

Learning from Negative Samples in Biomedical Generative Entity Linking

Anonymous ACL submission

Abstract

Generative models have become widely used in biomedical entity linking (BioEL) due to their excellent performance and efficient memory usage. However, these models are usually trained only with positive samples—entities that match the input mention’s identifier—and do not explicitly learn from hard negative samples, which are entities that look similar but have different meanings. To address this limitation, we introduce ANGEL (Learning from Negative Samples in Biomedical Generative Entity Linking), the first framework that trains generative BioEL models using negative samples. Specifically, a generative model is initially trained to generate positive entity names from the knowledge base for given input entities. Subsequently, both correct and incorrect outputs are gathered from the model’s top-k predictions. Finally, the model is updated to prioritize the correct predictions through preference optimization. Our models fine-tuned with ANGEL outperform the previous best baseline models by up to an average top-1 accuracy of 1.4% on five benchmarks. When incorporating our framework into pre-training, the performance improvement further increases to 1.7%, demonstrating its effectiveness in both the pre-training and fine-tuning stages. We will make our models and code publicly available upon acceptance.

1 Introduction

Biomedical entity linking (BioEL) involves aligning entity mentions in text with standardized concepts from biomedical knowledge bases (KB) such as UMLS (Bodenreider, 2004) or MeSH (Lipscomb, 2000). BioEL encounters significant challenges due to the diverse and ambiguous nature of biomedical terminology, including synonyms, abbreviations, and terms that look similar but have different meanings. For instance, ‘ADHD’ (CUI:C1263846, where CUI stands for Concept

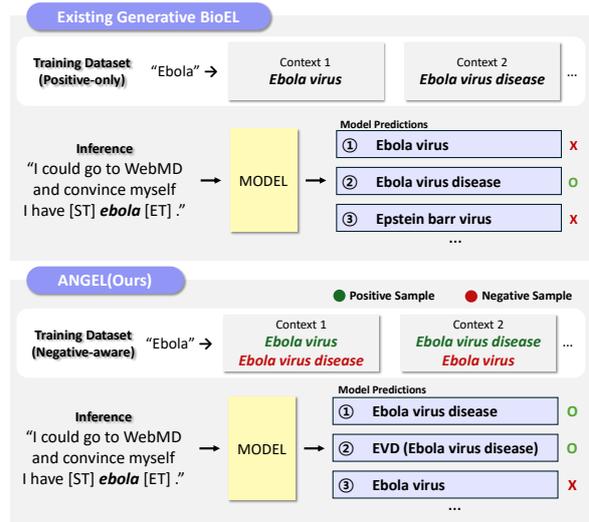


Figure 1: Comparison of training approaches between existing generative BioEL models and our ANGEL method. The main limitation of current generative BioEL methods is that they are trained only on positive samples. This restricts their ability to distinguish between entity names that look similar but different meanings depending on the context. Our ANGEL framework addresses this issue by training the model to prefer positive samples over negative ones.

Unique ID) has synonyms such as hyperkinetic disorder and attention deficit hyperactivity disorder. Additionally, ‘ADA’ can be mapped to either adenosine deaminase (CUI:C1412179) or American Diabetes Association (CUI:C1705019) depending on the context in which the entity appears.

Recent studies have focused on addressing these challenges, broadly categorized into two approaches: similarity-based and generative BioEL. Similarity-based models (Sung et al., 2020; Liu et al., 2021; Lai et al., 2021; Bhowmik et al., 2021; Agarwal et al., 2022) encode input mentions and entities from KBs into the same vector space using embedding models. They then calculate similarity scores to identify the most similar entities for each input entity. Although these approaches have

058 achieved remarkable improvements, they require
059 significant space to index and load embedding vec-
060 tors for all candidate entities (De Cao et al., 2020).
061 Furthermore, representing both the input and candi-
062 date entities as single vectors using a bi-encoder
063 can limit the quality of their representations, mak-
064 ing it difficult to handle challenging cases.

065 On the other hand, generative models (De Cao
066 et al., 2020; Yuan et al., 2022b), built upon an
067 encoder-decoder structure (Lewis et al., 2020; Raf-
068 fel et al., 2020), directly generate the most likely
069 entity name from the KB for the input entity. The
070 output space is dynamically controlled through a
071 constrained decoding strategy, ensuring that only
072 entities from the target KB are generated. Genera-
073 tive models offer several advantages over similarity-
074 based models, including greater memory efficiency
075 and higher performance. They eliminate the need
076 to index large external embedding vectors, and
077 their auto-regressive formulation effectively cross-
078 encodes the input document and candidate entities.

079 However, existing generative models are trained
080 solely on positive samples and do not explicitly
081 learn from negative samples. Despite their high
082 performance, they encounter limitations when dis-
083 tinguishing between biomedical entities with sim-
084 ilar surface forms but different meanings. Al-
085 though similarity-based models address this issue
086 by incorporating negative samples through syn-
087 onym marginalization (Sung et al., 2020) or con-
088 trastive learning (Liu et al., 2021), applying these
089 approaches to generative models is not straightfor-
090 ward. Consequently, generative models may overfit
091 to surface-level features, reducing the models’ abil-
092 ity to generalize effectively across varied contexts,
093 as illustrated in Figure 1.

094 To harness the benefits of generative approaches
095 while overcoming their limitation of not using neg-
096 ative samples, we introduce a novel training frame-
097 work, ANGEL. Our framework operates in two
098 stages: positive-only training and negative-aware
099 training. In the first stage, a generative model is
100 trained to generate biomedical terms from the KB
101 that share the same identifier as the given input
102 entity. In the second stage, we gather both correct
103 and incorrect outputs from the model’s top-k pre-
104 dictions. The model is then updated to prioritize the
105 correct predictions using the direct preference opti-
106 mization (DPO) algorithm (Rafailov et al., 2024).
107 Models trained on our ANGEL framework sig-
108 nificantly outperform the previous best similarity-
109 based and generative BioEL models, achieving an

average accuracy improvement of 1.7% across five
110 datasets. Our contributions are as follows: 111

- We introduce ANGEL, the first-of-its-kind 112
training framework that utilizes negative sam- 113
ples in generative entity linking. ANGEL 114
overcomes the limitations of existing genera- 115
tive approaches by effectively employing neg- 116
ative samples during training. 117
- ANGEL is a versatile framework, demonstrat- 118
ing its applicability in both the pre-training 119
and fine-tuning phases, leading to perfor- 120
mance improvements at each stage. Addition- 121
ally, our method is model-agnostic, consis- 122
tently improving results across various back- 123
bone language models, with gains ranging 124
from 0.9% to 1.7%. 125
- Our best model, pre-trained and fine-tuned 126
with our framework, outperforms the previ- 127
ous best baseline model by 1.7% across five 128
benchmark datasets. 129

2 Related Work 130

2.1 Biomedical Entity Linking 131

132 Biomedical entity linking (BioEL), also known
133 as biomedical entity normalization, is a crucial
134 task because of its application in several down-
135 stream tasks in the biomedical domain, such as
136 literature search (Lee et al., 2016), knowledge ex-
137 traction (Li et al., 2016a; Xiang et al., 2021; Zhang
138 et al., 2023), knowledge graph alignment (Cohen
139 and Hersh, 2005; Lin et al., 2022), and automatic di-
140 agnosis (Shi et al., 2021; Yuan and Yu, 2024). Typi-
141 cally, it is assumed that the target mention is already
142 provided, and the task is solely to link this mention
143 to the appropriate entity name from the KB. End-
144 to-end BioEL (Zhou et al., 2021; Ujji et al., 2021),
145 which also involves identifying mentions within a
146 sentence, is being actively researched, but this is
147 not our focus and will not be discussed in detail.

148 Traditional classification-based approaches
149 (Limsopatham and Collier, 2016a; Miftahutdinov
150 et al., 2019) employed a softmax layer for classifi-
151 cation, treating concepts as categorical variables
152 and thereby losing the detailed information of
153 concept names. Similarity-based (Sung et al.,
154 2020; Liu et al., 2021; Lai et al., 2021; Zhang
155 et al., 2022) models have significantly improved
156 BioEL performance, which encodes mentions
157 and candidate entity names in the same vector

space. They are characterized by high memory consumption due to the need to encode entities into pre-computed embeddings, posing scalability challenges with large datasets (De Cao et al., 2020). Several studies have integrated the concept of clustering into BioEL (Angell et al., 2021; Agarwal et al., 2022).

2.2 Generative Entity Linking

Generative models have become a powerful method for entity linking by overcoming the limitations of similarity-based models. The GENRE framework (De Cao et al., 2020) was the first to demonstrate this approach. To enhance precision and reduce memory usage, GENRE introduced a constrained decoding method (Hokamp and Liu, 2017) using a prefix tree (trie), which restricts the output space to valid entity names. This technique also facilitates easy updates to the set of entities, making the system highly adaptable to changes in the KB. In the biomedical field, notable examples of generative models include GenBioEL (Yuan et al., 2022b) and BioBART (Yuan et al., 2022a). GenBioEL, in particular, is the first model to apply a generative model BART (Lewis et al., 2020) to BioEL, after pre-training it using UMLS. Additionally, several hybrid approaches, known as retrieve-and-generate methods, have been proposed (Xu et al., 2023; Lin et al., 2024). In these methods, a similarity-based model first retrieves the top-k candidates, which are then reranked using a generative model. Although generative approaches have shown high performance, their training has typically been limited to positive samples, as discussed in the introduction section. In this study, we introduce the use of negative samples during training and demonstrate that this approach can significantly enhance the performance of generative models.

3 Method

3.1 Task Formulation

Let $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ be a human-labeled dataset, where \mathbf{x}_n represents an input text and y_n is the gold identifier defined in a KB denoted by \mathcal{E} . Each \mathbf{x}_n contains a target entity mention \mathbf{m}_n along with its surrounding contextual information \mathbf{c}_n^- and \mathbf{c}_n^+ , which represents the tokens before and after the entity mention \mathbf{x}_n , respectively. For simplicity, we will omit the subscript n . Our goal is to map each mention \mathbf{x} to its corresponding identifier y^*

from the set of entity names \mathcal{E} as follows:

$$y^* = \mathcal{F}(\operatorname{argmax}_{e \in \mathcal{E}} p_\theta(e|\mathbf{x})), \quad (1)$$

where \mathcal{F} is a mapping function that converts entities to their identifiers, and θ represents the model parameters.

Previous generative BioEL approaches train the model to generate a synonym $s \in S$ for the given mention in an autoregressive manner as follows:

$$p_\theta(\mathbf{s} | \mathbf{x}, \mathbf{v}) = \prod_{t=1}^T p_\theta(s_t | s_{<t}, \mathbf{x}, \mathbf{v}), \quad (2)$$

where $S \subset \mathcal{E}$ is the set of entity names (i.e., synonyms) corresponding to the identifier y^* , and T is the number of tokens of the synonym \mathbf{s} and s_t indicates the t -th token of the synonym. The prefix prompt \mathbf{v} to the decoder, represented as ‘[BOS] \mathbf{m} is’, is designed to make the decoder’s output resemble a natural language sentence, which helps to minimize discrepancies between language modeling and fine-tuning on the BioEL task. The target mention in the input is surrounded by the special tokens, [ST] and [ET], as follows:

$$[\text{BOS}] \mathbf{c}^- [\text{ST}] \mathbf{m} [\text{ET}] \mathbf{c}^+ [\text{EOS}],$$

where the special tokens [BOS] and [EOS] represent the ‘Begin Of Sentence’ and ‘End Of Sentence,’ respectively.

As shown in Equation 2, existing models are trained only to output synonyms corresponding to the given mention (i.e., positive samples), without learning from negative samples. In contrast, we introduce a new method called negative-aware training, which allows the model to learn by comparing both positive and negative samples, enhancing the model’s generalizability. We will describe our framework in detail in the following sections.

3.2 ANGEL Framework

Our framework consists of two main stages: positive-only training, which warms up the model on target datasets, and negative-aware training, which continuously improves the model by learning from negative samples (see Figure 2).

Positive-only Training In this initial stage, the goal is to learn the morphological similarities among synonyms. To achieve this, we train the model to generate synonyms (i.e., positive samples) that are predefined in the KB, similar to

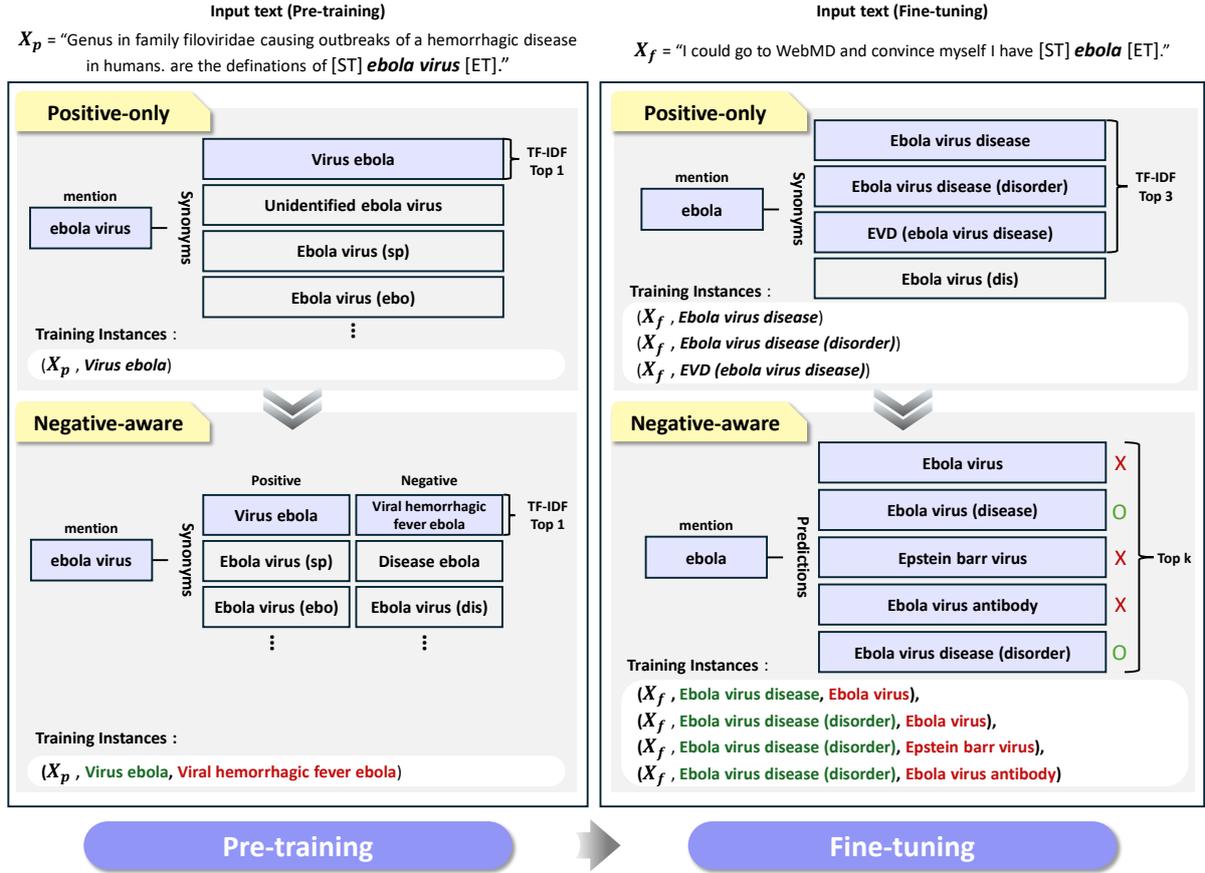


Figure 2: Overview of our method ANGEL (Learning from Negative Samples in Biomedical Generative Entity Linking). The core idea of ANGEL is to enhance both pre-training and fine-tuning by incorporating negative samples, which are obtained either through TF-IDF similarity or the model’s top-k predictions. This approach helps the model distinguish subtle differences between correct and incorrect entities.

traditional generative methods (see Equation 2). Our preliminary study indicated that using all synonyms, $s \in \mathcal{S}$, was not effective. Also, relying solely on the top-1 synonym, as done in a previous study (Yuan et al., 2022b), may limit generalizability. Therefore, we select several of the most similar synonyms to the mention based on their TF-IDF similarity, which is calculated as follows:

$$\hat{\mathcal{S}} = \text{argsort}_{s \in \mathcal{S}}(\text{TFIDF}(\mathbf{m}, s)), \quad (3)$$

where $\text{TFIDF}(\cdot)$ returns similarity scores. We use the top-k subset $\hat{\mathcal{S}}_k = \hat{\mathcal{S}}[:k] = \{\hat{s}_1, \dots, \hat{s}_k\}$ as training instances for each mention.

Negative-aware Training Although surface similarities are a useful feature for BioEL, over-relying on them can limit the model’s generalization ability. To address this issue, we update the model using negative-aware training. First, we obtain the top-k predictions of the model for mentions in the training dataset. We then automatically construct

a training dataset consisting of triplets: a mention with context (if it exists), a correct prediction, and an incorrect prediction. Among all pairs of correct and incorrect predictions, we select only those pairs where the incorrect prediction’s rank is higher than the correct prediction’s rank. Particularly, when the highest ranked entity is the correct one, we pair this entity with the highest ranked incorrect entity to preserve the model’s prior learning. Finally, using this dataset \mathcal{D}' , we then optimize the model with the DPO algorithm. This maximizes the likelihood of generating the correct prediction, e_w , over the incorrect prediction, e_l , defined as follows:

$$\mathcal{L}(p_\theta; p_{\text{ref}}) = -\mathbb{E}_{(\mathbf{x}, e_w, e_l) \sim \mathcal{D}'} \left[\log \sigma \left(\beta \log \frac{p_\theta(e_w | \mathbf{x})}{p_{\text{ref}}(e_w | \mathbf{x})} - \beta \log \frac{p_\theta(e_l | \mathbf{x})}{p_{\text{ref}}(e_l | \mathbf{x})} \right) \right], \quad (4)$$

where p_θ is the generative model to be trained, p_{ref} is the generative model that was trained in the pre-

vious stage using positive-only training, σ is a sigmoid function, and β is a hyperparameter.

Applying ANGEL in Pre-training Our framework is versatile, supporting not only fine-tuning with labeled datasets but also pre-training with the KB itself. Initially, we conduct positive-only training using synonym lists defined in UMLS, the most extensive KB in the biomedical field. We generate its surrounding contextual information automatically for each entity, using clause templates or definitions, as outlined in GenBioEL (Yuan et al., 2022b).¹ We then use the TF-IDF similarity to identify the top synonym and set it as the target (Equation 3). For negative-aware training, using the model’s top-k predictions is impractical due to UMLS’s vastness, which includes over 3 million entities. Instead, we use entities with the highest TF-IDF similarity but different identifiers from the input mentions as negative samples. This approach significantly reduces computation, enabling effective training across the entire UMLS.

4 Experiments

4.1 Datasets

We utilized five popular BioEL benchmark datasets: NCBI-disease (Doğan et al., 2014), BC5CDR (Li et al., 2016b), COMETA (Basaldella et al., 2020), AskAPatient (Limsopatham and Collier, 2016b), and MedMentions (Mohan and Li, 2019), with the ST21pv subset used for MedMentions. Due to the lack of a test set in the AskAPatient dataset, we adhere to the 10-fold evaluation protocol outlined by Limsopatham and Collier (2016b). Also, this dataset does not include context for the mentions. In the following tables, NCBI-disease, AskAPatient, and MedMentions are denoted as NCBI, AAP, and MM-ST21pv, respectively. Please refer to Appendix B for detailed dataset descriptions and statistics.

4.2 Baseline Models

We use top-performing similarity-based models (Sung et al., 2020; Liu et al., 2021; Lai et al., 2021; Zhang et al., 2022) and generative models (Lewis et al., 2020; Yuan et al., 2022a,b) as our baselines. To ensure a fair comparison with our model under identical experimental conditions, we replicate the following generative models and then

apply our ANGEL framework to them: (1) BART-large (Lewis et al., 2020) is an encoder-decoder language model pre-trained on a general-domain corpus. (2) BioBART-large (Yuan et al., 2022a) is the BART-large model continuously pre-trained on a biomedical-domain corpus. (3) GenBioEL (Yuan et al., 2022b) is initialized with the weights of the BART-large model and then pre-trained specifically for BioEL using UMLS.

In BioEL, several studies utilize retrieve-and-generate methods (Xu et al., 2023; Lin et al., 2024). This involves a similarity-based model retrieving the top-k candidates from the KB, followed by a generative model reranking these candidates. We exclude these methods from our experiments to focus on introducing a novel training method that uses negative samples for generative entity linking and demonstrating its effectiveness in a single generative model. Future research could apply our methodology to enhance the reranking process of generative models.

4.3 Implementation Details

Our framework is applied to each of these models during fine-tuning, referred to as ANGEL_{FT}, and during both pre-training and fine-tuning, referred to as ANGEL_{PT+FT}. For pre-training, we utilized the 2020AA version of the UMLS database,² which comprises 3.09M entities, of which 199K concepts contain definitions. In the pre-training phase, the model was trained for 5 epochs, with checkpoints created every 500 steps. We selected the best checkpoints based on the validation sets. During fine-tuning on MM-ST21pv, we also used the 2020AA version of UMLS because the 2017AA version was not directly accessible. Please note that the reported scores of baseline models were measured based on the 2017AA version of UMLS. For determining the best hyperparameters in positive-only training, we searched for the optimal learning rate within the range of 1e-5 to 3e-7 and adjusted the batch size between 8 and 16, following the approach of Yuan et al. (2022b). For the negative-aware training, we searched for the optimal learning rate from 2e-5 to 1e-6 and experimented with batch sizes ranging from 16 to 64. In pre-processing, following Yuan et al. (2022b), we expanded abbreviations using AB3P (Sohn et al., 2008), lowercase texts, mark mention boundaries with special tokens [ST] and [ET], and discard

¹Detailed descriptions of data generation during the pre-training stage are provided in Appendix A.

²<https://www.nlm.nih.gov/research/umls/licensedcontent/umlsarchives04.html>

Model	NCBI	BC5CDR	COMETA	AAP	MM-ST21pv	Average
<i>Similarity-based BioEL</i>						
BioSYN (Sung et al., 2020)	91.1	-	71.3	82.6	-	-
SapBERT (Liu et al., 2021)	92.3	-	75.1	89.0	-	-
ResCNN (Lai et al., 2021)	92.4	-	80.1	-	55.0	-
KRISSBERT (Zhang et al., 2022)	91.3	-	-	-	72.2	-
<i>Generative BioEL (reported)</i>						
BART (Lewis et al., 2020)	90.2	92.5	80.7	88.8	71.5	84.7
BioBART (Yuan et al., 2022a)	89.9	93.3	81.8	89.4	71.8	85.2
GenBioEL (Yuan et al., 2022b)	91.9	93.3	81.4	89.3	-	-
<i>Generative BioEL (reproduced)</i>						
BART [†] (Lewis et al., 2020)	90.3	93.0	80.4	88.7	70.1	84.5
+ ANGEL _{FT} (Ours)	91.4 (+1.1)	93.6 (+0.6)	81.3 (+0.9)	89.5 (+0.8)	71.2 (+1.1)	85.4 (+0.9)
BioBART [†] (Yuan et al., 2022a)	89.4	93.5	81.3	89.3	71.3	85.0
+ ANGEL _{FT} (Ours)	91.9 (+2.5)	94.7 (+1.2)	82.2 (+0.9)	<u>89.9</u> (+0.6)	73.4 (+2.1)	<u>86.4</u> (+1.4)
GenBioEL [†] (Yuan et al., 2022b)	91.0	93.1	80.9	89.3	70.7	85.0
+ ANGEL _{FT} (Ours)	<u>92.5</u> (+1.5)	94.4 (+1.3)	<u>82.4</u> (+1.5)	<u>89.9</u> (+0.6)	71.9 (+1.2)	86.2 (+1.2)
+ ANGEL _{PT+FT} (Ours)	92.8 (+1.8)	<u>94.5</u> (+1.4)	82.8 (+1.9)	90.2 (+0.9)	<u>73.3</u> (+2.6)	86.7 (+1.7)

Table 1: The top-1 accuracy of the models across the five BioEL datasets. The [†] symbol indicates that the results have been reproduced. Our ANGEL framework is applied to generative BioEL models during fine-tuning (ANGEL_{FT}) and both pre-training and fine-tuning (ANGEL_{PT+FT}). We exclude the performance of similarity-based models on BC5CDR, as they were evaluated separately on the chemical and disease subsets, differing from our settings.

mentions that overlap or are missing from the target KB. The best hyperparameter configurations are detailed in Appendix C. We used the source codes provided by Yuan et al. (2022b)³ and alignment handbook (Tunstall et al., 2023)⁴. During pre-training stage, we trained our model using eight 80G A100 GPUs for 12 hours. During fine-tuning stage, we used a single A100 GPU.

4.4 Results

Consistent with previous studies (Sung et al., 2020; Liu et al., 2021), we use accuracy at top-1 (Acc@1) as our evaluation metric. This metric measures the percentage of mentions where the model correctly ranks the gold standard identifier as the top choice. Table 1 demonstrates that our framework consistently improves the performance of generative models. Specifically, our fine-tuning method (i.e., ANGEL_{FT}) improves the Acc@1 scores of BART, BioBART, and GenBioEL by 0.9%, 1.4%, and 1.2%, respectively. When pre-training is also applied (i.e., ANGEL_{PT+FT}) to GenBioEL, the improvement increases to 1.7%, underscoring the effectiveness of both pre-training and fine-tuning in ANGEL. Our best model (i.e., GenBioEL with ANGEL_{PT+FT}) outperforms all baseline models, whether they are

³<https://github.com/Yuanhy1997/GenBioEL>

⁴<https://github.com/huggingface/alignment-handbook>

similarity-based or generative BioEL models.

5 Analysis

5.1 Ablation Study

We conducted detailed analyses on our negative-aware training. Additionally, the impact of the number of synonyms in positive-only training is provided in Appendix D.

Effect of Pre-training Table 2 highlights the effectiveness of ANGEL’s pre-training by comparing other pre-training methods. BART, pre-trained using a standard language modeling objective but not specifically tailored for BioEL tasks, shows the lowest performance. In contrast, GenBioEL, pre-trained using synonyms from UMLS in a similar manner to our positive-only training, initially demonstrates a substantial performance advantage over BART. However, this gap narrows considerably after fine-tuning, to the point where it is no longer statistically significant. When ANGEL’s negative-aware training is applied to GenBioEL, its performance improves significantly, achieving gains of 16.6% on BC5CDR and 10.9% on AAP. Even after fine-tuning, the performance gap remains noticeable, with a difference of 1.4% on BC5CDR and 0.9% on AAP.

Model	FT	BC5CDR	AAP
BART	✗	0.8	15.6
GenBioEL	✗	33.1	50.6
+ ANGEL (Ours)	✗	49.7	61.5
BART	✓	93.0	88.7
GenBioEL	✓	93.1	89.3
+ ANGEL (Ours)	✓	94.5	90.2

Table 2: The top-1 accuracy of models with different pre-training strategies, along with the fine-tuned scores. ‘FT’ denotes fine-tuning, with ✗ representing pre-trained models without fine-tuning, and ✓ indicating models fine-tuned on human-annotated training sets.

Method	BC5CDR	AAP
ANGEL (Ours)	94.5	90.2
GenBioEL	93.1 (-1.4)	89.3 (-0.9)
Prediction-based e_l ⇒ TF-IDF-based e_l	94.4 (-0.1)	90.0 (-0.2)
$p_\theta(e_l) > p_\theta(e_w)$ Pairs ⇒ All Possible Pairs	94.0 (-0.5)	90.0 (-0.2)
e_l within Top-5 ⇒ Top-10 Predictions	94.4 (-0.1)	89.9 (-0.3)

Table 3: The ablation study on positive (e_w) and negative (e_l) pair selection during negative-aware fine-tuning. ‘⇒’ indicates a modification in our method, specifically in the selection of either negative samples or pairs.

Selection of Positive and Negative Pairs Table 3 presents various methods for constructing positive-negative pairs in negative-aware training. We investigated the effects of three different aspects: negative sampling techniques, the ranking of negative samples, and top-k selection, on the BC5CDR and AAP datasets. Notably, all three model variants significantly improved performance compared to the baseline GenBioEL model, highlighting the effectiveness of our negative-aware training, irrespective of the specific techniques used for selecting positive-negative pairs. (1) First, when we modified our approach to extract negative samples based on TF-IDF, similar to the method used during pre-training, performance declined by 0.1% on the BC5CDR and by 0.2% on the AAP. This indicates that our approach, which allows the model to learn from its errors, is more effective than relying on TF-IDF similarity. While the TF-IDF-based method tends to select pairs where the positive and negative examples have similar surface forms, our error-driven approach enables the selection of more diverse negative samples without such constraints. (2) Additionally, when we substituted our

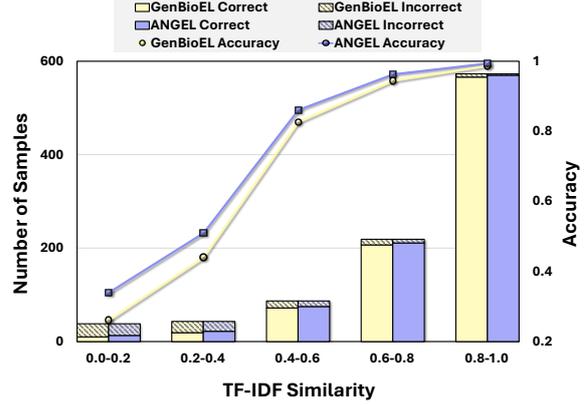


Figure 3: In-depth evaluation of GenBioEL and our ANGEL models based on the TF-IDF similarity between the input mentions and gold-standard entities. The NCBI-disease dataset was used. Further analysis is provided in Appendix E.

negative selection method—where negative samples are ranked higher than positive ones—with an approach that includes all positive-negative pairs regardless of rank, the performance dropped by 0.5% on the BC5CDR and by 0.2% on the AAP. (3) Finally, increasing the top-k selection from the top-5 to the top-10 predictions resulted in a performance decline of 0.1% on the BC5CDR and 0.3% on the AAP. Increasing the top-k can indeed gather more diverse negative samples; however, it may also increase the collection of typical samples that the model finds less confusing, potentially degrading performance. Therefore, maintaining a proper balance between diversity and difficulty in sample selection is crucial.

5.2 Error Analysis

We conducted an in-depth evaluation of the models based on the similarity between the input mentions and the gold-standard entities. Similarity was calculated using tri-gram TF-IDF, with the gold-standard entity determined as the candidate synonym with the highest similarity score to the input mention. The similarity scores, ranging from 0 to 1.0, were divided into five bins, and accuracy was measured for each bin. As shown in Figure 3, errors predominantly occurred in the 0.0-0.2 and 0.2-0.4 bins. This suggests that models tend to struggle when the surface forms of the input mentions are not closely aligned with those of the gold-standard entities. Our method improves the generalizability of the model, leading to an overall reduction in GenBioEL’s errors across all bins, with particularly notable improvements in cases of low similarity.

Rank	SapBERT	GenBioEL	ANGEL (Ours)
... aggressive the same way someone with [ST] <i>ASPD</i> [ET] would be, except teenagers ... (SNOMED CT:26665006)			
1	ASP	Anankastic personality disorder	Antisocial personality disorder (disorder)*
2	Acquired immune deficiency syndrome (disorder)	Borderline personality disorder	Antisocial personality disorder*
3	Acquired immune deficiency syndrome	Oppositional defiant disorder	Borderline personality disorder (disorder)
4	Mesalazine	Antisocial personality disorder*	Obsessive compulsive disorder (disorder)
5	Cryopyrin associated periodic syndrome (disorder)	Oppositional defiant disorder (disorder)	Dissocial personality disorder*
... I switched from lantus to [ST] <i>basaglar</i> [ET] in january and ... (SNOMED CT:411529005)			
1	Beagle	Linagliptin substance	Insulin glargine substance*
2	Basiliximab sodium	Benzodiazepine substance	Insulin glargine*
3	Basiliximab substance	Carisoprodol substance	Insulin glulisine substance
4	Albiglutide	Cariprazine	Ulipristal substance
5	Albiglutide substance	Benzocaine containing product	Lansoprazole
... effects on amino acid (r-aminobutyric acid (GABA), [ST] <i>glutamine</i> [ET], aspartate and glutathione) levels ... (MeSH:D018698)			
1	Glutamine	Glutamine	L-glutamine
2	Glutamic acid*	Glutamic acid*	Glutamine
3	L-glutamine	Glutamylmethionine	D-glutamine
4	L-glutamic acid*	Glutamylalanine	Glutamic acids
5	Glutamic acids	Glutaminic acids	Glutamic acid*

Figure 4: The top 5 predictions from different BioEL models are presented. Entity names with correct identifiers are highlighted in boldface with an asterisk. The first and second examples highlight the strengths of our model, while the final example illustrates its limitations. For a detailed explanation, please refer to the main text.

487 However, significant challenges remain, as the ac- 488 curacy of our model is only 34.2% in the 0–0.2 bin, 489 highlighting the need for further improvement.

490 5.3 Case Study

491 Figure 4 illustrates the predictions of SapBERT, 492 GenBioEL, and ANGEL. In the first example, the 493 mention ‘ASPD’ is an abbreviation for ‘antisocial 494 personality disorder’ (also known as ‘dissocial per- 495 sonality disorder’). SapBERT incorrectly predicts 496 ‘ASP’ due to the similarity in surface form. Gen- 497 BioEL struggles to distinguish between correct en- 498 tity names and those containing the words ‘person- 499 ality disorder’. In contrast, our model successfully 500 identifies the correct entities, without being mis- 501 led by false entity names that contain overlapping 502 terms. The second example involves the mention 503 ‘basaglar,’ a biosimilar medication that contains in- 504 sulin glargine, a long-acting insulin. The challenge 505 here arises from the fact that product names can 506 differ significantly from the biomedical terms used 507 to describe their active ingredients. This discrep- 508 ancy leads to failures in both SapBERT and Gen- 509 BioEL, as they struggle to connect the brand name 510 to its corresponding biomedical entity. Neverthe- 511 less, our model successfully identifies the correct 512 entity, showcasing its ability to handle such com- 513 plex cases effectively. In the final example, our

method was less effective. For the mention of ‘glu- 514 tamine,’ neither SapBERT nor GenBioEL identi- 515 fied the correct answer, but they did rank ‘Glutamic 516 acid,’ the correct entity, within the top 5 candidates. 517 Our model, however, ranked the correct answer 518 slightly lower. Consequently, while our model 519 shows a notable improvement in top-1 accuracy, 520 the increase in top-5 accuracy is relatively modest 521 in some datasets. The effectiveness of our method 522 also varies across different datasets. We discuss 523 this limitation in more detail in Appendix F, noting 524 that such cases are an area for further exploration. 525

526 6 Conclusions

527 In this study, we discussed the importance of nega- 528 tive samples in training generative BioEL models 529 and introduced ANGEL, the first framework in this 530 field to effectively incorporate negative-aware train- 531 ing into a generative model. Our models demon- 532 strated the ability to learn subtle distinctions be- 533 tween entities with similar surface forms and con- 534 texts. Experimental results showed that ANGEL 535 outperformed existing similarity-based and gener- 536 ative models, with notable performance improve- 537 ments of 0.9%, 1.4%, and 1.7% for BART, Bio- 538 BART, and GenBioEL, respectively, while achiev- 539 ing the best performance across five public BioEL 540 datasets. 540

541 Limitations

542 Our method is versatile and applicable to any
543 generative model, but it has only been tested on
544 encoder-decoder models and not on decoder-only
545 models such as BioGPT (Luo et al., 2022). We
546 plan to further investigate the effect of our method
547 on these models. Additionally, it has not been
548 tested on recent open-source large language models
549 (LLMs) (Touvron et al., 2023; Chen et al., 2023).
550 While we acknowledge that incorporating compar-
551 isons with LLMs and further assessing the effec-
552 tiveness of our approach would be an interesting
553 direction, using LLMs for entity linking presents
554 new challenges. The primary concern with larger
555 models is their inefficiency, particularly regarding
556 slower inference speeds and higher memory re-
557 quirements, which may render them unsuitable for
558 most real-world applications. This issue becomes
559 particularly problematic in fields such as biomed-
560 ical information extraction, where processing mil-
561 lions of publications to extract meaningful insights
562 is essential.

563 Our negative-aware training method may not
564 be limited to a specific domain, yet we have only
565 evaluated it on biomedical-domain datasets, which
566 restricts the demonstration of its broad applicabil-
567 ity. Nevertheless, we would like to emphasize the
568 reasons for focusing on the biomedical domain.
569 Biomedical entity linking has unique characteris-
570 tics that differentiate it from other domains, mak-
571 ing this problem both challenging and interesting.
572 In general domains, ambiguity typically arises be-
573 tween different types of entities (e.g., whether “Liv-
574 erpool” refers to a city or a sports club). Simi-
575 larly, in the biomedical domain, ambiguity exists
576 between different types, such as whether “Ebola”
577 in Figure 1 refers to a disease or a virus. Addi-
578 tionally, biomedical entities often exhibit signifi-
579 cant variations in their surface forms, even when
580 they share the same identifier, i.e., they refer to the
581 same entity. As shown in Figure 4, “Basaglar” can
582 be expressed as other variations such as “insulin
583 glargine substance” or “insulin glargine.” Further-
584 more, terms like “substance” in the entity “insulin
585 glargine substance” overlap with many other en-
586 tities (e.g., “Basiliximab substance,” “Linagliptin
587 substance,” “Benzodiazepine substance”), making
588 the task even more complex. Therefore, distin-
589 guishing between numerous candidates with simi-
590 lar surface forms is especially crucial in biomedical
591 entity linking. We believe that our method, which

592 trains the model using negative samples with simi-
593 lar structures, is particularly well-suited to tackle
594 this challenge. However, exploring the applica-
595 tion of our approach in other domains would be a
596 valuable direction for future research.

Ethical Considerations

597 This study complies with ethical standards, ensur-
598 ing that all datasets and models adhere to their
599 respective licenses and usage terms. Biomedical
600 examples are included to illustrate the methodol-
601 ogy; however, they serve explanatory purposes and
602 may not fully represent real-world scenarios. While
603 the model achieves notable improvements, its limi-
604 tations in handling low-similarity cases underscore
605 the importance of rigorous validation prior to de-
606 ployment, especially in sensitive applications. To
607 minimize risks, the model is recommended as a re-
608 ference tool rather than for direct decision-making
609 in critical contexts. 610

References

- 611 Dhruv Agarwal, Rico Angell, Nicholas Monath, and
612 Andrew McCallum. 2022. Entity linking via explicit
613 mention-mention coreference modeling. In *Proceed-
614 ings of the 2022 Conference of the North Ameri-
615 can Chapter of the Association for Computational
616 Linguistics: Human Language Technologies*, pages
617 4644–4658. 618
- Rico Angell, Nicholas Monath, Sunil Mohan, Nishant
619 Yadav, and Andrew McCallum. 2021. Clustering-
620 based inference for biomedical entity linking. In
621 *Proceedings of the 2021 Conference of the North
622 American Chapter of the Association for Computa-
623 tional Linguistics: Human Language Technologies*,
624 pages 2598–2608. 625
- Marco Basaldella, Fangyu Liu, Ehsan Shareghi, and
626 Nigel Collier. 2020. COMETA: A corpus for medical
627 entity linking in the social media. In *Proceedings
628 of the 2020 Conference on Empirical Methods in
629 Natural Language Processing (EMNLP)*, pages 3122–
630 3137. Association for Computational Linguistics. 631
- Rajarshi Bhowmik, Karl Stratos, and Gerard De Melo.
632 2021. Fast and effective biomedical entity linking
633 using a dual encoder. In *Proceedings of the 12th
634 International Workshop on Health Text Mining and
635 Information Analysis*, pages 28–37. 636
- Olivier Bodenreider. 2004. The unified medical lan-
637 guage system (umls): integrating biomedical termi-
638 nology. *Nucleic acids research*, 32(suppl_1):D267–
639 D270. 640
- Eunsuk Chang and Javed Mostafa. 2021. The use of
641 snomed ct, 2013-2020: a literature review. *Journal*
642

643	<i>of the American Medical Informatics Association</i> ,	Jiao Li, Yueping Sun, Robin J Johnson, Daniela Sci-	696
644	28(9):2017–2026.	aky, Chih-Hsuan Wei, Robert Leaman, Allan Peter	697
645	Zeming Chen, Alejandro Hernández Cano, Angelika	Davis, Carolyn J Mattingly, Thomas C Wieggers, and	698
646	Romanou, Antoine Bonnet, Kyle Matoba, Francesco	Zhiyong Lu. 2016a. Biocreative v cdr task corpus:	699
647	Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf,	a resource for chemical disease relation extraction.	700
648	Amirkeivan Mohtashami, et al. 2023. Meditron-70b:	<i>Database</i> , 2016.	701
649	Scaling medical pretraining for large language mod-	Jiao Li, Yueping Sun, Robin J Johnson, Daniela Sci-	702
650	els. <i>arXiv preprint arXiv:2311.16079</i> .	aky, Chih-Hsuan Wei, Robert Leaman, Allan Peter	703
651	Aaron M Cohen and William R Hersh. 2005. A survey	Davis, Carolyn J Mattingly, Thomas C Wieggers, and	704
652	of current work in biomedical text mining. <i>Briefings</i>	Zhiyong Lu. 2016b. Biocreative v cdr task corpus:	705
653	<i>in bioinformatics</i> , 6(1):57–71.	a resource for chemical disease relation extraction.	706
654	Allan Peter Davis, Thomas C Wieggers, Michael C	<i>Database</i> , 2016.	707
655	Rosenstein, and Carolyn J Mattingly. 2012. Medic: a	Nut Limsopatham and Nigel Collier. 2016a. Normalis-	708
656	practical disease vocabulary used at the comparative	ing medical concepts in social media texts by learn-	709
657	toxicogenomics database. <i>Database</i> , 2012:bar065.	ing semantic representation. In <i>Proceedings of the</i>	710
658	N De Cao, G Izacard, S Riedel, and F Petroni. 2020.	<i>54th annual meeting of the association for compu-</i>	711
659	Autoregressive entity retrieval. In <i>ICLR 2021-9th In-</i>	<i>tational linguistics (volume 1: long papers)</i> , pages	712
660	<i>ternational Conference on Learning Representations</i> ,	1014–1023.	713
661	volume 2021. ICLR.	Nut Limsopatham and Nigel Collier. 2016b. Normalis-	714
662	Rezarta Islamaj Doğan, Robert Leaman, and Zhiyong	ing medical concepts in social media texts by learn-	715
663	Lu. 2014. Ncbi disease corpus: a resource for dis-	ing semantic representation. In <i>Proceedings of the</i>	716
664	ease name recognition and concept normalization.	<i>54th annual meeting of the association for compu-</i>	717
665	<i>Journal of biomedical informatics</i> , 47:1–10.	<i>tational linguistics (volume 1: long papers)</i> , pages	718
666	Ada Hamosh, Alan F. Scott, Joanna S. Amberger,	1014–1023.	719
667	Carol A. Bocchini, and Victor A. McKusick. 2004.	Zhenxi Lin, Ziheng Zhang, Meng Wang, Yinghui Shi,	720
668	Online mendelian inheritance in man (omim), a	Xian Wu, and Yefeng Zheng. 2022. Multi-modal con-	721
669	knowledgebase of human genes and genetic disor-	trastive representation learning for entity alignment.	722
670	ders. <i>Nucleic Acids Research</i> , 33:D514 – D517.	In <i>Proceedings of the 29th International Conference</i>	723
671	Chris Hokamp and Qun Liu. 2017. Lexically con-	<i>on Computational Linguistics</i> , pages 2572–2584.	724
672	strained decoding for sequence generation using grid	Zhenxi Lin, Ziheng Zhang, Xian Wu, and Yefeng Zheng.	725
673	beam search. In <i>Proceedings of the 55th Annual</i>	2024. Biomedical entity linking as multiple choice	726
674	<i>Meeting of the Association for Computational Lin-</i>	question answering. In <i>International Conference on</i>	727
675	<i>guistics (Volume 1: Long Papers)</i> , pages 1535–1546.	<i>Language Resources and Evaluation</i> .	728
676	Tuan Lai, Heng Ji, and ChengXiang Zhai. 2021. Bert	Carolyn E Lipscomb. 2000. Medical subject headings	729
677	might be overkill: A tiny but effective biomedical	(mesh). <i>Bulletin of the Medical Library Association</i> ,	730
678	entity linker based on residual convolutional neural	88(3):265.	731
679	networks. In <i>Findings of the Association for Com-</i>	Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco	732
680	<i>putational Linguistics: EMNLP 2021</i> , pages 1631–	Basaldella, and Nigel Collier. 2021. Self-alignment	733
681	1639.	pretraining for biomedical entity representations. In	734
682	Sunwon Lee, Donghyeon Kim, Kyubum Lee, Jaehoon	<i>Proceedings of the 2021 Conference of the North</i>	735
683	Choi, Seongsoon Kim, Minji Jeon, Sangrak Lim,	<i>American Chapter of the Association for Computa-</i>	736
684	Donghee Choi, Sunkyu Kim, Aik-Choon Tan, et al.	<i>tional Linguistics: Human Language Technologies</i> ,	737
685	2016. Best: next-generation biomedical entity search	pages 4228–4238.	738
686	tool for knowledge discovery from biomedical litera-	Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng	739
687	ture. <i>PloS one</i> , 11(10):e0164680.	Zhang, Hoifung Poon, and Tie-Yan Liu. 2022.	740
688	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	Biogpt: generative pre-trained transformer for	741
689	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	biomedical text generation and mining. <i>Briefings</i>	742
690	Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart:	<i>in bioinformatics</i> , 23(6):bbac409.	743
691	Denoising sequence-to-sequence pre-training for nat-	Zulfat Miftahutdinov, Elena Tutubalina, Samsung-	744
692	ural language generation, translation, and comprehen-	PDMI Joint AI Center, and PDMI RAS. 2019. Deep	745
693	sion. In <i>Proceedings of the 58th Annual Meeting of</i>	neural models for medical concept normalization in	746
694	<i>the Association for Computational Linguistics</i> , pages	user-generated texts. <i>ACL 2019</i> , page 393.	747
695	7871–7880.	Sunil Mohan and Donghui Li. 2019. Medmentions: A	748
		large biomedical corpus annotated with umls con-	749
		cepts. <i>arXiv preprint arXiv:1902.09476</i> .	750

A Details of Pre-training

Our pre-training process follows the KB-Guided Pre-training strategy outlined in the GenBioEL framework (Yuan et al., 2022b). We define clause templates to generate synthetic training examples by incorporating synonyms and definitions from the KB. Specifically, we select pairs of synonyms s_a and s_b from the set $\mathcal{S} \subset \mathcal{E}$, along with a definition d_y corresponding to the identifier y . The synonyms and definitions are integrated into one of the two predefined clause templates as follows:

[BOS] [ST] s_a [ET] **is defined as** d_y [EOS]
 or
 [BOS] d_y **describes** [ST] s_a [ET] [EOS].

The input for the decoder is “[BOS] s_a **is**” and the output should be “ s_b [EOS]”. When no definitions are available in the KB, we construct d_y using alternative synonyms. For concepts with only two synonyms, s_a and s_b are used as the synonyms, with d_y being the same as s_b . For concepts with only one synonym, s_a , s_b , and d_y are the same, resulting in a straightforward sentence. This strategy effectively creates simulated contexts for the model to learn from, improving its ability to generalize to new entities not included in the downstream datasets. The complete list of templates used in this process is provided in Table A.

B Datasets

Table B presents the statistics of the five datasets used, along with their corresponding target knowledge bases.

NCBI-disease (Doğan et al., 2014) The NCBI-disease dataset contains 793 PubMed abstracts annotated with 6,892 disease mentions that are mapped to 790 unique disease concepts using the MEDIC ontology (Davis et al., 2012). MEDIC is a medical dictionary that integrates disease concepts, synonyms, and definitions from both MeSH (Lipscomb, 2000) and OMIM (Hamosh et al., 2004), encompassing a total of 9,700 unique disease entities. This dataset is primarily used for disease recognition and concept normalization tasks.

BC5CDR (Li et al., 2016b) The BC5CDR dataset includes 1,500 PubMed abstracts with 4,409 chemical entities, 5,818 disease entities, and 3,116 chemical-disease interactions. All annotated

Template	
Encoder Side	Decoder Side
$s_a < \text{is defined as} > d_y$	
$s_a < \text{is described as} > d_y$	
$d_y < \text{are the definitions of} > s_a$	
$d_y < \text{describe} > s_a$	
$d_y < \text{define} > s_a$	
$d_y < \text{are the synonyms of} > s_a$	
$d_y < \text{indicate the same concept as} > s_a$	$s_a \text{ is } s_b$
$s_a < \text{has synonyms such as} > d_y$	
$s_a < \text{refers to the same concepts as} > d_y$	
$d_y < \text{is} > s_a$	
$d_y < \text{is the same as} > s_a$	
$s_a < \text{is} > d_y$	
$s_a < \text{is the same as} > d_y$	

Table A: The templates used for constructing pre-training samples. s_a is the input synonym, and s_b is the decoding target. d_y includes contextual information such as definitions and synonyms. Template words are enclosed in $< >$.

entities are mapped to the MeSH ontology (Lipscomb, 2000), which is a subset of UMLS (Bodenreider, 2004). This dataset is widely used for biomedical entity recognition and interaction studies. To fit the purpose of our study, we use only the chemical and disease annotations and discard the interaction annotations.

COMETA (Basaldella et al., 2020) COMETA focuses on layman medical terminology, compiled from four years of content across 68 health-related subreddits. This dataset consists of 20,000 biomedical entity mentions annotated with concepts from SNOMED CT (Chang and Mostafa, 2021). It is utilized for the normalization of consumer health expressions into standardized medical terminologies.

AskAPatient (AAP) (Limsopatham and Collier, 2016b) The AskAPatient dataset contains 8,662 phrases from social media language, each mapped to medical concepts from SNOMED CT (Chang and Mostafa, 2021). This dataset does not include contextual information, meaning that mentions are disambiguated solely based on the phrases themselves. Since the AskAPatient dataset lacks a test set, we employed a 10-fold cross-validation approach as outlined in the original paper by Limsopatham and Collier (2016a). The statistics reported are the averages across these folds.

Dataset	NCBI	BC5CDR	COMETA	AAP	MM-ST21pv
Entity types	disease	disease/chemical	medical concepts	medical concepts	21 UMLS types
<i>Data Examples</i>					
Training	5,784	9,285	13,489	15,665	121,498
Validation	787	9,515	2,176	793	40,600
Test	960	9,654	4,350	866	39,922
<i>KB statistics</i>					
Entity names	108,071	809,929	904,798	3,381	203,282
Identifiers	14,967	268,162	350,830	1,036	25,419

Table B: The statistics of the benchmark datasets and their corresponding KBs.

MM-ST21pv (Mohan and Li, 2019) The MedMentions dataset is a large-scale resource for biomedical entity recognition. The ST21pv subset includes 4,392 PubMed abstracts with over 200,000 entity mentions linked to 21 selected UMLS semantic types. This dataset provides a comprehensive resource for training and evaluating biomedical entity recognition systems. Unlike the original dataset, we use the 2020AA version of UMLS as the KBs because the 2017AA version of UMLS is not directly accessible. This leads to some differences after preprocessing due to variations between versions. Specifically, our dataset deviates from the original MedMentions dataset by 741 training samples (0.6%), 284 validation samples (0.7%), and 235 test samples (0.6%).

C Hyperparameter Configurations

Table C details the hyperparameters used for positive-only training and negative-aware training across the BioEL benchmark datasets. We search for the hyperparameter settings that are optimized for each dataset. We refer to the study of Yuan et al. (2022b) to determine the range of the hyperparameters. During pre-training, we use the same hyperparameters as in GenBioEL. For positive-only training, we explore a range of training steps between 20K and 40K, a learning rate between $2e-5$ and $3e-7$, and batch sizes from 8 to 16, except during pre-training. During negative-aware training, we fix the β at 0.1, in accordance with the basic configuration of DPO, and search the hyperparameter space using a learning rate between $1e-5$ and $1e-6$ and batch sizes ranging from 8 to 64.

D The Number of Synonyms

To evaluate the impact of incorporating multiple synonyms during fine-tuning, we conducted experiments by varying the number of synonyms associated with each mention, testing with 1, 3, and

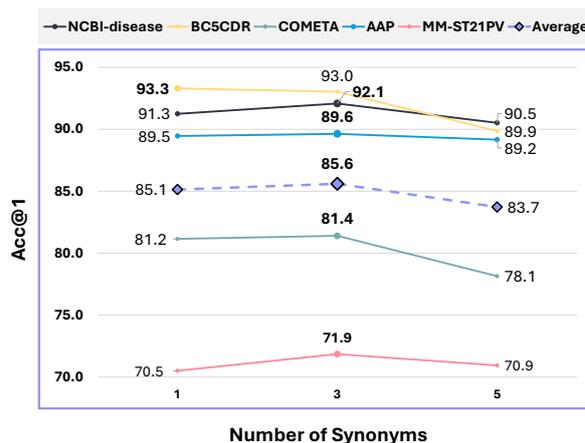


Figure A: The ablation study to determine the optimal number of synonyms. GenBioEL with ANGEL_{PT} was fine-tuned in this experiment. The scores are generally the highest when $k = 3$.

5 synonyms. As shown in Figure A, the average performance improves when using three synonyms compared to just one, but it declines when expanding to five synonyms. When the number of available synonyms is less than the specified number (e.g., fewer than 3 or 5 synonyms), all available synonyms are used to ensure maximum diversity in learning. Therefore, unlike GenBioEL, which utilized only the top-1 synonym, our approach incorporated up to the top-3 synonyms per mention, which proved to be optimal.

E Error Cases on COMETA

Similar to the analysis conducted on the NCBI-disease dataset (Figure 3), Figure B shows that the models in COMETA predominantly made the most errors in the 0.0-0.2 bin, where the similarity between input mentions and gold-standard entities is low. Our ANGEL framework improved GenBioEL’s performance across all bins, resulting in overall enhancement. Future work will necessitate the development of more advanced methods to

	Pre-training	Fine-tuning				
		NCBI	BC5CDR	COMETA	AAP	MM-ST21pv
<i>Positive-only Training</i>						
Training Steps	80K	20K	30K	40K	30K	40K
Learning Rate	4e-5	3e-7	5e-6	5e-6	5e-6	2e-5
Weight Decay	0.01	0.01	0.01	0.01	0.01	0.01
Batch Size	384	16	16	16	16	16
Adam ϵ	1e-8	1e-8	1e-8	1e-8	1e-8	1e-8
Adam β	(0.9,0.999)	(0.9,0.999)	(0.9,0.999)	(0.9,0.999)	(0.9,0.999)	(0.9,0.999)
Warmup Steps	1,600	0	500	500	0	1,000
Attention Dropout	0.1	0.1	0.1	0.1	0.1	0.1
Clipping Grad	0.1	0.1	0.1	0.1	0.1	0.1
Label Smoothing	0.1	0.1	0.1	0.1	0.1	0.1
<i>Negative-aware Training</i>						
Epochs	5	1	1	1	1	1
Learning Rate	1e-5	1e-5	1e-5	5e-6	5e-6	5e-6
β (DPO)	0.1	0.1	0.1	0.1	0.1	0.1
Weight Decay	0.01	0.01	0.01	0.01	0.01	0.01
Batch Size	1,024	64	16	64	8	64
Warmup Steps	100	100	100	100	100	100

Table C: The hyperparameters for positive-only training and negative-aware training.

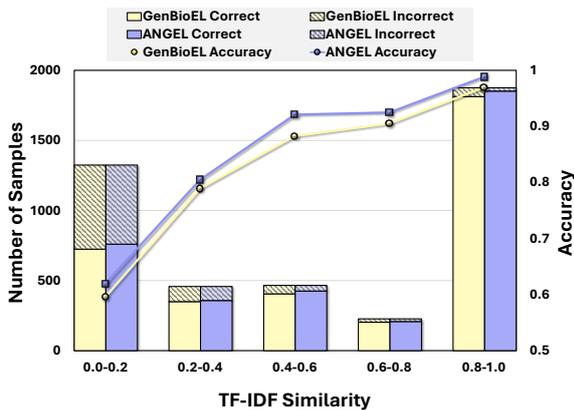


Figure B: In-depth evaluation of GenBioEL and our ANGEL models based on the TF-IDF similarity between the input mentions and gold-standard entities. The COMETA dataset was used.

specifically address errors in low similarity ranges.

F Top-5 Accuracy

Table D presents our model’s top-1 and top-5 accuracy on the BC5CDR and AAP datasets. It compares the performance of our model in its baseline form (GenBioEL) and after fine-tuning (ANGEL_{FT}) and combined pre-training and fine-tuning (ANGEL_{PT+FT}). Our approach consistently boosts top-1 accuracy across all datasets, though the trends in top-5 accuracy are less uniform. In BC5CDR, both top-1 and top-5 accuracy show significant improvements: top-1 accuracy rises by 1.4 percentage points (from 93.1% to

Model	BC5CDR		AAP	
	Acc@1	Acc@5	Acc@1	Acc@5
GenBioEL	93.1	95.7	89.3	95.4
+ ANGEL _{FT}	94.4	96.5	89.5	94.7
+ ANGEL _{PT+FT}	94.5	96.8	90.2	95.2

Table D: Comparison of top-1 and top-5 accuracy between the baseline model and models trained with ANGEL method after fine-tuning and pre-training on the BC5CDR and AAP datasets.

94.5%), and top-5 accuracy increases by 1.1 percentage points (from 95.7% to 96.8%). However, the AAP dataset exhibits a different pattern. While top-1 accuracy improves by 0.9 percentage points (from 89.3% to 90.2%), top-5 accuracy slightly declines: there is a 0.7 percentage points drop (from 95.4% to 94.7%) after fine-tuning and a 0.2 percentage points decrease (from 95.4% to 95.2%) after combined pre-training and fine-tuning. This decline in top-5 accuracy may be due to the AAP dataset’s limited contextual information, forcing the model to rely predominantly on the mention form, making it more challenging to maintain high accuracy across multiple predictions. Additionally, the negative sampling strategy could unintentionally bias the model toward optimizing top-1 accuracy, thereby impacting top-5 performance.

In conclusion, while our method consistently improves top-1 accuracy, the occasional slight decreases in top-5 accuracy, as observed in the AAP dataset, underscore the need for further refinement

1007 to maintain balanced accuracy across different rank-
1008 ing levels. Future work should focus on training
1009 strategies that preserve or enhance top-5 accuracy
1010 alongside top-1 improvements.