

BioCreAtIvE task 1A: gene mention finding evaluation

Alexander Yeh ^{1§}, Lynette Hirschman¹, Alexander Morgan¹, Marc Colosimo¹

¹The MITRE Corporation, 202 Burlington Road, Bedford, MA 01730, USA

[§]Corresponding author

Email addresses:

AY: asy@mitre.org

LH: lynette@mitre.org

AM: amorgan@mitre.org

MC: mcolosimo@mitre.org

Abstract

Background

The biological research literature is a major repository of knowledge. As the amount of literature increases, it will get harder to find the information of interest on a particular topic. There has been an increasing amount of work on text mining this literature, but comparing this work is hard because of a lack of standards for making comparisons.

Results

We took part in running BioCreAtIvE (Critical Assessment for Information Extraction in Biology), an open common evaluation of systems on a number of biological text mining tasks. We report here on task 1A, which deals with finding mentions of genes and related entities in text. The task makes use of data and evaluation software provided by the (US) National Center for Biotechnology Information (NCBI). 15 teams took part in task 1A.

Conclusion

A number of teams achieved scores over 80% F-measure (balanced precision and recall). This is good, but still somewhat lags the best scores achieved in some other domains such as newswire, due in part to the complexity and length of gene names, compared to person or organization names in newswire. "Finding mentions" is a basic task, which can be used as a building block for other text mining tasks, but the teams that tried to use their task 1A systems to help on other BioCreAtIvE tasks report mixed results.

Background

The biological research literature is a major repository of knowledge. Unfortunately, the amount of literature has gotten so large that it is often hard to find the information of interest on a particular topic. There has been an increasing amount of work on text mining this literature, but currently, there is no way to compare the systems developed because they are run on different data sets to perform different tasks [Hirschman 2002A]. Challenge evaluations have been successful in making such comparisons. Examples include the ongoing CASP evaluations (Critical Assessment of Techniques for Protein Structure Prediction) for protein structure prediction [CASP], the series of Message Understanding Conferences (MUCs) for information extraction on newswire text [Hirschman 1998], and the ongoing Text Retrieval Conferences (TREC) for information retrieval [TREC][Voorhees 2002]. Also, in 2002, we ran the first challenge evaluation of text mining for biology; this was an evaluation for classifying papers and genes with experimental evidence for gene products [Yeh 2003].

As mentioned in [Yeh 2003], the idea behind these series of open evaluations has been to attract teams to work on a problem by providing them with real (or realistic) training and test data, as well as objective evaluation metrics. These data sets are often hard to obtain, and the open evaluation makes it much easier for groups to build systems and compare performance on a common problem. If many teams are involved, the results are a measure of the state-of-the-art for that task. In addition, when the teams share information about their approaches and the evaluations are

repeated over time, then the research community can demonstrate measurable forward progress in a field.

To further the field of biological text mining, the “BioCreAtIvE” evaluation was run in 2003, with a workshop in March 2004 to discuss the results [BioCreAtIvE 2004]. The evaluation consisted of two tasks: Task 1 focused on extraction of gene names (Task 1A) and normalization of genes (Task 1B) from PubMed abstracts. Task 2 was a more advanced task focused on functional annotation, using full text information to classify a protein as to its molecular function, biological process and/or location within a cell. Task 1 consisted of two parts: This paper reports on task 1A, entity mention extraction. This extraction is a basic text mining operation. Its output is the input text, annotated with the mentions of interest; this can be used as a building block for other tasks, such as task 1B and task 2.

The gene mention task presents a number of difficulties. One difficulty is that gene (or protein) mentions are often English common nouns (as opposed to proper nouns, which, in English, are the nouns normally associated with names) and are often descriptions. In fact, many entities are named with ordinary words, for example some *Drosophila* (fruit fly) gene names are *blistery*, *inflated*, *period*, *punt*, *vein*, *dorsal*, *kayak*, *canoe* and *midget*. In addition, new entities are constantly being discovered and/or renamed with these common nouns. Also, many names originate as descriptions and can be quite complex, e.g., *hereditary non-polyposis colorectal cancer (hnpcc) tumor suppressor genes*.

Task and data

The data and evaluation software for task 1A were provided by W. John Wilbur and Lorraine Tanabe at the National Center for Biotechnology Information (NCBI). Every mention of interest is marked, so this task corresponds to the “named entity” task used in the natural language processing community. The data is marked for mentions of “names” related to genes, including binding sites, motifs, domains, proteins, promoters, etc. The data comes with a particular tokenization (word segmentation), and this tokenization determines the boundaries of what is marked. A token is either entirely part of a markable or not. A token *cannot* be split between a marked part and an unmarked part. For testing, the systems take as input the tokenized unannotated sentences; the output is the list of gene names for each sentence, with the start and stop token offsets. For evaluation, the system output is then compared to the “gold standard” hand-annotated answer key.

The data consists of sentences from Medline [Medline] abstracts that have been manually annotated for mentions of names of genes and related entities. Half of the sentences were chosen from abstracts likely to contain such names. The other half were chosen from abstracts likely *not* to contain such names. See [Tanabe 2004] (also in this volume) for further detail on the construction of the Task 1A data. The approximate sizes of the various data sets are given in Table 1.

The (final) test set is also known as “round1”, and for the evaluation, its sentences were renumbered to give no indication of what Medline abstracts they came from. The original sentence numbers were derived in part from the Medline/Pubmed id number of the abstract from which the sentence was drawn.

There is no detailed, multi-page explicit set of guidelines describing what is markable. Instead, there is a description provided with the data that gives a page or two listing of the types of entities that are and are not markable. Examples of markables are mutants (e.g., *p53 mutant*) and words like *codon*, *antibody*, etc. when combined with a gene name. Examples of non-markables include generic terms (e.g., the term *zinc finger* by itself) and mutations (e.g., *p53 mutations*).

Here are 2 excerpts from the training corpus (sentences 110312525757 and 13393732909):

The LMW FGF-2 up-regulated the PKC epsilon levels by 1.6 fold; by contrast the HMW isoform down-regulated the level...
...a protein related to ***SNF1 protein kinase***.

The underlines indicate the markable entities. The ***italic boldface*** indicates what alternative mentions can substitute for a markable. Note that for “SNF1” and “protein kinase”, an allowed alternative is “SNF1 protein kinase”, which includes both of them.

Correspondingly, the answer file contains the following mentions: “LMW FGF-2”, “PKC epsilon”, “HMW isoform”, “SNF1” and “protein kinase”.

Stored in another file are the alternative mentions that can be tagged and still count as being correct. For the answers mentioned above, here are the allowed alternative mentions: “FGF-2”, “PKC”, “HMW”, and “SNF1 protein kinase”.

When scoring, an exact match to an answer or an allowed alternative is needed to get credit for finding an answer. So for example, if for the answer LMW FGF-2, a system instead returns “The LMW FGF-2”, that system would get both a false negative (not matching the answer or its alternative) and also a false positive (the returned item does not match an answer or any alternative).

Results

15 teams entered submissions for this evaluation. Submissions were classified as either “open” or “closed”.

Closed: The system producing the submission is only trained on the task 1A “train” and “(development) test” (devtest) data sets, with no additional lexical resources

Open: The system producing the submission can make use of external word lists, other data sets, etc.

Most teams provided an “open/closed” classification for their submissions. If they did not, we made a classification based on a short system description that the teams provided. When we were not sure, we assumed “open”.

Teams were allowed to send up to 4 submissions each, as long as they included a “closed” submission. Teams only sending “open” submissions were allowed to send up to 3 submissions. We received a total of 21 “closed” submissions (plus 2 more

received late and deemed “unofficial”) and 19 “open” submissions (also plus 2 more received late and deemed “unofficial”).

The submissions were measured by their balanced F-score, recall and precision.

- Balanced F-score is the harmonic mean between recall and precision.
Balanced F-score = $2 * \text{Recall} * \text{Precision} / (\text{Recall} + \text{Precision})$
- Recall is the fraction or percentage of the answers in the answer key that were found by a submission.
- Precision is the fraction or percentage of the answers returned by a submission that turn out to be correct.

Scores achieved by the submissions

Many of the high performing submissions achieved scores quite close together. For example, with balanced F-score, the first and second highest teams were only 0.6% apart, and the second and third highest teams were even closer at 0.2% apart. This is close enough to be possibly affected by the disagreements in annotation that arise with just about any task on finding entity mentions. An example is that with this particular task, a partial review of the test set changed 0.4% (25 of 6000) of the answers.

It would be good to test the results for statistical significance, but this has not been done yet. It is quite possible that these differences are small enough to not be statistically significant.

Table 2 shows the high and low scores, as well as the 1st, 2nd (median) and 3rd quartile balanced-F, recall and precision scores for the 40 official submissions. One can see a compression or skew of the scores towards the high end.

- The high, 1st, 2nd and even 3rd quartile scores are relatively close to each other compared to the low scores
- With F-score, the top 3 teams have scores within 1% of each other
- With recall, the top 2 teams are separated by about 2%
- With precision, the top 2 teams are separated by about 1%

Generally, the open submissions did better than the closed submissions. An exception is that for the highest recall score, the top closed score is actually better than the top open score. The compression at the high scores also occurred for the gap between the open and closed submissions.

- For the higher scores (like high and 1st quartile), there was little difference (2% or less) between the open and closed submission scores
- For the lower scores (like low and 3rd quartile), the open submissions scores are measurably better than the closed scores

Figure 1 shows the balanced-F scores of the 40 official plus 4 unofficial submissions. The open submissions are in a dark solid and the closed submissions are in white with an outline. The submissions are labeled with the user (u) number of the team. 13 official submissions from 4 different teams achieved an F-score of 80% or higher (in the figure, this appears as 0.8 or higher). For most teams, their open submission(s) scored higher than their closed submission(s). Team 11 was an exception, as was team 5 to some extent. The gap between a team’s open and closed submissions was small

compared to the gap between the submissions from different users. However, team 21 had a large gap between their open and closed submissions.

Figure 2 shows a plot of the precision versus recall of the 40 official plus 4 unofficial submissions. The official open submissions are shown with diamonds, the official closed submissions are shown with squares. Unofficial submissions are shown with gray outlines (and clear centers) of diamonds and squares, respectively. Eight official submissions (from 3 different teams) achieved both a recall and precision of 80% or higher (appears in the figure as 0.8 or higher). As a set, the submissions with both a recall and precision of 60% or more seem to have a fairly balanced precision and recall. But for the most part, submissions which had a recall or submission below 60% tended to have a better recall than precision.

Some observations

Most teams made use of the training data in their system development. However, in reading the task 1A participants' system descriptions team u27 did not [BioCreAtIvE 2004]. Also, as far as we can tell, neither did team u18 (based on a short description not in the reference). Like many similar tasks, task 1A has its own unique features. This is probably a reason why, relative to the other teams, these two teams did not get very good results: u27's submission had a 61% balanced F-score, while u18's submission had 55% (both in the 4th quartile range for official open submissions). One indication of these unique features comes from Tamames [Tamames 2004](Discussion of the results - task 1A), whose system had not considered entities like domains, regions and mutants as "gene names" that should be marked, where as task 1A did include such entities.

A common comment from several Task 1A participants (for example, see the post-processing descriptions in Dingare [Dingare 2004](sec. 2.3) and Kinoshita [Kinoshita 2004](sec. 3)) was that one of the more difficult aspects of task 1A was determining the starting and ending boundaries of the gene-or-protein names. The requirement for an exact match to the answer key (or alternative) increased the difficulty.

As has been mentioned, many of the open and closed submissions achieved fairly close results. One possible reason for this is that, to the extent that this task is unique, outside sources will not help performance that much. Another possible reason is that for the most part, we relied on the teams themselves to classify their submissions as being "open or closed". In viewing the task 1A system descriptions [BioCreAtIvE 2004], one can see that the different teams varied in what resources they thought were allowed in a closed submission. As an example, when using a sub-system that generates part-of-speech (POS) tags, some (but not all) teams use such a POS sub-system for a "closed" submission even when the sub-system itself was trained on another annotated corpus. This is an indirect reliance on an outside corpus. Some teams treated this indirect reliance as permissible for a closed submission (for example, Dingare [Dingare 2004](sec. 2.1) and Zhou [Zhou 2004](sec. 1)), some teams did not.

Summary of System Descriptions

For task 1A, the teams tended to use one of the three following approaches at the top level of their system (see the participants' system descriptions [BioCreAtIvE 2004]):

1. Some type of Markov modelling.

2. Support vector machine (SVM). Typically, the input information on the word being classified would come from a small window of the words near that word of interest.
3. Rules. As far as we could tell, the rules were usually manually generated.

Many of the systems had pre- and/or post-processing stages in addition to the main approach taken. One system combined several other systems via a voting scheme [Zhou 2004].

These teams used a variety of features in their systems. For some features, many of the systems used entire other sub-systems to find the feature value. An example is using a part-of-speech (POS) tagger to find a word's part-of-speech. These sub-systems often used an approach that differed from the overall system's approach.

The four teams with 80% or higher F-scores had post-processing stages in addition to the main approach taken, and made use of many different features. All four of these teams performed some type of Markov modelling at the system's top level [Dingare 2004][Kinoshita 2004][Zhou 2004][McDonald 2004]. However, the teams used different techniques on their Markov models: maximum entropy, hidden Markov models (HMM) and conditional random fields. In addition, one team [Zhou 2004], also had an SVM system at the top level: decisions were made by having two HMMs and an SVM system vote. Also, note that when comparing different systems, the choice of features used is often at least as important as the approach/algorithm used. Yeh [Yeh 2000] gives an example of this.

Task 1A as a Building Block

One of reasons for having task 1A is that a task 1A system can serve as a building block for other tasks, like task 1B or task 2 of the BioCreAtIvE evaluation. The task 1B evaluation focused on finding the list of the distinct genes (of a particular species) mentioned in a Medline abstract, where the list contained the normalized, canonical names for those genes. Task 2 focused on classifying what a protein does and where in a cell it is found, and on returning text passages as evidence to support these classifications.

To what extent was it viable to use task 1A systems as a building block for more advanced capabilities? It turns out that three of the teams taking part in task 1A also took part in task 1B. In addition, one of the three teams also took part in a portion of task 2. So an interesting question is how useful these three teams found their task 1A systems to be when working on task 1B or 2.

One team with a high precision (80%+) task 1A system used the mentions found by their 1A system as the input for their 1B system [Tamames 2004](Task 1B): their 1B system then tried to find the normalized version of the mentions found by their task 1A system.

The story was more complicated for two other teams with both high precision (80%+) and high recall (78%+) task 1A systems. One team was from Pennsylvania (1A: [McDonald 2004], 1B: [Crim 2004]). The other team was from Edinburgh and Stanford (1A: [Dingare 2004], 1B: [Hachey 2004]). Both these two teams looked at some version of finding mentions with their task 1A system and then compared the

found mentions against the synonym lists for the genes of interest for task 1B. One complication that both teams found was this approach could easily produce a low precision for 1B, due to many genes sharing many of the same synonyms.

The Pennsylvania team also found that for genes from two (fly and yeast) of the three organisms of interest in task 1B (mouse was the 3rd organism), the task 1A tagger was not that accurate. A possible explanation given was that the task 1A training data did not have enough examples from these two organisms. For task 1B, the Pennsylvania team in the end did not use their task 1A tagger.

The Edinburgh/Stanford team found that using the original task 1A training set and lots of features tended to lower their recall of the 1B genes. They raised the recall by training their 1A system using the noisy task 1B training data and a reduced set of the possible features.

The Edinburgh/Stanford team also took part in task 2.1. In this task, a system is given an article, a protein mentioned in that article, and a classification of that protein that a person made based on that article. The system's job is to find a passage of text in that article that supports the classification made for that protein. The description for the team's task 2.1 system [Krymolowski 2004] made no mention of using their task 1A system or trying it on some part of task 2.1.

Discussion

One unique aspect of the data: enforcing a particular tokenization

As mentioned before, every entity mention task such as task 1A will have some features that are more or less unique to it. For task 1A, one such feature is that the data comes with a particular tokenization (word segmentation). Furthermore, this tokenization affects what counts as a mention, because either all of a token is tagged as part of a mention, or none of that token is tagged. This can cause problems when one just wants to tag part of a token as part of a mention. An example is the phrase

... a *protein kinase* a-mediated pathway ...

where the odd word tokens are underlined, but the even tokens are not. Here the token "a-mediated" is not useful, as the mention that one would really like to tag is "protein kinase a".

This tokenization is important because it affects what counts as a mention. Here are some rules (Lorraine Tanabe, personal communication):

1. If "X" is a token which is a gene name, then "X" is usually marked. An example is "CBF1" in the phrase "... of CBF1 in yeast ..." (in training data's sentence 90233781202).
2. If a token is of the form "X-" or "X-Y", where "X" is a gene name and "Y" is an adjective or verb, then the token is usually NOT marked. An example is "EGF-induced" in the phrase "... block EGF-induced mitogenesis and ..." (in training data's sentence 94547351603).

3. An exception to (2): when the Y in "X-Y" is "like", then "X-Y" is usually marked. Also, if the form is "X-Y Z", where "X-Y" is as in (2), and "Z" is a token like "domain", then "X-Y" is usually marked as part of the mention "X-Y Z". An example is "SH2-binding domain".

Disagreements in the data

In tasks like task 1A, small disagreements usually exist on what to annotate and what not to annotate. An example in task 1A is phrases of the form "X pathway(s)", where X is a phrase that is marked as part of a gene mention. An initial review of the test set found the following annotation variations (afterwards, all test set cases were changed to have "X" and "X pathway(s)" be both allowed alternative answers):

- 4 cases where "X pathway(s)" is NOT an allowed alternative to "X". An example is X= "Mek-Erk1/2" in the phrase "... the *Mek-Erk1/2* pathway by ..." (sentence 14076).
- 10 cases where "X" and "X pathway(s)" are both allowed alternatives. An example is X= "Ras/Raf/MAPK" in the phrase "... the *Ras/Raf/MAPK* pathway." (sentence 10544).

Similarly, the training set has

- 12 cases where "X pathway(s)" is NOT an allowed alternative to "X".
- 11 cases where "X pathway(s)" and "X" are allowed alternative answers.

Such variation in annotation makes it more difficult to learn or to formulate a rule for how to handle these kinds of constructions.

Lessons learned for future evaluations

If and when a future task 1A evaluation is run, we list the following issues to consider:

1. Tokenization is non-trivial for biological terms. Perhaps one should *not* enforce a fixed tokenization of the data. This non-enforcement will be expensive because one will need to change both how the data is annotated and how the system results are compared against the gold standard.
2. On a related matter, because of the difficulties in exactly determining a mention's boundaries, there is interest in also counting inexact matches to answers as being correct. One needs to be careful how this is done. For example, if missing either the first or last token still counts as correct, then just returning "epsilon" would count as finding "PKC epsilon".
3. For *open* versus *closed* submissions, one should either remove the distinction, or be more explicit as to what is allowed for a *closed* submission.
4. A suggestion was made to pad the test set with a lot of extra material that would not be scored, which will make it harder to "cheat" by manually examining the test set. If one were to do this, one would need to announce this ahead of time. One reason is that some automated approaches need more processing time than others. Another reason is that some automated approaches, such as transductive support vector machines [Joachims 1999], make use of statistics derived from the entire un-annotated test set.
5. At least one team [Dingare 2004] automatically searched for the PubMed/Medline abstract associated with each test set sentence. They used

the abstract as a surrounding context, and it seemed to be helpful. In many “real uses” of a task 1A system, a system will probably have such surrounding text. So one should probably just give these abstracts to every participant in the future.

6. There is also a question of what is a permissible resource to use:
 - One example is that with PubMed/Medline, a system could also look-up MESH terms, etc. associated with the Medline abstract for each sentence. If one plans to use such a tagging system before an abstract is assigned MESH labels (assignment is done manually), then such information will not be available in real usage, and such information should not be permitted.
 - Given a possible entity “X”, at least one team [Dingare 2004] did web searches for contexts like “X gene”, which support “X” being a possible entity. This seemed to be of limited help. Should this be permitted in the future? This probably depends on the anticipated “real” uses for such a feature. When tagging older material (such as the task 1A test set), the web will have relevant material. When tagging new text that describes new gene(s), the web will probably not have much, if any material.

Conclusions

For the BioCreAtIvE task 1A of gene mention finding, a number of teams achieved an 80-83% balanced F-score. This is similar to results for some other similar biological mention finding tasks, and is somewhat behind the 90%+ balanced F-scores achieved on English newswire named entity tasks [Hirschman 2002B]. Based on an observation offered by Kevin Cohen (who was on one of the teams [Kinoshita 2004]), a hypothesis for discrepancy is that gene names tend to be longer than comparable newswire names. To investigate this, we compared the length distribution of gene names in the test set for task 1A; this distribution is shown in Figure 3, and is compared to the distribution for name length of organizations in a newswire task. The newswire results are computed from the MUC-6 data, which is available from the Linguistic Data Consortium [LDC]. The average length of the task 1A gene names was 2.09, compared to 1.69 for ORGANIZATION names in the MUC-6 data. Given this distribution, we fitted a simple logistic regression model to both data sets, assuming that the success rate would be multiplicative based on the number of words in a name. This allowed us to extrapolate back to a single-word error rate for both tasks, allowing us to factor out differences in name length. For gene names, a 92% success rate on a single word gene name gave an overall task performance of 83%, the observed high score. For the MUC-6 organization names, a 95.5% single word success rate yielded a 93% success rate overall, which was the highest recorded result for MUC-6. In using this simple model, we recognize that it is not mathematically valid to use F-measure in place of accuracy. However, it does provide a crude approximation for how much of the task difficulty can be attributed to difference in name lengths among different tasks. This comparison leaves a residual 3-4% discrepancy between performance on the tasks for the single-word case. We hypothesize that this may be due to interannotator variability, leading to “noise” in the training and test data. For the MUC-7 task [Marsh 1998], interannotator agreement was measured at 97%, which is almost certainly significantly higher than for the gene mention task, which has not yet been formally measured.

In terms of successful approaches, the teams that achieved an 80% or more balanced F-score tended to use some type of Markov modelling at the top system level. However, these teams also had post-processing stages in addition to the main approach taken, and the different teams made use of different features. These stages and features can have just as much an effect on performance as the main approach taken.

One of the reasons to have task 1A is that it should be a useful building block to work on other tasks, like BioCreAtIvE task 1B. Three teams tried using their task 1A system for task 1B. Their experiences are mixed, with two of the three teams finding that a task 1A system trained on the task 1A training data often does not work so well on task 1B. One of these two teams improved things by retraining their 1A system using the noisy task 1B data.

A 2nd test set is available for task 1A, so it would be straightforward to run a task 1A evaluation in the future using this 2nd test set. Three questions to think about in any future evaluation are the following:

- What will it take to improve task 1A performance?
- How much will improving task 1A performance help with other tasks (like tasks 1B and 2)?
- How can one make a task 1A system be a more useful building block for other tasks?

Acknowledgements

This paper reports on work done in part at the MITRE Corporation under the support of the MITRE Sponsored Research Program and the National Science Foundation (contract number EIA-0326404). W. John Wilbur and Lorraine Tanabe of the National Center for Biotechnology Information provided the data and evaluation software used in the evaluation. Copyright © 2004 the MITRE Corporation. All rights reserved.

References (convert []'s to numbers later)

- [Hirschman 2002A]. Hirschman L, Park JC, Tsujii J, Wong L, Wu CH: **Accomplishments and challenges in literature data mining for biology.** *Bioinformatics* 2002, **18**:1553-1561.
- [CASP]. CASP [<http://predictioncenter.llnl.gov/>]
- [Hirschman 1998]. Hirschman L: **The evolution of evaluation: lessons from the message understanding conferences.** *Computer Speech and Language* 1998, **12**:281-305.
- [TREC]. TREC [<http://trec.nist.gov/>]
- [Voorhees 2002]. Voorhees EM, Buckland LP (Ed): *J. The Eleventh Text Retrieval Conference (TREC 2002): 2002.* NIST Special Publication 500-XXX, Gaithersburg, Maryland, [http://trec.nist.gov/pubs/trec11/t11_proceedings.html]

- [Yeh 2003]. Yeh AS, Hirschman L, Morgan AA: **The Evaluation of text data mining for database curation: lessons learned from the KDD challenge cup.** *Bioinformatics* 2003, **19**:i331-i339.
- [BioCreAtIvE 2004]. *BioCreAtIvE Workshop Handouts*, Granada, Spain, March 2004.
[\[http://www.pdg.cnb.uam.es/BioLINK/workshop_BioCreative_04/handout/index.html\]](http://www.pdg.cnb.uam.es/BioLINK/workshop_BioCreative_04/handout/index.html)
- [Medline]. **Medline** [<http://www.ncbi.nlm.nih.gov/PubMed/>]
- [Tanabe 2004] Tanabe L, Xie N Thom LH, Matten W, Wilbur WJ: **GENETAG: A Tagged Corpus for Gene/Protein Named Entity Recognition.** This volume.
- [Tamames 2004]. Tamames J: **Text Detective: BioAlma's gene annotation tool.** *BioCreAtIvE Workshop Handouts*, Granada, Spain, March 2004.
- [Dingare 2004]. Dingare S, Finkel J, Manning C, Nissim M, Alex B: **Exploring the Boundaries: Gene and Protein Identification in Biomedical Text.** *BioCreAtIvE Workshop Handouts*, Granada, Spain, March 2004.
- [Kinoshita 2004]. Kinoshita S, Ogren P, Cohen KB, Hunter L: **Entity identification in the molecular biology domain with a stochastic POS tagger: the BioCreative task.** *BioCreAtIvE Workshop Handouts*, Granada, Spain, March 2004.
- [Zhou 2004]. Zhou GD, Shen D, Zhang J, Su J, Tan SH, Tan CL: **Recognition of Protein/Gene Names from Text using an Ensemble of Classifiers and Effective Abbreviation Resolution.** *BioCreAtIvE Workshop Handouts*, Granada, Spain, March 2004.
- [McDonald 2004]. McDonald R, Pereira F: **Identifying Gene and Protein Mentions in Text Using Conditional Random Fields.** *BioCreAtIvE Workshop Handouts*, Granada, Spain, March 2004.
- [Yeh 2000]. Yeh A: **Comparing two trainable grammatical relations finders.** In *Proceedings of the 18th International Conference on Computational Linguistics (COLING 2000): 31 July – 4 August 2000; Saarbrueken.* 2000: 1146-1150.
- [Crim 2004]. Crim J, McDonald R, Pereira F: **Automatically Annotating documents with Normalized Gene Lists.** *BioCreAtIvE Workshop Handouts*, Granada, Spain, March 2004.
- [Hachey 2004]. Hachey B, Nguyen H, Nissim M, Alex B, Grover C: **Grounding Gene Mentions with Respect to Gene Database Identifiers.** *BioCreAtIvE Workshop Handouts*, Granada, Spain, March 2004.
- [Krymolowski 2004]. Krymolowski Y, Alex B, Leidner JL: **BioCreative Task 2.1: The Edinburgh-Stanford system.** *BioCreAtIvE Workshop Handouts*, Granada, Spain, March 2004.
- [Joachims 1999]. Joachims T: **Transductive Inference for Text Classification using Support Vector Machines.** In *Proceedings of the 16th International Conference on Machine Learning (ICML-99).* 1999.
- [Hirschman 2002B]. Hirschman L, Morgan A, Yeh A: **Rutabaga by any other name: extracting biological names.** *J. of Biomedical Informatics* 2002, **35**:247-259.
- [LDC]. **Linguistic Data Consortium** [<http://wave ldc.upenn.edu>]
- [Marsh 1998]. Marsh E, Perzanowski D: **MUC-7 Evaluation of IE Technology: Overview of Results.**
[\[http://www.itl.nist.gov/iaui/894.02/related_projects/muc/\]](http://www.itl.nist.gov/iaui/894.02/related_projects/muc/)

Figures

Figure 1 - Balanced F-scores of the 40+4 submissions

Figure 2 - Precision versus recall of the 40+4 submissions

Figure 3 - Percent of names of a given length for BioCreAtIvE task 1A gene names and MUC-6 organization names

Tables

Data Set	Sentences	Gene Mentions
training	7500	9000
(development) test	2500	3000
(final) test	5000	6000

Table 1 - Data set size

	Balanced F-score		Recall		Precision	
	open	closed	open	closed	open	closed
High	83%	83%	84%	85%	86%	86%
Quartile 1	81%	80%	81%	79%	83%	81%
Median (Q2)	78%	74%	74%	72%	80%	72%
Quartile 3	67%	59%	70%	62%	72%	53%
Low	25%	16%	42%	36%	17%	11%

Table 2 - F-score, recall and precision quartiles for the 40 official submissions