

# 1 ModEx: A text mining system for extracting 2 mode of regulation of Transcription Factor- 3 gene regulatory interaction

4  
5 Saman Farahmand<sup>1,2</sup>, Todd Riley<sup>1,2</sup> and Kourosh Zarringhalam<sup>1,3\*</sup>

6 <sup>1</sup> Computational Sciences PhD program, University of Massachusetts Boston, Boston, USA

7 <sup>2</sup> Department of Biology, University of Massachusetts Boston, Boston, USA;

8 <sup>3</sup> Department of Mathematics, University of Massachusetts Boston, Boston, USA

9  
10 \* Correspondence: [kourosh.zarringhalam@umb.edu](mailto:kourosh.zarringhalam@umb.edu), Department of Mathematics, University of  
11 Massachusetts Boston, 100 Morrissey Boulevard, 02125 Boston, USA

12  
13 **Abstract**— Transcription factors (TFs) are proteins that are fundamental to transcription and regulation of  
14 gene expression. Each TF may regulate multiple genes and each gene may be regulated by multiple TFs.  
15 TFs can act as either activator or repressor of gene expression. This complex network of interactions  
16 between TFs and genes underlies many developmental and biological processes and is implicated in several  
17 human diseases such as cancer. Hence deciphering the network of TF-gene interactions with information  
18 on mode of regulation (activation vs. repression) is an important step toward understanding the regulatory  
19 pathways that underlie complex traits. There are many experimental, computational, and manually curated  
20 databases of TF-gene interactions. In particular, high-throughput ChIP-seq datasets provide a large-scale  
21 map of transcriptional regulatory interactions. However, these interactions are not annotated with  
22 information on context and mode of regulation. Such information is crucial to gain a global picture of gene  
23 regulatory mechanisms and can aid in developing machine learning models for applications such as  
24 biomarker discovery, prediction of response to therapy, and precision medicine. In this work, we introduce  
25 a text-mining system to annotate ChIP-seq derived interaction with such meta data through mining PubMed  
26 articles. We evaluate the performance of our system using the gold standard small scale manually curated  
27 TRUSST database. Our results show that the method is able to accurately extract mode of regulation with

28 F-score 0.77 on TRRUST curated interaction and F-score 0.96 on intersection of TRUSST and ChIP-  
29 network. We provide a HTTP REST API for our code to facilitate usage.

30 **Availability:** Source code and datasets are available for download on GitHub:

31 <https://github.com/samanfrm/modex>

32 **HTTP REST API:** <https://watson.math.umb.edu/modex/>

33

34 **Index Terms**— Text mining, information extraction, name entity recognition, biological NLP, biomedical  
35 literature, gene regulatory network, mode of regulation.

36

## 37 1. INTRODUCTION

38 Gene regulatory networks are essential in many cellular processes, including metabolism, signal  
39 transduction, development, and cell fate [1]. At the transcriptional level, regulations of genes are  
40 orchestrated by concerted action between Transcription Factors (TFs), histone modifies, and distal *cis*-  
41 regulatory elements to finely tune and modulate expression of genes. Sequence-specific Transcription  
42 Factors play a key role in regulating gene transcription at the transcriptional level. They bind specific DNA  
43 motifs to regulate promoter activity and either enhance (activate) or repress (inhibit) expression of the  
44 genes. Deciphering transcriptional regulatory networks is crucial for understanding cellular mechanisms  
45 and response at a molecular level and can shed light on molecular basis of complex human diseases [2–5].  
46 Moreover, knowledge on interactions between genes and biomolecules is an essential building block in  
47 several pathway inference and gene enrichment analysis methods that aim to annotate an altered set of  
48 transcripts with biological function [6,7]. There are several sources of publicly available biomolecular  
49 interactions, including, signaling pathways, metabolic pathways, and protein-protein interactions [8–10].  
50 Databases of transcriptional regulatory network include JASPAR [11], the Open Regulatory Annotation  
51 database (ORegAnno) [12], Swissregulon [13], the Transcriptional Regulatory Element Database (TRED)  
52 [14], the Transcription Regulatory Regions Database (TRRD) [15], TFactS [16], TRRUST [17], and the  
53 Human Transcriptional Regulation Interactions database (HTRIdb) [17]. These databases include  
54 biologically validated and computationally inferred gene regulatory interactions and have been assembled  
55 with a variety of approaches, including reverse engineering approaches based on high-throughput gene  
56 expression experiments [18–20], text mining approaches [21], and manual curation [22]. These databases  
57 are a valuable source of gene regulatory information, however, there are several constraints that limit their  
58 usability. For example, databases of computationally predicted and expression-driven interactions are

59 typically very noisy. Importantly, the majority of the databases do not report the direction of regulation (up  
60 or down) - which is crucial to understanding the functional behavior of the cell.

61 A high-throughput experimental approach for identifying regulatory interaction is chromatin  
62 immunoprecipitation followed by sequencing (ChIP-seq). In ChIP-seq methodologies, antibodies that  
63 recognizes a specific TF are used to pull down attached DNA for sequencing. The ENCODE (Encyclopedia  
64 of DNA Elements) consortium [23] has produced vast amount of publicly available high-throughput ChIP-  
65 seq experiments that are processed and deposited into databases such as GTRD [24] and ChIP-Atlas [25]  
66 (>35,000 experiments). These databases can be utilized to construct a high coverage transcriptional  
67 regulatory network. Although these interactions are experimentally derived, they are still very noisy as the  
68 experiments are performed under different conditions and in different cell lines. Moreover, ChIP-seq does  
69 not provide direct information on mode of regulation.

70 The most reliable source of regulatory information is obtained by manual curation of peer-reviewed  
71 biomedical literature by domain experts and can be considered as the gold standard. Commercial vendors  
72 such as Ingenuity ([www.ingenuity.com](http://www.ingenuity.com)) offer pathway inference analysis algorithms on such manually  
73 curated networks of regulatory interactions. There are also public sources of curated causal gene regulatory  
74 interactions, such as TRRD, TRED, TFactS, and TRRUST. However, these databases are very small in  
75 their scope, covering only a fraction of TF-gene interactions. Overall, manual curation of biomedical  
76 literature is very time consuming, requires extensive resources, and does not scale to the pace at which  
77 biomedical literature is expanding [26]. For this reason, biomedical text mining has been extensively used  
78 to automate the process of biomolecular relation extraction from the literature. As literature lacks  
79 standardized representation of text, automatic routines for information extraction from textual context is  
80 very challenging [27].

81 There is a vast amount of literature on text-mining for various application [27]. Essential text mining steps  
82 for biomedical relation extraction can be divided into 3 steps: (1) information retrieval (IR), (2) entity  
83 recognition (ER), and (3) information extraction (IE). Together, they can be utilized to identify and extract  
84 specific biological knowledge from literature [28,29].

85 IR tools retrieve relevant text information from articles, abstracts, paragraphs, and sentences corresponding  
86 to subject of interest. A popular IR approach for biomedical application is the use of Boolean model logic  
87 (AND/OR) for extracting relevant information containing specific biological terms [27]. Prominent IR tools  
88 that use the Boolean logic model are iHOP [30] and PubMed. PubMed utilizes human-indexed MeSH terms  
89 to reduce the search space and retrieve relevant abstracts containing user specified keywords. iHOP builds  
90 on PubMed and is able to detect co-occurrence of terms. A limitation of iHOP is that the terms must occur  
91 in the same sentence.

92 After the IR step, ER must be used to identify relation between biological entities. This is a challenging  
93 step as entity names are not unique. Therefore, ER tools must take textual context into consideration to  
94 accurately detect entities. For example, gene names may have different variations in ortho-graphical  
95 structure (e.g. ABL1, Abl1, Abl-1) or multiple synonyms (e.g. ABL1, ABL, CHDSKM, Abelson tyrosine-  
96 protein kinase 1). ER methods, typically divide the task into two steps, (1) identify the entities and their  
97 location in the context, and (2) assign unique identifiers to the entities [27]. Fortunately, multiple  
98 terminological databases, such as Gene Ontology [31], UMBLS [32], BioLexicon [32], and Biothesaurus  
99 [32] provide information on biological entities and name variations and can be used to detect biological  
100 entities such as genes or proteins [33–35].

101 Lastly, Relation Extraction (RE) is an IR tasks for extracting pre-defined facts relating to an entity or entities  
102 in the text [36]. In biomedical domain, multiple RE methods have been developed to extract information  
103 relating to genes [17], such as Mutation-Disease associations, protein-protein interaction [37,38], pathway  
104 curation [39], gene methylation and cancer relation [40], biomolecular events [41], metabolic reactions [42]  
105 and gene-gene interactions [43]. For gene regulatory networks, which is the focus of this paper, the RE  
106 system must detect and extract a *causal* relation between a protein and a gene (e.g., A regulated B). This  
107 task is very complex, even for human experts [44]. To illustrate, consider the causal relation “aatf  
108 upregulates c-myc” that should be deduced from the following sentence: “down-regulation of c-myc gene  
109 was accompanied by decreased expressions of c-myc effector genes coding for htert, bcl-2, and aatf” [45].  
110 Extracting a positive regulatory interaction between aatf and c-myc is quite challenging using simple RE  
111 methods. For example, the RE method, may naively annotated the interaction as negative because of the  
112 keyword “decreased”. However, by taking “down-regulation” into account, the RE method would able to  
113 correctly extract a positive regulation from this sentence.

114 Construction of a causal transcriptional regulatory network by traditional means of text mining is hampered  
115 by these challenges and as a result, fully automated text-mining based models are limited in their scope and  
116 accuracy [27]. Combining experimentally-derived regulatory interactions from high-throughput sources  
117 with text-mining approaches can bridge the gap between the two approaches and address their  
118 shortcomings.

119 In this work, we present a hybrid model ModEx, to mine the biomedical literature to extract and annotate  
120 causal transcriptional regulatory interactions derived from high-throughput ChIP-seq datasets. Specifically,  
121 we applied our text-mining method to extract experimental TF-gene relations reported in ChIP-Atlas  
122 (assembled from all publicly available ChIP-seq experiments) from biomedical literature and annotated the  
123 retrieved interaction with meta-data, such as full supporting sentences, PubMed ID, and importantly mode  
124 of regulation, which is missing from ChIP-seq data. It is important to note that our approach bypasses

125 several of the challenges of fully automated text-mining methods, including query translation for a  
126 particular interaction, relevant citation retrieval, entities recognition and regulatory annotation. We  
127 evaluated the performance of our model using the TRRUST network [22], which contains 9,396 manually  
128 curated regulatory interactions. Our model was able to achieve an F-score 0.77 in retrieving and annotating  
129 the TRUSST network. When applied to TRRUST reported interaction that are also present in ChIP-seq  
130 data, the method achieved an F-score of 0.96.

131

## 132 **2. MATERIALS AND METHODS**

133 We begin by a brief overview of our text mining approach for extracting and annotating ChIP-seq derived  
134 TF-gene interactions with meta-data. We acquired all citations from PubMed abstracts by submitting  
135 queries to the database with appropriate Boolean logic regarding entities and their synonyms. State-of-the-  
136 art external ER systems such as PubTator [46] and beCAS [47] along with our ER system were utilized to  
137 obtain a list of biological entities in the abstract. We then used the Stanford dependency parser [48] to  
138 extract dependencies on different sentences and merge the parse trees into a parse graph. The major  
139 advantage of this parse graph is its potential to identify long-range dependency relations across sentence  
140 boundaries. Candidate relations were created by extracting subtrees connecting pairs of entities from the  
141 dependency graph. Finally, we extracted the mode of regulation based on two sets of manually-annotated  
142 positive and negative causal categories (consisting > 100 verbs and their inflections). In the subsequent  
143 sections we describe the details of our text-mining system.

144

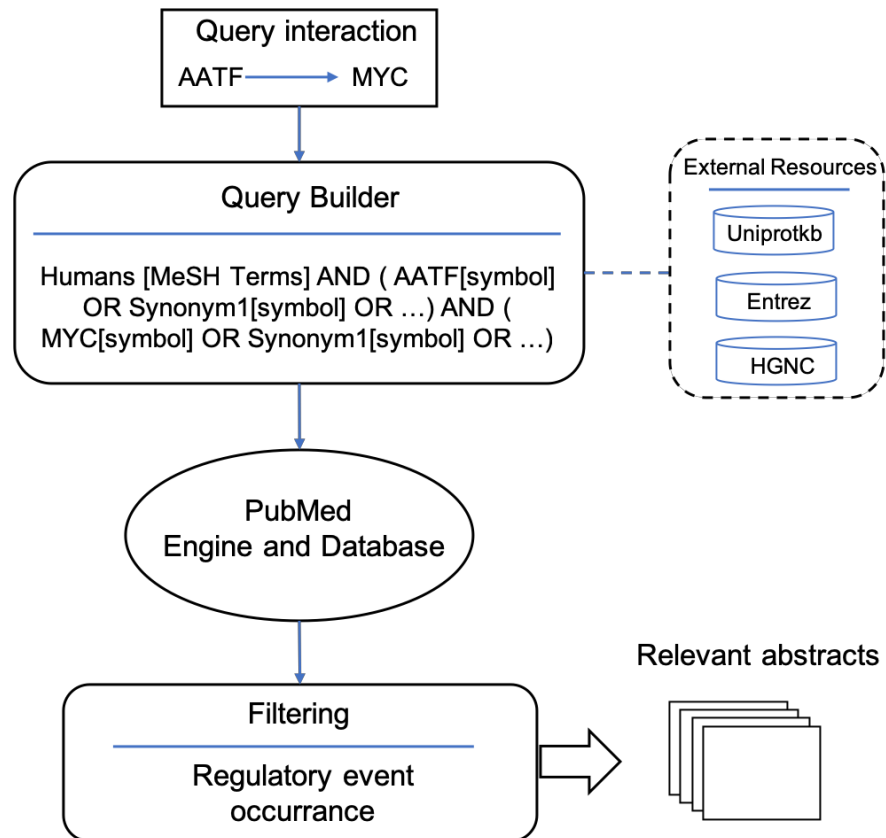
### 145 **2.1 Data sets**

146 PubMed database was used to query the entities relating to interaction in ChIP-Atlas. PubMed provides  
147 more than 25 million biomedical and life sciences journal articles. TRRUST regulatory network [49] was  
148 utilized as gold standard to evaluate the performance of ModEx. TRRUST is a manually-curated database  
149 of human transcriptional regulatory network with partial information on mode of regulation. It contains  
150 9,396 regulatory interactions of 800 human transcription factors, 5,066 of which are annotated with  
151 information on mode of regulation (3,148 repression and 1,918 activation). We also obtained TF-gene  
152 interaction data from ChIP-seq experiments, deposited on the ChIP-Atlas database [25]. This database  
153 contains all publicly available high-throughput ChIP-seq experiments. We assembled regulatory networks  
154 from these interactions using various cutoff criteria for ChIP-seq peak signal score and distance to the TSS.  
155 The least stringent criterion results in a network with 4 million interactions between 758 TFs and 18,874  
156 target genes. There is no reported mode of regulation in this database.

## 157 2.2 Extraction of relevant citations

158 For each regulatory interaction in our assembled ChIP-derived network, we developed an IR system to  
159 retrieve the information from the literature. Figure 1 illustrates the overall workflow of our IR component  
160 to fetch relevant citations associated with the regulatory interaction. We built a query based on the entities  
161 participated in the interactions to retrieve abstracts from PubMed database.

162



163

164

165 **Fig. 1.** The Information Retrieval framework. The steps are as follows: first, a Boolean query is built according  
166 to the associated entities in the regulatory interaction. It uses several external databases to complement the query  
167 with more synonyms and aliases. Then, the query is submitted to the PubMed database and abstracts are  
168 retrieved for processing. Abstracts with no regulatory events are excluded for further analysis.

169

170 Each query was supplemented with synonyms acquired from several external resources, including HGNC,  
171 Entrez, and UniprotKB to fetch more relevant abstracts. All related citations were acquired by submitting  
172 a query with appropriate Boolean logic (AND/OR) on entities and their synonyms. A MeSH descriptor term

173 (e.g. Humans) was also incorporated in the query to reduce the search space. For examples the query for  
174 AATF and MYC regulatory interaction is, “humans[mesh terms] AND (AATF[sym] OR BFR2[sym] OR  
175 CHE-1[sym] OR CHE1[sym] OR DED[sym]) AND (MYC[sym] OR MRTL[sym] OR MYCC[sym] OR  
176 BHLHE39[sym] OR C-MYC[sym])”.

177 The queries were submitted via the PubMed engine, a search engine that provides access to the MEDLINE  
178 database of references and abstracts on life sciences and biomedical articles. In our implementation we  
179 utilized Biopython [50] to run the queries through PubMed engine. We applied a filter on retrieved abstracts  
180 and retained only those containing expert-generated “regulatory events” as presented in Table 1. Each  
181 category contains more than 50 verbs and their inflections. For example, the AATF-MYC query outlined  
182 above, resulted in 4 relevant abstracts (PMIDs: 20549547, 17006618, 17006618, 20924650).

183

184 **TABLE 1.** Regulatory events categories

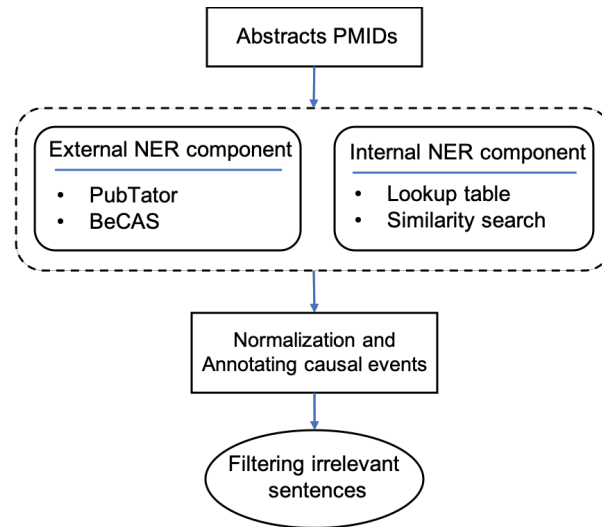
Category	NO. Events	Examples
Positive	500	increase, induce, activate, enhance, up-regulate, ...
Negative	511	reduce, decrease, suppress, block, down-regulate, decrease, ...

185

### 186 **2.3 Gene and regulatory event recognition**

187 The next step in the pipeline is to identify biological entities within the abstracts. Figure 2 shows the NER  
188 component of our system. Two external state-of-the-art NER systems were utilized to annotate the retrieved  
189 abstracts with an accurate and complete list of biological entities. The first system is PubTator [51], a web-  
190 based system for assisting biocuration. PubTator utilizes a HTTP REST interface, equipped with multiple  
191 state-of-the-art text mining algorithms to run queries. Using this system, we queried the retrieved PMIDs  
192 and obtained entity annotations in a JSON encoded text. Additionally, we utilized BeCAS [51] (another  
193 online NER tool) to improve the coverage of the entities. BeCAS, like PubTator, provides a RESTful API  
194 for biomedical name identification. It can run queries directly on provided text or PMIDs and returns  
195 associated annotations as an XML document.

196



197  
198

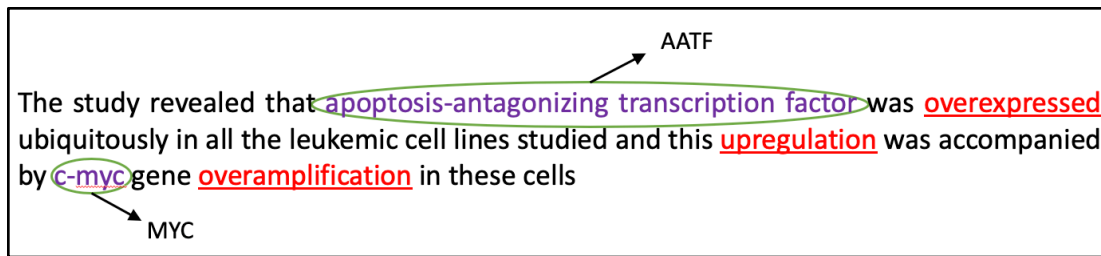
199 **Fig. 2.** The Gene and regulatory event recognition workflow. Each PubMed ID retrieved by IR component are  
200 submitted to the external NER tools (PubTator and BeCAS) for annotating genes in the abstracts. It follows  
201 complementary annotations using our internal NER component including a lookup table for covering acronyms, and  
202 a similarity search to identify lexical variations for gene names.

203 To further enhance the NER system, we implemented and added an additional NER component as follows.  
204 Abstracts were normalized to uppercase format and searched for gene acronyms using a manually-curated  
205 lookup table [52]. This table includes long term / short term pair association to recognize entities, which  
206 were missed by the external NER tools. For instance, AR is a short term for “Androgen Receptor” and was  
207 only detected as an entity (transcription factor) using this lookup table. Furthermore, we utilized a name  
208 similarity metric to identify strings with lexical variations such as whitespace and punctuations. For  
209 instance, “IL-12” and “IL12” are two lexical variations of “Interleukin 12”. The former version was not  
210 identified by the External NER systems. In our implementation, we set the entity detection threshold based  
211 on Jaro similarity [53] of 0.9 or larger between the query entity and the string in the abstract.

212 Next, we normalized the annotated word or a group of words corresponding to a gene to their HGNC symbol  
213 for simplification of downstream analysis. Regulatory events were also annotated using our expert-  
214 generated categories (Table 1). Figure 3 illustrates the normalization of gene names and annotation of  
215 regulatory events. Sentences that contained no regulatory event were excluded from further analysis.

216





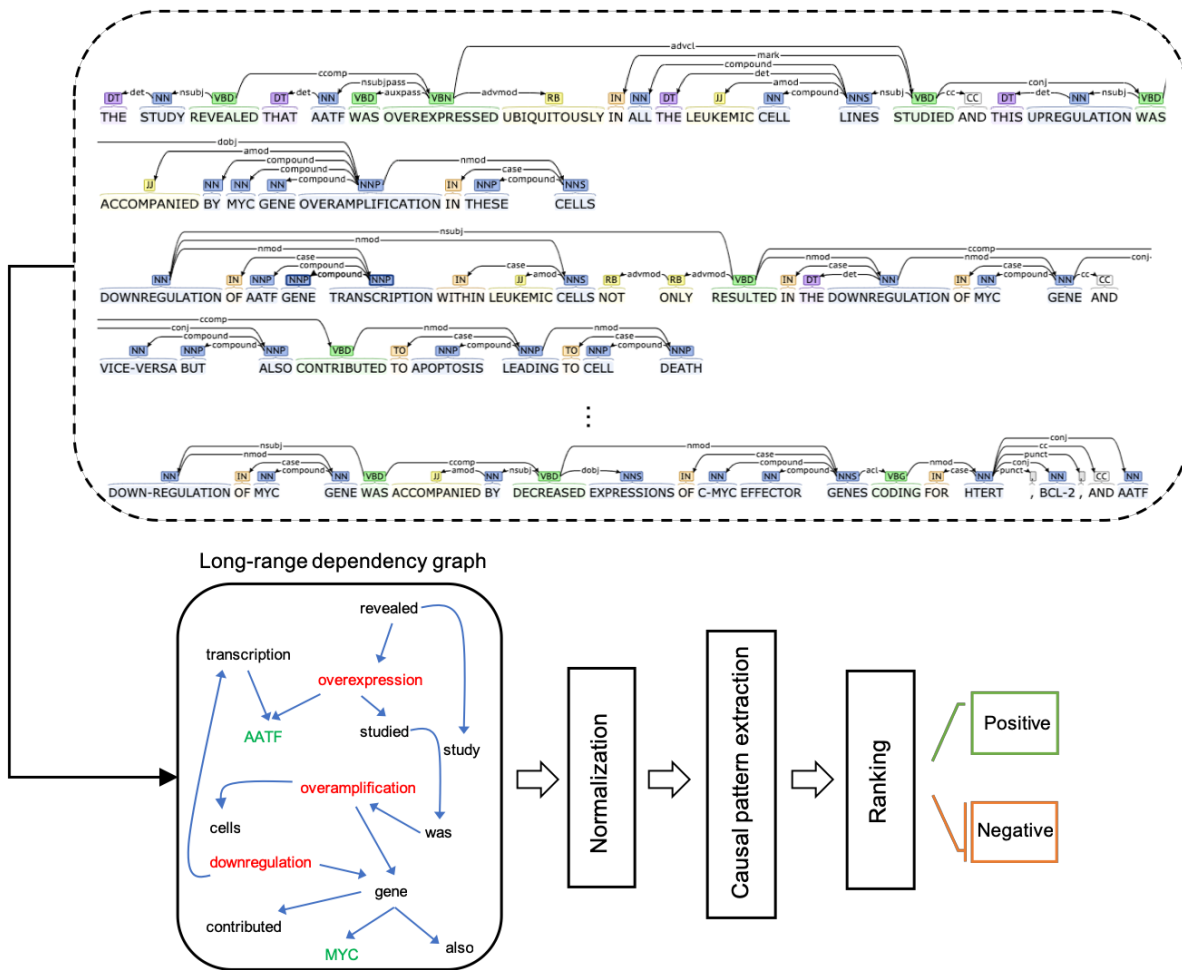
217

218 **Fig. 3.** An example of gene entity normalization and regulatory events annotation. All of the words or group of  
219 words associated to target entities (purple color) are normalized to their HGNC symbol for simplification. Causal  
220 regulatory events also are annotated according to their categories, and sentences with no regulatory event are  
221 excluded for further consideration.

#### 222 **2.4. Extracting mode of regulation**

223 For each causal interaction, its associated annotated sentences within relevant abstracts were submitted to  
224 the Stanford dependency parser [48] and a dependency parse tree was generated. Dependency trees  
225 extracted from different sentences were merged into a single large graph. The merging process is  
226 straightforward; Each dependency relation includes one head word/node and one dependent word/node.  
227 Nodes from different dependency relations representing the same word were. PMID was recorded for each  
228 edge in the parse tree to indicate its source. Each edge in the parse tree was assigned a weight based on the  
229 number of occurrences of dependency relations. The rationale for using this weighted parse tree is that it can  
230 be used to identify long-range dependency relations across sentence boundaries that would otherwise be  
231 missed. Figure 4 shows the relation extraction workflow of our method. Absolute frequency of a  
232 dependency relation obtained from the merging step can somewhat reflect the semantic relation of the head  
233 word and the dependent word.

234



235

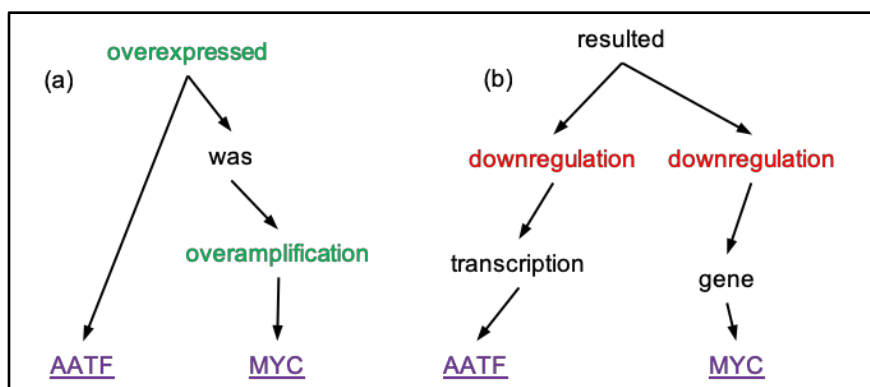
236

237 **Fig. 4.** Regulatory sign extraction workflow. A long-range dependency graph is constructed by merging all of  
 238 dependency trees corresponding to the evidence sentences. The weights of the graph are the number of occurrences  
 239 of dependency relations. Candidate regulatory signs are identified using common subtrees with at least one regulatory  
 240 event in the graph. Finally, a sign of regulation is assigned to the query interaction through the ranking task.

241 Our system, ModEx, creates candidate relations by extracting subtrees with common ancestors connecting  
 242 the pair of query genes as leaves. These subtrees must contain at least one causal event describing the  
 243 candidate relation between the given pair of genes. Subtrees were extracted by applying a Depth First  
 244 Search along with a boolean visited array to avoid possible loops. Nodes with two paths to the entities were  
 245 considered as a root of the subtree. Next, we utilized a rule based approach to describe relations using three  
 246 commonly used language constructs [36]. The first rule is effector-relation-effectee (e.g. A activates B).  
 247 The second rule is relation-of-effectee-by-effector (e.g. Activation of A by B). These rules were applied to  
 248 both paths from root to query entities to identify their regulatory dependency. Figure 5a illustrates the  
 249 regulatory relation extraction using these rules. Some sentences in the literature have complex structures,

250 which cannot be captured by these language constructs. To address this, we also incorporated a negation  
251 rule to increase the performance of the RE system. For example, consider the following sentence: “LMP1  
252 suppresses the transcriptional repressor ATF3, possibly leading to the TGFβ-induced ID1 upregulation”  
253 [54]. In the first pass the system assigns a positive mode to the interaction between ATF3 and ID1.  
254 However, there is a negative interaction between the TF and target gene. The negation rule considers the  
255 negative event “suppresses” related to ATF3 and switch the positive mode to negative. Figure 5b shows a  
256 subtree reflecting the negation rule.

257



258

259

260 **Fig. 5.** Examples to illustrate the rules for finding the regulatory sign. panel (a) shows an example for simple rule  
261 (effector-relation-effectee) in which the RE system can assign a positive sign to this candidate pattern. In panel (b),  
262 we can see the impact of the negation rule to extract accurate sign to this pattern. Two paths from root to query entities  
263 contain negative regulatory events which carries an activation/positive sign for the pattern.

264 We then ranked each subtree based on the sign of regulatory interaction between the query genes. The  
265 weights of the graph encode repetition of regulatory relations across sentences and abstracts. we considered  
266 the weights when there were more than one regulatory event associated with the target gene. In this case,  
267 an event with higher weight was selected for ranking the subtree. We also considered distance of events to  
268 the target gene when the weights in the subtree were equal. The closest event to the target entity will take  
269 the highest priority for determining the interaction sign. Finally, we investigated signs in every candidate  
270 subtree and assigned a total sign of regulation to the interaction using a voting scheme.

271

## 272 2.5. ModEx HTTP Interface

273 We implemented an HTTP REST server for users to programmatically annotate gene regulatory networks  
274 using ModEx. Clients should make HTTP requests to the server with a particular format, specifying the  
275 query entities and optional MeSH term to annotate. The query has to be requested in the following format:  
276 *TFEntrezID\_TargetEntrezID\_MeSHterm[optional]*. For instance, a query to the server for AATF-MYC  
277 should be formatted as “/signex/26574\_4609\_humans”. The server returns extracted annotation along with  
278 associated citations and sentences from PubMed database if any evidence exist. The server can be accessed  
279 at: <https://watson.math.umb.edu/modex/>

280

281

## 282 **3 RESULTS**

### 283 **3.1 Classification Performance**

284 We evaluated the performance of our method using the TRRUST database, a manually-curated network or  
285 regulatory interaction with partial information on mode of regulation. TRRUST is a high-quality database  
286 and can be considered as gold standard for our benchmark. We applied our method to 5,066 regulatory  
287 interactions in TRRUST for which information on mode of regulation was available. Table 2 shows the  
288 summary statistic of the results.

289

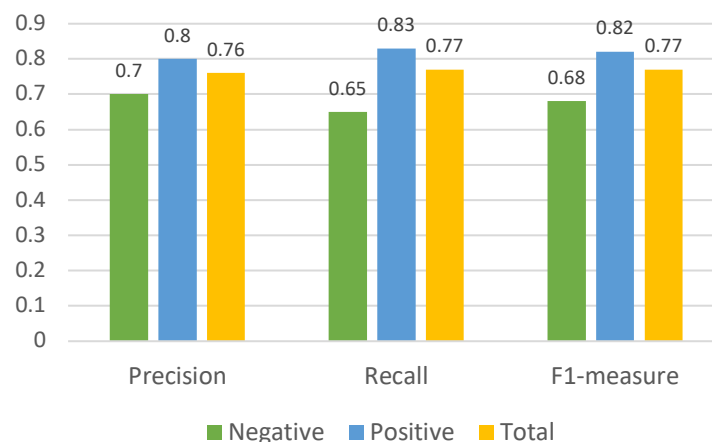
290 **TABLE 2.** Summary statistics of performing ModEx on TRRUST

System	Without evidence	With evidence	Detected with sign	
			4225	
modEx	182	4884	Positive	Negative
			2659	1566

291

292 Our method did not detect any PubMed abstracts for 182 of interactions. ModEx detected 4,225 signs  
293 corresponding to 4,225 regulatory interactions including 2,659 positive and 1,566 negative interactions. We  
294 compare the identified signs by ModEx with reported signs in the TRRUST database. Our system correctly  
295 extracted 2,216 positive and 1,024 negative signs with overall accuracy of 0.76. Figure 6 shows the

296 classification result of ModEx on the TRRUST database using various metrics (Precision, Recall and F-  
297 score).  
298



299

300 **Fig. 6.** Classification results of ModEx on TRRUST

301 We also compared the citations of the 4,884 retrieved by our system with the citations reported in the  
302 TRRUST database. All citations match, with accuracy of 1.0. This validation using the gold standard  
303 demonstrates the ability of our system to correctly identify relevant citations, extract causal interactions  
304 between TF and gene, and detect mode of regulation. We used our system to annotate the remaining  
305 interactions in the TRUSST database for which no mode of regulation is reported.

### 306 3.2 ChIP-Atlas Analysis

307 We next sought to extract and annotate ChIP-seq derived TF-gene causal regulatory interactions from  
308 literature using our system. Such meta-data and evidence from literature can increase the confidence in the  
309 TF-gene interactions identified by ChIP-seq experiments and further shed light on the mechanism of  
310 interaction. Information on mode of regulation in particular can be helpful to enhance the accuracy of  
311 enrichment algorithms for regulatory pathway inference [55].

312 We applied ModEx to ChIP-seq interactions, with moderately stringency criteria, i.e., binding distance  
313 within 1k of the TSS and ChIP peak score > 950, resulting in 43,444 interactions. The system was able to  
314 detect and annotate 1,592 of interactions in PubMed database. Table 3 outlines the results.

315

316 **TABLE 3.** Summary statistics of performing ModEx on ChIP-Atlas

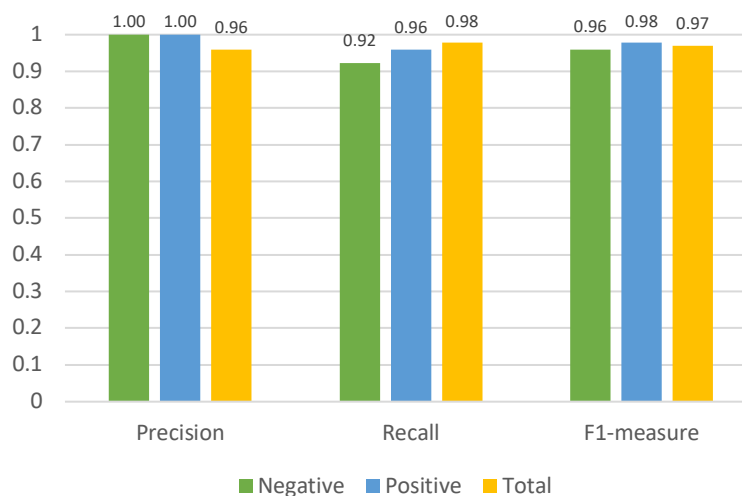
System	Overall	With evidence	Detected with sign	
			Positive	Negative
			1,592	
ModEx	43,444	5,133	1,421	171

317

318 Some of the retrieved annotated ChIP-seq interactions also appear in the TRRUST database (69 total). We  
 319 compared the identified mode of regulations of ChIP-seq interactions with the reported signs in the  
 320 TRRUST database. Figure 7 summarizes the classification results. As can be seen the agreement is very  
 321 high, indicating that our method can reliably identify and annotate ChIP interaction when they are reported  
 322 in literature. Additionally, we compared our acquired evidence (PMIDs) by ModEx with citations reported  
 323 in TRRUST. Our IR module was able to fetch the relevant evidence from PubMed database with accuracy  
 324 0.88.

325

326



327

328

329 **Fig. 7.** Classification results of ModEx on intersctio on TRRUST and ChIP-Atlas

### 330 3.3 Directional enrichment analysis

331 To demonstrate the utility of our annotated network, we used our network in conjunction with a directional  
 332 enrichment analysis algorithm [55] to identify drivers of differential expressed genes. We utilized  
 333 quaternaryProd, a gene set enrichment algorithm that can take advantage of direction of regulation on causal  
 334 biological interaction graphs to identify regulators of differential gene expression. The algorithm can take  
 335 a signed transcriptional regulatory network, such as TRRUST or our annotated ChIP-network along with a  
 336 differential expression profile as input and outputs a set of candidate active regulators. The ability of the  
 337 algorithm to identify regulators of differential gene expression relies heavily on the quality and the coverage  
 338 of the regulatory network on which the queries are performed. To test the utility of our network, we used  
 339 this algorithm along with differential expression profiles from controlled over-expression experiments used  
 340 in the original study. The over-expression experiments consist of differential gene expression profile from  
 341 a controlled in vitro E2F3 over expression [56] and c-Myc [56]. We inputted three networks into the  
 342 algorithm (1) the original TRUSST network, (2) annotated TRUSST network, and (3) annotated TRRUST  
 343 augmented with annotated ChIP-Atlas. By annotated TRRUST, we refer to the TRRUST network where  
 344 interaction with no reported mode of regulation were annotated using our system. Differential gene  
 345 expression analysis of these data sets resulted in 272, and 220 differentially expressed genes  
 346 respectively. Table 4 outlines the top 5 regulators predicted by the algorithm on E2F3 experiment sorted  
 347 by the FDR corrected *p-values* of the scoring scheme. For the E2F3 experiment, E2F1 is returned as the top  
 348 hypothesis regulator by the algorithm incorporating our annotated networks. E2F1 and E2F3 are close  
 349 family members and have a very similar role as transcription factors that function to control the cell cycle  
 350 and are similarly implicated in cancer [57]. It is interesting to note that original TRRUST database does not  
 351 include enough information for algorithm to recover E2F1, however the signal strengthens when TRUSST  
 352 is annotated with our system and a much more significant p-value is obtained when TRRUST is augmented  
 353 with annotated ChIP-Atlas. This shows that annotating ChIP-seq data provides significant additional power  
 354 to identify upstream regulators in conjunction with freely available causal networks.

355

356 **TABLE 4.** Directional enrichment analysis results on E2F1 expression signatures

TRRUST			Annotated TRRUST			Annotated TRRUST with ChIP-Atlas		
Name	Regulation	Adj. P-Val	Name	Regulation	Adj. P-Val	Name	Regulation	Adj. P-Val
REL	Down	1.5e-3	<b>E2F1</b>	<b>Up</b>	<b>1.2e-5</b>	<b>E2F1</b>	<b>Up</b>	<b>1.2e-7</b>
PROX1	Up	2.9e-3	PROX1	Up	2.1e-3	PROX1	Up	2.1e-3
SUGP1	Down	3.2e-3	SUGP1	Down	2.4e-3	SUGP1	Down	2.4e-3

NFIL3	Up	6.3e-3	RELA	Down	3.5e-3	RELA	Down	3.5e-3
TFDP1	Up	7e-3	TFDP1	Up	3.9e-3	TFDP1	Up	3.9e-3

357

358 Application of the method to c-Myc differential expression profile shows the same pattern. The annotated  
 359 TRRUST with ChIP-Atlas recovered MAX as one of the top 20 regulators with low p-value compared to  
 360 TRRUST. It has been demonstrated that oncogenic activity of c-Myc requires dimerization with MAX  
 361 [58].

362

363 **TABLE 5.** Directional enrichment analysis results on c-Myc expression signatures

TRRUST			Annotated TRRUST			Annotated TRRUST with ChIP-Atlas		
Name	Regulation	Adj. P-Val	Name	Regulation	Adj. P-Val	Name	Regulation	Adj. P-Val
MYBL2	Up	3.6e-3	MAFA	Down	1.3e-2	USF2	Up	8.8e-3
MXI1	Down	4.0e-3	MKL1	Down	1.3e-2	MAFA	Down	1.3e-2
AATF	Up	4.0e-3	GLI3	Down	1.7e-2	MKL1	Down	1.3e-2
ENO1	Down	4.0e-3	KAT2B	Up	1.9e-2	GLI3	Down	1.7e-2
NR1D1	Up	6.4e-3	SOX6	Down	1.9e-2	KAT2B	Up	1.9e-2
TLE3	Up	6.4e-3	HDAC1	Down	2.5e-2	SOX6	Down	1.9e-2
TOP2B	Down	6.4e-3	MYBL2	Down	2.6e-2	HDAC1	Down	2.5e-2
L3MBTL1	Up	6.4e-3	HDAC7	Up	3.0e-2	MYBL2	Down	2.6e-2
MAFA	Down	6.4e-3	ILF3	Down	3.0e-2	HDAC7	Up	3.0e-3
TLX1	Up	7.9e-3	ELK1	Up	3.6e-2	ILF3	Down	3.0e-2
HLF	Up	1.0e-2	GATA6	Up	3.9e-2	ELK1	Up	3.6e-2
DLX5	Up	1.1e-2	PPARG	Down	4.3e-2	GATA6	Up	3.9e-2
HOXA1	Up	1.1e-2	SATB1	Up	4.4e-2	PLAG1	Up	4.0e-2
MKL1	Down	1.2e-2	ARNT	Up	4.8e-2	CREB1	Up	4.1e-2
IFI16	Down	1.4e-2	RXRA	Up	5.3e-2	PPARG	Down	4.3e-2
PRDM1	Down	1.5e-2	ZEB1	Down	5.5e-2	SATB1	Up	4.4e-2
CEBPPE	Down	1.6e-2	E2F4	Down	5.5e-2	ZNF143	Up	4.6e-2
HDAC1	Down	1.6e-2	PPARD	Down	5.6e-2	ARNT	UP	4.8e-2
GLI3	Down	1.6e-2	XBP1	Up	6.5e-2	<b>MAX</b>	<b>Up</b>	<b>5.3e-2</b>
STAT4	Up	1.8e-2	DDIT3	Down	6.5e-2	RXRA	Up	5.4e-2

364



## 365 **Conclusion**

366 In this work we presented a fully automated text-mining system to extract and annotate causal regulatory  
367 interaction between transcription factors and genes from the biomedical literature. As a starting point, our  
368 method uses putative TF-gene interactions derived from high-throughput ChIP-seq or other experiments  
369 and seeks to collect evidence and meta-data in the biomedical literature to support the interaction. It should  
370 be noted that annotating a priori known interactions differs significantly in scope and complexity from  
371 general text-mining approaches for biomedical relation extraction. The later attempts to extract the causal  
372 relation from biomedical text directly, without prior knowledge of the entities and the interaction, whereas  
373 in our method the relation is know from biological experiments and curated databases a priori, thereby  
374 reducing the complexity significantly. This approach bridges the gap between data-driven methods and  
375 text-mining methods for constructing causal transcriptional gene regulatory networks and overcomes some  
376 of the drawbacks of either approach. With the rapid increase in high-throughput experiments and  
377 biomedical literature, hybrid method such as the one proposed can make a significant impact in biological  
378 knowledge retrieval.

379 We used a gold-standard manually curated dataset and demonstrated that our approach can reliably identify  
380 the relevant literature and extract the correct interaction and meta-data. We applied our method to high-  
381 throughput ChIP-seq data and provided literature support for ~1,500 interactions. Our annotated ChIP-  
382 derived transcriptional regulatory interaction can be used in conjunction with directional enrichment  
383 methods that aim to identify regulators of differential gene expression. Moreover, we use our system to  
384 annotate the interactions in the TRRUST database for which more of regulation is not reported. Our system  
385 can also be used as a tool to mine the literature for investigate interactions in newly performed ChIP-seq  
386 experiments, where researchers are interested to investigate a specific interaction between a protein and a  
387 gene. To facilitate usage, we implemented an HTTP REST server for users to programmatically annotate  
388 gene regulatory networks using ModEx available to download at: <https://watson.math.umb.edu/modex/>.  
389 The annotated ChIP-network as well as annotated TRRUST can be obtained from:  
390 <https://doi.org/10.6084/m9.figshare.8251502.v1>

## 391 **ACKNOWLEDGMENT**

392 The authors are grateful to Dr. Nurit Haspel for helpful discussions and comments.

## 393 **REFERENCES**

394

- 395 [1] H. de Jong, Modeling and Simulation of Genetic Regulatory Systems: A Literature Review,  
396 J. Comput. Biol. 9 (2002) 67–103. doi:10.1089/10665270252833208.

- 397 [2] G. Karlebach, R. Shamir, Modelling and analysis of gene regulatory networks, *Nat. Rev.*  
398 *Mol. Cell Biol.* 9 (2008) 770–780. doi:10.1038/nrm2503.
- 399 [3] S. Farahmand, S. Goliaei, N. Ansari-Pour, Z. Razaghi-Moghadam, GTA: a game theoretic  
400 approach to identifying cancer subnetwork markers, *Mol. Biosyst.* 12 (2016) 818–825.  
401 doi:10.1039/C5MB00684H.
- 402 [4] S. Farahmand, M.H. Foroughmand-Araabi, S. Goliaei, Z. Razaghi-Moghadam, CytoGTA: A  
403 cytoscape plugin for identifying discriminative subnetwork markers using a game theoretic  
404 approach, *PLoS One.* 12 (2017) e0185016. doi:10.1371/journal.pone.0185016.
- 405 [5] Z. Razaghi-Moghadam, A. Namipashaki, S. Farahmand, N. Ansari-Pour, Systems genetics of  
406 nonsyndromic orofacial clefting provides insights into its complex aetiology, *Eur. J. Hum.*  
407 *Genet.* 27 (2019) 226–234. doi:10.1038/s41431-018-0263-7.
- 408 [6] M. Miwa, T. Ohta, R. Rak, A. Rowley, D.B. Kell, S. Pyysalo, S. Ananiadou, A method for  
409 integrating and ranking the evidence for biochemical pathways by mining reactions from  
410 text, *Bioinformatics.* 29 (2013) i44–i52. doi:10.1093/bioinformatics/btt227.
- 411 [7] S. Bagewadi, T. Bobić, M. Hofmann-Apitius, J. Fluck, R. Klinger, Detecting miRNA Mentions  
412 and Relations in Biomedical Literature, *F1000Research.* 3 (2015) 205.  
413 doi:10.12688/f1000research.4591.3.
- 414 [8] D. Szklarczyk, A. Franceschini, S. Wyder, K. Forslund, D. Heller, J. Huerta-Cepas, M.  
415 Simonovic, A. Roth, A. Santos, K.P. Tsafou, M. Kuhn, P. Bork, L.J. Jensen, C. von Mering,  
416 STRING v10: protein–protein interaction networks, integrated over the tree of life, *Nucleic*  
417 *Acids Res.* 43 (2015) D447–D452. doi:10.1093/nar/gku1003.
- 418 [9] M. Kanehisa, S. Goto, KEGG: Kyoto Encyclopedia of Genes and Genomes, *Nucleic Acids Res.*  
419 28 (2000) 27–30. doi:10.1093/nar/28.1.27.
- 420 [10] E.G. Cerami, B.E. Gross, E. Demir, I. Rodchenkov, O. Babur, N. Anwar, N. Schultz, G.D.  
421 Bader, C. Sander, Pathway Commons, a web resource for biological pathway data, *Nucleic*  
422 *Acids Res.* 39 (2011) D685–D690. doi:10.1093/nar/gkq1039.
- 423 [11] A. Sandelin, W. Alkema, P. Engström, W.W. Wasserman, B. Lenhard, JASPAR: an open-  
424 access database for eukaryotic transcription factor binding profiles, *Nucleic Acids Res.* 32  
425 (2004) 91D – 94. doi:10.1093/nar/gkh012.
- 426 [12] O.L. Griffith, S.B. Montgomery, B. Bernier, B. Chu, K. Kasaian, S. Aerts, S. Mahony, M.C.  
427 Sleumer, M. Bilenky, M. Haeussler, M. Griffith, S.M. Gallo, B. Giardine, B. Hooghe, P. Van  
428 Loo, E. Blanco, A. Ticoll, S. Lithwick, E. Portales-Casamar, I.J. Donaldson, G. Robertson, C.  
429 Wadelius, P. De Bleser, D. Vlieghe, M.S. Halfon, W. Wasserman, R. Hardison, C.M.  
430 Bergman, S.J.M. Jones, ORegAnno: an open-access community-driven resource for  
431 regulatory annotation, *Nucleic Acids Res.* 36 (2007) D107–D113. doi:10.1093/nar/gkm967.
- 432 [13] M. Pachkov, I. Erb, N. Molina, E. van Nimwegen, SwissRegulon: a database of genome-  
433 wide annotations of regulatory sites, *Nucleic Acids Res.* 35 (2007) D127–D131.  
434 doi:10.1093/nar/gkl857.
- 435 [14] C. Jiang, Z. Xuan, F. Zhao, M.Q. Zhang, TRED: a transcriptional regulatory element  
436 database, new entries and other development, *Nucleic Acids Res.* 35 (2007) D137–D140.  
437 doi:10.1093/nar/gkl1041.
- 438 [15] N.A. Kolchanov, O.A. Podkolodnaya, E.A. Ananko, E. V. Ignatieva, I.L. Stepanenko, O. V. Kel-  
439 Margoulis, A.E. Kel, T.I. Merkulova, T.N. Goryachkovskaya, T. V. Busygina, F.A. Kolpakov,  
440 N.L. Podkolodny, A.N. Naumochkin, I.M. Korostishevskaya, A.G. Romashchenko, G.C.

- 441 Overton, Transcription Regulatory Regions Database (TRRD): its status in 2000, *Nucleic*  
442 *Acids Res.* 28 (2000) 298–301. doi:10.1093/nar/28.1.298.
- 443 [16] A. Essaghir, F. Toffalini, L. Knoops, A. Kallin, J. van Helden, J.-B. Demoulin, Transcription  
444 factor regulation can be accurately predicted from the presence of target gene signatures  
445 in microarray gene expression data, *Nucleic Acids Res.* 38 (2010) e120–e120.  
446 doi:10.1093/nar/gkq149.
- 447 [17] H. Han, H. Shim, D. Shin, J.E. Shim, Y. Ko, J. Shin, H. Kim, A. Cho, E. Kim, T. Lee, H. Kim, K.  
448 Kim, S. Yang, D. Bae, A. Yun, S. Kim, C.Y. Kim, H.J. Cho, B. Kang, S. Shin, I. Lee, TRRUST: a  
449 reference database of human transcriptional regulatory interactions, *Sci. Rep.* 5 (2015)  
450 11432. doi:10.1038/srep11432.
- 451 [18] E. Segal, M. Shapira, A. Regev, D. Pe’er, D. Botstein, D. Koller, N. Friedman, Module  
452 networks: identifying regulatory modules and their condition-specific regulators from  
453 gene expression data, *Nat. Genet.* 34 (2003) 166–176. doi:10.1038/ng1165.
- 454 [19] S. Farahmand, S. Goliaei, Z.R.M. Kashani, S. Farahmand, Identifying Cancer Subnetwork  
455 Markers Using Game Theory Method, in: Springer, Singapore, 2019: pp. 105–109.  
456 doi:10.1007/978-981-10-4505-9\_17.
- 457 [20] A. Vajdi, N. Haspel, H. Banaee, A new DP algorithm for comparing gene expression data  
458 using geometric similarity, in: 2015 IEEE Int. Conf. Bioinforma. Biomed., IEEE, 2015: pp.  
459 1157–1161. doi:10.1109/BIBM.2015.7359846.
- 460 [21] M. Krallinger, A. Valencia, Text-mining and information-retrieval services for molecular  
461 biology, *Genome Biol.* 6 (2005) 224. doi:10.1186/gb-2005-6-7-224.
- 462 [22] H. Han, J.-W. Cho, S. Lee, A. Yun, H. Kim, D. Bae, S. Yang, C.Y. Kim, M. Lee, E. Kim, S. Lee, B.  
463 Kang, D. Jeong, Y. Kim, H.-N. Jeon, H. Jung, S. Nam, M. Chung, J.-H. Kim, I. Lee, TRRUST v2:  
464 an expanded reference database of human and mouse transcriptional regulatory  
465 interactions, *Nucleic Acids Res.* 46 (2018) D380–D386. doi:10.1093/nar/gkx1013.
- 466 [23] T.E.P. ENCODE Project Consortium, The ENCODE (ENCyclopedia Of DNA Elements) Project.,  
467 *Science.* 306 (2004) 636–40. doi:10.1126/science.1105136.
- 468 [24] I. Yevshin, R. Sharipov, T. Valeev, A. Kel, F. Kolpakov, GTRD: a database of transcription  
469 factor binding sites identified by ChIP-seq experiments, *Nucleic Acids Res.* 45 (2017) D61–  
470 D67. doi:10.1093/nar/gkw951.
- 471 [25] S. Oki, T. Ohta, G. Shioi, H. Hatanaka, O. Ogasawara, Y. Okuda, H. Kawaji, R. Nakaki, J. Sese,  
472 C. Meno, ChIP-Atlas: a data-mining suite powered by full integration of public ChIP-seq  
473 data, *EMBO Rep.* 19 (2018) e46255. doi:10.15252/embr.201846255.
- 474 [26] S. Gebel, R.B. Lichtner, B. Frushour, W.K. Schlage, V. Hoang, M. Talikka, A. Hengstermann,  
475 C. Mathis, E. Veljkovic, M. Peck, M.C. Peitsch, R. Deehan, J. Hoeng, J.W. Westra,  
476 Construction of a computable network model for DNA damage, autophagy, cell death, and  
477 senescence., *Bioinform. Biol. Insights.* 7 (2013) 97–117. doi:10.4137/BBI.S11154.
- 478 [27] L.J. Jensen, J. Saric, P. Bork, Literature mining for the biologist: from information retrieval  
479 to biological discovery, *Nat. Rev. Genet.* 7 (2006) 119–129. doi:10.1038/nrg1768.
- 480 [28] A.M. Cohen, W.R. Hersh, A survey of current work in biomedical text mining, *Brief.*  
481 *Bioinform.* 6 (2005) 57–71. doi:10.1093/bib/6.1.57.
- 482 [29] P. Zweigenbaum, D. Demner-Fushman, H. Yu, K.B. Cohen, Frontiers of biomedical text  
483 mining: current progress, *Brief. Bioinform.* 8 (2007) 358–375. doi:10.1093/bib/bbm045.
- 484 [30] R. Hoffmann, A. Valencia, Implementing the iHOP concept for navigation of biomedical

- 485 literature, *Bioinformatics*. 21 (2005) ii252–ii258. doi:10.1093/bioinformatics/bti1142.
- 486 [31] M. Ashburner, C.A. Ball, J.A. Blake, D. Botstein, H. Butler, J.M. Cherry, A.P. Davis, K.  
487 Dolinski, S.S. Dwight, J.T. Eppig, M.A. Harris, D.P. Hill, L. Issel-Tarver, A. Kasarskis, S. Lewis,  
488 J.C. Matese, J.E. Richardson, M. Ringwald, G.M. Rubin, G. Sherlock, Gene Ontology: tool  
489 for the unification of biology, *Nat. Genet.* 2000 251. (2000).
- 490 [32] O. Bodenreider, The Unified Medical Language System (UMLS): integrating biomedical  
491 terminology., *Nucleic Acids Res.* 32 (2004) D267-70. doi:10.1093/nar/gkh061.
- 492 [33] A. Yeh, A. Morgan, M. Colosimo, L. Hirschman, BioCreAtIvE Task 1A: gene mention finding  
493 evaluation, *BMC Bioinformatics*. 6 (2005) S2. doi:10.1186/1471-2105-6-S1-S2.
- 494 [34] Y. Mao, K. Van Auken, D. Li, C.N. Arighi, P. McQuilton, G.T. Hayman, S. Tweedie, M.L.  
495 Schaeffer, S.J.F. Laulederkind, S.-J. Wang, J. Gobeill, P. Ruch, A.T. Luu, J.-j. Kim, J.-H. Chiang,  
496 Y.-D. Chen, C.-J. Yang, H. Liu, D. Zhu, Y. Li, H. Yu, E. Emadzadeh, G. Gonzalez, J.-M. Chen,  
497 H.-J. Dai, Z. Lu, Overview of the gene ontology task at BioCreative IV, Database. 2014  
498 (2014) bau086–bau086. doi:10.1093/database/bau086.
- 499 [35] A.A. Morgan, Z. Lu, X. Wang, A.M. Cohen, J. Fluck, P. Ruch, A. Divoli, K. Fundel, R. Leaman,  
500 J. Hakenberg, C. Sun, H. Liu, R. Torres, M. Krauthammer, W.W. Lau, H. Liu, C.-N. Hsu, M.  
501 Schuemie, K.B. Cohen, L. Hirschman, Overview of BioCreative II gene normalization.,  
502 *Genome Biol.* 9 Suppl 2 (2008) S3. doi:10.1186/gb-2008-9-s2-s3.
- 503 [36] K. Fundel, R. Kuffner, R. Zimmer, RelEx--Relation extraction using dependency parse trees,  
504 *Bioinformatics*. 23 (2007) 365–371. doi:10.1093/bioinformatics/btl616.
- 505 [37] L. Qian, G. Zhou, Tree kernel-based protein–protein interaction extraction from biomedical  
506 literature, *J. Biomed. Inform.* 45 (2012) 535–543. doi:10.1016/J.JBI.2012.02.004.
- 507 [38] A. Vajdi, K. Zarringhalam, N. Haspel, Patch-DCA: Improved Protein Interface Prediction by  
508 utilizing Structural Information and Clustering DCA scores, *BioRxiv*. (2019) 656074.  
509 doi:10.1101/656074.
- 510 [39] K.E. Ravikumar, K.B. Waghlikar, D. Li, J.-P. Kocher, H. Liu, Text mining facilitates database  
511 curation - extraction of mutation-disease associations from Bio-medical literature, *BMC*  
512 *Bioinformatics*. 16 (2015) 185. doi:10.1186/s12859-015-0609-x.
- 513 [40] Y.-C. Fang, P.-T. Lai, H.-J. Dai, W.-L. Hsu, MeInfoText 2.0: gene methylation and cancer  
514 relation extraction from biomedical literature, *BMC Bioinformatics*. 12 (2011) 471.  
515 doi:10.1186/1471-2105-12-471.
- 516 [41] M. Gerner, F. Sarafranz, C.M. Bergman, G. Nenadic, BioContext: an integrated text mining  
517 system for large-scale extraction and contextualization of biomolecular events,  
518 *Bioinformatics*. 28 (2012) 2154–2161. doi:10.1093/bioinformatics/bts332.
- 519 [42] J. Czarnecki, I. Nobeli, A.M. Smith, A.J. Shepherd, A text-mining system for extracting  
520 metabolic reactions from full-text articles, *BMC Bioinformatics*. 13 (2012) 172.  
521 doi:10.1186/1471-2105-13-172.
- 522 [43] E.K. Mallory, C. Zhang, C. Ré, R.B. Altman, Large-scale extraction of gene interactions from  
523 full-text literature using DeepDive, *Bioinformatics*. 32 (2015) btv476.  
524 doi:10.1093/bioinformatics/btv476.
- 525 [44] R.J. Mooney, R. Bunescu, Mining knowledge from text using information extraction, *ACM*  
526 *SIGKDD Explor. Newsl.* 7 (2005) 3–10. doi:10.1145/1089815.1089817.
- 527 [45] S. Bhatia, D. Kaul, N. Varma, Potential tumor suppressive function of miR-196b in B-cell  
528 lineage acute lymphoblastic leukemia, *Mol. Cell. Biochem.* 340 (2010) 97–106.

- 529 doi:10.1007/s11010-010-0406-9.
- 530 [46] C.-H. Wei, H.-Y. Kao, Z. Lu, PubTator: a web-based text mining tool for assisting biocuration,  
531 *Nucleic Acids Res.* 41 (2013) W518–W522. doi:10.1093/nar/gkt441.
- 532 [47] T. Nunes, D. Campos, S. Matos, J.L. Oliveira, BeCAS: biomedical concept recognition  
533 services and visualization, *Bioinformatics.* 29 (2013) 1915–1916.  
534 doi:10.1093/bioinformatics/btt317.
- 535 [48] D. Klein, C.D. Manning, Fast Exact Inference with a Factored Model for Natural Language  
536 Parsing, in: S. Becker, S. Thrun, K. Obermayer (Eds.), *Adv. Neural Inf. Process. Syst.* 15, MIT  
537 Press, 2003: pp. 3–10. [http://papers.nips.cc/paper/2325-fast-exact-inference-with-a-](http://papers.nips.cc/paper/2325-fast-exact-inference-with-a-factored-model-for-natural-language-parsing.pdf)  
538 [factored-model-for-natural-language-parsing.pdf](http://papers.nips.cc/paper/2325-fast-exact-inference-with-a-factored-model-for-natural-language-parsing.pdf).
- 539 [49] H. Han, J.-W. Cho, S. Lee, A. Yun, H. Kim, D. Bae, S. Yang, C.Y. Kim, M. Lee, E. Kim, S. Lee, B.  
540 Kang, D. Jeong, Y. Kim, H.-N. Jeon, H. Jung, S. Nam, M. Chung, J.-H. Kim, I. Lee, TRRUST v2:  
541 an expanded reference database of human and mouse transcriptional regulatory  
542 interactions, *Nucleic Acids Res.* 46 (2018) D380–D386. doi:10.1093/nar/gkx1013.
- 543 [50] P.J.A. Cock, T. Antao, J.T. Chang, B.A. Chapman, C.J. Cox, A. Dalke, I. Friedberg, T.  
544 Hamelryck, F. Kauff, B. Wilczynski, M.J.L. de Hoon, Biopython: freely available Python tools  
545 for computational molecular biology and bioinformatics, *Bioinformatics.* 25 (2009) 1422–  
546 1423. doi:10.1093/bioinformatics/btp163.
- 547 [51] C.-H. Wei, H.-Y. Kao, Z. Lu, PubTator: a web-based text mining tool for assisting biocuration,  
548 *Nucleic Acids Res.* 41 (2013) W518–W522. doi:10.1093/nar/gkt441.
- 549 [52] K.E. Ravikumar, M. Rastegar-Mojarad, H. Liu, BELMiner: adapting a rule-based relation  
550 extraction system to extract biological expression language statements from bio-medical  
551 literature evidence sentences, *Database.* 2017 (2017). doi:10.1093/database/baw156.
- 552 [53] M.A. Jaro, Advances in Record-Linkage Methodology as Applied to Matching the 1985  
553 Census of Tampa, Florida, *J. Am. Stat. Assoc.* 84 (1989) 414. doi:10.2307/2289924.
- 554 [54] A.K. Lo, C.W. Dawson, K.W. Lo, Y. Yu, L.S. Young, Upregulation of Id1 by Epstein-Barr Virus-  
555 encoded LMP1 confers resistance to TGF $\beta$ -mediated growth inhibition, *Mol. Cancer.* 9  
556 (2010) 155. doi:10.1186/1476-4598-9-155.
- 557 [55] C.T. Fakhry, P. Choudhary, A. Gutteridge, B. Sidders, P. Chen, D. Ziemek, K. Zarringhalam,  
558 Interpreting transcriptional changes using causal graphs: new methods and their practical  
559 utility on public networks, *BMC Bioinformatics.* 17 (2016) 318. doi:10.1186/s12859-016-  
560 1181-8.
- 561 [56] A.H. Bild, G. Yao, J.T. Chang, Q. Wang, A. Potti, D. Chasse, M.-B. Joshi, D. Harpole, J.M.  
562 Lancaster, A. Berchuck, J.A. Olson, J.R. Marks, H.K. Dressman, M. West, J.R. Nevins,  
563 Oncogenic pathway signatures in human cancers as a guide to targeted therapies, *Nature.*  
564 439 (2006) 353–357. doi:10.1038/nature04296.
- 565 [57] H.-Z. Chen, S.-Y. Tsai, G. Leone, Emerging roles of E2Fs in cancer: an exit from cell cycle  
566 control, *Nat. Rev. Cancer.* 9 (2009) 785–797. doi:10.1038/nrc2696.
- 567 [58] B. Amati, M.W. Brooks, N. Levy, T.D. Littlewood, G.I. Evan, H. Land, Oncogenic activity of  
568 the c-Myc protein requires dimerization with Max, *Cell.* 72 (1993) 233–245.  
569 doi:10.1016/0092-8674(93)90663-B.
- 570