A High-Precision Health-Relatedness Score for Phrases to Mine Cause–Effect Statements from the Web

Anonymous ACL submission

Abstract

The measurement of the health-relatedness of a phrase is important when mining the 003 web at scale for health information, e.g., 004 when building a search engine or when carrying out health-sociological analyses. We propose a new termhood scoring scheme that allows for the prediction of the health-800 relatedness of phrases at high precision. An evaluation on several corpora of cause-effect statements (heuristically and professionally labeled) yields about 60% recall at over 90% precision, outperforming state-of-theart vocabulary-based approaches and per-014 forming on par with BERT while being less resource-demanding. A new resource of over 4 million health-related cause-effect statements is compiled, such as "Studies 017 show that stress induces insomnia.", which explicitly connect symptoms ('stress') as 019 claimed causes for conditions ('insomnia'). It consists of over 4 million sentences from more than 2 million unique web pages and 234,000 unique websites.¹

1 Introduction

025

029

037

Health sociology investigates society's interaction with health, where an important subject of interest is how consumers obtain and perceive health-related information. The web, as a main source (Sbaffi and Rowley, 2017), has been frequently studied in this regard throughout the past two decades. Three systematic reviews summarize the outcomes of 79, 157, and 165 studies, respectively (see Table 1): The studies typically focus on a single medical domain and range in size from a handpicked single page to up to 1,524 pages, with averages of 100.5, 78.5, and 50.3 pages per study.

Virtually all the aforementioned studies have been carried out manually. In order to enable

Rev.	Stu	idies	V	ges			
	Year	Count	Min	Max	Mean	Stddev	Sum
a	2001	79	3	1,147	100.5	157.7	7,796
b	2013	165	3	388	78.5	73.4	12,870
с	2017	157	1	1,524	50.3	133.9	7,891

Table 1: Key statistics of the number of websites or web pages analyzed in studies of online health information as reviewed by (a) Eysenbach et al. (2002), (b) Zhang et al. (2015), (c) Daraz et al. (2018). Some studies are part of more than one review; most do not differentiate websites from web pages.

040

041

042

043

044

045

047

048

049

051

052

053

054

056

058

059

060

061

062

063

064

065

066

067

068

069

scaling up such studies, further automation of various prerequisite tasks is required: (1) the discovery and acquisition of websites and web pages with relevance to health, (2) the extraction of specific health-related statements, and (3) the attribution of health-related statements to authoritative sources (e.g., for fact checking). While the first and third step have been and are subject to ongoing research and development, the second step has received much less attention thus far, especially given the requirement of reaching a high precision so as to minimize noise in subsequent analyses.

Since a substantial portion of the information need of health consumers relates to causes and effects, be it the etiology of a condition or the effect of a treatment, we focus on this specific case and contribute towards automating the aforementioned second step as follows: (1) A new approach for measuring the health-relatedness of phrases with high precision is introduced (Section 3). (2) Based on our approach, a new resource compiles healthrelated cause-effect statements at web scale (Section 4). (3) In an in-depth evaluation, the approach is compared to several state-of-theart approaches, outperforming state-of-the-art baselines for medical entity linking, while performing on par with BERT while requiring significantly less resources (Section 5).

¹Code and data will be published alongside the paper; an excerpt is found as supplementary material.

070

072

073

075

077

078

079

081

083

084

086

087

090

096

097

098

100

101

102

103

104

106

107

108

109

110

111

112

113

114

115

116

117

118 119

120

2 Related Work

The impact that online information can have on a consumer's health has sparked the interest of the health-sociological research community ever since the web established itself as an information source in society. For example, user surveys investigate consumers' perceptions of online health information (Diaz et al., 2002), e-health services (Andreassen et al., 2007), as well as the criteria by which consumers judge the quality of a website (Sun et al., 2019).

Information quality appears to be the mostinvestigated characteristic. Numerous studies systematically reviewed the quality of websites with respect to specific topics like orthodontics (Jiang, 2000) and performance-enhancing drugs (Brennan et al., 2013). Apart from specific topics, restrictions to particular portions of the web are also common. Examples include smallscale studies of dietary advice (Cooper et al., 2012) and the misinterpretation (Yavchitz et al., 2012) or exaggeration (Summer et al., 2014) of clinical trial results in online news. Recent research focused on social media (Suarez-Lledo and Alvarez-Galvez, 2021), particularly health misinformation on Twitter (Broniatowski et al., 2018; Bal et al., 2020). The accuracy of health information in search result snippets has also been investigated (Bondarenko et al., 2021).

Besides the mostly manual analyses, some quality assessment tasks have been automated, such as the detection of websites listing unproven cancer treatments (Aphinyanaphongs and Aliferis, 2007) as well as fake medical websites (Abbasi et al., 2012), and determining if a website conforms to the HON Code (Boyer and Dolamic, 2015; Boyer et al., 2017), a health information quality standard for websites.

In terms of discriminating between health and non-health-related content, most previous work has focused on classifying entire articles or pages. For example, medical vocabularies are used to detect news articles related to health (Watters et al., 2002; Zheng et al., 2002), and convolutional neural networks to detect mental health-related Reddit posts (Gkotsis et al., 2017). Little previous work exists on classifying phrases or terms as health-related, preventing the automation fact-checking. Further afield, keyword extraction and automatic ontology creation are related, where the goal is to extract prototypical words for a particular domain. For example, the C-value/NC-value method extracts multi-word domain terms from a corpus using term frequencies (Frantzi et al., 2000). Its reliance on the syntactic structure of extracted candidate words render it inapplicable to arbitrary phrases. 121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

More straightforwardly applied are the family of contrastive termhood scores. which relate term frequencies from a domain corpus to term frequencies from one or more out-of-domain corpora. These include tf-idf-inspired measures (Basili et al., 2001; Kim et al., 2009), measures estimating how exclusive a term is for a domain (Khurshid et al., 2000; Park et al., 2008), and combinations or extensions thereof (Wong et al., 2007; Bonin et al., 2010). We transfer contrastive termhood scoring to measuring health-relatedness and compare it with the state-of-the-art medical entity linker, Quick-UMLS (Soldaini and Goharian, 2016). Unlike classical medical entity linking algorithms, like MetaMap (Aronson, 2001) and cTakes (Savova et al., 2010), QuickUMLS is faster, achieves higher F1 and recall on several benchmarks, and can be tuned to prioritize precision or recall. Whereas entity linkers using neural language models (Neumann et al., 2019; Nejadgholi et al., 2019) are trained on entire abstracts and require additional context to extract entity candidates, rendering them inapplicable to phrases.

3 Measuring Health-Relatedness

Determining if a phrase is health-related is an issue of ambiguity. Homonomy (same surface form, different meaning) and polysemy (same surface form, different sense) render this decision difficult.² This section revisits so-called termhood scores, which measure the degree to which a given word is specific to a certain domain (Kageura and Umino, 1996). We introduce a new generalized score for phrases, and show how to tailor it to the health domain. Underlying our generalized termhood score are contrastive weight (CW) (Basili et al., 2001), term domain-specificity (TDS) (Khurshid et al., 2000; Park et al., 2008), and discriminative weight (DW) (Wong et al., 2007).

²The word 'cancer' can refer to a clearly health-related malignant tumor, but also to the zodiac sign, which is less likely to appear in a health-related context.

168 169

170

171

172

173

174

175

176

177

178

179

180

181 182

183

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

201

204

207

209

210

211

3.1 Contrastive Termhood Scores

All three considered contrastive termhood scores rely on a corpus of domain-specific text and a contrastive corpus of out-of-domain text. Formally, the health corpus H and the contrastive general corpus G are each represented by the multisets of all words in their texts. The corpus frequency $cf_C(c)$ of a word w in a corpus C denotes the absolute number of w's occurrences in C, while the relative corpus frequency $rf_C(w)$ denotes $cf_C(w)/|C|$ and the inverse corpus frequency icf(w) denotes

$$icf(w) = \log\left(\frac{|H| + |G|}{cf_H(w) + cf_G(w)}\right).$$

The contrastive weight CW of a word w is defined as

$$CW(w) = \log \left(cf_H(w) + 1 \right) \cdot icf(w).$$

It is strongly related to $tf \cdot idf$, but instead of term and inverse document frequency, it uses corpus and inverse corpora frequency. The term domain-specificity TDS measures the domain exclusivity of a word w:

$$TDS(w) = \log\left(\frac{rf_H(w) + 1}{rf_G(w) + 1} + 1\right).$$

The discriminative weight DW was originally defined as the product of CW and an unnormalized version of TDS, which used the corpus frequency *cf* instead of the relative frequency *rf*. Since varying corpora sizes heavily affect the unnormalized TDS score, we replace it with its normalized version and simply define DW as

 $\mathrm{DW}(w) = \mathrm{CW}(w) \cdot \mathrm{TDS}(w).$

3.2 Generalized Phrase Termhood

To calculate termhood scores for phrases instead of words, it appears straightforward to average a phrase's individual word termhood scores: However, this does not work well for health-related phrases with many out-ofdomain or stop words, like "unnecessary plastic surgery". Even though 'surgery' has a high termhood score, the overall average is rather low, due to the out-of-domain words 'unnecessary' and 'plastic'. We propose two schemes that avoid the issues of the simple average.

The first uses a weighted average to boost a phrase's words with high termhood. The idea

Co	rpus	Language 1	Documents	Words
G	Wikipedia	mixed, layp.	$12,\!265,\!374$	$3.0 \cdot 10^{9}$
$H_1 \\ H_2 \\ H_3 \\ H_4$	PubMed PubMed Centr. Textbooks Encyclopedias	scientific scientific clinical, educ mixed, laype	$\begin{array}{r} 31,847,923\\ 3,611,361\\ \text{ational} 434\\ \text{rson} 67,967 \end{array}$	$\begin{array}{c} 3.8\!\cdot\!10^9 \\ 5.4\!\cdot\!10^9 \\ 1.4\!\cdot\!10^7 \\ 9.3\!\cdot\!10^6 \end{array}$

Table 2: Overview of our evaluation corpora.

is that a single highly health-related word is able to dictate the termhood score of a phrase, thereby increasing recall. Phrases with a high (unweighted) average termhood will still be ranked high so that precision is not affected. We calculate the weighted average of the termhood scores x_1, \ldots, x_m of an *m*-word phrase as the generalized mean

$$M_p(x_1,\ldots,x_m) = \left(\frac{1}{m}\sum_{i=1}^m x_i^p\right)^{\frac{1}{p}}$$
 22

212

213

214

215

216

217

218

219

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

with the non-zero real-valued parameter p. For p = 1, the generalized mean corresponds to the arithmetic mean. By increasing p, the mean is biased towards the higher-valued term-hood scores; in the extreme case of $p = \infty$, the largest x_i is returned.

As the second scheme, we propose to also compute the weighted average termhood over the n-grams of a phrase. While the unigram 'plastic' is relatively unrelated to health, the bigram 'plastic surgery' certainly is healthrelated. Though the above generalized mean already increases the bigram's termhood compared to a simple average, the high occurrence frequency of the bigram itself is an even better indicator for its health-relatedness. Due to the sparsity of larger n-grams, especially prevalent in smaller corpora, we average the termhood scores of a phrase over multiple ngrams. Let s denote the phrase w_1, \ldots, w_m and let $s_{i,k}$ denote the subphrase w_i, \ldots, w_{i+k} of s $(0 \le k \le m - 1 \text{ and } i \in \{1, \dots, m - k\})$. Let the above termhood scores t(.) (i.e., CW, TDS, and DW) all be pre-calculated up to *n*-grams. The phrase termhood $PT_{t,n,p}(s)$ for phrase s is then defined as

$$\operatorname{PT}_{t,n,p}(s) = \frac{1}{n} \sum_{k=0}^{n-1} \left(\frac{1}{m-k} \sum_{i=1}^{m-k} t(s_{i,k})^p \right)^{\frac{1}{p}}.$$
 247

248 249

252

254

255

257

258

259

260

261

262

264

266

267

269

270

273

274

275

276

277

278

279

281

284

286

287

289

3.3 Adaptation to the Health Domain

We select Wikipedia³ as our contrastive corpus G because of its domain variety, relatively uniform language, and accessibility to the general public. As candidates for a health corpus H, we consider and evaluate four alternatives, each with its own (dis)advantages (see Table 2 for an overview).

The first three corpora use documents provided by the National Library of Medicine: a dump of over 30 million abstracts from PubMed⁴, a subset of over 3 million fulltext publications from PubMed Central⁵ and 434 textbooks of the textbook and monograph category from the NCBI Bookshelf⁶. While both PubMed based corpora are large scale, their language is mainly scientific. The textbook corpus contains more clinical language, which we hypothesize to more closely match the expected proficiency level of web language.

Finally, we also crawled the entries of five consumer-oriented medical encyclopedias (Appendix A). Because the encyclopedias are purposefully written in layperson's terms, its joint language distribution is assumed to be most similar to the target language distribution.

3.4 Pilot Experiments

Comparing the three scores, Figures 1a-c show that, for the PubMed corpus H_1 , all scores rank out-of-domain words and stop words lower than health-related words. However, the distributions of CW and TDS differ substantially, with the DW striking a balance between both. While the TDS ranks comparably few words as extremely health-related, the CW has a more even distribution with less extreme differences. Especially, 'ward' has a large difference in ranking between both scores. While it occurs frequently within the health domain so that CW attributes a high health-relatedness, it also occurs frequently in the general domain. Its lacking exclusiveness leads to the TDS scoring it comparably low.

To gain an intuition into the effect of using the different health corpora H_1 to H_4 with the termhood scores, Figures 1c and d compare

Subset	Opt. Measur	e Size	\mathbf{HR}	Р	R
Full	F1	2,968,345	25.6%	0.78	0.83
Full	Prec.	1,623,968	14.0%	0.90	0.73
Support	F1	111,406	61.8%	0.88	0.93
Support	Prec.	103,792	57.6%	0.91	0.89

Table 3: Descriptive statistics of four health-related cause–effect networks extracted from the CauseNet. The number of statements in each dataset, proportion of health-related statements as well as estimated precision and recall are listed.

294

295

296

297

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

330

331

332

the DW using the PubMed and Encyclopedia corpora. Here, the difference between scientific and non-scientific language between the two corpora is evident. The word 'experiment' has a comparably high termhood score using the PubMed corpus, but is similarly ranked to stop words using the Encyclopedia corpus. Otherwise, the shape of the distributions and location of examples are fairly similar.

4 Case Study: Cause–Effect Statements

To evaluate and demonstrate our approach, we apply it to a large graph of cause-effect statements. The CauseNet (Heindorf et al., 2020) is a graph of over 11 million pairs of cause and effect phrases extracted from all sentences found in the web pages of the ClueWeb12 web crawl.⁷ It is important to note that these statements are *claimed* cause-effect statements, i.e., statements that have been made on some web page. Therefore, it contains of cause-effect statements for which empirical evidence can be found ("earthquake \rightarrow tsunami"), but also many for which this is not the case ("jupiter opposing mars \rightarrow bad luck on the job"). The statements were extracted using a linguistic pattern matching, achieving an estimated precision of 83%. Precision can be further increased to an estimated 93% by only considering statements with high support, i.e. statements which were extracted more than once using different linguistic patterns. The increase in precision of course takes a toll on recall. Only about 1.6%of statements have high support.

We evaluate our termhood approach (see Section 5.3) on several manually labeled subsets of the CauseNet. Based on this evaluation, we extract four different health-related cause–effect networks, one maximizing the F1-

 $^{^3 {\}rm Specifically}$ a dump of English Wikipedia articles from June 1st, 2021.

⁴https://pubmed.ncbi.nlm.nih.gov/

⁵https://www.ncbi.nlm.nih.gov/pmc/

 $^{^{6} \}rm https://www.ncbi.nlm.nih.gov/books$

 $^{^{7}} https://www.lemurproject.org/clueweb12/$



Figure 1: Histograms of termhood frequencies for the (a) CW, (b) TDS, and (c) DW on the PubMed corpus, and for (d) DW on the Encyclopeida corpus. Example words are highlighted.

measure, and one maximizing the F1-measure with at least 90% precision for both the full and the high-precision CauseNet with high support, and release these to the public for further analyses. Table 3 summarizes the descriptive statistics of each resource. Optimizing for high precision on the full CauseNet yields 1,623,968 health-related statements. With high precision, an estimated 14% and an estimated 58% of all statements within the full support subsets of the CauseNet are respectively health-related.

5 Evaluation

333

334

337

338

339

341

342

343

345

346

347

354

358

364

366

368

This section reports on an in-depth evaluation of our health-relatedness score compared to state-of-the-art entity linking and BERT-based baselines and different parameterizations on four labeled cause–effect statement datasets.

5.1 Baselines

Vocabulary-Based Approach. A precision oriented approach for determining healthrelatedness of a phrase is to check if it, or subphrases of it, are part of a medical vocabulary. To judge the performance of our healthrelatedness score, we therefore compare it to a vocabulary-based approach based on Quick-UMLS (Soldaini and Goharian, 2016), a stateof-the-art medical entity linker. The proportion of words within a phrase which could be matched to UMLS concepts is then used as a health-relatedness score.

In more detail, given a phrase s, first all medical entity mentions E are extracted. Overlapping entity mentions, e.g. 'cancer' is contained in "breast cancer", are handled by taking only the longest mentioned entity, resulting in a subset of non-overlapping entity mentions \hat{E} . Next, stop words⁸ not contained in any entity mentions are removed from s, yielding \hat{s} . The vocabulary health-relatedness score V(s) is then computed as

$$V(s) = \frac{|\hat{s}| - \sum_{e_i \in \hat{E}} |e_i|}{|\hat{s}|}.$$
373

370

371

372

374

375

376

377

378

379

381

382

383

384

387

388

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

Entity mentions are linked to the UMLS Metathesaurus (Humphreys and Lindberg, 1993), which is a mix of medical vocabularies of varying specificity. We investigate three decreasingly specific vocabulary subsets in an attempt to increase precision: the MeSH hierarchy⁹, the MeSH hierarchy with additional synonyms (MeSH Syn) and the entire UMLS Metathesaurus (see Appendix B for further details). For all three variants, we also consider restricting the set of concepts to a set of medically specific semantic types (ST21pv) as proposed in the MedMentions entity linking dataset (Mohan and Li, 2019). Finally, several different string similarity thresholds (Jaccard similarity was used in this work) were tested to allow for fuzzy string matching and increase recall.

BERT-Based Approach.

As a second baseline, we use a BERT-based sequence classifier which is trained to predict if a sequence of tokens originates from a healthrelated corpus. Starting from a pretrained SciBERT (Beltagy et al., 2019) model, we finetune two different models. One model each is trained to predict if a noun phrase originates from the PubMed H_1 or the Encyclopedia H_4 corpus. Noun phrases from the Wikipedia corpus G serve as negative samples. Further details about the training procedure can be found in Appendix C.

⁹https://www.nlm.nih.gov/mesh/meshhome.html

⁸The English nltk stop words list is used.

5.2 Reference Datasets

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

We apply three labeling strategies across four reference datasets. The first reference dataset is collected from Wikidata and labeled using a heuristic. With the help of a medical practitioner (a practicing orthopedist and professor), we gather nine general root concepts that include the majority of health-related concepts. Then all 9,317 Wikidata relations with the has cause (P828) and/or has effect (P1542) predicates are extracted. All relations for which both the cause and effect concepts are direct or indirect children of the root concepts are considered health-related. See Appendix D for a full list of root concepts and further details. As this dataset propagates labels heuristically, we consider it a silver standard.

Next, we manually classified two different sets of CauseNet statements; 1,000 randomly sampled statements from each, the full CauseNet (Full), and the high-support subset (Support). A subset of 100 statements from the Full dataset was labeled by 3 separate annotators, achieving a Cohen's Kappa of 0.77. These datasets are considered as gold standard.

Finally, after evaluating on the aforementioned datasets, we sampled 1,000 statements from the full CauseNet that were closest to the decision threshold of the termhood classifier with the highest F1 score and at least 90% precision on the Full dataset. The aforementioned practitioner labeled the dataset, with the additional option to label statements with unsure. For lack of a better term, we call this dataset a platinum standard, as it is professionally labeled and specifically focuses on a difficult subset of cause–effect statements. The best approaches are evaluated on the full dataset (Practitioner-Full) and the confidence splits (Practitioner-Sure, Practitioner-Unsure).

Table 4 gives an overview of dataset statistics. Interestingly, the proportion of health statements within the CauseNet datasets varies substantially. The higher the support, the more likely a statement is health-related. Additionally, the high proportion of statements marked as unsure by the practitioner shows the difficulty of the task. While some statements were marked unsure because of unknown terminology, most were borderline decisions because of ambiguous concepts

Dataset	Size	Health-related	Length
Wikidata	9317	31.0%	4.90
Full	1000	19.7%	7.21
Support	1000	50.3%	3.40
Practitioner-Full	1000	77.2%	5.99
Practitioner-Sure	594	82.0%	5.84
Practitioner-Unsure	406	70.2%	6.21

Table 4: Overview of the evaluation datasets, including the proportion of true health-related cause– effect statements, and the average number of words per cause/effect phrase.

(e.g., poor treatment \rightarrow problems), or the difficulty to delineate the health domain from other related domains (e.g., biological processes cold air \rightarrow bronchoconstriction).

5.3 Results

To combine the individual termhood scores of the cause and effect phrases we use "and" (both cause and effect scores need to exceed the decision threshold) as an upper bound precisionoriented operator. Following the rationale of using the generalized mean for increasing recall in health-phrase detection, we also test the generalized mean for combining cause and effect phrase termhood. To differentiate between parameters, we denote p for averaging n-gram termhood by p_n and for averaging phrase termhood by p_p . By setting $p_p = \infty$, the maximum phrase termhood score of either cause and effect is used. Thereby, $p_p = \infty$ is the same as using "or" (one of cause or effect phrase termhood scores need to exceed the decision threshold) and acts as the complement to the "and" operator and as a recall-oriented upper bound.

To evaluate the different approaches we run a grid search over the parameters and thresholds on the silver and gold-standard datasets and test for significance using a bootstrap test with 5,000 permutations. See Appendix E for a full description of parameters. Table 5 gives an overview of the best variants for each approach in terms of F1 measure.

While no approach is able to statistically significantly outperform all others, all termhood scores and the BERT-based approach are able to statistically significantly (p < 0.05) outperform the vocabulary approaches across all datasets. The vocabularies are unable to achieve high precision. While it is usually possible to tune the decision threshold to achieve perfect precision, the binary classification, i.e. 468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

456

457

458

459

460

461

	Approach	Parameters	Operator	Р	R	F1
	MeSH	Jacc.=0.9	$p_p{=}1$	0.57	0.74	0.65
ą	MeSH Syn	Jacc.=0.9	$p_p=2$	0.55	0.78	0.65
lat	UMLS	Jacc.=1.0	AND	0.49	0.83	0.62
kić	BERT	PubMed	$p_p{=}\infty$	0.72	0.82	0.77
N1	CW Encyc.,	$n=1, p_n=10$	$p_p=5$	0.70	0.74	0.72
-	TDS Encyc.,	$n{=}1, p_n{=}10$	$p_p=2$	0.70	0.88	0.78
	DW Encyc.,	$n{=}1,\ p_n{=}2$	$p_p=2$	0.74	0.83	0.78
	MeSH	Jacc.=1.0	$p_p=10$	0.44	0.77	0.56
	MeSH Syn	Jacc.=1.0	$p_p=1$	0.57	0.60	0.59
	UMLS	Jacc.=1.0	AND	0.40	0.69	0.51
[[IJ	BERT	PubMed	$p_p=10$	0.84	0.79	0.81
ц	CW Encyc.,	$n{=}1, p_n{=}5$	$p_p=5$	0.77	0.79	0.78
	TDS Encyc.,	$n=3, p_n=1$	$p_p=2$	0.74	0.86	0.79
	DW Encyc.,	$n{=}3,\ p_n{=}1$	$p_p=1$	0.78	0.83	0.80
	MeSH	Jacc.=1.0	$p_p=1$	0.62	0.87	0.72
	MeSH Syn	Jacc.=0.9	$p_p=1$	0.74	0.76	0.75
ort	UMLS	Jacc.=1.0	AND	0.60	0.91	0.72
pc	BERT	PubMed	$p_p=10$	0.92	0.87	0.90
jug	CW Encyc.,	$n=3, p_n=10$	$p_p=2$	0.82	0.88	0.85
01	TDS Encyc.,	$n=2, p_n=1$	$p_p=1$	0.88	0.93	0.90
	DW Encyc.,	$n{=}3, p_n{=}5$	$p_p=1$	0.86	0.93	0.90

Table 5: Parameterizations of each approach optimized for F1 on the silver- and gold-standard evaluation datasets.

a concept is either health-related or not, creates an upper bound for the vocabulary based approaches. Even setting the threshold such that only fully matched phrases are included leads to some false positive predictions.

Termhood Score Comparison. Compared to the vocabulary approaches, the termhood scores have more granular distributions and the decision threshold can therefore be tuned to achieve high precision. Table 6 lists the best performing approaches with at least 90%precision in terms of F1-measure (none of the vocabulary approaches were able to achieve more than 90% precision). When high precision is required, the term domain-specificity outperforms the contrastive weight on the two CauseNet evaluation datasets, featuring substantially higher recall at similar precision. However, the exact opposite relationship can be seen for the Wikidata dataset. By combining both CW and TDS, the DW achieves the best or only marginally worse F1 scores on all evaluation datasets.

519Taking a closer look at the effect the vari-520ous health corpora have on classification per-521formance shows that the Encyclopedia corpus522always leads to the best performance. Irrespec-523tive of dataset, every termhood score is able to524achieve the highest overall F1 score in both un-

	Approach	Parameters	Operato	r P	\mathbf{R}	$\mathbf{F1}$
kidata	BERT CW Encyc. TDS Encyc	$, n=1, p_n=10$ n=2, n=1	$p_p = 5$ $n_r = 1$	- 0.90 0.90	- 0.39 0.17	0.54 0.28
Wi	DW Encyc.	, $n=2$, $p_n=1$, $n=1$, $p_n=5$	$p_p = 1$	0.91	0.50	0.64
Full	BERT CW Encyc. TDS Encyc. DW Encyc.	PubMed , $n=1, p_n=5$, $n=2, p_n=1$, $n=3, p_n=2$	$\begin{array}{c} \texttt{AND} \\ p_p = 2 \\ p_p = 1 \\ p_p = 1 \end{array}$	0.92 0.91 0.90 0.92	0.50 0.52 0.62 0.58	0.65 0.66 0.73 0.71
Support	BERT CW Encyc. TDS Encyc. DW Encyc.	PubMed , $n=1, p_n=10$, $n=3, p_n=1$, $n=2, p_n=1$	$p_p = 10$ $p_p = 5$ $p_p = 1$ $p_p = 1$	0.92 0.91 0.91 0.90	0.87 0.75 0.89 0.87	0.90 0.82 0.90 0.89

Table 6: Parameterizations of each approach with at least 90% precision optimized for F1 on the silverand gold-standard datasets.

constrained (Table 5) and constrained precision scenarios (Table 6) using the Encyclopedia corpus. This effect does not translate to the BERT model. The models trained on the PubMed H_1 corpus outperform the Encyclopedia H_4 corpus trained models on all datasets. 525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

In contrast, the effects of the n-gram and generalized mean parameters on performance are more subtle. The CW and the TDS each prefer high and low p_n values respectively. Especially the contrastive weight profits from the possibility to increase the p_n value and subsequently increase recall. The term domain-specificity is already precision-oriented and therefore performs best with lower p_n values. The DW again strikes a balance between both and prefers higher or lower p_n values depending on the proportion of health-related labels in the dataset.

Finally, the *n*-gram variants have the smallest impact on performance. When switching to n = 3, the CW loses 10% points in recall on the Full dataset, while the TDS and DW have their largest drops at 9% and 5% points for n = 1 respectively. Over all other datasets the drop in performance is negligible. This most likely stems from the fact that the Full dataset has by far the longest average event length at 7.21 words per event, which enables the termhood scores to take full advantage of longer n-grams.

Practitioner Evaluation. We finally evaluate the termhood scores and BERT-based approach on the difficult subset of relations labeled by a medical practitioner. As a reminder, based on the results of the evaluation on the silverand gold-standard datasets, we use the discrim-

518

496

497

	Approach	Parameters	Operator	Р	R	F1
	BERT	PM	AND	0.91	0.18	0.30
Ш	CW Encyc.,	$n{=}3, p_n{=}10$	AND	0.91	0.19	0.31
Ē	TDS Encyc.,	$n{=}3, p_n{=}1$	AND	0.90	0.06	0.11
	DW Encyc.,	$n{=}2,\ p_n{=}1$	AND	0.91	0.15	0.25
	BERT	PM	AND	0.90	0.82	0.86
ure	CW Encyc.,	$n{=}2, p_n{=}2$	AND	0.90	0.59	0.72
$\mathbf{S}_{\mathbf{U}}$	TDS Encyc.,	$n{=}1, p_n{=}1$	AND	0.90	0.34	0.50
	DW Encyc.,	$n{=}3,p_n{=}1$	AND	0.90	0.66	0.76
e	BERT	PM	AND	1.00	0.02	0.05
sur	CW Encyc.,	$n{=}1, p_n{=}5$	AND	0.94	0.05	0.10
Jns	TDS Textbo	ok, $n=1, p_n=2$	$p_p{=}\infty$	0.90	0.10	0.18
	DW Textbo	ok, $n=2, p_n=2$	$p_p{=}5$	0.91	0.07	0.14

Table 7: Parameterizations of each statistical approach with at least 90% precision optimized for F1 on the practitioner evaluation datasets.

inative weight (Encyclopedia corpus, n = 3, $p_n = 2$, $p_p = 1$) and sample the 1,000 relations closest to the decision threshold (120) tuned for high F1 with precision over 90%.

Table 7 gives an overview of the best paramemeterizations. Again requiring high precision, we find that performance drastically drops on the Practitioner-Full dataset. All scores have to set a high decision threshold and use the "and" operator to reach the high precision and requirement and thereby sacrifice recall. Considering the split of into Practitioner-Sure and Practitioner-Unsure relations however shows that the approaches especially struggle with the relations the practitioner was not confident about. The performance on the Practitioner-Sure subset on the other hand is on par with the best approaches on the Full dataset, with the BERT approach significantly outperforming the termhood scores.

To gain insight into the performance drop of the unsure dataset, we sample the highest scored true negative relations of the best discriminative weight approach, i.e. the relations which the approach considers to likely be healthrelated that the practitioner labeled as health unrelated. See Table 8 for the top four examples. We find that the two issues, health domain demarcation and handling general concepts mentioned in Section 5.2, reflect themselves in the classification. A "stroke" causing a "reduced cell count", as in the second to last example of Table 8, might have medical relevance, but could just as well be the plain description of a biological process. Second, while it is unclear which "small amounts" are meant

$\mathbf{Cause} \rightarrow \mathbf{Effect}$	DW		
	Cause	Effect	
$\begin{tabular}{lllllllllllllllllllllllllllllllllll$	$\begin{array}{c} 155.95 \\ 164.37 \\ 179.32 \\ 128.57 \\ 108.52 \end{array}$	$140.32 \\ 129.70 \\ 110.98 \\ 160.20 \\ 176.77$	

Table 8: Highest scored true negative statements by the best-performing discriminative-weight approach on the Practitioner-Unsure dataset. The first effect contains a spelling error not handled by the web crawl extraction.

in final example of Table 8, the discriminative weight nonetheless considers the concept as likely health-related because of its frequent usage. 597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

6 Conclusions

We develop a novel approach to determine the health-relatedness of arbitrary short phrases with a high precision, developing a new generalized termhood score. To demonstrate and evaluate our approach in a realistic setting, we apply it, among other datasets, to a web-scale graph of cause–effect statements. In comparison to state-of-art entity linking approaches, our approach is the only one capable of achieving the high precision required for practical purposes, outperforming the baseline approaches on all evaluation datasets. Combined with our generalization, the discriminative weight score proves to be most robust, with the term domainspecificity performing slightly better in highprecision scenarios.

We apply the best precision and F1-oriented approaches to the full CauseNet graph, and a precision oriented subset of the graph. The result is a new resource of high-precision healthrelated statements at an unprecedented scale, suitable for investigating health-sociological questions automatically. At an estimated precision of 0.9 and estimated recall of 0.73, the precision-oriented extraction on the full graph contains 1,623,968 health-related statements from 4,420,897 statements as well as 234,355 and 2,139,563 unique websites and web pages. This opens up new possibilities for the quantitative analysis of health-related information on the web.

592

593

595

596

Ethical Considerations

633

Research on health-related tasks can be often sensitive as its haphazard transfer into prac-635 tice may cause significant harm. Though our 636 research is not aimed at supporting medical 637 treatments, its envisioned application in healthsociological analyses may cause these analyses to include or exclude pieces of information in error. This is why we explicitly aim for a high pre-641 cision: ensuring that a phrase that achieves a 642 high health-relatedness score is actually healthrelated protects both the time and effort of 644 health sociologists tasked with analyzing them, as well as the privacy of people whose web content is ambiguous or otherwise close to health, but not quite crossing the line from being subject to critical interpretation by health experts. 649 Nevertheless, for some applications, achieving 650 a high recall, and thus be inclusive of all that is health-related at the expense of false positives might also be important. The limitations of 653 our approach in this regard are clearly outlined, yet we do see potential of shifting its operating 655 point toward that end.

References

663

666

Ahmed Abbasi, Fatemeh "Mariam" Zahedi, and Siddharth Kaza. 2012. Detecting Fake Medical Web Sites Using Recursive Trust Labeling. ACM Transactions on Information Systems 30, 4 (Nov. 2012), 22:1-22:36. https://doi.org/10.1145/2382438.2382441

Hege K. Andreassen, Maria M. Bujnowska-Fedak, Catherine E. Chronaki, Roxana C. Dumitru, Iveta Pudule, Silvina Santana, Henning Voss, and Rolf Wynn. 2007. European Citizens' Use of E-Health Services: A Study of Seven Countries. BMC Public Health 7, 1 (April 2007), 53.

https://doi.org/10.1186/1471-2458-7-53

Yin Aphinyanaphongs and Constantin Aliferis. 671 2007. Text Categorization Models for Identifying Unproven Cancer Treatments on the Web. Studies 673 in Health Technology and Informatics 129, Pt 2 (2007), 968-972.

A. R. Aronson. 2001. Effective Mapping of Biomedical Text to the UMLS Metathesaurus: The MetaMap Program. Proceedings of the AMIA 679 Symposium (2001), 17–21.

Rakesh Bal, Sayan Sinha, Swastika Dutta, Rishabh Joshi, Sayan Ghosh, and Ritam Dutt. 681 682 2020. Analysing the Extent of Misinformation in 683 Cancer Related Tweets. Proceedings of the

International AAAI Conference on Web and Social	684
Media 14 (May 2020), 924–928.	685
Roberto Basili, Alessandro Moschitti,	686
Maria Teresa Pazienza, and Fabio Massimo	687
Zanzotto. 2001. A Contrastive Approach to Term	688
Extraction. In <i>Proceedings of the TIA 2001:</i>	689
<i>Terminologie et Intelligence Artificielle.</i> 119–128.	690
Iz Beltagy, Kyle Lo, and Arman Cohan. 2019.	691
SciBERT: A Pretrained Language Model for	692
Scientific Text. In Proceedings of the 2019	693
Conference on Empirical Methods in Natural	694
Language Processing and the 9th International	695
Joint Conference on Natural Language Processing	696
(EMNLP-IJCNLP). Association for Computational	697
Linguistics, Hong Kong, China, 3615–3620.	698
https://doi.org/10.18653/v1/D19-1371	699
Alexander Bondarenko, Ekaterina Shirshakova,	700
Marina Driker, Matthias Hagen, and Pavel	701
Braslavski. 2021. Misbeliefs and Biases in	702
Health-Related Searches. In 30th ACM	703
International Conference on Information and	704
Knowledge Management (CIKM 2021). ACM.	705
https://doi.org/10.1145/3459637.3482141	706
Francesca Bonin, Felice Dell'Orletta, Giulia	707
Venturi, and Simonetta Montemagni. 2010. A	708
Contrastive Approach to Multi-Word Extraction	709
from Domain-Specific Corpora. In <i>Proceedings of</i>	710
the International Conference on Language	711
Resources and Evaluation. European Language	712
Resources Association, Valletta, Malta, 19–21.	713
Célia Boyer and Ljiljana Dolamic. 2015.	714
Automated Detection of HONcode Website	715
Conformity Compared to Manual Detection: An	716
Evaluation. Journal of Medical Internet Research	717
17, 6 (June 2015), e3831.	718
https://doi.org/10.2196/jmir.3831	719
Célia Boyer, Cédric Frossard, Arnaud Gaudinat,	720
Allan Hanbury, and Gilles Falquetd. 2017. How to	721
Sort Trustworthy Health Online Information?	722
Improvements of the Automated Detection of	723
HONcode Criteria. <i>Procedia Computer Science</i> 121	724
(Jan. 2017), 940–949.	725
https://doi.org/10.1016/j.procs.2017.11.122	726
Brian P. Brennan, Gen Kanayama, and	727
Harrison G. Pope. 2013. Performance-Enhancing	728
Drugs on the Web: A Growing Public-Health	729
Issue. <i>The American Journal on Addictions</i> 22, 2	730
(2013), 158–161.	731
https://doi.org/10.1111/j.1521-0391.2013.00311.x	732
David A. Broniatowski, Amelia M. Jamison, SiHua	733
Qi, Lulwah AlKulaib, Tao Chen, Adrian Benton,	734
Sandra C. Quinn, and Mark Dredze. 2018.	735
Weaponized Health Communication: Twitter Bots	736
and Russian Trolls Amplify the Vaccine Debate.	737
<i>American Journal of Public Health</i> 108, 10 (Oct.	738
2018), 1378–1384.	739
https://doi.org/10.2105/AJPH.2018.304567	740

- Benjamin E. J. Cooper, William E. Lee, Ben M. 741 Goldacre, and Thomas A. B. Sanders. 2012. The 742 Quality of the Evidence for Dietary Advice given 743 744 in UK National Newspapers. Public Understanding of Science 21, 6 (Aug. 2012), 664-673. 745 https://doi.org/10.1177/0963662511401782
- Lubna Daraz, Allison S. Morrow, Oscar J. Ponce, 747 Wigdan Farah, Abdulrahman Katabi, Abdul 748 Majzoub, Mohamed O. Seisa, Raed Benkhadra, 749 Mouaz Alsawas, Prokop Larry, and M. Hassan 750 Murad. 2018. Readability of Online Health 751 Information: A Meta-Narrative Systematic Review. American Journal of Medical Quality: The Official Journal of the American College of Medical 754 Quality 33, 5 (2018), 487–492. https://doi.org/10.1177/1062860617751639

Joseph A. Diaz, Rebecca A. Griffith, James J. Ng, Steven E. Reinert, Peter D. Friedmann, and 758 Anne W. Moulton. 2002. Patients' Use of the 759 Internet for Medical Information. Journal of 760 761 General Internal Medicine 17, 3 (2002), 180–185. https://doi.org/10.1046/j.1525-1497.2002.10603.x

Gunther Eysenbach, John Powell, Oliver Kuss, and Eun-Ryoung Sa. 2002. Empirical Studies Assessing the Quality of Health Information for Consumers on the World Wide WebA Systematic Review. JAMA 287, 20 (May 2002), 2691-2700. https://doi.org/10.1001/jama.287.20.2691

Katerina Frantzi, Sophia Ananiadou, and Hideki Mima. 2000. Automatic Recognition of Multi-Word Terms:. The c-Value/Nc-Value 772 Method. International journal on digital libraries 3, 2 (2000), 115–130.

George Gkotsis, Anika Oellrich, Sumithra Velupillai, Maria Liakata, Tim J. P. Hubbard, 775 776 Richard J. B. Dobson, and Rina Dutta. 2017. Characterisation of Mental Health Conditions in Social Media Using Informed Deep Learning. 778 Scientific Reports 7, 1 (March 2017), 45141. 779 780 https://doi.org/10.1038/srep45141

774

777

784

786

787

790

792

793

Stefan Heindorf, Yan Scholten, Henning Wachsmuth, Axel-Cyrille Ngonga Ngomo, and Martin Potthast. 2020. CauseNet: Towards a Causality Graph Extracted from the Web. In Proceedings of the 29th ACM International Conference on Information & Knowledge Management. Association for Computing Machinery, New York, NY, USA, 3023–3030.

B L Humphreys and D A Lindberg. 1993. The UMLS Project: Making the Conceptual Connection between Users and the Information They Need. Bulletin of the Medical Library Association 81, 2 (April 1993), 170-177.

You-Ling Jiang. 2000. Quality Evaluation of 794 Orthodontic Information on the World Wide Web. American Journal of Orthodontics and Dentofacial Orthopedics 118, 1 (July 2000), 4–9. https://doi.org/10.1067/mod.2000.104492

Kyo Kageura and Bin Umino. 1996. Methods of Automatic Term Recognition: A Review. Terminology. International Journal of Theoretical and Applied Issues in Specialized Communication 3, 2 (Jan. 1996), 259–289. https://doi.org/10.1075/term.3.2.03kag

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

Ahmad Khurshid, Lee Gillman, and Lena Tostevin. 2000. Weirdness Indexing for Logical Document Extrapolation and Retrieval. In Proceedings of the Eighth Text Retrieval Conference (TREC-8). Gaithersburg, USA, 1–8.

Su Nam Kim, Timothy Baldwin, and Min-Yen Kan. 2009. An unsupervised approach to domain-specific term extraction. In Proceedings of the Australasian Language Technology Association Workshop. Sydney, Australia, 94–98. https://www.aclweb.org/anthology/U09-1013

Sunil Mohan and Donghui Li. 2019. MedMentions: A Large Biomedical Corpus Annotated with UMLS Concepts. arXiv:1902.09476 [cs] (Feb. 2019). arXiv:1902.09476 [cs]

Isar Nejadgholi, Kathleen C. Fraser, Berry De Bruijn, Muqun Li, Astha LaPlante, and Khaldoun Zine El Abidine. 2019. Recognizing UMLS Semantic Types with Deep Learning. In Proceedings of the Tenth International Workshop on Health Text Mining and Information Analysis (LOUHI 2019). Association for Computational Linguistics, Hong Kong, 157–167. https://doi.org/10.18653/v1/D19-6219

Mark Neumann, Daniel King, Iz Beltagy, and Waleed Ammar. 2019. ScispaCy: Fast and Robust Models for Biomedical Natural Language Processing. In Proceedings of the 18th BioNLP Workshop and Shared Task. Association for Computational Linguistics, Florence, Italy, 319-327. https://doi.org/10.18653/v1/W19-5034 arXiv:arXiv:1902.07669

Youngja Park, Siddharth Patwardhan, Karthik Visweswariah, and Stephen C Gates. 2008. An Empirical Analysis of Word Error Rate and Keyword Error Rate. In Proceedings of the Ninth Annual Conference of the International Speech Communication Association. Brisbane, Australia, 2070-2073.

Guergana K Savova, James J Masanz, Philip V Ogren, Jiaping Zheng, Sunghwan Sohn, Karin C Kipper-Schuler, and Christopher G Chute. 2010. Mayo Clinical Text Analysis and Knowledge Extraction System (cTAKES): Architecture, Component Evaluation and Applications. Journal of the American Medical Informatics Association 17, 5 (Sept. 2010), 507–513. https://doi.org/10.1136/jamia.2009.001560

Laura Sbaffi and Jennifer Rowley. 2017. Trust and Credibility in Web-Based Health Information: A Review and Agenda for Future Research. Journal

- of Medical Internet Research 19, 6 (2017), e218.
 https://doi.org/10.2196/jmir.7579
- Luca Soldaini and Nazli Goharian. 2016.
 Quickumls: A Fast, Unsupervised Approach for
 Medical Concept Extraction. In *MedIR Workshop*, *SIGIR*. Association for Computing Machinery,
 New York, NY, USA, 1–4.

Victor Suarez-Lledo and Javier Alvarez-Galvez.
2021. Prevalence of Health Misinformation on
Social Media: Systematic Review. Journal of
Medical Internet Research 23, 1 (Jan. 2021),
e17187. https://doi.org/10.2196/17187

Petroc Sumner, Solveiga Vivian-Griffiths, Jacky
Boivin, Andy Williams, Christos A. Venetis,
Aimée Davies, Jack Ogden, Leanne Whelan,
Bethan Hughes, Bethan Dalton, Fred Boy, and
Christopher D. Chambers. 2014. The Association
between Exaggeration in Health Related Science
News and Academic Press Releases: Retrospective
Observational Study. *BMJ* 349 (Dec. 2014), g7015.
https://doi.org/10.1136/bmj.g7015

Yalin Sun, Yan Zhang, Jacek Gwizdka, and
Ciaran B. Trace. 2019. Consumer Evaluation of
the Quality of Online Health Information:
Systematic Literature Review of Relevant Criteria
and Indicators. *Journal of Medical Internet Research* 21, 5 (May 2019), e12522.
https://doi.org/10.2196/12522

Carolyn Watters, Wanhong Zheng, and Evangelos Milios. 2002. Filtering for Medical News Items. Proceedings of the American Society for Information Science and Technology 39, 1 (2002), 284–291.

889 https://doi.org/10.1002/meet.1450390131

884

885

890

894

Wilson Wong, Wei Liu, and Mohammed
Bennamoun. 2007. Determining Termhood for
Learning Domain Ontologies Using Domain
Prevalence and Tendency. In Proceedings of the 6th
Australasian Conference on Data Mining. Citeseer,
47–54.

Amélie Yavchitz, Isabelle Boutron, Aida Bafeta,
Ibrahim Marroun, Pierre Charles, Jean Mantz,
and Philippe Ravaud. 2012. Misrepresentation of
Randomized Controlled Trials in Press Releases
and News Coverage: A Cohort Study. *PLOS Medicine* 9, 9 (Sept. 2012), e1001308.
https://doi.org/10.1371/journal.pmed.1001308

Yan Zhang, Yalin Sun, and Bo Xie. 2015. Quality
of Health Information for Consumers on the Web:
A Systematic Review of Indicators, Criteria, Tools,
and Evaluation Results. Journal of the
Association for Information Science and
Technology 66, 10 (2015), 2071–2084.
https://doi.org/10.1002/asi.23311

910 Wanhong Zheng, Evangelos Milios, and Carolyn
911 Watters. 2002. Filtering for Medical News Items
912 Using a Machine Learning Approach. Proceedings
913 of the AMIA Symposium (2002), 949–953.

A Encyclopedia Links

http://health.am/encyclopedia	915
https://medlineplus.gov/encyclopedia.html	916
https://merriam-webster.com/medical	917
https://ucsfhealth.org (var. sub pages)	918
https://www.rxlist.com/	919
drug-medical-dictionary/article.htm	920

914

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

B UMLS Vocabularies

For all vocabulary subsets we use the 2020AB revision of the UMLS Metathesaurus. For the MeSH subset we gather all concepts contained in the MeSH vocabulary and filter out all atoms contained in MeSH. The MeSH Syn. subset also includes all MeSH concepts, but keeps all atoms linked to those concepts irrespective of the vocabulary that atom is from. For the full UMLS subset we use all concepts and atoms from every Category 0 (no additional restrictions or license terms apply) vocabulary.

C BERT Training

The BERT approach was trained by finehuggingface¹⁰ the transformers tuning allenai/scibert_scivocab_uncased check-PyTorch¹¹ and PyTorchLightning¹² point. were used to train the model using a batch size of 32 and learning rate of 0.000005. The input text was split into sentences using nltk¹³ and noun phrases extracted using spacy.¹⁴. Due to the large corpora sizes fine-tuning converged before a single complete epoch was reached. Therefore, training was halted after no decrease in training loss was reached for 15 consecutive training loss samples, where a sample was taken every 1,000 steps.

D Wikidata Details

Root wikidata concepts: fungus (Q764), protein (Q8054), microorganism (Q39833), biogenic substance (Q289472), medical procedure (Q796194), disease causative agent (Q2826767), etiology (Q5850078), physiological condition (Q7189713) and medicinal product (Q86746756).

¹³https://www.nltk.org/

¹⁰https://huggingface.co/

¹¹https://pytorch.org/

¹²https://www.pytorchlightning.ai/

 $^{^{14} \}rm https://spacy.io/$

All relations with a chain of predicates starting at one of the root concepts, and consisting of only subclass of (P279), parent taxon (P171), risk factor (P5642), and optionally ending with instance of (P31) to both the cause and effect concepts are considered as related to health.

956

957

958

961

962

963

964

965

966

967

968

969

In total, 11,160 relations were extracted from Wikidata. We removed 799 invalid relations with missing concept labels. We additionally removed all relations pertaining to COVID-19 (1,044 in total), because these are severely overrepresented and COVID-19 was not yet found in the Textbook corpus, nor CauseNet.

E Grid Search Parameters

For the three vocabulary approaches (MeSH, 970 MeSH Syn. and UMLS) seven Jaccard distance 971 thresholds $(0.4, 0.5, \ldots, 0.9, 1.0)$ were tested. 972 For the three termhood scores (CW, TDS and 973 DW) four health corpora (PubMed, PubMed 974 Central, Textbook, Encyclopedias), three n-975 gram sizes n (uni-, bi- and trigrams) and five 976 values for p_n (1, 2, 5, 10, ∞) for averaging 977 n-gram termhood scores using the generalized 978 mean are tested. Finally, for the final relation 979 classification, the same set of values for p_p are 980 981 tested for averaging cause and effect scores and in addition to the "and" operator. 982