# Biblio-MetReS for user-friendly mining of genes and biological processes in scientific documents

Abstract One way to initiate the reconstruction of molecular circuits is by using automated text-mining techniques. Developing more efficient methods for such reconstruction is a topic of active research, and those methods are typically included by bioinformaticians in pipelines used to mine and curate large literature datasets. Nevertheless, experimental biologists have a limited number of available user-friendly tools that use text-mining for network reconstruction and require no programming skills to use. One of these tools is Biblio-MetReS. Originally, this tool permitted an on-the-fly analysis of documents contained in a number of web-based literature databases to identify co-occurrence of proteins/genes. This approach ensured results that were always up-to-date with the latest live version of the databases. However, this "up-to-dateness" came at the cost of large execution times. Here we report an evolution of the application Biblio-MetReS that permits constructing co-occurrence networks for genes, GO process, Pathways, or any combination of the three types of entities and graphically represent those entities. We show that the performance of Biblio-MetReS in identifying gene co-occurrence is as least as good as that of other comparable applications (STRING and iHOP). In addition, we also show that the identification of GO processes is on par to that reported in the latest BioCreAtIvE challenge. Finally, we also report the implementation of a new strategy that combines on-the-fly analysis of new documents with preprocessed information from documents that were encountered in previous analyses. This combination simultaneously decreases program run time and maintains "up-to-dateness" of the results. **Availability:** http://metres.udl.cat/index.php/downloads **Contact: metres.cmb@gmail.com**

2   Anabel Usie[1,2], Hiren Karathia[1], Ivan Teixidó[2], Rui Alves[1,*] and Francesc Solsona[2,*]

3   [1]Department of Basic Medical Sciences, Edifici Recerca Biomedica I, Universitat de
4   Lleida & IRBLleida, Av Rovira Roure 80, 25198, Lleida, Spain.

5   [2]Department of Industrial & Informatics Engineering, Escola Politècnica Superior,
6   Universitat de Lleida, Av Jaume II 69, 25001, Lleida, Spain.

7   [*] Corresponding authors
8   e-mails: ralves@cmb.udl.cat, francesc@diei.udl.cat

## Introduction

10  The reconstruction of molecular circuits is an important research goal in the biological

11  sciences. One of the ways to achieve that reconstruction starts with the use of

12  automated text-mining techniques to identify networks of genes that co-occur in

13  scientific documents. Subsequent human curation of these co-occurrence networks can

14  then lead to accurate circuit reconstruction.

15      In this process, the automated identification of co-occurrence gene networks is

16  crucial because databases of relevant scientific documents contain many more entries

17  than those that can be manually analyzed (Alves and Sorribas 2007; Markowetz and

18  Spang 2007; Arighi, Chohen et al. 2013; Krallinger, Leitner et al. 2013). A gold standard

19  of these databases, MEDLINE, contains more than $19 \times 10^6$ records, with 2000-4000

20  new entries being added each day (NCBI, 2013). Extracting biological information from

21  such large databases requires text-mining methods and tools that are able to

22  automatically integrate and summarize useful biological information across the database

23  records.

24      The development of text-mining methods that enable circuit reconstruction from

25  scientific documents is an area of active development (Camon, Barrell et al. 2005;

26  Hoffmann and Valencia 2005; Huang, Ding et al. 2008; Chen, Liu et al. 2010; Arighi, Lu

27  et al. 2011; Kano, Bjorne et al. 2011; Kim, Pyysalo et al. 2011; Szklarczyk, Franceschini

28  et al. 2011; Usié, Karathia et al. 2011; Bossy, Jourde et al. 2012; Kim, Nguyen et al.

29  2012; Pyysalo, Ohta et al. 2012; Arighi, Chohen et al. 2013; Krallinger, Leitner et al.

30  2013). The performance of those methods for automated identification of the circuits

31  (Camon et al., 2005), of their components (genes/proteins), and of the inter-component

32  relationships, has been systematically evaluated over the last few years, for example

33  through the BioNLP (Kano, Bjorne et al. 2011; Kim, Pyysalo et al. 2011; Bossy, Jourde

34  et al. 2012; Kim, Nguyen et al. 2012; Pyysalo, Ohta et al. 2012) and BioCreAtIvE

35  initiatives (Huang, Ding et al. 2008; Chen, Liu et al. 2010; Arighi, Lu et al. 2011; Wu,

36  Arighi et al. 2012; Arighi, Chohen et al. 2013; Krallinger, Leitner et al. 2013).

37  To briefly summarize, there are three large classes of methods that have been

38  employed for the reconstruction of co-occurrence gene networks: dictionary-based

39  methods, morphology-based methods, and context-based methods (Vazquez, Krallinger

40  et al. 2011). Dictionary-based methods rely on matching compiled lists of terms to their

41  appearances in the text of documents (Yang, Lin, & Li, 2008). Morphology-based

42  methods rely on the morphological structure of specific classes of words to single them

43  out in documents (Malouf, 2002; Peng & Schuurmans, 2003). Finally, context-based

44  methods can be divided into Machine Learning or Natural Language Processing

45  techniques: The former identify patterns in the structure of the text that help to recognize

46  the presence of the relevant entities in documents; the later draw on our knowledge

47  about the grammar and syntax rules of natural languages to recognize those entities.

48  These three general approaches can be combined in order to improve NER (for

49  example see (Arighi, Chohen et al. 2013; Krallinger, Leitner et al. 2013) and references

50  therein).

51  In general, methods participating in evaluations such as BioCreAtIvE or BioNLP are

52  implemented in tools that can be included in web-services and assist curators in the

53  maintenance of large databases of biological knowledge. Examples of this are given in

54  (Arighi, Chohen et al. 2013; Krallinger, Leitner et al. 2013). In most cases, using these

55  methods and tools requires that one becomes an expert computer user and learns how

56  to program.

57  Experimental scientists that are interested in being users of, without becoming experts

58  in, text-mining methods to directly reconstruct networks of gene co-occurrence for their

59  genes of interest in scientific documents have a much smaller set of available user-

60  friendly tools. The first that became available was iHOP (Hoffmann and Valencia 2005),

61  which was later joined by STRING (Franceschini, Szklarczyk et al. 2013). These user

62  friendly and intuitive web applications allow anyone to reconstruct the network of co-

63  occurrences contained in Medline abstracts and/or Pubmed documents. Users of these

64  applications face two important limitations. First, the applications rely on preprocessed

65  versions of the Medline/Pubmed databases, which means that searches are fast but

66  results are always out of date. Second, the coverage of full text documents by the

67  applications is, at best, limited.

68  Recognizing these limitations, Biblio-MetReS (**Biblio**metric **Met**abolic network

69  **Re**construction **S**erver (Usie, Karathia et al. 2011)) was implemented for the same target

70  audience as STRING or iHOP, but relying on two differential features with respect to

71  those applications. The first was that it would search databases and analyze documents

72  on the run, thus providing the users with the most up-to-date results available on the

73  web. The second was that full text documents were also analyzed, as were other

74  databases besides Medline/Pubmed. These two features made Biblio-MetReS

75  significantly slower than STRING and iHOP.

76  Here, we report an evolution of the application Biblio-MetReS that permits constructing

77  co-occurrence networks for genes, GO process, Pathways, or any combination of the

78  three types of entities and graphically represent those entities. No other user-friendly

79  application that we are aware of simultaneously allows the type of mixed analysis and

80   graphical representation afforded by Biblio-MetReS. We show in a comparative analysis

81   that the performance of Biblio-MetReS in identifying gene co-occurrence is as least as

82   good as that of other comparable applications (STRING and iHOP). In addition, we also

83   show that the identification of GO processes is on par to that reported in the latest

84   BioCreAtIvE challenge (Arighi, Chohen et al. 2013). Finally, we also report the

85   implementation of a new strategy that combines on-the-fly analysis of new documents

86   with preprocessed information from documents that were encountered in previous

87   analyses. This combination simultaneously decreases program run time and maintains

88   "up-to-dateness" of the results.


89   ## Methods

90   ### Organism Selection

91   Biblio-MetReS is organism centric. Users must select their organism of interest from a

92   list of more than 1200 organisms with fully sequenced and annotated genomes before

93   starting any search. They must also decide whether they want to perform GO and/or

94   pathway term co-occurrence analysis. After these decisions are made, the program

95   loads the genes for the organism from the program's central database. If selected, terms

96   from the GO classification, KEGG and/or Panther pathways are also loaded into the

97   application's front end.

98   ### Document Analysis

99   To analyze networks of co-occurrence Biblio-MetReS needs users to select at least one

100   gene from their organism of interest and at least one database in which to perform

101   document analysis. The seed list of genes is then used to identify relevant documents in

102   the selected database(s) of documents. The text in the flagged documents is then

103   analyzed to identify additional genes from the organism of interest, as well as GO and/or

104    Pathway terms. This allows users to identify co-occurrence among GO/Pathway entities

105    and between GO/Pathway entities and gene/protein entities in sentences, paragraphs or

106    documents. For further description of this procedure, please see section 1 of the

107    supplementary materials.

108    Biblio-MetReS uses exact matching of gene names to an internal dictionary of

109    synonyms to identify co-mentions of genes/proteins in the text of scientific records. The

110    gene synonyms are those officially defined by NCBI. For any given gene, all synonyms

111    are searched for in the text of flagged documents. Similarly, Biblio-MetReS uses exact

112    matching of GO terms to an internal dictionary of synonyms defined by the gene

113    ontology consortium (Gene Ontology 2013) to identify co-mentions of GO terms in the

114    text of scientific records. The same exact matching is done to identify mentions to

115    entities from the complete joint sets of KEGG (Kotera, Hirakawa et al. 2012) and

116    Panther pathways (Mi, Lazareva-Ulitsky et al. 2005).

117    Co-occurrence of any two terms is analyzed at three levels. First, all possible

118    pairs of terms are searched for in each document as a whole. Then, each document is

119    divided into paragraphs and the pairs of terms identified in the document are searched

120    for within each paragraph. Finally, each paragraph is divided into single sentences and

121    the pairs of terms identified in that paragraph are searched for within each sentence.

122    Within each level, the distance between each term in the pair is not taken into account.

123    The database containing the organisms and their gene names, as well as the GO

124    terms, is updated every three months using information compiled automatically from

125    NCBI and GO.

126    **Calculating the significance of term co-occurrences**

127 To attribute statistical significance to the co-occurrence of a pair of genes, pathway

128 terms, GO terms, gene-pathway terms, gene-GO terms or GO-pathway terms, we

129 calculate several metrics. First, we measure how frequently the different pairs co-occur

130 in sentences, paragraphs and/or documents. We then take the odds ratio of the

131 frequency of occurrences in the first two categories with respect to that of the third. The

132 closer to one these odds ratios are, the more frequent it is that both genes are

133 mentioned only in the same sentences or paragraphs of a document, rather than

134 appearing haphazardly in different sections of the text.

135       Second, we calculate how much information we gain by having two terms, $T_i$ and

136 $T_j$, co-occur, when compared to the individual occurrences of the terms. To estimate this

137 we use information theory. The individual probability of occurrence of a term is denoted

138 as p(Ti) and it is formally defined as p(Ti)=a/n, where a is the number of documents

139 where Term $T_i$ appears, and n is the total number of documents. The joint probability of

140 co-occurrence of two terms, $p(T_i, T_j)$, is defined as $p(T_i,T_j)$=b/n, where b is the number of

141 documents where terms $T_i$ and $T_j$ simultaneously appear, and n is the total number of

142 documents. The mutual information, $MI(T_i, T_j)$, is then calculated as follows:

143 $MI(T_i,T_j)=p(T_i,T_j)\times\log(p(T_i,T_j)/(p(T_i)p(T_j)))$

144       Finally, and in order to attribute some form of statistical significance to the co-

145 occurrence of a pair of terms, we do as follows. Consider a set of n sentences

146 (paragraphs, documents) [1 ..., n]. For a given term k define

147 $y_{i,k}=\begin{cases} 1 & term\ k\ occurs \in sentence\ (paragraph,\ document)\ i \\ 0 & otherwise \end{cases}$   Now, for terms k1 and k2 define

148 $\varphi_{k1,k2}=y_{i,k1}\times y_{i,k2}$

149    which has value 1 when both terms co-occur and 0 otherwise.

150       Both these variables have a Bernoulli distribution. If the occurrence of terms k1

151    and k2 is independent, then $p(\varphi_{k1,\ k2}) = p(yk1)\ p(yk2)$ would be expected, where $p(y_{k,\cdot})$ is

152    the relative frequency of occurrence of term $y_{k,\cdot}$ and $p(\varphi_{k1,\ k2})$ is the relative frequency of

153    co-occurrence of terms $k_1$ and $k_2$ in the total number n of sentences (paragraphs,

154    documents). Then, a Pearson statistic can be used to test for independence of

155    occurrence between k1 and k2 by comparing the observed frequencies, $n1 = n \times p(\varphi_{k1,\ k2})$

156    and $n2 = (1-p(\varphi_{k1,\ k2})) \times n$, with the expected frequencies under the null hypothesis of

157    independence, which would be $m1 = n \times p(yk1) \times p(yk2)$ and $m2 = n \times (1-p(yk1) \times p(yk2))$.

158    The Pearson statistic is computed as follows

$$X^2 = \sum_{i=1}^{2} \frac{(\text{\textit{¿}} - mi)^2}{mi}$$

159    This statistics follows a chi-square distribution with one degree of freedom, i.e. $\square^2_1$

160    $\sim X^2$; hence, the p-value can be calculated as $p = Pr\ \square^2_1 > \chi^2)$ to assess whether the

161    observed co-occurrence is higher than the one expected by pure chance.

162 **Precompilation strategy**

163    Biblio-MetReS v2.0 implements a precompilation approach that works in the following

164    way. Any search done will identify a given number of documents in the database(s)

165    selected by the user(s), for an organism of interest. If a given document has not been

166    found in any previous search by any user, in the context of that organism, Biblio-MetReS

167    will analyze it as described in section 1 of supplementary materials and all information

168 contained in that document and relevant for the analysis will be stored in a central

169 database (see section 3 of supplementary materials for detailed information). If a given

170 document has been previously found by any user, its information will be directly

171 accessed from our central database, and the document will not be reanalyzed. This

172 means that newly found documents are mined on the fly by the program to find and

173 count mentions of relevant entities, while mentions in documents that have been

174 previously found are simply looked up in our central database.

## Results
### Biblio-MetReS and Biblio-MetReS Player

177 Biblio-MetReS v2.0 can be used to identify genes/proteins from more than 1200 different

178 organisms in records stored in a variety of databases. Users download the application

179 and run it locally. A functioning internet connection and a local copy of JAVA are required

180 for the program to work. Upon starting the program, users login to the central Biblio-

181 MetReS database and choose which organism they are interested in and whether they

182 want to search only for co-occurrence of genes/proteins or if they also want to include

183 biological pathways and/or GO biological processes in the analysis. Once this choice is

184 made, the program loads the necessary information from the central database. Taking

185 this approach, instead of including all the data locally in the program installation, permits

186 making an application that is much smaller in size and needs less RAM to function

187 properly. Subsequently, users select the source of documents that they want to analyze,

188 as well as the genes/proteins and/or pathways/GO biological processes that they want

189 to search for. They can also include their own handpicked list of processes to be

190 searched. Once the search is launched, Biblio-MetReS will identify documents in the

191 relevant databases that contain mentions to the relevant search items. After identifying

192    these documents, the application fully analyzes them to identify mentions for any

193    additional gene or Pathway/GO biological process via a dictionary matching approach.

194    The co-occurrence of the different entities is analyzed at the level of the whole

195    document, of individual paragraphs and of individual sentences, and the significance of

196    this co-occurrence is calculated as described in (Usié et al., 2011). The information that

197    is relevant for the co-occurrence calculations is stored in the central Biblio-MetReS

198    database. Any subsequent searches that identify the same document will not reanalyze

199    it; instead, these numbers are directly retrieved from that database. Once the analysis is

200    complete, the users can visualize it in graphical and textual form. Links to the

201    documents and sentences where co-occurrences are found are provided. Graphical

202    visualization of the results can be done in different ways. Users can create graphs for

203    the global co-occurrences network and for gene- or pathway/process-centric co-

204    occurrences at the document/paragraph and sentence levels. Significance and Mutual

205    Information of each co-occurrence is also provided in tabular form. The graphical

206    representation of the networks is automatically stored in local xml files. These files can

207    be opened using a small app, Biblio-MetReS Player, which can be downloaded from the

208    Biblio-MetReS website. This permits reviewing previously obtained networks without

209    having to redo the search. All this process is summarized in Figure 1.

210    **Comparative performance and benchmarking of new types of entities**

211    The benchmarking of the program and its improvements with respect to version 1.0 was

212        carried out using four organisms of interest: *Saccharomyces cerevisiae, Homo*

213        *sapiens, Escherichia coli* and *Drosophila melanogaster*. For each organism we

214        used as a search seed a set of genes belonging to Glycolysis, Lysine metabolism

215        and RNA processing pathways. The genes were chosen to reproduce the

216        experiments reported in (Usié, Karathia et al. 2011). Details are shown in

217        Supplementary Table S1. We benchmarked three different aspects of Biblio-

218        MetReS. First, we benchmarked the comparative identification of genes between

219        Biblio-MetReS, iHOP, and STRING, given that these three applications have

220        similar target audiences. Second, we benchmarked the ability of Biblio-MetReS to

221        identify biological processes/pathways. Finally, we benchmarked the

222        improvements in Biblio-MetReS run time made by implementing the combined

223        pre-processing/live analysis strategy.

224        Benchmarking the comparative identification of genes between Biblio-MetReS,

225  iHOP, and STRING was done in the following way. We used the genes and organisms

226  described in Supplementary Table S1 to interrogate independently Biblio-MetReS, iHOP,

227  and STRING. For each of the three applications, the complete set of results for each

228  gene from the same pathway were pooled together for analysis. The results from

229  STRING were further filtered to eliminate all genes and interactions that were not

230  literature based. Supplementary Table S2 in supplementary materials compares the

231  results for the three applications. In summary, Biblio-MetReS find the largest number of

232  genes, followed by STRING, and iHOP. The number of genes found by STRING and

233  Biblio-MetReS are of the same order of magnitude, while iHOP finds between one and

234  two orders of magnitude less genes. This result derives from the fact that iHOP analyzes

235  only Medline abstracts, while Biblio-MetReS and STRING analyze the full text of

236  Pubmed publications, in addition to the Medline abstracts. This is confirmed by the fact

237  that, when Biblio-MetReS is run only to analyze Medline abstracts, it finds a similar

238  number of genes as iHOP (data not shown). As was observed in Supplementary Table 3,

239  the genes found by each of the three applications for the same experiments only

240  partially overlap and this is explained by the different datasets analyzed by each of the

241  programs and by partially different dictionaries of gene synonyms (Usié et al., 2011).

242    Neither STRING nor iHOP permit identifying GO terms and their associations to

243    genes. Therefore we cannot perform experiments that are similar to the comparative

244    benchmarking experiments described above. In light of this, benchmarking of Biblio-

245    MetReS' ability to identify biological processes/pathways was done in the following way.

246    To perform the GO identification benchmark we used the test and development sets of

247    the BioCreAtIvE IV GO task corpus [3]. We used Biblio-MetReS dictionary matching

248    approach to identify GO terms in the non-annotated documents and then analyzed the

249    corresponding annotated documents. We found that Biblio-MetReS identified 100% of

250    the annotated GO terms in both sets (2963 terms in the training set and 2243 terms in

251    the development set). Biblio-MetReS also identifies 2259 additional GO terms in the

252    development set and 2119 additional GO terms in the training set. For the purpose of

253    our testing these terms must be considered false positive. Taking this into account, the

254    precision in the development set is 50%, while in the training set it increases to 58%.

255    The F-score performance of Biblio-MetReS is 33% in the development set and 37% in

256    the training set, which is on par with the best approaches presented in the lattest

257    BioCreAtIvE IV challenge (Mao et al., 2013).

258    Benchmarking run time was done in two ways. First, we search only the Pubmed

259    database. Second, we search by selecting all literature databases available in Biblio-

260    MetReS. In both tests we used all the seed genes from Supplementary Table S1. Each

261    seed is used by Biblio-MetReS as a query search. This query search is launched twice.

262    When the first search is done there are no preprocessed documents in Biblio-MetReS'

263    database. The information in documents is analyzed on-the-fly and stored. Then the

264    searches are repeated, now with the documents stored in Biblio-MetReS' database. This

265    allows us to estimate the percentage of run-time saved by preprocessing the

266     documents. The results are shown in Figure 2 and in Supplementary Table S2 of the

267     supplementary materials. On average we get decreases in run time of more than 90%.


268     ## Discussion

269     Here we present the new version of Biblio-MetReS, a user friendly tool for the

270     identification of gene/protein co-occurrence networks in scientific documents. The

271     major changes with respect to version 1.0 have to do with the search and

272     analysis process of the documents, which can now be up to 95% faster than in

273     the previous version. In addition, the tool now also searches for co-occurrences

274     of biological processes and pathways, to help users to more easily establish the

275     biological circuits in which their genes of interest may be involved in.

276     The methods used by the application to identify genes and proteins, as well as

277     biological processes and pathways, in the documents are dictionary-based. These

278     methods perform on par with iHOP and STRING for gene and protein identification and

279     with the best BioCreAtIvE methods for biological process identification (see

280     Supplementary Materials).

281     Taken together, the new application further facilitates the identification of

282     functional relationships between proteins and aids in identifying the biological processes

283     and circuits in which those proteins may be involved. Although GO term search has

284     been implemented in several literature search tools (see for example (Doms and

285     Schroeder 2005; Plake, Royer et al. 2009), among others), no other user-friendly tool

286     permits simultaneous graphical reconstruction of networks of co-occurrence between

287     genes, GO terms, and Pathway terms.

288    As is demonstrated by the BioCreAtIvE challenge (Mao et al., 2013), the problem

289    of identifying entities in scientific texts is far from solved. Although Biblio-MetReS aims at

290    giving non-expert users the possibility of performing such identification and use that

291    identification to extract biological knowledge, there is much room for improvement. We

292    are implementing an offline system to automatically search, analyze, and store

293    information about gene/protein and pathway/biological processes co-occurrences in the

294    documents. This will contribute to decrease the dependence of Biblio-MetReS on the

295    users and their searches to preprocess information and make searches faster.

296    ## Funding

302    ## References

303    Alves, R. and A. Sorribas (2007). "In silico pathway reconstruction: Iron-sulfur cluster biogenesis
304        in Saccharomyces cerevisiae." BMC systems biology **1**.

305    Arighi, C. N., K. B. Chohen, et al. (2013). Proceedings of the Fourth BioCreative Challenge
306        Evaluation Workshop. Washington, Biocreative Challenge.

307    Arighi, C. N., Z. Lu, et al. (2011). "Overview of the BioCreative III Workshop." BMC
308        Bioinformatics **12**(Suppl 8).

309    Bossy, R., J. Jourde, et al. (2012). "BioNLP Shared Task--The Bacteria Track." BMC
310        Bioinformatics **13 Suppl 11**: S3.

311    Camon, E. B., D. G. Barrell, et al. (2005). "An evaluation of GO annotation retrieval for
312        BioCreAtIvE and GOA." BMC Bioinformatics **6 Suppl 1**.

313 Chen, Y., F. Liu, et al. (2010). "BioLMiner and the BioCreative II.5 challenge." BMC
314     Bioinformatics **11**(Suppl 5).

315 Doms, A. and M. Schroeder (2005). "GoPubMed: exploring PubMed with the Gene Ontology."
316     Nucleic Acids Res **33**(Web Server issue): W783-W786.

317 Franceschini, A., D. Szklarczyk, et al. (2013). "STRING v9.1: protein-protein interaction
318     networks, with increased coverage and integration." Nucleic Acids Res **41**(Database
319     issue): D808-815.

320 Gene Ontology, C. (2013). "Gene Ontology annotations and resources." Nucleic Acids Res
321     **41**(Database issue): D530-535.

322 Hoffmann, R. and A. Valencia (2005). "Implementing the iHOP concept for navigation of
323     biomedical literature." Bioinformatics **21**: 252-258.

324 Huang, M., S. Ding, et al. (2008). "Mining physical protein-protein interactions from the
325     literature." Genome Biol **9**(Suppl 2).

326 Kano, Y., J. Bjorne, et al. (2011). "U-Compare bio-event meta-service: compatible BioNLP event
327     extraction services." BMC Bioinformatics **12**: 481.

328 Kim, J.-D., S. Pyysalo, et al. (2011). Overview of BioNLP Shared Task.

329 Kim, J. D., N. Nguyen, et al. (2012). "The Genia Event and Protein Coreference tasks of the
330     BioNLP Shared Task 2011." BMC Bioinformatics **13 Suppl 11**: S1.

331 Kotera, M., M. Hirakawa, et al. (2012). "The KEGG databases and tools facilitating omics
332     analysis: latest developments involving human diseases and pharmaceuticals." Methods
333     in molecular biology (Clifton, N.J.) **802**: 19-39.

334 Krallinger, M., F. Leitner, et al. (2013). Proceedings of the Fourth BioCreative Challenge
335     Evaluation Workshop. Washington, BioCreative Challenge.

336 Markowetz, F. and R. Spang (2007). "Inferring cellular networks--a review." BMC Bioinformatics
337     **8 Suppl 6**.

338 Mi, H., B. Lazareva-Ulitsky, et al. (2005). "The PANTHER database of protein families,
339     subfamilies, functions and pathways." Nucleic Acids Res **33**(Database issue): D284-288.

340 Plake, C., L. Royer, et al. (2009). "GoGene: gene annotation in the fast lane." Nucleic Acids Res
341     **37**(Web Server issue): W300-304.

342 Pyysalo, S., T. Ohta, et al. (2012). "Overview of the ID, EPI and REL tasks of BioNLP Shared
343     Task 2011." BMC Bioinformatics **13 Suppl 11**: S2.

344 Szklarczyk, D., A. Franceschini, et al. (2011). "The STRING database in 2011: functional
345     interaction networks of proteins, globally integrated and scored." Nucleic Acids Res
346     **39**(Database issue): D561-568.

347 Usie, A., H. Karathia, et al. (2011). "Biblio-MetReS: a bibliometric network reconstruction
348     application and server." BMC Bioinformatics **12**: 387.

349    Usié, A., H. Karathia, et al. (2011). "Biblio-MetReS: A bibliometric network reconstruction
350        application and server." BMC Bioinformatics **12**(1).

351    Vazquez, M., M. Krallinger, et al. (2011). "Text Mining for Drugs and Chemical Compounds:
352        Methods, Tools and Applications." Molecular Informatics **30**(6-7): 506-519.

353    Wu, C. H., C. N. Arighi, et al. (2012). "BioCreative-2012 Virtual Issue." Database (Oxford) **2012**:
354        bas049.
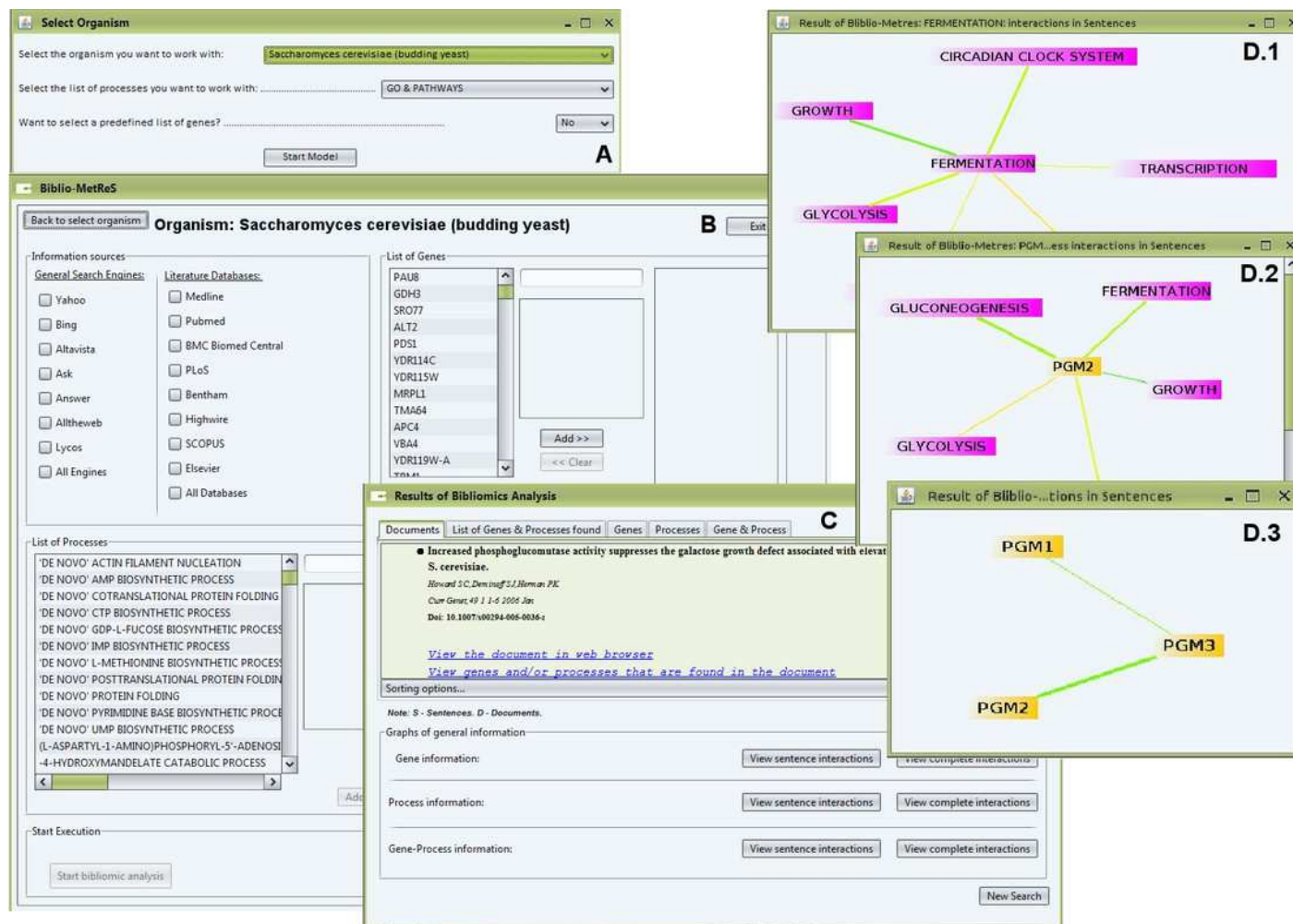
**Figures.**

**Figure 1 - User workflow for Biblio-MetReS.**


**Figure 2 - Effect of preprocessing documents on Biblio-MetReS' run time.**

In brief, genes from three KEGG-defined pathways are used for this test. Panels A.x show experimental results for glycolysis genes. Panels B.x show experimental results for Lysine metabolism genes. Panels C.x show experimental results for RNA processing genes. Three organisms are used in this benchmark. Panels Y.1 show results for *Homo sapiens*, panels Y.2 show results for *Drosophila melanogaster*, panels Y.3 show results for *Escherichia coli,* and panels Y.4 show results for *Saccharomyces cerevisiae*. These pathways and organisms were chosen to remain consistent with the tests performed in (Usié et al., 2011). Searches were done selecting all the databases in the application. Graphs can be interpreted as follows. Light gray bars indicate the run time for Biblio-MetReS when the corresponding gene is searched for the first time. In this case the program has to do a full document analysis on the fly and no information has been preprocessed. Darker gray bars indicate the run time for Biblio-MetReS when the search for the corresponding gene is repeated, and preprocessed information is already present in Biblio-MetReS' central database. The column "All" indicates the run-time for searching all genes in the graph simultaneously, after individual searches for each gene had already been done and results preprocessed and stored.

# Figure 1

User workflow for Biblio-MetReS.

**Figure 1 - User workflow for Biblio-MetReS.**

# Figure 2

Effect of preprocessing documents on Biblio-MetReS' run time.

**Figure 2 - Effect of preprocessing documents on Biblio-MetReS' run time.** In brief, genes from three KEGG-defined pathways are used for this test. Panels A.x show experimental results for glycolysis genes. Panels B.x show experimental results for Lysine biosynthesis genes. Panels C.x show experimental results for RNA degradation genes. Three organisms are used in this benchmark. Panels Y.1 show results for *Homo sapiens*, panels Y.2 show results for *Drosophila melanogaster*, panels Y.3 show results for *Escherichia coli,* and panels Y.4 show results for *Saccharomyces cerevisiae*. These pathways and organisms were chosen to remain consistent with the tests performed in (Usié et al., 2011) . Searches were done selecting all the databases in the application. Graphs can be interpreted as follows. Light gray bars indicate the run time for Biblio-MetReS when the corresponding gene is searched for the first time. In this case the program has to do a full document analysis on the fly and no information has been preprocessed. Darker gray bars indicate the run time for Biblio-MetReS when the search for the corresponding gene is repeated, and preprocessed information is already present in Biblio-MetReS' central database. The column "All" indicates the run-time for searching all genes in the graph simultaneously, after individual searches for each gene had already been done and results preprocessed and stored.