

Context-Aware Word Sense Disambiguation in Narrative Texts: A Hybrid BERT-LSTM Approach

Muha Ioana

January 10, 2026

Abstract

Word Sense Disambiguation (WSD) is a fundamental challenge in Natural Language Processing, requiring models to determine the correct meaning of a polysemous word based on its context. This report presents a solution for SemEval 2026 Task 5: Rating Plausibility of Word Senses in Ambiguous Sentences through Narrative Understanding. The task involves predicting human judgments of the plausibility of a specific word sense within short, five-sentence stories (the "AmBiStory" dataset). We investigate the limitations of static word embeddings (GloVe) in narrative contexts, demonstrating their susceptibility to "mean collapse" and inability to resolve fine-grained ambiguities. To address this, we propose a hybrid Cross-Encoder architecture that utilizes a pre-trained Transformer (BERT) for contextual feature extraction, coupled with a Bidirectional LSTM for sequence modeling. Our ablation study reveals that replacing static embeddings with contextualized representations results in a statistically significant improvement in Spearman Correlation (from 0.0 to 0.15), highlighting the necessity of deep contextual awareness for resolving narrative ambiguity.

1 Introduction

Natural Language Processing (NLP) aims to bridge the gap between human communication and computational understanding. A central problem in this field is Word sense disambiguation (WSD) [6], defined as the ability to identify the correct meaning of an ambiguous word within a specific sentence or context. The difficulty of WSD stems from the inherent complexity of language, particularly phenomena such as homonymy (words with the same form but unrelated meanings) and polysemy (words with multiple related meanings).

Standard computational models often struggle with these distinctions because a word's meaning is heavily dependent on "contextual dependency". For example, the word "track" carries entirely different semantic values in the context of "detectives following a trail" versus "recording a song." Accurate disambiguation requires models to look beyond the target word itself and integrate information from the entire surrounding document.

This report addresses SemEval 2026 Task 5, which focuses on rating the plausibility of word senses in the "AmBiStory" dataset. Unlike traditional classification tasks that assign a single label, this task requires predicting a continuous plausibility score (1-5) that correlates with human judgment. The dataset consists of short stories where a specific homonym is used in an ambiguous sentence, surrounded by a pre-context and an ending that clarifies the intended meaning.

Our objective was to develop a neural network architecture capable of capturing the narrative flow required to disambiguate these homonyms. We hypothesize that traditional static embeddings (like GloVe) are insufficient for this task because they conflate multiple senses into a single vector representation. To test this, we implemented and compared several architectures, culminating in a Hybrid Cross-Encoder BERT-LSTM. This model leverages the state-of-the-art contextual understanding of Transformers while retaining the sequence modeling capabilities of Recurrent Neural Networks (RNNs).

This report details our methodology, beginning with a review of related WSD strategies, followed by a description of our proposed hybrid architecture. We present a comprehensive ablation study demonstrating the impact of attention mechanisms and contextual embeddings on performance, and conclude with an analysis of the model's behavior on the validation set.

2 Related Work

Research into Word Sense Disambiguation (WSD) has a long history, dating back to the foundational work of Warren Weaver in the 1940s, who first identified the statistical nature of language ambiguity. Over the decades, approaches to WSD have evolved through three distinct phases: Knowledge-Based approaches, Supervised Machine Learning, and current State-of-the-Art Deep Learning methods.

2.1 Knowledge-Based Approaches

Early WSD systems relied heavily on external lexical resources, such as dictionaries and thesauri. A seminal example is the Lesk Algorithm (1986) [4], which determines the correct sense of a word by measuring the overlap between the word’s dictionary definition and the words in its surrounding context. Subsequent methods utilized structured knowledge bases like WordNet [5] to map words to specific ”synsets” (sets of synonyms). While these methods have the advantage of being interpretable and not requiring large labeled datasets, they often suffer from limited coverage and an inability to handle new or evolving word senses.

2.2 Supervised Machine Learning

In the 1990s, the availability of annotated corpora led to the rise of supervised learning methods. Algorithms such as Naive Bayes, Support Vector Machines (SVM), and Decision Trees were trained on sentences where target words were manually tagged with their correct senses. SVMs, in particular, were found to yield high results compared to other traditional classifiers. However, these approaches are constrained by the ”knowledge acquisition bottleneck” - the high cost and difficulty of producing large-scale, manually annotated datasets required for training.

2.3 Deep Learning and Contextual Embedding

The most recent advancements in WSD have been driven by Deep Learning. Neural networks, specifically Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks [3], became prominent due to their ability to process sequences and capture long-range dependencies in text. Studies have shown that Bi-LSTMs can achieve high accuracy (up to 96% in some language-specific tasks) by learning complex patterns in the data).

A critical limitation of early neural approaches was the use of static word embeddings (e.g., GloVe [7] or Word2Vec). These models assign a fixed vector to every word, meaning the word ”bank” receives the same representation regardless of whether it appears in a financial or geographical context. This conflation of meanings makes it difficult for the model to resolve polysemy.

To address this, the field has shifted toward pre-trained Transformer Models like BERT [1] (Bidirectional Encoder Representations from Transformers). Unlike static embeddings, BERT generates ”contextualized” word representations, where the vector for a word is dynamically derived from its specific surrounding sentence. This allows the model to inherently distinguish between fine-grained senses before classification even begins. Recent surveys [2] indicate that hybrid approaches-combining the feature-extraction power of Transformers with the sequence modeling of RNNs-are a promising direction for resolving ambiguity in complex narrative tasks.

3 Methodology

To address the challenge of rating word sense plausibility in narrative contexts, we developed a hybrid neural architecture that combines the feature-extraction capabilities of Pre-trained Transformers with the sequence modeling strengths of Recurrent Neural Networks. Our final model, the Cross-Encoder BERT-LSTM, was designed to overcome the limitations of static embeddings (such as GloVe) by leveraging deep contextual representations.

3.1 Architecture Overview

The proposed model operates in a four-stage pipeline:

1. Cross-Encoder Input Processing: The definition and story are concatenated into a single sequence to enable early interaction between the target word and its definition.
2. Contextual Feature Extraction: A pre-trained BERT model (frozen) extracts dense vector representations for each token, resolving polysemy at the embedding level.
3. Dimensionality Reduction (Bottleneck): A projection layer compresses the high-dimensional BERT vectors to reduce model capacity and prevent overfitting.
4. Sequence Modeling & Regression: A Bidirectional LSTM processes the compressed sequence, followed by a regression head that predicts the plausibility score (1–5).

3.2 Input Representation: The Cross-Encoder Strategy

In our preliminary experiments using a Siamese architecture (processing the Definition and Story in parallel branches), we observed that the model struggled to identify which specific word in the story was the target homonym. This resulted in a high generalization error.

To correct this, we adopted a Cross-Encoder strategy. Instead of feeding the definition and story as separate inputs, we concatenate them into a single input string separated by the special separator token [SEP]:

$$\text{Input} = [\text{CLS}] \text{ Definition } [\text{SEP}] \text{ Story } [\text{SEP}]$$

This configuration forces the Self-Attention mechanism within the Transformer to explicitly compare every word in the definition against every word in the story before the embeddings are generated. This effectively "highlights" the relevant homonym in the story based on the provided definition, addressing the "Contextual Dependency" challenge inherent in WSD tasks.

3.3 Feature Extraction: Contextual Embeddings

We utilize BERT-base-uncased as the embedding backbone. Unlike static embeddings (e.g., GloVe), which assign a fixed vector to the word "track" regardless of context, BERT generates dynamic representations based on the surrounding sentence structure². Given the limited size of the training dataset (~2,000 samples) relative to the number of parameters in BERT (110 million), fine-tuning the entire Transformer was computationally prohibitive and prone to catastrophic overfitting. Therefore, we froze the weights of the BERT layers, using the model purely as a feature extractor. The output is a sequence of vectors of dimension $d_{bert} = 768$.

3.4 Dimensionality Reduction: The Projection Bottleneck

A significant technical challenge encountered during training was the disparity between the high dimensionality of BERT outputs (768) and the small dataset size. Direct feeding of 768-dimensional vectors into the LSTM led to severe overfitting, where Training Loss approached 0.4 while Validation Loss stagnated at 1.7.

To mitigate this, we introduced a Projection "Bottleneck" Layer before the LSTM. This layer consists of a dense linear transformation followed by a LeakyReLU activation and Dropout:

$$h_{proj} = \text{Dropout}(\text{LeakyReLU}(W_{proj} \cdot x_{bert} + b))$$

We project the embeddings from $d_{bert} = 768$ down to $d_{lstm_in} = 256$. This reduces the number of parameters in the subsequent LSTM layer by approximately 66%, forcing the model to filter out noise and retain only the most salient semantic features required for disambiguation.

3.5 Sequence Modeling: Bidirectional LSTM

The compressed feature sequence is passed to a Bidirectional LSTM (BiLSTM). While Transformers are excellent at capturing global context, LSTMs are effective at modeling the specific sequential flow of a narrative. We use a hidden dimension of 128 for the LSTM.

The BiLSTM processes the sequence in both forward and backward directions, producing two hidden states for each token. We represent the final "summary" of the story by concatenating the last hidden state of the forward pass (\vec{h}_T) and the last hidden state of the backward pass (\overleftarrow{h}_0):

$$v_{story} = \text{Concat}(\vec{h}_T, \overleftarrow{h}_0)$$

3.6 Regression Head and Optimization

The final narrative vector v_{story} is passed through a Multi-Layer Perceptron (MLP) regressor. To ensure the predictions remain within the valid range of the task, we apply a scaled Sigmoid activation function to the final output:

$$\hat{y} = 4.0 \cdot \sigma(z) + 1.0$$

This constrains the output strictly between 1.0 and 5.0. The model is trained using Mean Squared Error (MSE) loss. We employ the Adam optimizer with a learning rate of $5e - 4$ and L2 Weight Decay ($1e - 4$) to further regularize the weights and improve generalization on the unseen development set.

4 Experiments & Results

To evaluate the effectiveness of the proposed architecture, we conducted a systematic ablation study. We compared our final Cross-Encoder BERT-LSTM against several baseline configurations to isolate the contributions of specific components, such as the embedding type (Static vs. Contextual) and the input interaction strategy (Siamese vs. Cross-Encoder).

4.1 Experimental Setup

All models were trained on the train.json split and evaluated on the dev.json split of the AmbiStory dataset. We utilized two primary evaluation metrics:

- Spearman Correlation (r): Measures the monotonic relationship between the predicted scores and the human-annotated averages. This is the primary metric for assessing whether the model captures the ranking of plausibility.
- Accuracy (within Standard Deviation): Measures the percentage of predictions falling within the range of $[Average - \sigma, Average + \sigma]$. This metric accounts for the subjectivity and variance in human annotations.

4.2 Ablation Study

Table 1 summarizes the performance of the different architectures developed during this study.

4.3 Analysis of Result

1. Failure of Static Embeddings (The Baseline): The baseline model utilizing GloVe embeddings exhibited a phenomenon known as "Mean Collapse." With a Spearman correlation near zero, the model failed to learn any meaningful relationship between the story and the rating. Instead, it converged to predicting the global average rating (approx. 3.0) for all samples to minimize Mean Squared Error. This confirms that static embeddings, which conflate multiple meanings of a homonym into a single vector, lack the necessary granularity for this task.

Table 1: Ablation Study: Comparison of Model Architectures and Input Strategies

Model Architecture	Embedding	Input Strategy	Val Loss	Spearman	Accuracy
Baseline LSTM	GloVe (Static)	Siamese	1.60	-0.01	52.5%
Attention BiLSTM	GloVe (Static)	Siamese + Attn	1.70	0.06	53.5%
Hybrid LSTM	BERT (Frozen)	Siamese	1.76	0.10	57.3%
Cross-Encoder LSTM	BERT (Frozen)	Concatenated	1.72	0.15	59.0%

Note: "Siamese" indicates the Definition and Story were processed separately. "Concatenated" indicates the Cross-Encoder strategy ([CLS] Def [SEP] Story). Accuracy refers to predictions within one standard deviation of the human average.

2. Impact of Contextual Embeddings: Replacing GloVe with BERT (Hybrid LSTM) yielded the single largest performance jump, raising the Spearman correlation to 0.10. This validates our hypothesis that resolving polysemy requires dynamic, context-aware representations. However, the Siamese configuration-where the definition and story were encoded independently-limited the model’s ability to focus on the specific target word within the narrative.
3. Superiority of the Cross-Encoder: The final Cross-Encoder architecture achieved the highest performance (Spearman 0.15, Accuracy 59%). By concatenating the definition and story into a single sequence, the BERT self-attention mechanism could directly compare the definition against the homonym in the text. This "early interaction" allowed the subsequent LSTM to process a sequence where the relevant context was already highlighted, leading to statistically significant improvements ($p < 0.001$).
4. Dimensionality and Regularization: We observed that the high dimensionality of BERT features (768) initially caused rapid overfitting. Introducing the "Bottleneck" projection layer (reducing dimensions to 256) and increasing Dropout to 0.5 were decisive factors in stabilizing the training. This allowed the model to generalize better to unseen stories, as evidenced by the consistent improvement in validation metrics compared to the earlier high-capacity models.

5 Error Analysis & Discussion

While the Cross-Encoder BERT-LSTM achieved the best performance among the tested architectures, the accuracy of 5% and Spearman correlation of 0.15 indicate that the task remains challenging. A deeper analysis of the results reveals inherent difficulties in the AmbiStory dataset and specific limitations of our approach.

5.1 Interpretation of "Accuracy within Standard Deviation"

The primary metric, accuracy within one standard deviation 1, accounts for the subjectivity of the task. Unlike binary classification, human annotators often disagree on the exact plausibility score for a given story. A standard deviation of ≈ 1.2 suggests that a story rated as a "3" by one human might be rated as a "2" or "4" by another. Therefore, our model’s accuracy of 59% does not imply it is "wrong" 41% of the time; rather, it indicates that in nearly 60% of cases, the model’s prediction fell within the acceptable range of human consensus. This aligns with findings in the literature that WSD is often subjective and context-dependent.

5.2 Failure Modes

Qualitative analysis of the predictions reveals two primary sources of error:

1. Subtle Narrative Cues: The model struggles when the "ending" sentence relies on world knowledge or subtle inference rather than explicit keywords. For example, if the ending implies a specific meaning through sarcasm or a complex metaphor, the LSTM-even with BERT features-may fail to detect the shift in plausibility.

2. The Regression To Mean Effect: Although the Cross-Encoder mitigated the "mean collapse" observed in the baseline, the model remains conservative. It rarely predicts extreme scores (1.0 or 5.0), preferring the safer range of 2.5–3.5. This limits the Spearman correlation, as the model fails to commit strongly to highly plausible or highly implausible scenarios.

5.3 Computational Constraints

The reliance on a "frozen" BERT model (feature extraction) rather than full fine-tuning was a necessary trade-off due to computational constraints and the small dataset size. While this prevented catastrophic forgetting, it also meant the embeddings were not optimized specifically for the "AmbiStory" domain. A fully fine-tuned Transformer might better capture the specific stylistic nuances of these short stories.

6 Conclusion

This report presented a comparative study of neural architectures for the task of Rating Plausibility of Word Senses in Ambiguous Sentences. We investigated the transition from static word embeddings to contextualized deep learning representations.

Our experiments demonstrated that traditional static embeddings (GloVe) are fundamentally ill-suited for resolving polysemy in narrative contexts, resulting in models that fail to capture semantic rankings. By integrating Contextual Embeddings (BERT), we significantly improved performance. Furthermore, we showed that a Cross-Encoder strategy-concatenating the definition and story into a single input sequence-outperforms the Siamese architecture by enabling the model to attend to the specific target homonym effectively.

The final Cross-Encoder BERT-LSTM architecture achieved a Spearman correlation of 0.15 and an accuracy of 59%, establishing a statistically significant baseline over random chance. Future work could improve upon these results by employing larger pre-trained models (e.g., RoBERTa-Large), implementing full fine-tuning with careful regularization, or augmenting the dataset with additional synthetic examples to better train the regression head.

References

- [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [2] Robbel Habtamu and Beakal Gizachew. State-of-the-art approaches to word sense disambiguation: A multilingual investigation. *School of Information Technology and Engineering, Addis Ababa Institute of Technology*, 2024.
- [3] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [4] Michael Lesk. Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone. In *Proceedings of the 5th Annual International Conference on Systems Documentation*, pages 24–26, 1986.
- [5] George A Miller. Wordnet: A lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [6] Roberto Navigli. Word sense disambiguation: A survey. *ACM Computing Surveys (CSUR)*, 41(2):1–69, 2009.
- [7] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.