

Enhancing Data Curation for Clinical Trial Registries: Application of Language Models for Drug and Disease Recognition and Normalization

Anonymous ACL submission

Abstract

Clinical trial registry reviews can reveal crucial insights into medical research quality and scope. The current process for generating reports from these registries relies heavily on manual data curation, which includes categorizing trials by disease type and classifying drugs. These tasks are time-consuming and prone to human error. In the present work, we explore the use of automated techniques for extracting drug and disease information, as well as their linking to a medical ontology. By improving the data capture and curation, our aim is to contribute to the development of new systems for reviewing and monitoring clinical trial registries. All resources are available on GitHub¹.

1 Introduction

Public clinical trial registries, such as ClinicalTrials.gov, are essential resources that enable stakeholders—including researchers, patients, health-care providers, and policymakers—to navigate the landscape of drug development. These registries allow for the monitoring of emerging therapeutic targets and substances, and ensuring that new treatments meet safety and efficacy standards. Additionally, they can facilitate the tracking of adverse drug reactions and support the evaluation of clinical trial design quality (Saberwal, 2021).

However, extracting information from these resources is challenging due to large data volume, incomplete and unstructured reporting and variability in terminology (Tse et al., 2018; Pillamarapu et al., 2019; Shi and Du, 2024). While the Aggregate Analysis of ClinicalTrials.gov database (AACT)² has been released in 2011 to enhance access to the data, it provides little automated validation and harmonization of data elements (Tasneem et al., 2012). For example, a recent study found that the

interventions section of ClinicalTrials.gov included non-drug-related terms, hindering comprehensive drug trend analysis (Namiot et al., 2023). Therefore, the current process of evidence synthesis from trial registries relies heavily on manual data curation, including the tasks of categorizing trials by disease type and classifying drugs (Hirsch et al., 2013; Liu et al., 2018). This approach is time-consuming and prone to human error, which might result in inconsistencies and missed information.

Computational methods, especially natural language processing (NLP), can support clinical evidence synthesis by structuring, standardizing, and semantically analyzing data (Marshall et al., 2017; Thomas et al., 2017). Techniques like Named Entity Recognition (NER) identify and categorize text elements such as drug and disease names (Wang et al., 2018). Complementary, Entity Linking (EL) matches these identified elements to unique identifiers in knowledge bases and enables entity normalization, i.e., their uniform representation (Shen et al., 2015; Shi et al., 2023).

Thus, we used NER to explore ways of enhancing the existing condition and intervention fields in ClinicalTrials.gov. We compared neural NER outputs with the existing AACT manual annotations and evaluated a state-of-the-art method for linking entities to the Systematized Nomenclature of Medicine Clinical Terminology (SNOMED CT) (Cornet and de Keizer, 2008).

2 Methods

2.1 Reference Corpus

We worked with a dataset of annotated trials from ClinicalTrials.gov (NeuroTrialNER)³. This dataset includes entity-level annotations in trial titles and summaries, identifying entities like disease names (called “conditions”) and drugs. We analyzed the

¹<https://anonymous.4open.science/r/NeuroTrialDataCuration-3F46/>

²<https://aact.ctti-clinicaltrials.org/>

³Developed within our group, the work is currently under anonymized review. See [Anonymous GitHub](#).

test set of 153 trials, focusing on those with condition annotations (144 trials, 345 annotations) and drug annotations (50 trials, 100 annotations).

2.2 Named Entity Recognition

2.2.1 Model

We used BioLinkBERT as the model reported to achieve the best results on the test set for condition (F1 0.85) and drug name recognition (F1 0.90) (Yasunaga et al., 2022). In previous work the model was fine-tuned on the NeuroTrialNER train set, and we ran inference on a local CPU.

2.2.2 Evaluation

We were interested to understand what are the differences between the entities extracted from the text using BioLinkBERT, and the existing values provided in the ClinicalTrials.gov (AACT) records.

For each clinical trial, we aggregated the token-wise NER extractions into unique entities at the abstract level to enable comparison with AACT. We then determined whether each unique entity from AACT and BioLinkBERT appeared in one or both annotations. Overlaps were identified based on exact or partial token matches, with partial matches defined by significant character overlap, as described in the Appendix A.

To better understand cases where entities were present only in AACT or the BioLinkBERT extractions, we sampled 20 instances where an entity was returned by only one system. Each instance was manually reviewed and classified either as a synonym, false positive, or as a unique true positive for one system, thus a false negative for the other.

2.3 Named Entities Linking

2.3.1 Manual Annotation

NeuroTrialNER did not include annotations for linking named entities to SNOMED nomenclature. To assess performance, two annotators independently linked each manually annotated condition and drug entity from the test set to the ontology entries using the SNOMED CT web browser⁴. They identified the most accurate matches, extracting the concept name and concept IDs. The process is detailed in Appendix B. Inter-annotator agreement (IAA) was measured using Cohen’s kappa statistic, and we report the 95% confidence intervals (CI) (Cohen, 1960).

⁴SNOMED CT Browser

2.3.2 Dictionary Lookup

We used a names dictionary based technique as a simple baseline for the entity linking task. We combined reference terminology from multiple knowledge bases, detailed in Appendix C. This resulted in a dictionary of 25,933 unique drug names and 18,458 unique condition names, including synonyms and lexical variations.

Following the method outlined in Wood (2023), we linked entity words that matched entries from our dictionary. This approach did not accommodate misspellings.

2.3.3 SapBERT and SNOMED

We utilized the Self-alignment Pretraining for BERT (SapBERT) model from the Huggingface library, pre-trained on PubMedBERT full texts, without further fine-tuning or change to the hyperparameters⁵. Inference with the model was performed on local CPU.

We acquired SNOMED CT data from NIH⁶, isolating concepts and synonyms in the categories disorder, finding, procedure and medicinal product. SapBERT vector representations were created for each SNOMED concept and synonym.

For each named entity from the test set, we generated a SapBERT embedding and used it to match the closest SNOMED concepts based on Euclidean distance (Huang et al., 2008). Note that his setup did not take the mention’s context into account (Kartchner et al., 2023). The top five closest matches and their distances (cdist) were retrieved.

2.3.4 Evaluation

The assessment of the linking quality was performed in terms of precision, recall and the F1-measure, as defined in (Shen et al., 2015) and shown in Appendix D.

2.3.5 Experiments

The “cdist” value can be interpreted as an indicator of the match’s accuracy, with larger distances suggesting lower confidence in the match. We aimed to determine an optimal “cdist” threshold, above which entities should not be linked to SNOMED due to a high likelihood of being false positives. To achieve this, we explored various threshold values.

Moreover, it is possible that the manually annotated SNOMED target is not the top match returned by the system but falls within the top k closest

⁵cambridgeltl/SapBERT-from-PubMedBERT-fulltext

⁶NIH SNOMED CT International Edition, April 1, 2024

168 matches ($k=2,3,4,5$). We therefore analyzed how
169 performance varies when considering whether the
170 target entity is among these closely matched enti-
171 ties.

172 3 Results

173 3.1 Named Entity Recognition

174 The UpSet plot in **Figure 1 A** shows the intersec-
175 tion of condition entities extracted by AACT and
176 BioLinkBERT. Around 35% (254) of the entities
177 were recognized by both methods. AACT and Bi-
178 oLinkBERT also uniquely returned 218 and 228
179 conditions, respectively. For drug entities, 52 drug
180 names were overlapping (**Figure 1 B**). Addition-
181 ally, BioLinkBERT uniquely identified 54 drugs,
182 while AACT uniquely identified 32 drugs.

183 To understand the discrepancies between the
184 methods, we manually reviewed a random sam-
185 ple of the entities recognized by only one of the
186 systems and identified the following patterns (see
187 also Table 1 in Appendix E):

- 188 • **Different entity surface forms:** The methods
189 identified the same entity, but they had lexical
190 variations. This was more frequent with drug
191 entities (57%) than conditions (42%).
- 192 • **Unique entity by one method:** BERT de-
193 tected more detailed conditions and interven-
194 tion information. AACT contained entities
195 that BERT could not extract because they were
196 not mentioned in the trial descriptions.
- 197 • **False Positives:** AACT had only 2% false
198 positives, while BERT had 15% for conditions
199 and 5% for drugs.

200 3.2 Named Entity Linking

201 3.2.1 Manual Linking

202 The Cohen’s kappa score between the two anno-
203 tators for linking drug entities was 0.85 (CI: 0.78,
204 0.92), and for linking conditions, it was 0.79 (CI:
205 0.75, 0.84). The two annotators manually reviewed
206 the disagreements and reached a consensus on the
207 final target SNOMED entity, which was then used
208 for model evaluation.

209 3.2.2 Dictionary Lookup

210 Of the 100 drug mentions, 52% were successfully
211 linked using the exact string matching dictionary
212 lookup strategy. This method also successfully
213 linked 123 (36%) of the annotated condition en-
214 tities. Linking conditions was more challenging

215 because the annotations included extra disease char-
216 acteristics such as stage and severity, which were
217 not present in our target disease knowledge bases.

218 3.2.3 Optimal Entity Linking Performance

219 The highest F1 scores were obtained at a cdist
220 threshold of 7.73 for conditions, achieving an F1
221 score of 0.76 (**Figure 2 A**). For drug entities, the
222 highest F1 score of 0.92 was achieved at a cdist
223 threshold of 8.18 (**Figure 3 A** in Appendix F).

224 As seen in **Figure 2 B**, at lower cdist thresh-
225 olds, the model was more stringent, accepting only
226 very close matches. This resulted in higher preci-
227 sion but lower recall, as the model missed some
228 true matches that had a higher Euclidean distance.
229 Conversely, at higher cdist thresholds, the model
230 was less strict, which increased recall by including
231 more true matches, but also decreased precision
232 due to the inclusion of more false positives.

233 3.2.4 Performance at different k

234 **Figure 2 B** demonstrates the relationship between
235 the number of included closest entities (k) and
236 the performance of the entity linking model. The
237 results showed that while precision and recall in-
238 crease with the number of closest entities consid-
239 ered. This indicates that the correct entity is fre-
240 quently found within the top 5 closest entities, sug-
241 gesting that these entities are closely related. Simi-
242 lar results were obtained for drug entities, see
243 **Figure 3 B** in Appendix F.

244 4 Discussion

245 Analysis of entities unique to either AACT or
246 BERT revealed that the same entity often appeared
247 in both extractions with different surface forms,
248 highlighting the challenge of handling extensive
249 synonyms in the biomedical domain (Kartchner
250 et al., 2023). The neural NER approach offered
251 more detailed and standardized annotations. For
252 example it included disease stages and severity
253 grades, while excluding drug dosage information.
254 This suggests the potential for using this technique
255 to automatically extract and standardize entities
256 from trial descriptions, enhancing the granularity
257 and completeness of the data.

258 We also showed that a neural entity linker to
259 a standard medical vocabulary could address the
260 challenge of different entity surface forms. This
261 would facilitate data aggregation across different
262 trials and enable analysis at various hierarchical

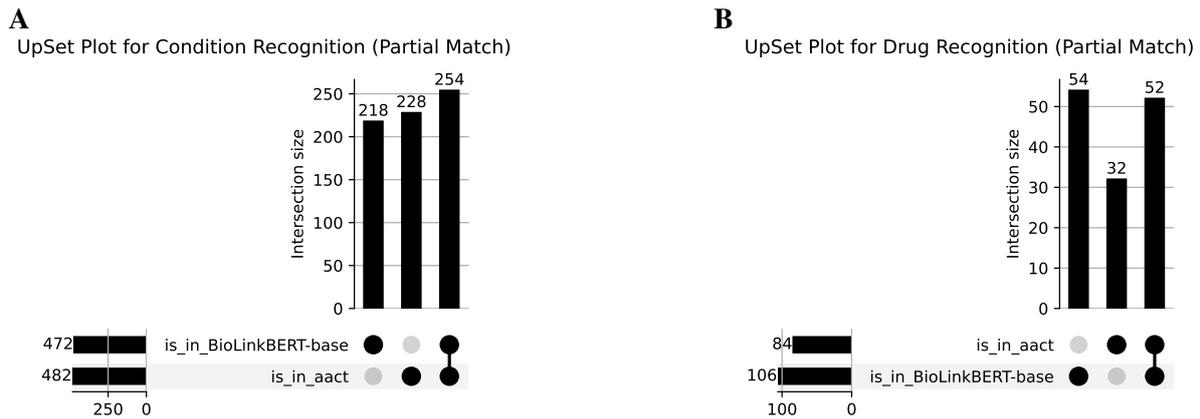


Figure 1: Named entities comparison. Horizontal bars show the total number of unique entities (sets) recognized by each method. The vertical bars indicate the size of intersections between sets. A single filled dot means the set is coming from only one of the outputs, and a connecting line indicates overlapping entities.

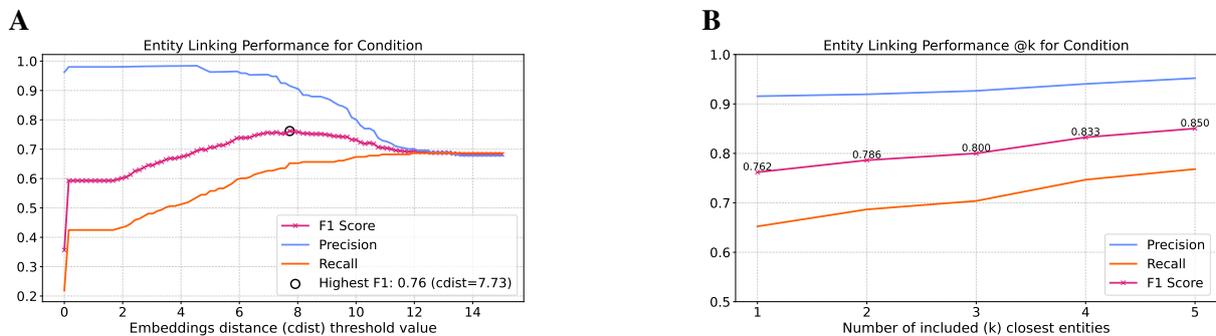


Figure 2: Entity linking experiments for condition. **A**: Impact of different Euclidean distance (cdist) threshold values. **B**: Performance change when considering a k number of closest entities.

263 levels. However, our experiments also uncovered 264
 265 two challenges. The first was the insufficient 266
 267 coverage of concepts in SNOMED. For example, 268
 269 while “Everolimus” is correctly categorized as a 270
 271 substance, its well-known brand name, “Afinitor”, 272
 273 is missing, as well as the drug “Priopidine”. The 274
 275 second challenge stemmed from the lack of con- 276
 277 textual information when using SapBERT, which 278
 279 made it difficult to determine the correct target of 280
 281 the entities. Specifically, for disease names, it was 282
 283 challenging to decide whether an entity referred 284
 285 to a disease or a symptom/finding without the con- 286
 287 text it appeared in. For example, human annotators 288
 289 preferred “Depressive disorder (disorder)” for the 290
 291 entity “depressive”, whereas SapBERT returned 292
 293 “Symptoms of depression (finding)”.

279 5 Conclusion

280 This study evaluated the impact and limitations of 281
 282 NLP-based techniques for automated data extrac- 283
 284 tion from clinical trial registry data. Our findings 285
 286

283 indicate that NER can retrieve entities from trial 284
 285 titles and summaries, potentially replacing or com- 286
 287 plementing the manually provided data in AACT. 288
 289 Additionally, we explored linking entities to struc- 290
 291 tured representations in an ontology and standard- 292
 293 ize variations, addressing a gap in AACT.

289 Future work could expand upon these findings 290
 291 in several directions. First, we identified the need 292
 293 for additional ontologies or knowledge bases to 294
 295 address the issue of missing entities. Second, while 296
 297 we tested a single entity linking approach, there is 298
 299 a need for a more comprehensive benchmarking of 300
 301 different methodologies.

296 A promising future application would be inte- 297
 298 grating these techniques into existing trial registry 299
 300 platforms. This could enhance the data capturing 301
 302 and curation process, making it more complete, 303
 304 standardized and less prone to human errors, thus 305
 306 enhancing the usability and interoperability of the 307
 308 data for downstream tasks, such as monitoring and 309
 310 evidence synthesis.

304 Limitations

305 Our research was limited to trials in neuroscience
306 from ClinicalTrials.gov. However, we believe that
307 the methodologies and approaches we employed
308 could be adapted for use with other clinical trial
309 registry platforms and medical domains.

310 The estimation of optimal parameters for entity
311 linking was exploratory and conducted on the
312 test set of the NeuroTrialNER corpus. However,
313 it would be more rigorous to annotate the validation
314 corpus, optimize the parameters based on that,
315 and then report the performance scores on the test
316 set. Furthermore, the manual annotations for entity
317 linking did not take into account the context in
318 which the entity appeared. It might be necessary
319 to refine the annotation guidelines and differentiate
320 more clearly between target SNOMED concepts
321 such as disorders and findings.

322 Finally, our research was conducted exclusively
323 using English-language data. Expanding this work
324 to include other languages could enrich the dataset
325 and offer more comprehensive insights into global
326 clinical practices.

327 References

328 Jacob Cohen. 1960. [A coefficient of agreement for
329 nominal scales](#). *Educational and psychological mea-
330 surement*, 20(1):37–46.

331 Ronald Cornet and Nicolette de Keizer. 2008. [Forty
332 years of SNOMED: a literature review](#). *BMC Medi-
333 cal Informatics and Decision Making*, 8(1):S2.

334 Bradford R. Hirsch, Robert M. Califf, Steven K.
335 Cheng, Asba Tasneem, John Horton, Karen Chiswell,
336 Kevin A. Schulman, David M. Dilts, and Amy P.
337 Abernethy. 2013. [Characteristics of Oncology Clin-
338 ical Trials: Insights From a Systematic Analysis
339 of ClinicalTrials.gov](#). *JAMA Internal Medicine*,
340 173(11):972–979.

341 Anna Huang et al. 2008. [Similarity measures for text
342 document clustering](#). In *Proceedings of the sixth new
343 zealand computer science research student confer-
344 ence (NZCSRSC2008)*, Christchurch, New Zealand,
345 volume 4, pages 9–56.

346 David Kartchner, Jennifer Deng, Shubham Lohiya, Te-
347 jasri Kopparthi, Prasanth Bathala, Daniel Domingo-
348 Fernández, and Cassie S. Mitchell. 2023. [A Com-
349 prehensive Evaluation of Biomedical Entity Linking
350 Models](#). *Proceedings of the Conference on Empirical
351 Methods in Natural Language Processing. Confer-
352 ence on Empirical Methods in Natural Language
353 Processing*, 2023:14462–14478.

Xu Liu, Yuan Zhang, Ling-Long Tang, Quynh Thu Le,
Melvin L. K. Chua, Joseph T. S. Wee, Nancy Y. Lee,
Brian O’Sullivan, Anne W. M. Lee, Ying Sun, and
Jun Ma. 2018. [Characteristics of Radiotherapy Trials
Compared With Other Oncological Clinical Trials in
the Past 10 Years](#). *JAMA Oncology*, 4(8):1073–1079.

Iain J Marshall, Joël Kuiper, Edward Banner, and By-
ron C Wallace. 2017. [Automating biomedical evi-
dence synthesis: RobotReviewer](#). In *Proceedings of
the conference. Association for Computational Lin-
guistics. Meeting*, volume 2017, page 7. NIH Public
Access.

Eugenia D. Namiot, Diana Smirnovová, Aleksandr V.
Sokolov, Vladimir N. Chubarev, Vadim V. Tarasov,
and Helgi B. Schiöth. 2023. [The international clin-
ical trials registry platform \(ICTRP\): data integrity
and the trends in clinical trials, diseases, and drugs](#).
Frontiers in Pharmacology, 14. Publisher: Frontiers.

Mounika Pillamarapu, Abhilash Mohan, and Gayatri
Saberwal. 2019. [An analysis of deficiencies in the
data of interventional drug trials registered with Clin-
ical Trials Registry - India](#). *Trials*, 20(1):535.

Gayatri Saberwal. 2021. [The Many Uses of Data in
Public Clinical Trial Registries](#). *Current Science*,
120(11):1686.

Wei Shen, Jianyong Wang, and Jiawei Han. 2015. [Entity
Linking with a Knowledge Base: Issues, Techniques,
and Solutions](#). *IEEE Transactions on Knowledge
and Data Engineering*, 27(2):443–460. Conference
Name: IEEE Transactions on Knowledge and Data
Engineering.

Jiyun Shi, Zhimeng Yuan, Wenxuan Guo, Chen Ma,
Jiehao Chen, and Meihui Zhang. 2023. [Knowledge-
graph-enabled biomedical entity linking: a survey](#).
World Wide Web, 26(5):2593–2622.

Xuanyu Shi and Jian Du. 2024. [Constructing a finer-
grained representation of clinical trial results from
ClinicalTrials.gov](#). *Scientific Data*, 11(1):41. Pub-
lisher: Nature Publishing Group.

Asba Tasneem, Laura Aberle, Hari Ananth, Swati
Chakraborty, Karen Chiswell, Brian J McCourt, and
Ricardo Pietrobon. 2012. [The database for aggregate
analysis of ClinicalTrials.gov \(AACT\) and subse-
quent regrouping by clinical specialty](#). *PloS one*,
7(3):e33677.

James Thomas, Anna Noel-Storr, Iain Marshall, Byron
Wallace, Steven McDonald, Chris Mavergames, Paul
Glasziou, Ian Shemilt, Anneliese Synnot, Tari Turner,
et al. 2017. [Living systematic reviews: 2. combin-
ing human and machine effort](#). *Journal of clinical
epidemiology*, 91:31–37.

Tony Tse, Kevin M Fain, and Deborah A Zarin. 2018.
[How to avoid common problems when using Clini-
calTrials.gov in research: 10 issues to consider](#). *Bmj*,
361.

409 Yanshan Wang, Liwei Wang, Majid Rastegar-Mojarad,
410 Sungrim Moon, Feichen Shen, Naveed Afzal, Sijia
411 Liu, Yuqun Zeng, Saeed Mehrabi, Sunghwan Sohn,
412 et al. 2018. [Clinical information extraction appli-](#)
413 [cations: a literature review](#). *Journal of biomedical*
414 *informatics*, 77:34–49.

415 Thomas A Wood. 2023. Drug named entity recognition
416 (computer software), version 1.0.1. To appear.

417 Michihiro Yasunaga, Jure Leskovec, and Percy Liang.
418 2022. [LinkBERT: Pretraining language models with](#)
419 [document links](#). In *Proceedings of the 60th Annual*
420 *Meeting of the Association for Computational Lin-*
421 *guistics (Volume 1: Long Papers)*, pages 8003–8016,
422 Dublin, Ireland. Association for Computational Lin-
423 guistics.

424 A NER Overlap Calculation Details

425 The partial match similarity assessment was cal-
426 culated considering both the number of matching
427 characters and their positions within the strings to
428 determine the closeness of the match⁷. For instance,
429 if the AACT annotation is “hemiplegic cerebral
430 palsy”, and the BioLinkBERT prediction is “cere-
431 bral palsy”, this qualifies as a partial match.

432 B Annotation Guideline for Entity 433 Linking

434 Guidelines

- 435 1. Read the entity from the list of extracted NERs
436 (e.g., column unique_condition_target).
- 437 2. If the entity is not clear, look up the clinical
438 trial from which it was extracted and read the
439 context in which it appears.
 - 440 2.1. If the linking would be possible only
441 through the context, add this as a flag in
442 the designated column of the annotations
443 file (column context_required).
- 444 3. If it is clear what concept is represented,
445 search for it in the SNOMED browser.
 - 446 3.1. If the concept is found:
 - 447 3.1.1. Preferably look for (disease) or (sub-
448 stance) main concept; e.g., for the
449 entity "tic", prefer "Tic disorder (dis-
450 order)" instead of "Tic (finding)".
 - 451 3.1.2. Copy the concept and the
452 concept ID into columns
453 target_snomed_concept and
454 target_snomed_concept_id.

⁷We used the get_close_matches function with cutoff=0.6
from: <https://docs.python.org/3/library/difflib.html>

3.1.3. Always keep the semantic tag, even
if the concept is of another semantic
tag, e.g., (procedure) - this will be
used to know which other semantic
concept to include in the SNOMED
graph.

3.1.4. If more than one entity is extracted,
add all the corresponding matches for
linking, separated with a comma.

3.2. If the concept is not found:

3.2.1. Try to reduce the entity to its main
components, e.g., if the entity was
"post-operative atrial fibrillation" and
this returns no hits from the database,
look for "atrial fibrillation" only;
also, if the entity is an adjective, try
with the noun form, e.g., if the entity
was "acromegalic", try looking for
"acromegaly".

3.2.2. Consider using a synonym.

3.2.3. If a more generic concept is returned,
then use this generic concept in-
stead, e.g., for the entity "autoim-
mune neurological diseases", the re-
sulting match is "Autoimmune dis-
ease (disorder)".

3.2.4. If there is still no meaning-
ful concept returned, write
n.a. for snomed_concept and
snomed_concept_id.

485 C Dictionary Sources for EL

486 For a comprehensive list of neurological and
487 psychiatric diseases, we combined two primary
488 sources: the International Classification of Dis-
489 eases 11th Revision⁸ (ICD-11) and the MeSH terms
490 list⁹. This integration resulted in a list of 18,458
491 unique disease names, including synonyms and
492 lexical variations, categorized under “Mental, be-
493 havioural or neurodevelopmental disorder” and
494 “Neurologic Manifestations”. For drug names, we
495 compiled data from DrugBank¹⁰, Wikipedia, Med-
496 linePlus, and MeSH terms¹¹.

497 D Linking Evaluation Measures

498 Following (Shen et al., 2015), we measured the
499 following metrics to assess the entity linking per-

⁸<https://icd.who.int/icdapi>

⁹Version 2023 obtained as an XML file from
<https://www.nlm.nih.gov/databases/download/mesh.html>

¹⁰<https://go.drugbank.com/>

¹¹<https://pypi.org/project/drug-named-entity-recognition/>

500 formance.

501 Precision measures the accuracy of the entity
502 linking system by evaluating the proportion of cor-
503 rectly linked entity mentions out of all mentions
504 linked by the system.

$$505 \text{ Precision} = \frac{|\text{correctly linked entity mentions}|}{|\text{linked mentions generated by system}|}$$

506 Recall assesses the completeness of the entity link-
507 ing system by evaluating the proportion of cor-
508 rectly linked entity mentions out of all mentions
509 that should have been linked.

$$510 \text{ Recall} = \frac{|\text{correctly linked entity mentions}|}{|\text{entity mentions that should be linked}|}$$

511 The F1 score is combines precision and recall to
512 provide a single score for evaluation:

$$513 F_1 = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

514 E Overview of NER Discrepancies

515 Table 1 presents various types of discrepancies ob-
516 served in the extraction of condition and drug en-
517 tities using BioLinkBERT and AACT. The evalua-
518 tion was based on 20 randomly sampled entities for
519 each type of disagreement, totaling 40 examples
520 per entity category (drug and disease).

521 F Linking Performance for DRUG

522 **Figure 3 A** illustrates the performance of the en-
523 tity linking model for drug entities at various em-
524 beddings distance (cdist) thresholds. As the cdist
525 threshold increases, the model becomes less strin-
526 gent, which impacts these metrics. The highest F1
527 score of 0.92 is achieved at a cdist threshold of 8.18,
528 indicating an optimal balance between precision
529 and recall at this point.

530 **Figure 3 B** shows the entity linking performance
531 for drug entities as a function of the number of
532 included closest entities (k). As the number of clos-
533 est entities increases from 1 to 5, both Recall and
534 F1 Score improve, reaching their peak at k=4 and
535 k=5 with an F1 Score of 0.964. Precision remains
536 high and stable throughout, indicating that includ-
537 ing more closest entities improves recall without
538 significantly compromising precision.

Difference Type	Conditions		Drug		Comment
	Frequency	Example (BERT vs AACT)	Frequency	Example (BERT vs AACT)	
Both extractions represent the same entity	17 (42%)	spine cancer vs spinal bone metastases	23 (57%)	dextrose vs dextrose 5% in water	Often the BERT extractions contained less noise and could be more easily aggregated.
Correct entity available only in AACT	7 (18%)	no annotation vs ovarian cancer	7 (18%)	migraine medications vs verapamil + paroxetine	In all cases, those entities were not available in the title of trial brief description.
Correct entity available only in BERT	8 (20%)	respiratory muscle dysfunction vs muscle weakness	6 (15%)	clozapine vs nmdac plus aifa	BERT's extractions contained more fine-grained details for conditions. Also, interventions tested together with a new intervention are annotated.
False positives AACT	1 (2%)	ketamine	1 (2%)	blood sampling	Observed entities that do not belong to the class.
False positives BERT	6 (15%)	post-, lack	2 (5%)	5 (from dextrose 5% in water)	Observed extraction errors from BERT, e.g., partial extractions of an entity (2 cases).

Table 1: Types of discrepancies for condition and drug entities extraction using BioLinkBERT and AACT. The evaluation of the results for each entity type was based on 20 randomly sampled entities for each disagreement type, i.e., 20 examples where entity extracted only by BERT and not by AACT, and 20 examples where only by AACT but not by BERT.

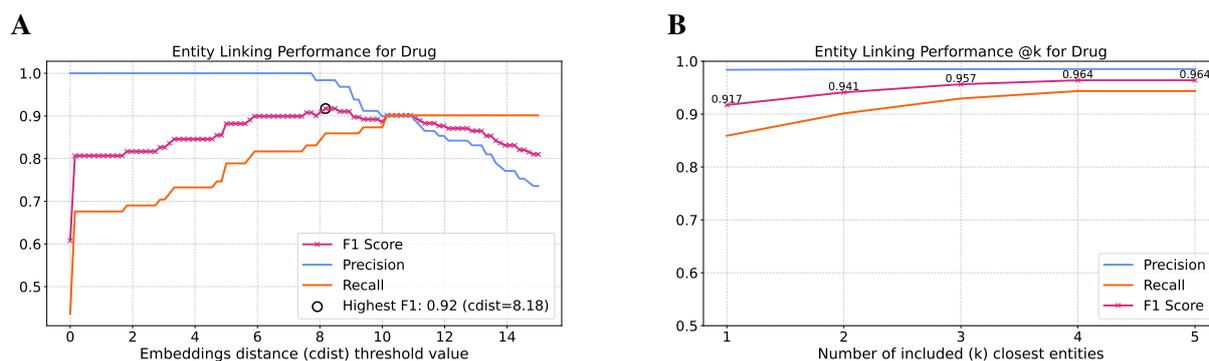


Figure 3: Entity linking experiments for drug. **A:** Impact of different Euclidean distance (cdist) threshold values. **B:** Performance change when considering a k number of closest entities.