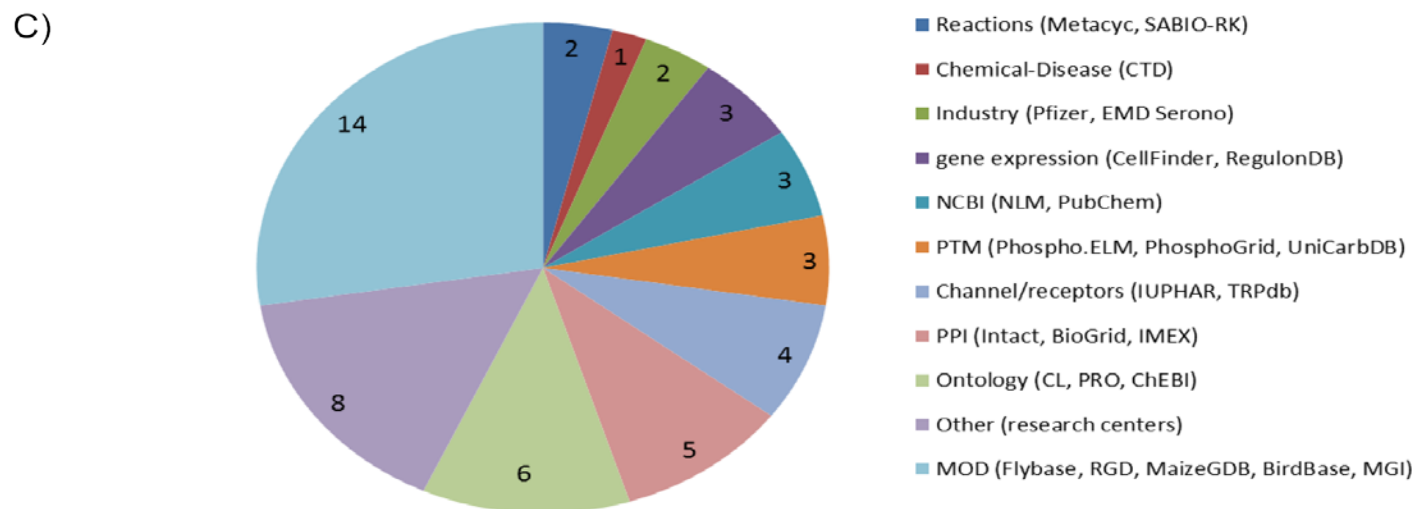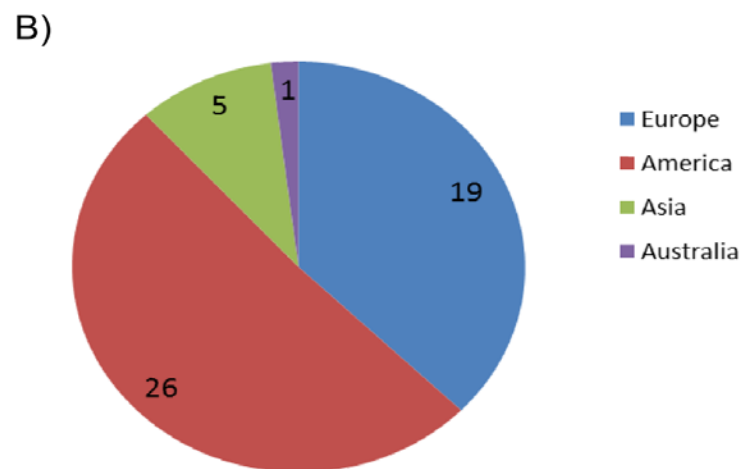| System | # Curators | |
|---|---|---|
| | Full | Partial |
| CellFinder | 4 | 2 |
| Ontogene | 2 | 3 |
| MarkerRIF | 3 | 3 |
| SciKnowMine | 1 | 5 |
| BioQRator | 6 | 6 |
| RLIMS-P | 3 | 5 |
| EGAS | 4 | 4 |
| Tagtog | 7 | 2 |
| Argo | 3 | 1 |
| Total | 33 | 31 |

B)



Europe — 19
America — 26
Asia — 5
Australia — 1

C)



- Reactions (Metacyc, SABIO-RK) — 2
- Chemical-Disease (CTD) — 1
- Industry (Pfizer, EMD Serono) — 2
- gene expression (CellFinder, RegulonDB) — 3
- NCBI (NLM, PubChem) — 3
- PTM (Phospho.ELM, PhosphoGrid, UniCarbDB) — 3
- Channel/receptors (IUPHAR, TRPdb) — 4
- PPI (Intact, BioGrid, IMEX) — 5
- Ontology (CL, PRO, ChEBI) — 6
- Other (research centers) — 8
- MOD (Flybase, RGD, MaizeGDB, BirdBase, MGI) — 14

# Time on task in full level participation curation task and curator experience level

| System | no TM (min) | TM (min) | TM (min) | $t_{noTM}/t_{TM}$ | $t_{noTM}/t_{TM}$ | Curation experience (years) |
|---|---|---|---|---|---|---|
| BioQRator | 275 | 195 | | 1.4 | | <1 |
| | 70 | 100 | | 0.7 | | >3 |
| | 160 | 180 | | 0.9 | | >3 |
| | 150 | 150 | | 1.0 | | 1-3 |
| Egas | 93.71 | 60.13 | | 1.6 | | <1 |
| | 184 | 120 | | 1.5 | | 1-3 |
| | 104.91 | 26.21 | | 4.0 | | <1 |
| | 64.48 | 60.86 | | 1.1 | | >3 |
| MarkerRIF | 212 | 90 | 145 | 2.4 | 1.5 | 1-3 |
| | 115 | 84 | 70 | 1.4 | 1.6 | <1 |
| | 170 | 95 | 103 | 1.8 | 1.7 | <1 |
| RLIMS-P | 585 | 560 | | 1.0 | | >3 |
| | 301 | 186 | | 1.6 | | 1-3 |
| | 164 | 161 | | 1.0 | | <1 |

# Subjective measure



**Overall Impression** (pool from 3 questions)

**Task Completion** (pool from 3 questions)

**Design** (pool from 4 questions)

**Learnability** (pool from 4 questions)

**Usability** (pool from 5 questions)

Key

- no response
- positive response
- neutral response
- negative response

## The curation time does not always go hand by hand with user overall system satisfaction

With BioQRator and RLIMS-P curators are satisfied with system even the time required in the no-TM-assisted versus TM-assisted curation was comparable for each.

**Some reasons**:

-system provides a nice interface with functionalities that in the long run makes the monotonous curation work more enjoyable

-some systems have both retrieval and extraction steps, the curators appreciated the retrieval step because it saves a lot of time in article selection. However, the task was measured on the extraction step and most time was spent on normalization.

*System accessibility*:  due to one of the following; firewalls, system temporarily down, or inability to log in.

*Error messages*: either no error message displayed or the error message did not satisfactorily explain the problem.

*Hidden functionality:* key functionality for executing the TM task not apparent to curators.

*Language and icons*: icons and names of sections/functionalities non-intuitive or used TM jargon.

*Look and feel:* Color choice for entity highlighting was not optimal for color blinded users for some of the systems

http://www.biocreative.org/tasks/biocreative-v/iat-task-biocurators/

Seven Systems for different tasks

Please select one or more systems from this list:

☐ **Argo** (Curation of phenotypes relevant to the chronic obstructive pulmonary disease (COPD) in the PhenomeNet database)

☐ **Egas** (Identification of clinical attributes associated with human inherited gene mutations, described in PubMed abstracts)

☐ **OntoGene** (Curators interested in bioconcepts currently supported by OntoGene, e.g., miRNA, gene, chemical, disease)

☐ **GenDisFinder** (Knowledge discovery of known/novel human gene-disease associations from biomedical literature)

☐ **MetastasisWay** (Look for the biomedical concepts and relations associated with metastasis and construct the metastasis pathway)

☐ **BELIEF** (A semi-automated curation interface which supports expert in relation extraction and encoding in the modelling language BEL (Biological Expression Language) )

☐ **EXTRACT** (List the environment type and organism name mentions identified in a given piece of text)
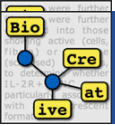
Evaluation Period: period June 22 to July 31
Flexible and remotely conducted
Total time commitment estimation over that period:
Full participation: 12h
Partial participation: 30min-1h

Sign In!

Manual curation is accurate, but does not scale. Text mining scales, but is not accurate

Interactive systems like those presented in IAT can provide curators with decision support:

- suggesting important papers to curate
- highlighting entities of relevance in text
- offering controlled vocabularies and ontologies
- on-the-fly error-correction
- removal of redundancy

o Tools developed in past BioCreative challenges have been integrated as modules in a subset of the participating systems, such as GenNorm in RLIMS-P; and PIE for protein-protein interaction article ranking and retrieval BioQRator.

o This demonstrates the importance of the traditional shared tasks to promote development of state-of-the-art text mining tasks that when mature these offer text mining solutions that can be integrated in a system framework.

**BioCreative Organizers**

**UAG Members**

**Text mining teams and biocurators that participated in the BioCreative tasks**