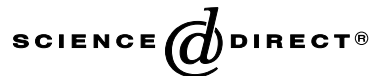


ACADEMIC
PRESS

Available online at www.sciencedirect.com



Journal of Biomedical Informatics xxx (2003) xxx-xxx

Journal of
Biomedical
Informatics

www.elsevier.com/locate/yjbin

2 Automatically identifying gene/protein terms in MEDLINE abstracts

3 Hong Yu,^{a,*} Vasileios Hatzivassiloglou,^{a,1} Andrey Rzhetsky,^b and W. John Wilbur^c

4 ^a Department of Computer Science, Columbia University, 1214 Amsterdam Avenue, New York, NY 10027, USA

5 ^b Department of Medical Informatics, Columbia Genome Center, Columbia University, 622 W. 168th St., VC-5, New York, NY 10032, USA

6 ^c National Center for Biotechnology Information, National Library of Medicine, NIH, Building 38A, Room 5S506, 8600 Rockville Pike,
7 Bethesda, MD 20894, USA

8 Received 4 January 2003

9 Abstract

10 *Motivation.* Natural language processing (NLP) techniques are used to extract information automatically from computer-
11 readable literature. In biology, the identification of terms corresponding to biological substances (e.g., genes and proteins) is a
12 necessary step that precedes the application of other NLP systems that extract biological information (e.g., protein-protein in-
13 teractions, gene regulation events, and biochemical pathways). We have developed GPmarkup (for “gene/protein-full name mark
14 up”), a software system that automatically identifies gene/protein terms (i.e., symbols or full names) in MEDLINE abstracts. As a
15 part of marking up process, we also generated automatically a knowledge source of paired gene/protein symbols and full names (e.g.,
16 *LARD* for *lymphocyte associated receptor of death*) from MEDLINE. We found that many of the pairs in our knowledge source do
17 not appear in the current GenBank database. Therefore our methods may also be used for automatic lexicon generation.

18 *Results.* GPmarkup has 73% recall and 93% precision in identifying and marking up gene/protein terms in MEDLINE abstracts.
19 *Availability:* A random sample of gene/protein symbols and full names and a sample set of marked up abstracts can be viewed at
20 <http://www.cpmc.columbia.edu/homepages/yuh9001/GPmarkup/>. *Contact.* hy52@columbia.edu. Voice: 718-796-2985; fax: 212-939-
21 7028.

22 © 2003 Published by Elsevier Science (USA).

23 *Keywords:* Automatic term recognition; Synonym; Mark up; Information extraction; Knowledge acquisition; Natural language processing

24 1. Introduction

25 The current MEDLINE database includes over 12
26 million computer-readable records in the biomedical
27 domain and is expanding rapidly; it is a rich resource for
28 biological knowledge including protein-protein interac-
29 tions [1], gene regulation events [2], sub-cellular locations
30 of proteins [3], and pathway discovery [4]. One way to
31 automatically extract information stored in MEDLINE
32 is to apply an information extraction system such as a

natural language processing (NLP) parser [5]. Identify- 33
ing gene/protein terms in MEDLINE abstracts is a nec- 34
essary step towards an information extraction system. 35

Genes and proteins are usually represented by sym- 36
bols and names in literature. The names usually are the 37
long forms of their symbols and describe the functions 38
of the genes or proteins. We hypothesize that authors 39
define gene/protein symbols in their articles when the 40
meanings are new in literature and the definitions can be 41
captured by a computer program. We also hypothesize 42
that if not all of the gene/protein symbols appearing in 43
an abstract are defined, the definition may appear in 44
other abstracts. Therefore literature redundancy (e.g., 45
the same genes or proteins are represented by different 46
authors in different articles) makes it plausible that we 47
may obtain automatically a relatively exhaustive gene/ 48
protein symbol and full name table from all of MED- 49
LINE. In this study, we empirically tested all of the 50
above hypotheses. 51

* Corresponding author. Present address: Department of Computer Science, Columbia University, Mudd 471, 116th, New York, NY 10029, USA. Fax: +212-666-0140.

E-mail addresses: Hongyu@cs.columbia.edu (H. Yu), vh@columbia.edu (V. Hatzivassiloglou), Andrey.Rzhetsky@dmi.columbia.edu (A. Rzhetsky), wilbur@ncbi.nlm.nih.gov (W. John Wilbur).

¹ Present address: Department of Computer Science, Columbia University, 1214 Amsterdam Avenue, New York, NY 10027, USA.

52 This study presents an algorithm and its implemen-
 53 tation for automatic identification of gene and protein
 54 terms (i.e., symbols or full names) in MEDLINE ab-
 55 stracts. As a part of the algorithm, we also present a
 56 method for automatically generating a knowledge
 57 source of paired gene/protein symbols (e.g., *LARD*) and
 58 full names (e.g., *lymphocyte associated receptor of death*)
 59 from MEDLINE. Our results show that a large number
 60 of the pairs in our knowledge source do not appear in
 61 LocusLink, a public database of gene/protein symbols
 62 and corresponding full names [6,7].

63 A key step in our marking up methodology is to pair
 64 gene/protein symbols to their names, so that we can use
 65 biological function keywords (e.g., kinase) to differen-
 66 tiate the symbols from other technical terms. For ex-
 67 ample, by mapping abbreviation *PKA* to full name
 68 *protein kinase A*, not to full form *path of the kinematic*
 69 *axis*, we are able to identify *PKA* is a protein term since
 70 keywords *protein* and *kinase* appear in the full form of
 71 *PKA*.

72 We previously have developed a method that auto-
 73 matically maps biomedical abbreviations to full forms.
 74 In this study, we incorporated biological domain
 75 knowledge into the method of mapping abbreviations to
 76 full forms to enhance the mapping between gene/protein
 77 symbols and full names. The biological domain knowl-
 78 edge was obtained from manually reviewing published
 79 guidelines of the nomenclature of genes and proteins.
 80 We then developed a method to differentiate paired
 81 gene/protein symbols and full names from other bio-
 82 medical abbreviations and full forms.

83 To mark up gene/protein terms in MEDLINE ab-
 84 stracts, we first mark up gene/protein symbols and full
 85 names when the full names are defined. We then look up
 86 a knowledge source to mark up the remaining gene/
 87 protein terms. We generate the knowledge source by
 88 extracting all pairs of gene/protein symbols and full
 89 names from over eleven million MEDLINE records
 90 (year 1966–2001).

91 2. Background

92 A number of rule-based, linguistic, statistical, ma-
 93 chine-learning, and hybrid approaches have been de-
 94 veloped to mark up gene/protein terms automatically in
 95 biological text. For example, Fukuda et al. (1998) ap-
 96 plied morphological cues to identify protein terms (e.g.,
 97 if a word contains uppercase letter(s) and special char-
 98 acter(s), the word is a protein term). Gaizauskas et al.
 99 (2000) identified protein terms through suffixes such as –
 100 *ase*. Proux et al. (1998) identified non-English words as
 101 gene terms. Linguistic approaches have mainly applied
 102 part-of-speech tagging [8] or shallow parsing [9] to
 103 identify noun phrases, from which gene/protein terms
 104 were obtained. Hybrid approaches have combined lin-

guistic with rule-based approaches for multi-word gene/
 protein term recognition. For example [8], applied Brill's
 tagger [10] in combination with rules such as “connect
 non-adjacent annotations if every word between them is
 either noun, adjective, or a numeral” to identify multi-
 word protein terms such as *ras guanine nucleotide ex-
 change factor SOS*. Tanabe and Wilbur [11] retrained
 Brill's tagger on the biomedical domain for gene/protein
 name-identification. Statistical approaches have clus-
 tered abstracts for keyword identification [12]. Machine-
 learning approaches have applied naïve Bayes [9], Hid-
 den Markov Models [13], and decision trees [14], to
 classify gene/protein terms. Other approaches include
 lookup in knowledge sources such as GenBank and
 SWISSPROT [15].

Our method of marking up gene/protein names is a
 mixture of pattern-recognition and knowledge-based
 approaches. We first map gene/protein symbols to full
 names when the full names are defined. Those gene/
 protein terms are then marked up. The rest of gene/
 protein terms are identified from the gene/protein sym-
 bol and full name knowledge source which we extracted
 automatically from MEDLINE.

2.1. Systems that automatically map gene and protein symbols to full names

A number of systems have been developed for auto-
 matic mapping between abbreviations and full names
 [16–23]. Those systems applied a variety of approaches
 including linguistic, rule, and statistical methods and
 reported precisions from 70–97%. Most of those systems
 tend to be domain independent and therefore may not
 perform ideally in a restricted domain such as biology.
 For example, most of pattern-recognition approaches
 [18,19] do not capture *ryk* (for *receptor tyrosine kinase*
related gene) since *y* represents *tyrosine* and *y* is not the
 first letter of *tyrosine*. In addition, most of the systems
 do not differentiate gene/protein symbols from other
 abbreviations and full names.

A system that was developed specifically for mapping
 protein symbols to full names is PNAD-CSS (for “pro-
 tein full name abbreviation dictionary construction
 support system”) [24]. PNAD-CSS used morphological
 features to recognize proper nouns as protein terms in
 biological abstracts [8]. Knowing a phrase may contain a
 protein symbol and full name, PNAD-CSS recognized
 parentheses and determined whether the parenthetical
 phrase was an abbreviation of the outer phrase. To map
 a protein symbol to its name, PNAD-CSS broke up
 words of the preceding phrase, and determined whether
 the parenthetical abbreviation candidate maps to the
 initial letters of the broken-up phrase. For example,
 consider the phrase “*megestrol acetate (megace)*.”
 PNAD-CSS parsed “*megestrol acetate*” as “*meges trol ac*
etate,” which is then matched to “*megace*.” For example,

Table 1

Guidelines that are useful for applying computational approaches to map a gene or a protein symbol to its full name

1. A gene symbol should stand for a description of a phenotype, a gene product or a gene function [26].
2. A gene symbol shall be short (between three to six characters) [26–32].
3. A gene symbol is an abbreviation of its full name [28].
4. If the symbol of a gene contains a character or property for which there is a recognized abbreviation, the abbreviation should be used; for example, the single-letter abbreviation for amino acids used in aminoacyl residues or approved biochemical Abbreviations such as GLC for glucose, GSH for glutathione [31] and *Bp* for *binding protein* [32].
5. The initial character should always be a letter [29–33].
6. All Greek symbols should be changed to letters in the Latin alphabet [31].
7. Amino acids have their special symbols [34].
8. The protein symbol is the same as the gene symbol [33].
9. The creator of a gene full name shall follow the guidelines and get consultation from curator of the guideline before journal publication [26].
10. Gene full names should be included in the abstracts of any relevant papers [26].

159 “*meg*,” “*ac*,” and “*e*” in “*megace*” match the initial
160 letter(s) of “*meges*,” “*ac*,” and “*etate*,” respectively.

161 We find that PNAD-CSS has some limitations: it
162 applies morphological cues for protein term recognition
163 and the morphological cues may falsely identify as
164 protein symbols other substances (e.g., *LSD-25* for *ly-*
165 *sergic acid diethylamide*), cell types (e.g., *BHK-21* for
166 *baby-hamster kidney-cell line*), procedures (e.g., *PCR* for
167 *polymerase chain reaction*) as well as clinical syndromes
168 and diseases (e.g., *CHF* for *congestive heart failure*). This
169 is because many abbreviations that are not gene/protein
170 symbols consist of upper-case letters and numbers. The
171 PNAD-CSS’ pattern-matching rules also did not contain
172 special rules for protein names (for example, *y* repre-
173 sents *tyrosine*).

174 Previously, we have developed a system, AbbRE (for
175 “abbreviation and full name recognition and extrac-
176 tion,” see [25]), that pairs biomedical abbreviations with
177 full names. AbbRE first selected parenthetical expres-
178 sions and the phrases preceding the parenthesis as can-
179 didate abbreviations and full names. It then applied a set
180 of the pattern-matching rules to map abbreviations to
181 full names. The rules were obtained from the common
182 conventions authors use to create abbreviations. The
183 following rules were included: (1) *the first letter of an*
184 *abbreviation matches the first letter of a meaningful word*
185 *of the full name;* (2) *the abbreviation matches the first*
186 *letter of each word in the full name;* (3) *the abbreviation*
187 *letter matches consecutive letters of a word in the full*
188 *name and* (4) *the abbreviation letter matches a middle*
189 *letter of a word in the full name if the first letter of the*
190 *word matches the abbreviation.* AbbRE had 70% recall
191 and 95% precision in identifying paired abbreviations
192 and full names in biomedical articles.

193 Though AbbRE’s pattern-matching rules did not
194 contain special rules for protein names, AbbRE is robust
195 and extensible. In this study (i.e., GPmarkup), we man-
196 ually examined the published guidelines of the nomen-
197 clature of genes and proteins and added to AbbRE special
198 rules to enhance its mapping gene/protein symbols to full
199 names. In addition, we added in rules for differentiating
200 gene/protein terms from other biomedical terms.

3. Methods and results

201

Our method section consists of six sub-sections: (1)
Mapping gene/protein symbols to full names as well as
abbreviations to full names. (2) Generating a knowledge
source of paired abbreviations and full names from
MEDLINE abstracts. (3) Filtering out other abbrevia-
tion-full name pairs to produce a knowledge source of
paired gene/protein symbols and full names. (4) Mark-
ing up gene/protein terms in MEDLINE abstracts. (5)
Evaluating GPmarkup. (6) Measuring the percentage of
defined gene/protein symbols in MEDLINE abstracts.

202
203
204
205
206
207
208
209
210
211

3.1. Mapping gene/protein symbols to full names

212

To understand how gene/protein abbreviation-full
name pairs are created in the first place, we examined a
number of published guidelines for the nomenclature of
genes and proteins. We found those guidelines are al-
most always species-specific (that is applicable only to
genes and proteins from, say, yeast, and not rat). Spe-
cies-specific may be caused by the fact that the com-
mittees for the nomenclature are formed by experts
specializing on a particular model organism. Table 1
lists guidelines that were useful for mapping abbrevia-
tions to full forms.

213
214
215
216
217
218
219
220
221
222
223

Analysis of the published guidelines allowed us to
identify some special abbreviations that are used for
gene/protein nomenclature (see Table 2) and to develop
the pattern-matching rules that map gene/protein sym-
bols to names.

224
225
226
227
228

3.1.1. Special abbreviations

229

see Table 2.

230

3.1.2. Pattern-matching rules

231

GPmarkup applies a set of pattern-matching rules to
map gene/protein symbols to full names when the full
names are defined within the documents. The pattern-
matching rules adapted AbbRE’s (as described in Sec-
tion 2.1) with the following modifications and exten-
sions:

232
233
234
235
236
237

- 238 *Rule 1: Any number and special character is ignored* 284
 239 *for mapping gene/protein symbols to full names.* 285
 240 We added in a rule to map letters only. We ignored
 241 numbers and special characters (e.g., “+”) due to the
 242 following two reasons:
 243 (1) Many numbers and special characters in a gene or a
 244 protein symbol do not appear in their full names.
 245 For example, *CYP2C19* for *cytochrome P450, sub-*
 246 *family IIC (mephenytoin 4-hydroxylase)*, where
 247 “19” is not represented and “2” is represented by
 248 “II.”
 249 (2) Many numbers in gene or protein symbols order dif-
 250 ferently in their full names (e.g., *ALOX12* for *ara-*
 251 *chidonate 12-lipoxygenase*, where “12” in the
 252 symbol “*ALOX12*” is after “*LOX*” that represents
 253 *lipoxygenase*, but before “*lipoxygenase*” in the full
 254 name “*arachidonate 12-lipoxygenase*”).
 255 *Rule 2: Special abbreviation substitutions* 299
 256 We substitute some nouns with their special abbrevi-
 257 ations when we apply the pattern-matching rules. For
 258 example, instead of mapping *DYRK1A* to *dual-specific-*
 259 *ity tyrosine phosphorylation regulated kinase 1A*, we at-
 260 tempt to map *DYRK1A* to *dual-specificity Y*
 261 *phosphorylation regulated kinase 1A*, where *tyrosine* has
 262 been replaced by *Y*. After the mapping, we recover the
 263 original terms. 300
 264 In reality, not all the authors use the special abbrevi-
 265 ations (listed in Table 2) for their nomenclature. An
 266 example is *PTK2B* for *protein tyrosine kinase 2 β*, where
 267 *tyrosine* is represented by its common abbreviation *T*
 268 instead of *Y*. Therefore, our algorithm considers both
 269 types of mapping (with and without substitution of a
 270 special noun with a shorthand) and selects the best
 271 matching version.
 272 For example, we attempt to map *PTK2B* to both
 273 *protein tyrosine kinase 2 β* and *protein Y kinase 2 β*; we
 274 map *DYRK1A* to both *dual-specificity tyrosine phos-*
 275 *phorylation regulated kinase 1A* and *dual-specificity Y*
 276 *phosphorylation regulated kinase 1A*.
 277 When a full name has more than one word that has
 278 many abbreviations, we include all of the combinations
 279 for substitution. For example, in case of *NK AIF* for
 280 *sodium-potassium ATPase inhibitory factor*, we attempted
 281 to map *NK AIF* to *sodium-potassium ATPase inhibitory*
 282 *factor*, *Na-potassium ATPase inhibitory factor*, *sodium-K*
 283 *ATPase inhibitory factor*, and *Na-K ATPase inhibitory*
 284 *factor*. We found that *Na-K ATPase inhibitory factor* was
 285 mapped and we recovered the original full name.
 3.1.3. *Parenthetic pattern* 286
 Prior to pattern-matching rules, GPmarkup selects
 candidate abbreviations and full names. For this task,
 GPmarkup recognizes special patterns such as “<ab-
 breviation>(<full name>)” or “<full name>(<abbrevi-
 ation>”. Recall AbbRE also recognized these patterns.
 However, AbbRE can not recognize gene/protein terms
 that incorporate nested parentheses. For example, Ab-
 bRE fails to map *acyl-coenzyme A (acyl-CoA) dehydro-*
genases to *ACD* from the following string extracted
 from [35] *the expression of various acyl-coenzyme A*
(acyl-CoA) dehydrogenases (ACD) since it parses into
 the following two components: 298
the expression of various acyl-coenzyme A (acyl-CoA) and dehy- 299
drogenases (ACD) 300
 To correct for this shortcoming, we introduced into
 the newer algorithm (GPmarkup) an additional rule to
 recognize gene/protein full names that incorporate pa-
 rentheses. It then parses the above string into the fol-
 lowing two components: 305
the expression of various acyl-coenzyme A (acyl-CoA) and the ex- 306
pression of various acyl-coenzyme A (acyl-CoA) dehydrogenases 307
(ACD) 308
 where the phrases preceding and within the parentheses
 in each component incorporate candidate abbreviations
 and full names, to which GPmarkup further applies its
 pattern-matching rules to map abbreviations to full
 names. 313
 3.2. *Generating a knowledge source of paired abbrevia-* 314
tions/full names from MEDLINE abstracts 315
 We applied GPmarkup to 11 million MEDLINE re-
 cords (1966–2001), which contain the same number of
 titles and over six million abstracts (note that not all
 MEDLINE records contain abstracts). We obtained a
 knowledge source that consisted of 574,327 unique pairs
 of abbreviations and full names. The most frequently
 defined abbreviations were *PCR* (*polymerase chain re-*
action, which appeared in 7988 abstracts) and *NO* (*nitric*
oxide, which appeared in 7855 abstracts). 324

Table 2
Special abbreviations that are used in gene/protein nomenclature

Type	
Amino acids	We use all one letter codes where these differ from the first letter of the amino acid. For example, <i>tyrosine</i> — <i>Y</i> (<i>SYK</i> for <i>spleen tyrosine kinase</i>)
Two chemical symbols used	<i>Sodium</i> — <i>Na</i> , <i>potassium</i> — <i>K</i> (<i>NK AIF</i> for <i>sodium-potassium ATPase inhibitory factor</i>)
Three other symbols used	<i>Inhibitor</i> — <i>N</i> or <i>NH</i> , <i>box</i> — <i>X</i> (<i>CDKN1A</i> for <i>cyclin-dependent kinase inhibitor 1A</i> (<i>p21</i> , <i>Cip1</i>), <i>CDX1</i> for <i>caudal type homeo box transcription factor 1</i>)

325 3.3. Filtering out other abbreviation-full name pairs to
 326 produce a knowledge source of paired gene/protein
 327 symbols and full names

328 The algorithm outlined above also identifies a large
 329 number of general abbreviations that are not gene/pro-
 330 tein symbols and full names. We therefore developed a
 331 rule-based approach to partition our knowledge source
 332 of abbreviation-full name pairs into gene/protein sym-
 333 bol-full name pairs and other abbreviation-full name
 334 pairs.

335 Our rule-based approach combines morphological
 336 cues, functional keywords, and position-functional
 337 keywords to filter out non-gene/protein terms. The ap-
 338 proach is described as follows:

339 If an abbreviation contains a number, the abbrevia-
 340 tion and full name is a gene/protein symbol-full name
 341 pair only if the full name contains one or more of the
 342 following keywords (denoted as set K1): *protein(s)*,
 343 *gene(s)*, *peptide(s)*, *molecule(s)*, *enzyme(s)*, *ligand(s)*,
 344 *compound(s)*, *receptor(s)*, *channel(s)*, *transcriptor(s)*,
 345 *regulator(s)*, *inhibitor(s)*, *antibody*, *antibodies*, *globu-
 346 lin(s)*, *factor(s)*, *motif*, *domain(s)*, *compound(s)*, *seg-
 347 ment(s)*, *subunit(s)*, *locus*, *loci*, *cassette(s)*, *chain*,
 348 *complex(es)*, *homeobox(es)*, *box(es)*, *member(s)*, *dele-
 349 tion*, *axon*, *family*, *families*, *chromosome(s)*, *sequence*,
 350 α , β , γ , *interleukin* and any words except for *disease*
 351 that ends in *-ase*.

352 If an abbreviation does not contain a number, the ab-
 353 breviation and full name is gene/protein symbol-full
 354 name pair only if the last word of the full name is a
 355 keyword in set K1.

356 We obtained functional keywords by manually ex-
 357 amining all of the entries in LocusLink. Note that some
 358 keywords (e.g., “gene”) in set K1 can appear as both the
 359 last word or the middle word of a gene/protein term
 360 (e.g., *Btg4* for *B-cell translocation gene 4* and *AFG3L1*
 361 for *AFG3 (ATPase family gene 3, yeast)-like 1*). On the
 362 other hand, some keywords (e.g., “chromosome”) do
 363 not appear as the last word of, but only within a gene/
 364 protein term (e.g., *C10ORF2* for *chromosome 10 open
 365 reading frame 2*).

366 We applied the rules to abbreviations and full names
 367 and generated a knowledge source of 86,767 unique
 368 pairs of gene/protein symbols and full names. The most
 369 frequently defined gene/protein symbols included *egf*
 370 (for *epidermal growth factor*, appears in 2023 ab-
 371 stracts), *il* (for *interleukin*, appears in 2183 abstracts),
 372 and *ldl* (for *low density lipoprotein*, appears in 2673
 373 abstracts).

374 3.4. Marking up gene/protein terms in MEDLINE
 375 abstracts

376 We further developed and implemented an algorithm
 377 to mark up gene/protein terms in MEDLINE abstracts.

GPmarkup first maps abbreviations to full names and 378
 then performs the markup for any abbreviation with an 379
 identified full name (details in Sections 3.2 and 3.3). For 380
 the remaining terms in abstracts, we looked up the 381
 knowledge sources of paired abbreviations and full 382
 names and paired gene/protein symbols and names. As 383
 an effort to achieve a higher precision, we only looked 384
 up multi-word gene/protein terms, since a single word 385
 term could be ambiguous (for example, *aap* denotes 386
antiarrhythmic peptide or *automatic action potential*, the 387
 former is a protein name, and the latter is not). 388

When a string can be mapped to several terms stored 389
 in our knowledge sources, GPmarkup favors longer 390
 term mapping and markup. It does not mark up a term 391
 which is used as a modifier of entity other than genes 392
 and proteins. For example, GPmarkup does not markup 393
 the protein term *amyloid β protein* in a string of *cerebral
 amyloid β protein angiopathy*, because the protein name 394
 is used as a modifier for the disease term *angiopath*. 395

GPmarkup applies direct matching (i.e., the string in 397
 text exactly appears in our knowledge sources) except 398
 that GPmarkup includes a word that immediately fol- 399
 lows a gene or a protein symbol or full name if the word 400
 either consists of a number or is a functional keyword 401
 including “gene,” “protein,” “homologue,” and “re- 402
 ceptor.” For example, knowing a β and *il12 p40* as gene 403
 or protein symbols, GPmarkup also identifies a $\beta 40$ and 404
il12 p40 homologue. 405

3.5. GPmarkup evaluation 406

We performed evaluation in the following three 407
 steps: (1) mapping abbreviations to full names, (2) fil- 408
 tering out other terms to produce a knowledge source 409
 of paired gene/protein symbols and names, and (3) 410
 marking up gene/protein terms in MEDLINE ab- 411
 stracts. We therefore evaluate GPmarkup phase by 412
 phase. We also compared the knowledge source of 413
 paired gene/protein symbols and full names with the 414
 ones in LocusLink. We evaluated by recall (i.e., num- 415
 ber of correct answers identified by our system divided 416
 the total number of correct answers) and precision (i.e., 417
 number of correct answers divided by the total number 418
 of answers specified by our system). We estimated 419
 confidence intervals for these measures based on the 420
 binomial distribution. 421

3.5.1. Mapping abbreviations to full names 422

We randomly (by time of publication) selected 30 423
 MEDLINE abstracts and asked three biomedical ex- 424
 perts (all with PhD or MD) to map abbreviations to full 425
 names when the full names are defined within the ab- 426
 stracts. The gold standard was determined by a majority 427
 vote of experts. GPmarkup correctly mapped 56 ab- 428
 breviations and full names out of a total of 59 pairs that 429
 were determined by experts. GPmarkup wrongly iden- 430

431 tified one pair that was not an abbreviation and full
432 name. GPmarkup's recall and precision in identifying
433 and extracting abbreviations and full names were, with
434 95% confidence intervals, 0.95 (0.86–0.99) and 0.98
435 (0.91–1.00), respectively.

436 3.5.2. Filtering out other terms

437 We then evaluated our rule-based approach for par-
438 titioning the knowledge source of abbreviation-full name
439 pairs into gene/protein symbol-full name pairs and other
440 abbreviation-full name pairs. We randomly selected 1000
441 pairs of gene/protein symbols and full names and 1000
442 pairs of other abbreviations and full names partitioned
443 by GPmarkup and evaluated recall and precision of the
444 partitioning. We asked experts (see 3.5.1) for help in
445 defining a gold standard. Table 3 lists the results of the
446 evaluation. Note that GPmarkup included some in-
447 complete-matches of abbreviations and full names (e.g.,
448 {*il-6*, *interleukin*}). Since the ratio of gene/protein sym-
449 bol-names to other abbreviation-full name pairs was
450 1:5.6 (86,767/[574,327–86,767]); the numbers were de-
451 scribed in Sections 3.2 and 3.3), GPmarkup had an ac-
452 curacy of 0.95 ± 0.02 , with 95% confidence. The figure
453 0.95 comes from the ratio $(982 + 949 * 5.6)/(1000 +$
454 $1000 * 5.6)$ which is based on the numbers in Table 3
455 and their relative frequencies as just computed.

456 3.5.3. Marking up gene/protein terms in MEDLINE 457 abstracts

458 We then evaluated GPmarkup in marking up gene/
459 protein terms in MEDLINE abstracts. We randomly (by
460 time of publication) selected 50 MEDLINE abstracts,
461 which consists of a total of 539 sentences (including the
462 title). Some selected abstracts did not cover biological
463 domain and therefore did not have gene/protein terms at

464 all. Therefore, we did not select only biological abstracts
465 for evaluation because we judge a false markup is as bad
466 as a missing markup. We therefore judged that a ran-
467 dom selection of abstracts best reflects our system's re-
468 call and precision.

469 Table 4 lists the evaluation results of the 50 abstracts.
470 GPmarkup applies XML format for term mark up. For
471 example, the tag “phr”(for “phrase”) has attributes in-
472 cluding “sem” (for “semantic category”) that has value
473 “gp” (for “gene and protein terms”) and “t” (for “tar-
474 get”) that represents gene/protein full names. We count
475 any appearance of gene/protein terms. For example, if
476 protein “*amyloid β protein*” appears three times in the
477 abstract, we count three instead of one for this case. We
478 posted a sample set of marked up abstracts at [http://](http://www.cpmc.columbia.edu/homepages/yuh9001/GPmarkup/)
479 [www.cpmc.columbia.edu/homepages/yuh9001/GPmark-](http://www.cpmc.columbia.edu/homepages/yuh9001/GPmarkup/)
480 [up/](http://www.cpmc.columbia.edu/homepages/yuh9001/GPmarkup/).

481 From Table 4, if we count a partial-matching as a
482 match, the recall and the precision of GPmarkup were,
483 with 95% confidence, 0.73 ± 0.05 $(222 + 15)/(222 +$
484 $15 + 88)$ and 0.93 ± 0.03 $(222 + 15)/(222 + 15 + 17)$,
485 respectively. We found all partial matches represent valid
486 proteins. However, if we do not include a partial-match-
487 ing as a match, the recall and precision of GPmarkup
488 were, with 95% confidence, 0.68 ± 0.05 $222/(222 + 15 +$
489 $88)$ and 0.87 ± 0.04 $(222/(222 + 15 + 17))$, respectively.

490 3.5.4. Comparing gene/protein symbols and full names 491 extracted from MEDLINE with LocusLink

492 We downloaded the knowledge source of paired gene/
493 protein symbols and full names from LocusLink [36].
494 LocusLink is maintained by the National Center for
495 Biotechnology Information. It presents information on
496 official nomenclature of genes and lists a total of 115,890
497 manually annotated paired gene symbols and full

Table 3

Evaluation results of GPmarkup in filtering the knowledge source of paired abbreviations and full names to produce a knowledge source of paired gene/protein symbols and full names

Evaluation cases	Expert judgments		
	Number of gene/protein symbol-full name pairs	Number of other abbreviation-full name pairs	Number of non abbreviation-full name pairs
1000 pairs of gene/protein symbols and full names as identified by GPmarkup	982	9 (e.g. <i>srg</i> for <i>spent restaurant grease</i>)	9 (e.g., <i>gene</i> for <i>genes</i>)
1000 pairs of other abbreviations and full names as identified by GPmarkup	1 (i.e., <i>A-Igg</i> for <i>Anti-human Igg</i>)	949	50 (e.g., <i>ph2</i> for <i>phages</i>)

Table 4

Evaluation results of GPmarkup

Type of category	GPmarkup identified
Complete-matching (e.g., $\langle \text{phr sem} = \text{"gp"} \text{ t} = \text{"signaling lymphocyte activation molecule"} \rangle \text{slam} \langle / \text{phr} \rangle$)	222
Partial-matching ^a (e.g., $\langle \text{phr sem} = \text{"gp"} \rangle \text{interleukin 1} \langle / \text{phr} \rangle \text{ receptor ii}$)	15
Missing (e.g., <i>2b4</i>)	88
False-matching ^b (e.g., $\langle \text{phr sem} = \text{"gp"} \rangle \text{acupuncture points and channels} \langle / \text{phr} \rangle$)	17

^a The correct full name is “interleukin 1 receptor ii”.

^b False-matching includes those non-gene and non-protein terms that are identified by GPmarkup.

498 names, though we found that only 65,987 entries have
499 both gene/protein symbols and full names.

500 We randomly selected 100 entries that incorporate
501 both symbols and full names from the LocusLink and
502 manually identify their existence in our knowledge
503 source of paired gene/protein symbols and full names.
504 We also randomly selected 100 unique gene/protein
505 symbol and full name pairs from our knowledge source
506 and manually identified their existence in LocusLink.

507 We found that 62 out of 100 selected pairs in our
508 knowledge source did not appear in LocusLink. Exam-
509 ples included {*ACY1-ACP*, *acyl-acyl carrier protein*},
510 {*GCDFP*, *gross cyst disease fluid protein*}, {*CCK-OP*,
511 *cholecystokinin octopeptide*} and {*l-PK*, *l* pyruvate ki-
512 *nase*} though some of the missing pairs represent protein
513 products instead of direct genes. For example, {*l-PK*, *l*
514 *pyruvate kinase*} is a spliced product of its gene {*PKLR*,
515 *pyruvate kinase*}² which appears in LocusLink and there
516 is no gene for {*CCK-OP*, *cholecystokinin octopeptide*}³.
517 Eight pairs partially matched to LocusLink. For exam-
518 ple, *PPI*, *peptide prolyl cis trans isomerase* appears in our
519 knowledge source. In LocusLink, we found {*PPIa*,
520 *peptidylprolyl isomerase a (cyclophilin a)*}.”

521 On the other hand, we found that only 40 LocusLink
522 entries could be found in our knowledge source (16 of
523 them have variations). We judged that four of those 60
524 failed entries are not gene/protein symbols and full
525 names (e.g., {*shs*, *sutherland-haan x-linked mental re-*
526 *tardation syndrome*}). To find whether the remaining 56
527 entries exist in MEDLINE, we searched 12 million
528 MEDLINE records (1966–2002). We applied direct
529 matching (case insensitive) and manually analyzed ab-
530 stracts that contained either the symbol or the full name
531 of those 56 failed entries. We failed to find the existence
532 of 50 of them in MEDLINE, either symbols or full
533 names. Examples include {*2700088m22rik*, *riken cDNA*
534 *2700088m22 gene*} and {*atp5b1l*, *atp synthase, h+*
535 *transporting, mitochondrial f1 complex, β polypeptide-*
536 *like 1*}. Of the rest of six entries, we could find symbols
537 in MEDLINE, but failed to find full names. Examples
538 include {*aspa*, *aspartoacylase (aminoacylase 2, canavan*
539 *disease*)} and {*assp6*, *argininosuccinate synthetase*
540 *pseudogene 6*}, for the former we found the full name
541 with variations, for the latter we found that the full
542 name did not exist in the MEDLINE record where the
543 symbol appeared.

544 3.6. The percentage of undefined gene/protein symbols and 545 full names

546 If all the gene/protein symbols and full names were
547 defined in MEDLINE abstracts, then GPmarkup would

also serve the purpose for disambiguation by assigning 548
full names to symbols. However, not all the gene/protein 549
symbols are defined in the abstracts. 550

551 We measured the percentage of defined gene/protein
552 symbols in MEDLINE abstracts. We randomly selected
553 100 abstracts (according to the time of publication) from
554 a total of 782,560 MEDLINE abstracts (1966–2001)
555 that were retrieved by the keyword “protein.” Those
556 abstracts contain 1069 sentences (including titles). We
557 measured the percentage of undefined gene/protein
558 symbols. We counted unique appearance of gene/protein
559 symbols within abstracts. Based on the authors’ judg-
560 ment, the numbers of defined and undefined gene/pro-
561 tein symbols were 92 and 27, respectively. The
562 percentage of defined gene/protein symbols and full
563 names was, with 95% confidence, 0.77 ± 0.08 .

564 4. Discussion

565 Many public databases such as GenBank have gene/
566 protein synonym knowledge sources. However, the da-
567 tabases are largely maintained manually and therefore
568 are not always up to date. GPmarkup can generate
569 automatically a knowledge source of paired gene/protein
570 symbols and full names from MEDLINE abstracts. The
571 automated fashion may reduce manual efforts. In addi-
572 tion, GPmarkup may capture the most up-to-date gene/
573 protein symbols and full names if the full names are
574 defined in abstracts and follow the guidelines of no-
575 menclature of genes and proteins.

576 We also found that a majority of gene/protein sym-
577 bols and full names extracted in our knowledge source
578 did not appear in LocusLink. Recall LocusLink consists
579 of a large number of mainly manually annotated paired
580 gene/protein symbols and full names. In addition, we
581 found a majority of pairs in LocusLink did not appear
582 in our knowledge source either; most of those pairs did
583 not even appear in MEDLINE by keyword search. The
584 results suggest that there is a gap between LocusLink
585 knowledge source and the actual text. This difference
586 may make it difficult to apply LocusLink directly for
587 looking up terms in MEDLINE. On the other hand,
588 since our knowledge source of paired gene/protein
589 symbols and names were directly extracted from
590 MEDLINE, they may be more useful as a knowledge-
591 based markup.

592 One limitation of GPmarkup is that not all the gene/
593 protein symbols and full names are defined in the ab-
594 stracts and therefore GPmarkup may not capture some
595 gene/protein symbols and full names. However, two
596 other factors alleviate this problem: authors are en-
597 couraged to define gene/protein full names in the ab-
598 stracts of any relevant papers [26], and the literature is
599 redundant. Therefore, applying GPmarkup to all of
600 MEDLINE abstracts is likely to capture a majority of

² GenBank Accession No. U47654.

³ For details see <http://arbl.cvmb.colostate.edu/hbooks/pathophys/endocrine/gi/cck.html>.

601 gene/protein symbols and full names that appear in the
602 text.

603 GPmarkup may also miss gene/protein symbols and
604 full names when authors do not follow the guidelines for
605 naming genes and proteins. To capture these gene/pro-
606 tein symbols and full names, we may integrate into
607 GPmarkup statistical approaches such as Hisamitsu and
608 Niwa's approach [18,20] of selecting phrases associated
609 with parentheses that were statistically significant. In
610 addition, GPmarkup may also miss abbreviations and
611 full names that are introduced through syntactic pat-
612 terns (e.g., appositions). In the near future we plan to
613 utilize the approaches of [37] that enumerated syntactic
614 patterns for abbreviation detection.

615 Other limitations include the ambiguity in usage of
616 gene/protein terms. For example, we do not differentiate
617 a gene term from a protein one. We do not differentiate
618 a general gene/protein term (e.g., *growth factors*) from a
619 specific one (e.g., *protein kinase A*). We also do not
620 identify to which organism, tissue, cell type, and sub-
621 location a gene/protein term refers. We propose to in-
622 tegrate the approach of [38] for disambiguating gene/
623 protein terms. We also hope to develop statistical NLP
624 approaches for further disambiguation.

625 Our study shows that many gene/protein symbols
626 (77%) are defined within the abstracts, GPmarkup can
627 map a majority of gene/protein symbols to full names.
628 GPmarkup does not mark up undefined gene/protein
629 symbols if the symbols have several full names in the
630 knowledge source of abbreviation-full name pairs. For
631 example, *aap* denotes *antiarrhythmic peptide*, *alkyl ac-*
632 *ceptor protein*, *alzheimer amyloid precursor protein*, *am-*
633 *inoantipyrine*, and *automatic action potential* in our
634 knowledge source and GPmarkup thus does not mark
635 up “*aap*” as a gene/protein term when it is not defined in
636 the abstract. We therefore sacrifice GPmarkup's recall
637 for high precision. In the future, we will integrate a
638 disambiguation method that assigns the full names from
639 our knowledge source to the ambiguous symbols. Once
640 a symbol is assigned to its full name, we can apply our
641 rule-based approach (see Section 3.3) determining whe-
642 ther the symbol is a gene/protein term.

643 Note that we recognized a gene/protein term if the
644 term actually represents a gene/protein in the abstract.
645 We described earlier that we did not mark up “*cerebral*
646 *amyloid β protein angiopathy*” as a protein name even
647 though “*cerebral amyloid β protein*” by itself is a protein
648 name. Other researchers may do differently [11].

649 5. Conclusion

650 This study shows that GPmarkup is efficient (73%
651 recall and 93% precision) in marking up gene/protein
652 terms in MEDLINE abstracts. Our results may provide
653 a useful supplement to manually curated resources such

as LocusLink (GenBank). A method to more accurately
654 identify the full names of undefined abbreviations would
655 increase the recall of GPmarkup and enhance its use-
656 fulness. 657

Acknowledgments

We want to thank Dr. Carol Friedman and Ivan
659 Iossifov for valuable discussions. This research was
660 supported in part by National Science Foundation In-
661 formation Technology Research Grant EIA-0121687
662 and National Institutes of Health Grant RO1
663 GM61372-01A2. 664

References

- [1] Blaschke C et al. Automatic extraction of biological information
666 from scientific text: protein–protein interactions. Proc Int Conf
667 Intell Syst Mol Biol 1999:60–7. 668
- [2] Stapley BJ, Benoit G. Biobibliometrics: information retrieval and
669 visualization from co- occurrences of gene names in Medline
670 abstracts. Pac Symp Biocomput 2000:529–40. 671
- [3] Stapley BJ, Kelley LA, Sternberg MJE. Predicting the sub-cellular
672 location of proteins from text using support vector machines. In:
673 PSB, Hawaii, 2002. 674
- [4] Ng SK, Wong M. Toward routine automatic pathway discovery
675 from on-line scientific text abstracts. Genome Inform Ser Work-
676 shop Genome Inform 1999;10:104–12. 677
- [5] Carol Friedman, P.K., Michael Krauthammer, Hong Yu, Andrey
678 Rzhetsky. GENIES: A Natural-Language Processing System for
679 the Extraction of Molecular Pathways from Complete Journal
680 Articles. In: ISMB, 2001. 681
- [6] Maglott DR et al. NCBI's LocusLink and RefSeq. Nucleic Acids
682 Res 2000;28(1):126–8. 683
- [7] Pruitt KD et al. Introducing RefSeq and LocusLink: curated
684 human genome resources at the NCBI. Trends Genet
685 2000;16(1):44–7. 686
- [8] Fukuda K et al. Toward information extraction: identifying
687 protein names from biological papers. Pac Symp Biocomput
688 1998:707–18. 689
- [9] Nobata C, C, CN, TJI. Automatic term identification and
690 classification in biology texts. In: Proceedings of the Natural
691 Language Pacific Rim Symposium (NLPRS'99), 1999. 692
- [10] Brill E. Transformation-based error-driven learning and natural
693 language processing: a case study in part of speech tagging.
694 Comput Linguistics 1995. 695
- [11] Tanabe L, Wilbur WJ. Tagging gene and protein names in
696 biomedical text. Bioinformatics 2002;18:1124–32. 697
- [12] Andrade MA, Valencia A. Automatic extraction of keywords
698 from scientific text: application to the knowledge domain of
699 protein families. Bioinformatics 1998;14(7):600–7. 700
- [13] Collier NH, Nobata C, Tshjii J. Extracting the names of genes and
701 gene products with a hidden markov model. In: Proceedings of the
702 18th International Conference on Computational Linguistics
703 (COLING'2000), 2000, p. 201–7. 704
- [14] Nobata C, Collier NH, Tsujii J. Comparison between tagged
705 corpora for the named entity task. In: Proceedings of the
706 workshop on comparing corpora (at ACL'2000), Kilgariff, A.,
707 Berber Sardinha, 2000. 708

- 709 [15] Krauthammer M et al. Using BLAST for identifying gene and
710 protein names in journal articles. *Gene* 2000;259(1–2):245–
711 52.
- 712 [16] Pakhomov S. Semi-supervised maximum entropy based approach
713 to acronym and abbreviation normalization in medical text. In:
714 Proc. 40th annual meeting of the association for computational
715 linguistics, Philadelphia, Pennsylvania, USA, 2002.
- 716 [17] Bowden PR, Eventt L, Halsted P. Automatic acronym acquisition
717 in a knowledge extraction program. In: *CompuTerm98*, Montreal,
718 Ontario, 1998.
- 719 [18] Hisamitsu T, Niwa Y. Extraction of useful terms from parenthetical
720 expression by using simple rules and statistical measures. In:
721 *CompuTerm98*, Montreal, Canada, 1998.
- 722 [19] Schwartz AS, Hearst MA. A simple algorithm for identifying
723 abbreviation definitions in biomedical text. In: *Pac Symp Bio-*
724 *comput*, 2003.
- 725 [20] Liu H, Friedman C. Mining terminological knowledge in large
726 biomedical corpora. In: *Pac Symp Biocomput*, 2003.
- 727 [21] Park Y, Byrd RJ. Hybrid text mining for finding abbreviations
728 and their definitions. In: *Proceedings of the 2001 Conference on*
729 *Empirical Methods in Natural Language Processing*, Pittsburgh,
730 PA, 2001.
- 731 [22] Larkey LS, Ogilvie P, Price MA. Acrophile: an automated
732 acronym extractor and server. In: *Proceedings of the Fifth*
733 *ACM International Conference on Digital Libraries*, 2000.
- 734 [23] Wren JD, Garner HR. Heuristics for identification of acronym-
735 definition patterns within text: towards an automated construction
736 of comprehensive acronym-definition dictionaries. *Methods*
737 *Inf Med*, 2002, in press.
- 738 [24] Yoshida M, Fukuda K, Takagi T. PNAD-CSS: a workbench for
739 constructing a protein name abbreviation dictionary. *Bioinform-*
740 *atics* 2000;16(2):169–75.
- 741 [25] Yu H, Hripcsak G, Friedman C. Mapping abbreviations to full forms
742 in biomedical articles. *J Am Med Inform Assoc* 2002;9(3):262–72.
- 743 [26] Kohli J. Genetic nomenclature and gene list of the fission yeast
744 *Schizosaccharomyces pombe*. *Curr Genet* 1987;11(8):575–89.
- [27] Nomenclature CPG. Chicken Poultry Genome Nomenclature
745 Web site. Available at <http://www.ri.bbsrc.ac.uk/chickmap/nomenclature.html>. Accessed May 1, 2001. 746
- [28] Zebrafish-Nomenclature, Zebrafish Nomenclature Committee and
747 Guidelines. Available at http://zfin.org/zf_info/nomen_comm.html. Accessed May 1, 2001. 748
- [29] Maltais LJ et al. Rules and guidelines for mouse gene nomen-
749 clature: a condensed version. International committee on standard-
750 ized genetic nomenclature for mice. *Genomics* 1997;45(2):471–6. 751
- [30] Antonarakis SE. Recommendations for a nomenclature system
752 for human gene mutations. Nomenclature working group. *Hum*
753 *Mutat* 1998;11(1):1–3. 754
- [31] Chicken-nomenclature, Nomenclature for naming loci, alleles,
755 linkage groups, and chromosomes to be used in poultry genome
756 publications and databases. Available at <http://www.ri.bbsrc.ac.uk/chickmap/nomenclature.html>. Accessed May 1, 2001. 757
- [32] Rat-nomenclature, Rat: Nomenclature Committee Guidelines.
758 Available at <http://ratmap.gen.gu.se/ratmap/WWWNomen/Brief.html>. Accessed May 1, 2001. 759
- [33] Horvitz HR et al. A uniform genetic nomenclature for the
760 nematode *Caenorhabditis elegans*. *Mol Gen Genet*
761 1979;175(2):129–33. 762
- [34] Aminoacid-nomenclature, Nomenclature and Symbolism for
763 Amino Acids and Peptides. Available at <http://www.chem.qmw.ac.uk/iupac/AminoAcid/>. Accessed May 1, 2001. 764
- [35] Kibayashi M, Nagao M, Chiba S. Influence of valproic acid on
765 the expression of various acyl-CoA dehydrogenases in rats.
766 *Pediatr Int* 1999;41(1):52–60. 767
- [36] GenBank, <http://www.ncbi.nlm.nih.gov/Genbank/GenbankOverview.html>. Accessed August 1, 2002. 768
- [37] Klavans J, Muresan S. Evaluation of the DEFINDER System for
769 Fully Automatic Glossary Construction. In: *Proceedings of the*
770 *AMIA Symposium*, 2001. 771
- [38] Hatzivassiloglou V, Duboue PA, Rzhetsky A. Disambiguating
772 proteins, genes, and RNA in text: a machine learning approach.
773 *Bioinformatics* 2001;17(Suppl 1):S97–S106. 774