

# Towards Simulating User Relevance Feedback

Jasmin Mousavi\*  
Oregon State University  
mousavij@oregonstate.edu

Nischal Aryal\*  
Oregon State University  
safarim@oregonstate.edu

Mahdis Safari\*  
Oregon State University  
safarim@oregonstate.edu

## ABSTRACT

This paper proposes a deep learning-based approach that simulates the user-feedback loop in entity matching and relevance modeling without the need for costly and time-consuming data collection. Our method formulates the relevancy modeling problem as a text/entity matching task and leverages NLP techniques to learn matching patterns from logged interactions. We evaluate our approach on datasets from three popular domains and show its effectiveness, achieving close human performance. Furthermore, our approach demonstrates robust performance on complex negative samples, indicating its potential to benefit any domain that relies on user feedback.

### PVLDB Reference Format:

Jasmin Mousavi\*, Nischal Aryal\*, and Mahdis Safari\*. Towards Simulating User Relevance Feedback. PVLDB, 16(1): XXX-XXX, 2023.  
doi:XX.XX/XXX.XX

## 1 INTRODUCTION

Human-in-the-loop learning is often employed when there is insufficient training data available to effectively train a model or when a model needs to demonstrate high performance with limited interactions, such as in online systems. The primary goal of human-in-the-loop learning is to enhance the model’s performance by integrating human knowledge. This is often achieved through a reinforcement learning approach that rewards the model based on user feedback. Applications of human-in-the-loop learning can be found in various tasks, such as data labeling, data analysis, data integration, entity extraction/resolution, and outlier/anomaly detection [29].

Nevertheless, depending on user feedback for model improvement can also present challenges for human-in-the-loop systems. Acquiring a sufficient amount of user data or maintaining long-term user engagement can impose substantial financial and cognitive burdens. To address this challenge, we propose learning a relevance model for simulating user feedback.

Relevance modeling involves identifying relevant documents or items based on the user’s query or intent [22]. A critical aspect of relevance modeling is entity matching, which aims to determine whether two different representations refer to the same real-world entity. The term “entity matching” also loosely refers to the broader problem of determining whether two heterogeneous representations of different entities should be associated together [16]. An

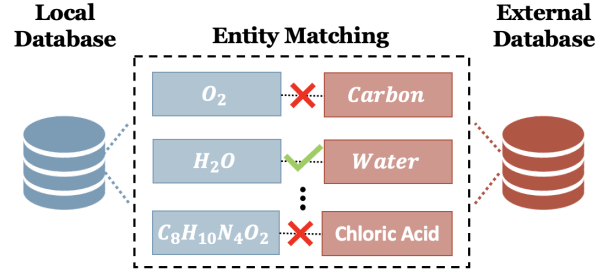


Figure 1: Entity matching for tuples from different databases: determine whether the pairs refer to the same entity.

example of entity matching is illustrated in Figure 2, where *Water* and *H<sub>2</sub>O* refer to the same entity, but have different representations.

In recent years, deep learning has made significant progress in semantic matching, largely attributed to its capacity to learn representations and generalize matching patterns from raw data, such as queries, documents, users, and items [11]. The successes achieved by deep learning in semantic matching can be extended to user relevance modeling by formulating this problem as a relevance matching task.

To effectively model user feedback in the context of relevance matching, it is insufficient to merely learn a similarity function between entities. Measuring similarity between two entities is more complex than identifying common characteristics since any two entities can be represented in an infinite number of ways. To accurately determine whether two representations refer to the same entity, the focus must be on identifying characteristics that are not only common but also *relevant* to the specific task at hand.

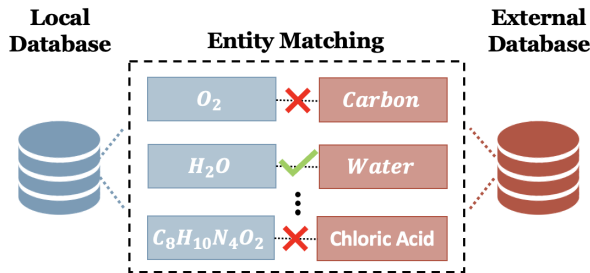
Determining relevance can be challenging since what is considered relevant is highly dependent on the context, priorities, and objectives of an individual user. Simply deriving the similarity between two entities is inadequate for modeling user feedback as it is not a unique relation. This is because relevant characteristics often differ based on the context. Therefore, when modeling user relevance feedback, it is crucial to have a comprehensive understanding of the context, the user’s goals and requirements, and the ability to identify the essential characteristics that are most relevant to the given task.

In this paper we aim to address the challenges in human-in-the-loop systems by introducing an approach that replaces user feedback. Our work relates to modeling user feedback in an entity-matching scenario [12, 19, 23, 27], where the user determines whether a tuple from the local database is a relevant match with a tuple from the external database. Given a local and external database, our approach eliminates the need for user feedback by using a learned relevance model to perform the relevance matching task.

\*These authors contributed equally to this work

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 16, No. 1 ISSN 2150-8097.  
doi:XX.XX/XXX.XX



**Figure 2: Entity matching for tuples from different databases: determine whether the pairs refer to the same entity.**

To summarize, this paper makes the following contributions:

- (1) This paper introduces a Siamese style network that leverages GPT-2 embeddings [21] for modeling user relevance feedback (Section 4)
- (2) We perform empirical experiments on three diverse datasets using logs from an online integration application [1] to assess the effectiveness of our proposed method (Section 6)
- (3) In our results, we demonstrate that our approach is robust towards complex negative samples (Section 6.4 & 6.5)

#### OLD INTRO:

Online Information Retrieval (IR) and Recommender Systems face the challenge of relevance matching, which involves ranking documents or items based on their relevance to a user’s query [22]. A critical aspect of relevance matching is entity matching.

Entity matching refers to the task of determining whether two different representations refer to the same real-world entity. The term “entity matching” also loosely refers to the broader problem of determining whether two heterogeneous representations of different entities should be associated together. [16]. An example of entity matching is illustrated in Figure 2, where *Water* and  $H_2O$  refer to the same entity, but have different representations.

Grounding this problem relies on user or expert feedback. However, gathering enough user data or ensuring continued long-term engagement from users can pose significant financial and cognitive challenges. In such scenarios, a learned model for user feedback can be beneficial for recommendation or retrieval problems.

Deep learning has made significant progress in recent years in matching through deep semantic models for search and neural collaborative filtering models for recommendation, primarily due to their ability to learn representations and generalize matching patterns from raw data like queries, documents, users, and items.[11] The successes of Deep learning in semantic matching can be extended to entity matching and relevance modeling, which is of great importance to various domains that depend on user feedback. The bottleneck of many online systems depends on user feedback to continue improving their models. The application of our approach goes beyond the IR domain and broadly impacts any domain that leverages user feedback, including Reinforcement Learning and numerous applications in Medical Sciences.

Modeling user feedback in the context of entity matching requires more than simply learning a similarity function between

entities. Because any two entities can be represented in an infinite number of ways, measuring similarity is more complex than finding common features between the two. Rather, we must determine which characteristics are relevant to the task at hand in order to more accurately measure whether two objects refer to the same entity. However, determining relevance can be difficult since what is relevant is highly dependent on the context, interests, and goals of a specific user. Simply deriving similarity between two entities is not sufficient for modeling user feedback, as it is not a unique relation. This is due to the fact that relevant characteristics often vary with respect to the user. Hence, modeling user feedback in the context of entity matching requires a thorough understanding of the user’s needs and goals, as well as the ability to determine which characteristics are most relevant to the task at hand.

This paper aims to address the challenges of entity matching and relevance modeling by introducing an approach that replaces user feedback. Our work relates to modeling user feedback in an entity-matching scenario[27][23][12], where the user determines whether a tuple from the local database is a relevant match with a tuple from the external database (Figure 3). Given a local and external database, our approach eliminates the need for user feedback by using a learned model to perform the matching task.

To summarize, this paper makes the following main contributions:

- (1) This paper introduces an approach to model user relevance feedback in an entity matching scenario, utilizing Natural Language Processing (NLP) techniques in order to simulate such feedback.
- (2) We performed empirical experiments on three diverse datasets using logs from an online entity matching system to assess the effectiveness of our proposed method.
- (3) In our results, we demonstrate that our approach is robust towards complex negative samples.

## 2 RELATED WORKS

The general matching problem has *traditionally* relied on techniques such as Support Vector Machines (SVMs) for identifying matches [2, 15]. A significant drawback of these conventional matching methods lies in their dependence on domain knowledge and feature engineering [16]. However, a promising solution to these challenges has emerged in the form of deep learning-based approaches. These approaches capitalize on the formidable representational capabilities of neural networks, allowing them to automatically extract features from raw data and learn complex matching patterns. This enables the establishment of robust connections between expressive representations.

Deep learning models have been widely adopted and formalized for the matching problem in diverse domains, including but not limited to NLP, computer vision, IR, and Recommender Systems. These models have demonstrated their effectiveness in addressing matching tasks related to text, images, documents, queries, users, and items [11]. By leveraging the representational power of deep learning, these models have shown great potential in achieving accurate and efficient matches across a wide range of applications.

**Information Retrieval.** In IR, the primary goal is to retrieve the top- $k$  matches for a given query. This can be formalized as a matching problem between a query and document.

**Recommender System.** The primary goal in Recommender Systems is to identify items that are most relevant to a user. This can be formalized as a matching problem between a user and items.

**Text Matching.** In text matching, primary goal is to find match for some given textual information. This can be formalized as a matching problem between textual data and another form of data (e.g. text, image, knowledge graph).

**Entity Matching** The entity matching problem is closely related to the general text matching problem. The goal of entity matching is to determine whether two different tuples refer to the same entity. This can be formalized as a matching problem between tuples.

#### OLD RELATED WORK:

The general matching problem has traditionally relied on techniques such as Support Vector Machines (SVMs) for identifying matches [2, 15]. However, recent advances in deep learning have yielded promising results in diverse domains, including natural language processing and computer vision.

The requirement of domain knowledge and feature engineering are the two biggest limitations of traditional matching methods [16]. Multiple deep learning-based approaches employing neural networks to automatically extract features from raw data and learn intricate matching patterns have been proposed to address these challenges.

This matching problem has also received significant attention across various domains, with commonalities observed in problems such as finding matches between **asymmetrical sequences, learning embedding spaces, relevance/matching scores, and architecture choices**. While our primary focus is on entity matching, we can still draw inspiration from related fields such as information retrieval, recommender systems and text matching.

**Information Retrieval.** In IR, the primary goal of matching is to retrieve the top- $k$  matches for a given query. This matching problem compares a short query against a long document, where the challenge lies in identifying the most relevant and useful information for the query from within the document.

A deep relevance matching model for ad-hoc retrieval is proposed in [11]. The model uses a deep neural network to learn a non-linear transformation of the query-document pairs into a joint embedding space, where the relevance score between the query and document is computed. The architecture of the model, which includes multiple layers of feedforward neural networks and a similarity measure based on cosine similarity.

**Recommender System.** In this domain, the entity matching problem involves comparing unstructured and asymmetric sequences, such as user data with item data, where the representations of user and item profiles may differ. Therefore, the architecture for this type of problem needs to take into account the possibility of varying representations between the two.

DeepCF [5] is a unified framework for learning representation and matching function in recommender systems. DeepCF adopts a deep neural network to learn the representations of users and items, and leverages a matching function to estimate the personalized preference score. The framework includes three modules: input module, representation learning module, and matching function

learning module. The input module preprocesses the raw data and feeds it into the representation learning module to learn the latent representations of users and items. The matching function learning module combines the learned representations and learns to predict the preference score.

**Text Matching.** In text matching, the problem of asymmetric sequence matching is similar to that in IR. However, due to the shorter length of both sequence in text matching, identifying a match becomes even more challenging. Unlike, in IR, the query is less likely to be a subset of the document, further adding to the complexity of the matching problem.

For text matching in asymmetrical domains, WD-Match [30] incorporate a Wasserstein distance regularization into sequence representation learning. Their approach involves two branches competing against each other: one estimates the Wasserstein distance using the projected features, while the other minimizes the Wasserstein distance regularized matching loss. The researchers demonstrate that regularizers helps WD-Match to generate feature vectors that are evenly distributed in the semantic space, making them more appropriate for matching.

Within the biomedical domain, SiameseCHEM [7] is a Siamese Recurrent Neural Network (SRNN) based on the BiLSTM with a self-attention mechanism for bioactivity prediction. The SRNN model uses two recurrent neural networks with shared weights to learn representations of two chemical compounds, and a self-attention mechanism to focus on the most informative parts of the representations. The model architecture also incorporates an attention-based pooling mechanism to further improve the performance of the model.

**Entity Matching** The entity matching problem is closely related to the general text matching problem, which has been extensively studied in the literature. Siamese architectures have emerged as a popular choice for text-matching tasks. In light of this, we draw motivation from a subset of recent state-of-the-art proposals for entity matching with these architectures.

The use of a Siamese hierarchical attention network for entity matching is proposed in [17]. This approach encodes the entity and description using a hierarchical attention mechanism to capture features at various levels of granularity. The Siamese architecture then generates a similarity score by comparing the encoded representations of the entity and description. However, the approach has limitations, including its reliance on high-quality entity descriptions and its inability to handle entities with multiple or ambiguous descriptions.

A Siamese-based BERT network (SiBERT) is proposed in [18] for the alignment of Chinese medical entities. The SiBERT network utilizes a Siamese neural network architecture to learn the similarity between two entities, where a shared BERT network encodes the two input entities. The output of the Siamese network is a binary classification indicating whether the two input entities match or not.

For multi-modal knowledge graphs, [3] present a Siamese network architecture for entity alignment between textual and visual information. The proposed architecture uses two sub-networks for the textual and visual modalities, which share weights and are trained jointly using a Siamese network approach to learn a similarity metric between entities. The authors leverage an attention

mechanism that selectively attends to different parts of the textual and visual features to improve the alignment accuracy.

TODO: [19]

This paper takes a different approach from those previously discussed above in that it focuses on the problem of simulating user feedback for systems that employ human-in-the-loop systems. Our objective in this paper is to introduce an approach for modeling user feedback in an entity matching scenario.

### 3 PROBLEM DEFINITION

Consider a system designed for relevance matching between two databases, where one is a local database consisting of  $N$  tuples and the other is an external database consisting of  $M$  tuples, as depicted in Figure 3.

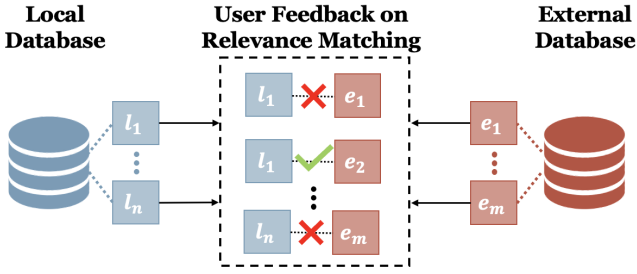


Figure 3: Incorporating user feedback for the relevance matching problem.

Given local tuple  $l_i$  and external tuple  $e_i$ , the user provides explicit relevance feedback to the system by identifying if  $l_i$  and  $e_i$  are relevant to each other (Equation 1). The system is dependent on user feedback for optimal performance.

$$relevance(l_i, e_i) = \begin{cases} 1, & \text{if } l_i \text{ and } e_i \text{ are relevant to each other} \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

Our objective is to mitigate this reliance on user data in online human-in-the-loop systems by learning a relevance model that can effectively simulate user feedback.

### 4 MODEL ARCHITECTURE

We utilize a Siamese network to model user feedback in the context of entity matching. The architecture of our model is illustrated in Figure 4.

#### 4.1 Encoder

Given a pair of tuples  $(l_i, e_i)$ , we encode local tuple  $l_i$  and external tuple  $e_i$  using an identical encoder with shared weights. The grey box in Figure 4 illustrates the encoder for our model.

**4.1.1 Embeddings.** Embeddings refer to vector representations of both external and local tuples. A straightforward approach to learning these embeddings is through static numbering of each word in these sequences, (*static word embeddings*). However, a major drawback of this approach is that all meanings of a word with multiple senses must share the same representation, which limits

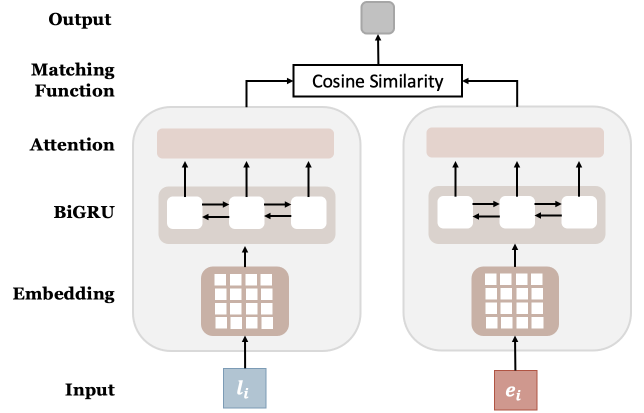


Figure 4: Model Architecture: Siamese network that encodes a pair of tuples  $(l_i, e_i)$  as input and measures the cosine similarity between them. The encoder, represented by the grey box, shares weights between both sides of the network.

its efficacy[6]. To overcome this limitation, we leverage *contextual word representations* generated by state-of-the-art language model, GPT-2 [21].

$$embed_{l_i} = \text{last-hidden-layer}(\text{GPT-2}(l_i)) \quad (2a)$$

$$embed_{e_i} = \text{last-hidden-layer}(\text{GPT-2}(e_i)) \quad (2b)$$

To obtain the embeddings, we feed a tuple (e.g.  $l_i$  or  $e_i$ ) as a sequence of words to GPT-2 and extract the last hidden layer as our embeddings. The final hidden layer of GPT-2 serves as a vectorized representation for each word in the input sequence.

**4.1.2 BiGRU.** Bidirectional Gated Recurrent Unit (BiGRU)[25] is a neural network architecture commonly used in NLP and sequence modeling tasks. BiGRU combines two Gated Recurrent Unit (GRU) layers, where one layer processes the input sequence in a forward direction, and the other processes the sequence in a backward direction. The architecture of BiGRU allows the model to capture contextual information from both the past and future context of each input token, thus enabling the model to better understand the meaning and structure of a sequence.

The BiGRU takes in a sequence embedding (Equation 2) as input and outputs a sequence representation  $b_t$  for tuple  $t \in \{l_i, e_i\}$ . Output for each side of the network is as follows:

$$b_{l_i} = \text{BiGRU}(embed_{l_i}) \quad (3a)$$

$$b_{e_i} = \text{BiGRU}(embed_{e_i}) \quad (3b)$$

For a single side of the network, we separately write the forward  $\overrightarrow{\text{BiGRU}}$  and backward  $\overleftarrow{\text{BiGRU}}$  directions as follows:

$$\overrightarrow{h_i} = \overrightarrow{\text{BiGRU}}(embed_i, \overrightarrow{h_{i-1}}) \quad (4)$$

$$\overleftarrow{h_i} = \overleftarrow{\text{BiGRU}}(embed_i, \overleftarrow{h_{i-1}}) \quad (5)$$

$$h_i = (\overrightarrow{h_i}, \overleftarrow{h_i}) \quad (6)$$

Hidden state  $h_i$  is constructed by concatenating the  $\overrightarrow{\text{BiGRU}}$  forward direction hidden state  $\overrightarrow{h_i}$  with the  $\overleftarrow{\text{BiGRU}}$  backward direction hidden state  $\overleftarrow{h_i}$ .

**4.1.3 Attention.** An attention function maps a query and a set of key-value pairs to an output. The output is computed as a weighted sum of the values, where each weight assigned to a value is computed by a dot product of the query with the corresponding key [28]. We employ attention to enhance the representation of a sequence by attending each token to every other token (*self-attention*).

Given  $b_t = (s_1, s_2, \dots, s_n)$  as a sequence output from our BiGRU layer, the Attention layer computes the weights  $\alpha_i$  for each token  $s_i$  based on its relationship with the other tokens in the same sequence. These attention weights not only reveal the relative importance of each token in the overall sequence semantics but also encode the similarities between tokens.

The attention mechanism takes in BiGRU hidden states matrix  $H$  and outputs matrix  $A$  containing the attention weights

$$H = (h_0, h_1, \dots, h_n) \quad (7)$$

$$A = \text{softmax}(W_2 \tanh(W_1 H^T)) \quad (8)$$

where  $W_1$  and  $W_2$  are trainable weight matrices.

The final representation for a tuple is constructed by first multiplying the attention weight matrix  $A$  with the BiGRU hidden states  $H$ , and then passing the result through a linear layer with a leaky rectified unit activation function:

$$t = \text{LeakyReLU}(W_3(AH)^T + b_3) \quad (9)$$

where  $W_3$  is a trainable weight matrix and  $b_3$  is the bias term.

## 4.2 Matching Function

Similarity metrics are mathematical functions that measure the degree of similarity or dissimilarity between two objects. Similarity metrics are widely used in data analysis and machine learning to compare data points.

One popular similarity metric used in entity matching is cosine similarity [26]. It measures the cosine of the angle between two vectors in a multi-dimensional space. The resulting value ranges from -1 to 1, where 1 indicates that the two vectors are identical and -1 indicates that they are completely dissimilar.

Given tuple representations for each branch of our network,  $t_l$  and  $t_e$  (Equation 9), we compute their pairwise cosine similarity as

$$T_l = \text{L2-normalize}(t_l) \quad (10)$$

$$T_e = \text{L2-normalize}(t_e) \quad (11)$$

$$\text{Sim}_C(T_l, T_e) = \frac{\langle T_l, T_e \rangle}{\|T_l\| \|T_e\|} \quad (12)$$

## 5 LOSS METRICS

Loss metrics are used in training neural networks to measure the error between the predicted outputs and the true target values. The choice of loss metrics is task dependent. Our relevance modeling task has a binary classification form as our models try to predict between *match(1)* and *non-match(0)*.

### 5.1 Cross Entropy Loss

The cross-entropy method was introduced as an adaptive approach for calculating probabilities of infrequent events and addressing combinatorial optimization problems[24]. The cross-entropy loss function measures the dissimilarity between predicted and true probability distributions and has been shown to effectively improve

model performance in binary classification tasks. [8]. The formula for cross-entropy loss can be expressed as follows:

$$\mathcal{H}(p, q) = - \sum_i p(i) \log q(i) \quad (13)$$

where  $p$  is the true probability distribution and  $q$  is the predicted class distribution.

### 5.2 Contrastive Loss

Contrastive loss is a popular loss function for evaluating the performance of siamese networks and distinguishing between two similar representations. It works by pulling together neighbors and pushing apart non-neighbors within a set of high-dimensional training vectors [13].

For example, given a batch of  $N$  local and external tuple matches  $(l_i, e_i)$ , contrastive loss leverages the  $N \times N$  possible tuple pairings by maximizing the cosine similarity of the  $N$  true matches while minimizing the cosine similarity of the  $N^2 - N$  non-matches. Given that our model outputs the cosine similarity  $\text{Sim}_C(T_l, T_e)$ , we compute contrastive loss as described in [20]:

$$\text{output} = \text{Sim}_C(T_l, T_e) \quad (14)$$

$$\text{labels} = [0 : \text{batchSize}] \quad (15)$$

$$\text{loss}_l = \mathcal{H}(\text{labels}, \text{output}) \quad (16)$$

$$\text{loss}_e = \mathcal{H}(\text{labels}, \text{output}^T) \quad (17)$$

$$\text{total\_loss} = \frac{\text{loss}_l + \text{loss}_e}{2} \quad (18)$$

where  $\mathcal{H}$  is cross entropy loss defined by Equation 13.

## 6 EXPERIMENTAL SETUP

### 6.1 Datasets

To evaluate the performance of our model, we leverage interaction logs from an online integration system [1] regarding datasets from three prominent areas: Products, Drugs, and News.

**Products.** The products dataset contains product information from e-commerce websites Amazon and Google [4]. Shared identifiers (ISBNs, SKUs, etc.) that could be used to easily match entities together were removed.

**Drugs.** The drug dataset used in our experiments consists of drug reviews collected from the website "Drugs.com" and corresponding descriptions of the same drugs scraped from "Wikipedia" [9]. While the language used in the reviews is informal and diverse, which may not always be directly related to the drug, Wikipedia articles are heavily edited and provide a more formal and standardized description of the drug.

**News.** The news dataset comprises articles from 38 prominent media organizations [10]. It includes both the article titles and summaries (local) as well as the articles themselves (external). The summaries were generated by different authors, using various techniques, resulting in varying degrees of overlap.

### 6.2 User Interaction

We analyze interactions between the online integration system and a user in the following fashion. During an interaction, the system presents one local tuple and 20 external tuples to the user. The

Dataset	Source	Attributes	Positive Samples	Negative Samples
<b>Products</b>	Local	name, description, manufacturer, price	632	12,797
	External	title, description, manufacturer, price		
<b>Drugs</b>	Local	drugName, condition, review	1791	34,439
	External	page_title, wikipedia_summary		
<b>News</b>	Local	title, article_summary	25,570	633,638
	External	article_content		

**Table 1: Details of datasets used in our evaluation with sampling distribution skewed towards negative (mismatch) samples**

user’s task is to identify which external tuple(s), if any, match the local tuple. For instance, if the user identifies one external tuple as a match for a given local tuple, it implies that the remaining 19 external tuples do not match the local tuple. We consider matches as positive samples and non-matches as negative samples.

In a given interaction, the presented external tuples are retrieved from an external database using the same query (i.e. external tuples share similar representations with respect to an interaction/query). Hence, negative samples are typically similar in representation to the positive sample(s) (i.e true match) for a given local tuple.

### 6.3 Models

**Baseline.** For our baseline, we construct a siamese style network similar to that in Figure 4. As described in Section 4, local and external tuples are processed by passing them through an embedding layer using GPT-2 *without fine-tuning*. The output of the embedding layer is passed to a fully connected layer for each branch of the network. The final layer takes the encoded representations from each branch and computes their cosine similarity. Due to the class imbalance in our datasets (i.e. more negative matches than positive matches), we use contrastive loss while training our model (5.2).

**BiGRU + Attention.** As shown in Figure 4, we first process the local and external tuples by embedding them with a pre-trained transformer-based model (e.g. GPT-2) *without fine-tuning*. These embeddings are then fed into a BiGRU. The output of the BiGRU is passed along to a ‘matching-layer’, which outputs a score of how relevant the external tuple is to the local tuple. To account for the imbalance in our data (i.e. more negative matches than positive matches), we use contrastive loss for training our model (5.2).

### 6.4 Experiment Details

We split the positive samples into *train* and *test* sets. The training set consists of 80% of the positive samples, while the test set consists of 20%.

To better simulate real-world user interactions and assess our model’s capacity to distinguish the correct match among similar negative samples, we develop a *hard-negative* test set. This test set is created by including positive samples from the *test* set and their corresponding negative samples extracted from the same interaction (Section 6.1). In other words, each batch will contain 1 positive sample and  $batchSize - 1$  corresponding negative samples. By including negative samples in a batch, we aim to better resemble a real-world user interaction, further testing our model’s ability to select the correct match among similar negative samples.

Due to time constraints, our models were trained for different durations depending on the size of the dataset. Specifically, we trained the models for 1000 epochs for the Products dataset, 800 epochs for the Drugs dataset, and 50 epochs for the News dataset. During training, we employed the Adam optimizer with weight decay [14]. Additionally, we set a batch size of 16 for both training and testing across all datasets.

### 6.5 Results & Analysis

**Train & Test Sets.** Our findings indicate that the *BiGRU + Attention* model achieved better prediction accuracy than the *Baseline* on the Products and Drugs datasets, however, we also observed that it had higher loss metrics on the Products dataset (Table 2). In contrast to this trend, we found that the *Baseline* outperformed the *BiGRU + Attention* model on the News dataset.

We find it intriguing that the *BiGRU + Attention* model fared worse than the *Baseline*, despite the News dataset featuring more coherent natural language text. It is possible that the embeddings from GPT-2 provided a better representation for this dataset. However, we believe that the *BiGRU + Attention* model will surpass the *Baseline* given adequate training time.

**Hard-Negative Test Set.** With the exception of the *Baseline* model on the News dataset, we observed that hard-negative test accuracy was higher than both the train and test accuracy, even though the dataset contained the same positive samples as the test set (Table 2). Despite the *Baseline* outperforming the *BiGRU + Attention* model on the test and train sets, we observe that *BiGRU + Attention* performs better on the hard-negative test set, which better resembles a real-world entity matching interaction.

We find the high accuracy trend observed on the hard-negative test set to be particularly noteworthy since it is the test set that most closely simulates real-world entity-matching interactions. We are uncertain as to why the hard-negative test set outperforms the standard test set, given that it includes the same positive samples. However, we hypothesize that our original assumption of the true match being more similar to the negative samples may require additional investigation.

## 7 CONCLUSION

Our proposed architecture effectively simulates user feedback in an entity matching scenario. We demonstrate that our model is robust against negative samples and can accurately identify the correct match in tests that closely resemble user interactions. Although our model did not outperform the baseline on the (largest) News



Dataset	Metric	Baseline			BiGRU + Attention		
		Train	Test	Hard-Negative	Train	Test	Hard-Negative
Products	Accuracy (%)	88.51	87.40	84.62	<b>89.90</b>	<b>90.55</b>	<b>1.00</b>
	Loss	0.36	0.33	-	0.38	0.45	-
Drugs	Accuracy (%)	82.40	79.39	86.01	<b>97.63</b>	<b>97.77</b>	<b>1.00</b>
	Loss	0.45	0.49	-	0.13	0.19	-
News	Accuracy (%)	<b>84.78</b>	<b>84.81</b>	77.94	76.81	76.53	<b>1.00</b>
	Loss	0.43	0.45	-	0.54	0.55	-

**Table 2: Experiment results. Train and Test set consist of only positive samples, while Hard-Negative Test set uses positive samples and their corresponding negative samples. Each Hard-Negative Test batch consists of 1 positive sample and (batch size - 1) negative samples. Prediction accuracy (Acc.) is determined by comparing the label with the highest softmax probability against the true prediction label.**

dataset, we believe that with additional training time, it will surpass the baseline’s performance.

## 7.1 Limitations & Future Work

A potential limitation of our approach pertains to the model settings and architecture choices. For instance, we did not incorporate fine-tuning for GPT-2 during training due to memory constraints. In future studies, we plan to investigate fine-tuning and compare our current recurrent-style architecture with a transformer-style architecture.

We also aim to improve our model’s performance by exploring the use of *hard-negatives* in the training process in order to better simulate real-world entity-matching scenarios. Moreover, a real-world scenario will include interactions where a batch of samples may only include negative samples. Therefore, we wish to explore the degree of certainty in our predictions by setting thresholds on the predicted probability distribution.

## REFERENCES

- [1] Christopher Buss, Jasmin Mosavi, Mikhail Tokarev, Arash Termehchy, Maier David, and Stefan Lee. 2023. *Effective Entity Augmentation By Querying External Data Sources*. Technical Report. <https://web.engr.oregonstate.edu/~termehca/papers/entityarg.pdf>
- [2] Surajit Chaudhuri, Bee-Chung Chen, Venkatesh Ganti, and Raghav Kaushik. 2007. Example-Driven Design of Efficient Record Matching Queries. In *Proceedings of the 33rd International Conference on Very Large Data Bases (Vienna, Austria) (VLDB '07)*. VLDB Endowment, 327–338.
- [3] Liyi Chen, Zhi Li, Tong Xu, Han Wu, Zhefeng Wang, Nicholas Jing Yuan, and Enhong Chen. 2022. Multi-Modal Siamese Network for Entity Alignment. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (Washington DC, USA) (KDD '22)*. Association for Computing Machinery, New York, NY, USA, 118–126. <https://doi.org/10.1145/3534678.3539244>
- [4] Sanjib Das, AnHai Doan, Paul Suganthan G. C., Chaitanya Gokhale, Pradap Konda, Yash Govind, and Derek Paulsen. [n.d.]. The Magellan Data Repository. <https://sites.google.com/site/anhaidgroup/projects/data>.
- [5] Zhi-Hong Deng, Ling Huang, Chang-Dong Wang, Jian-Huang Lai, and Philip S. Yu. 2019. DeepCF: A Unified Framework of Representation Learning and Matching Function Learning in Recommender System. *CoRR* abs/1901.04704 (2019). [arXiv:1901.04704](https://arxiv.org/abs/1901.04704) <http://arxiv.org/abs/1901.04704>
- [6] Kavin Ethayarajh. 2019. How Contextual are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings. *arXiv* (2019).
- [7] Gogishvili D Nittinger E Zhao H Tyrchan C Fernández-Llaneza D, Ulander S. 2021. Siamese Recurrent Neural Network with a Self-Attention Mechanism for Bioactivity Prediction. *ACS Omega* 3, 6 (2021), 11086–11094.
- [8] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT press. <https://www.deeplearningbook.org/>
- [9] Felix Gräßer, Surya Kallumadi, Hagen Malberg, and Sebastian Zaunseder. 2018. Aspect-based sentiment analysis of drug reviews applying cross-domain and cross-data learning. In *International Conference on Digital Health*. 121–125.
- [10] Max Grusky, Mor Naaman, and Yoav Artzi. 2018. Newsroom: A Dataset of 1.3 Million Summaries with Diverse Extractive Strategies. In *NAACL*.
- [11] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce Croft. 2016. A Deep Relevance Matching Model for Ad-Hoc Retrieval. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management (Indianapolis, Indiana, USA) (CIKM '16)*. Association for Computing Machinery, New York, NY, USA, 55–64. <https://doi.org/10.1145/2983323.2983769>
- [12] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce Croft. 2016. A Deep Relevance Matching Model for Ad-Hoc Retrieval. (2016), 55–64. <https://doi.org/10.1145/2983323.2983769>
- [13] R. Hadsell, S. Chopra, and Y. LeCun. 2006. Dimensionality Reduction by Learning an Invariant Mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, Vol. 2. 1735–1742. <https://doi.org/10.1109/CVPR.2006.100>
- [14] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [15] Hanna Köpcke, Andreas Thor, and Erhard Rahm. 2010. Learning-Based Approaches for Matching Web Data Entities. *IEEE Internet Computing* 14, 4 (2010), 23–31. <https://doi.org/10.1109/MIC.2010.58>
- [16] Yuliang Li, Jinfeng Li, Yoshihiko Suhara, Jin Wang, Wataru Hirota, and Wang-Chiew Tan. 2021. Deep Entity Matching: Challenges and Opportunities. *J. Data and Information Quality* 13, 1, Article 1 (jan 2021), 17 pages. <https://doi.org/10.1145/3431816>
- [17] Kai Ma, Liang Wu, Liufeng Tao, Wenjia Li, and Zhong Xie. 2018. Matching Descriptions to Spatial Entities Using a Siamese Hierarchical Attention Network. *IEEE Access* 6 (2018), 28064–28072. <https://doi.org/10.1109/ACCESS.2018.2837666>
- [18] Zerui Ma, Linna Zhao, Jianqiang Li, Xi Xu, and Jing Li. 2022. SiBERT: A Siamese-based BERT network for Chinese medical entities alignment. *Methods* 205 (2022), 133–139. <https://doi.org/10.1016/j.jymeth.2022.07.003>
- [19] Sidharth Mudgal, Han Li, Theodoros Rekatsinas, AnHai Doan, Youngchoon Park, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, and Vijay Raghavendra. 2018. Deep learning for entity matching: A design space exploration. In *Proceedings of the 2018 International Conference on Management of Data*. 19–34.
- [20] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [21] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language Models are Unsupervised Multitask Learners. (2019). <https://openai.com/blog/better-language-models/>
- [22] Jinfeng Rao, Lingling Liu, Yi Tay, Wei Yang, Peng Shi, and Jimmy Lin. 2019. Bridging the Gap between Relevance Matching and Semantic Matching for Short Text Similarity Modeling. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 5370–5381. <https://doi.org/10.18653/v1/D19-1540>
- [23] Jinfeng Rao, Lingling Liu, Yi Tay, Wei Yang, Peng Shi, and Jimmy Lin. 2019. Bridging the Gap between Relevance Matching and Semantic Matching for Short Text Similarity Modeling. (Nov. 2019), 5370–5381. <https://doi.org/10.18653/v1/D19-1540>

- [24] R. Rubinstein. 1997. Optimization of computer simulation models with rare events. *European Journal of Operational Research* 99, 1 (1997), 89–112.
- [25] M. Schuster and K.K. Paliwal. 1997. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing* 45, 11 (1997), 2673–2681. <https://doi.org/10.1109/78.650093>
- [26] Amit Singhal et al. 2001. Modern information retrieval: A brief overview. *IEEE Data Eng. Bull.* 24, 4 (2001), 35–43.
- [27] Saravanan Thirumuruganathan, Han Li, Nan Tang, Mourad Ouzzani, Yash Govind, Derek Paulsen, Glenn Fung, and AnHai Doan. 2021. Deep Learning for Blocking in Entity Matching: A Design Space Exploration. *Proc. VLDB Endow.* 14, 11 (oct 2021), 2459–2472. <https://doi.org/10.14778/3476249.3476294>
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. In *Advances in Neural Information Processing Systems*.
- [29] Xingjiao Wu, Luwei Xiao, Yixuan Sun, Junhang Zhang, Tianlong Ma, and Liang He. 2022. A survey of human-in-the-loop for machine learning. *Future Generation Computer Systems* (2022).
- [30] Weijie Yu, Chen Xu, Jun Xu, Liang Pang, Xiaopeng Gao, Xiaozhao Wang, and Ji-Rong Wen. 2020. Wasserstein distance regularized sequence representation for text matching in asymmetrical domains. *arXiv preprint arXiv:2010.07717* (2020).