

Tree Matching Networks for Natural Language Inference: Parameter-Efficient Semantic Understanding via Dependency Parse Trees

Jason Lunder
Eastern Washington University
jllunder@ewu.edu

December 2, 2025

Abstract

Transformer-based models like BERT achieve high accuracy on Natural Language Inference (NLI) but require hundreds of millions of parameters and extensive pretraining. We investigate whether explicit dependency tree structures can improve parameter efficiency by providing syntactic inductive bias that transformers must learn from scratch.

We adapt Graph Matching Networks to operate on dependency parse trees, creating Tree Matching Networks (TMN), and compare against BERT-based baselines with matched parameters on SNLI and SemEval tasks. TMN substantially outperforms BERT at small scales with reduced memory and training requirements, validating that dependency structure provides measurable benefits.

However, we identify a scaling plateau: increasing TMN parameters $600\times$ yields minimal improvement, indicating that simple aggregation methods create an architectural bottleneck. These findings motivate hybrid approaches that preserve structural benefits while addressing scalability limitations, such as using graph neural networks to pre-encode dependency trees before transformer-based aggregation with tree-aware positional encodings.

Introduction

Natural language understanding requires capturing both semantic content and syntactic structure. Modern transformer-based models like BERT [5] achieve impressive results on natural language inference (NLI) tasks but require hundreds of millions of parameters and extensive pretraining. We

investigate whether explicit structural representations can improve parameter efficiency while maintaining competitive performance.

Dependency parse trees explicitly encode syntactic relationships as graph edges, potentially providing inductive biases that reduce learning complexity. Consider the sentence “A dog runs in the park.” A dependency parser identifies “dog” as the nominal subject (nsubj) of “runs,” “in” as a prepositional modifier, and “park” as the object of the preposition. These relationships are encoded as labeled edges in the dependency tree. In contrast, transformer models must learn these relationships implicitly from token co-occurrence patterns across massive corpora.

Research Question

Can graph neural networks operating on dependency parse trees provide more parameter-efficient sentence embeddings than transformer-based models for natural language inference?

Specifically, we investigate whether Tree Matching Networks (TMN), our adaptation of Graph Matching Networks [11] to dependency trees, can outperform BERT at equivalent parameter counts on natural language inference.

Our Approach

We adapt Graph Matching Networks (GMN) [11] to operate on linguistic dependency trees with rich node and edge features. We compare two architectural paradigms: matching models that use cross-graph attention between sentence pairs

(TreeMatchingNet vs BertMatchingNet), and embedding models that process inputs independently with comparison only at the embedding level (TreeEmbeddingNet vs BertEmbeddingNet). Both TMN and BERT models use similar parameter counts ($\sim 36\text{M}$ for TMN, $\sim 41\text{M}$ for BERT) and undergo identical three-phase training: (1) contrastive pretraining on WikiQS and AmazonQA, (2) task-specific multi-objective InfoNCE training, and (3) direct classification fine-tuning.

Key Findings

TreeMatchingNet achieves 60.70% SNLI test accuracy compared to BertMatchingNet’s 35.38%, demonstrating $1.7\times$ superior performance with comparable parameter counts. This translates to improved parameter efficiency: TMN achieves $1.96\times$ better efficiency (1.69% vs 0.86% accuracy per million parameters). However, we observe a scaling plateau: increasing TMN parameters from 60K to 36M ($600\times$) yields minimal improvement, revealing an architectural bottleneck in the aggregation mechanism. Meanwhile, BertMatchingNet exhibits a failure mode where it predicts every test example as entailment, suggesting incompatibility between the cross-attention architecture and BERT’s design.

Contributions

We present an adaptation of Graph Matching Networks to linguistic dependency trees with rich features combining BERT embeddings, POS tags, and morphological annotations for natural language inference. Our evaluation compares tree-based GNN and BERT architectures at matched parameter counts ($\sim 36\text{M}$) with identical training protocols, enabling controlled analysis of structural inductive bias. We identify a scaling plateau where $600\times$ parameter increases yield minimal improvement, suggesting the aggregation mechanism as an architectural bottleneck. Through systematic comparison of matching versus embedding architectures, we isolate the contribution of cross-graph attention and observe compatibility differences between tree-based and transformer-based approaches.

Related Work

Natural Language Inference

The Stanford Natural Language Inference (SNLI) corpus [1] contains 570K sentence pairs labeled as entailment, contradiction, or neutral. This dataset has become a standard benchmark for evaluating semantic understanding in NLP systems. State-of-the-art models like EFL with RoBERTa-large [21] achieve approximately 90% accuracy but require 355M parameters. Our goal is to achieve competitive accuracy with substantially fewer parameters ($10\times$ reduction) by leveraging structural inductive biases rather than relying solely on model scale and extensive pretraining.

Graph Neural Networks and Graph Matching

Graph neural networks (GNNs) have emerged as effective models for learning on structured data [16]. Modern GNN architectures [12, 9, 20] compute node representations through iterative message passing along graph edges, aggregating information from local neighborhoods. These models are permutation-invariant by design and have been successfully applied to molecular property prediction, knowledge graph reasoning, and program analysis.

Graph Matching Networks [11] introduced a cross-graph attention mechanism for computing similarity between pairs of graphs. The GMN architecture consists of three key components: (1) graph propagation layers that update node representations via message passing, (2) cross-graph attention that computes attention-weighted matchings between nodes in different graphs, and (3) aggregation that produces graph-level embeddings. GMNs have been applied to graph edit distance learning, molecular similarity search, and binary function similarity analysis.

Our extension: We adapt GMN principles to linguistic dependency trees, incorporating domain-specific features (BERT embeddings, POS tags, morphology) and developing a multi-objective training protocol for multi-class NLI.

Dependency Parsing for NLU

Dependency trees have been used in NLP for semantic role labeling [7], relation extraction [3], and machine translation [6]. Early neural approaches

combined dependency structures with neural models: Recursive neural networks [17] and Tree-LSTMs [19] process parse trees in a bottom-up fashion, demonstrating that syntactic structure improves semantic representations. However, these approaches struggled with sequential processing bottlenecks and limited scalability.

Hybrid Structure-Transformer Approaches

Recent work explores combining dependency structures with transformers to preserve syntactic inductive bias while achieving scalability. SyntaxBERT [14] proposes late fusion (applying graph neural networks to transformer outputs) and joint fusion (interleaving GNN layers within transformer blocks). Dependency Transformer Grammars [23] modify transformer attention masks to simulate dependency transition systems, encoding syntactic constraints directly in attention patterns. Stack Attention [8] incorporates implicit syntactic structure through stack operations in attention mechanisms, learning syntax without explicit parse trees.

These approaches differ from our work in architectural positioning: they modify transformer attention mechanisms or post-process transformer outputs, while our proposed approach uses graph neural networks as a preprocessing step to create structurally-enriched node embeddings that serve as input to standard transformer aggregation. The current paper validates that dependency tree structure provides measurable benefits worth preserving in such hybrid architectures.

Parameter-Efficient NLP

Growing interest in parameter-efficient NLP has led to various approaches: knowledge distillation [15] transfers knowledge from large models to smaller ones, pruning [13] removes unnecessary parameters, and efficient architectures like ALBERT [10] reduce parameters through factorization and sharing. Our approach is complementary: we explore whether structural inductive biases can reduce parameter requirements from the ground up, rather than compressing existing large models.

Distance Metric Learning

Metric learning aims to learn distance functions that group similar examples together and separate dissimilar ones. Siamese networks [2] apply the same network to two inputs independently and compute similarity from the resulting embeddings. These architectures have achieved success in face verification [4, 18] and image matching [22].

Our graph matching models differ from standard Siamese networks: rather than processing inputs independently and comparing only at the output, we employ cross-graph attention throughout the propagation process. This enables early information fusion and allows the model to adjust representations based on what they are being compared to.

Tree-Based Similarity Learning for Natural Language Inference

Problem Formulation

Natural Language Inference (NLI) involves determining the relationship between a premise sentence P and hypothesis sentence H , with labels in $\{\text{Entailment (+1), Neutral (0), Contradiction (-1)}\}$. We use the Stanford Natural Language Inference (SNLI) corpus [1] as our primary evaluation benchmark. SNLI is a standard dataset for evaluating semantic understanding in NLP systems, containing sentence pairs derived from image captions with crowd-sourced annotations. Performance on SNLI demonstrates a model’s ability to understand entailment relationships, which requires capturing both semantic content and logical reasoning.

For generalization evaluation, we also use the SemEval Semantic Textual Similarity benchmark, which frames the task as a 2-class similarity problem (similar vs dissimilar). This tests whether structural benefits transfer beyond entailment to broader semantic similarity judgments.

Data Preparation. Raw text datasets were converted to dependency tree representations using TMN_DataGen (Section 3.1.1). For large-scale contrastive pretraining (WikiQS + AmazonQA), computational constraints limited full-dataset conversion and training: we converted 2.52M training pairs and randomly sample 600 batches per epoch (153,600 pairs, batch size 256) without replacement within each epoch. This sampling strategy processes $\sim 6\%$ of the contrastive corpus per epoch.

Table 1: Dataset statistics: original corpus sizes, converted tree counts, and training usage. Contrastive pretraining uses sampled batches (600/epoch) from the full corpus; primary and fine-tuning stages use complete datasets.

Dataset	Task	Original Corpus		Converted to Trees		Training Usage
		Train	Dev/Test	Train	Dev/Test	
<i>Contrastive Pretraining (Large-Scale Corpus)</i>						
WikiQS	Question Sim.	~28.8M	–	1,850,142	261,292 / 575,121	600 batches/epoch (153,600 samples)
AmazonQA	Q&A Pairs	~923K	–	669,943	83,138 / 85,089	
<i>Combined Contrastive Train:</i>					2,520,085	
<i>Primary Training and Fine-Tuning (Task-Specific)</i>						
SNLI	Entailment	550,152	10K / 10K	540,803	9,969 / 9,960	Complete dataset
SemEval	Similarity	3,000	750 / 6,750	2,967	750 / 6,638	

Each epoch uses a different random sample from the full pool. For task-specific datasets (SNLI, SemEval), we use complete converted datasets without sampling.

Data Processing Pipeline

We developed a custom data processing pipeline (TMN_DataGen) that transforms raw text into dependency trees with rich linguistic features.

Preprocessing

Text normalization includes Unicode handling, case preservation, and word boundary detection. We employ configurable strictness levels (0-3) to handle varying amounts of noise in different datasets.

Dependency Parsing

We use a multi-parser approach to combine the strengths of different parsing systems. DiaParser (based on Electra-base) provides accurate dependency tree structures and dependency relation labels, which form the graph edges in our model. SpaCy provides complementary linguistic annotations including part-of-speech tags, lemmas, and morphological features. This multi-parser strategy leverages DiaParser’s superior accuracy for syntactic structure while benefiting from SpaCy’s comprehensive linguistic feature extraction.

Feature Extraction

Each dependency tree node is represented by an 804-dimensional feature vector combining semantic and syntactic information. BERT-base-uncased

provides 768-dimensional contextual word embeddings that capture semantic content. Part-of-speech tags contribute 17 dimensions (one-hot encoded categories including NOUN, VERB, ADJ, ADV, etc.) extracted from SpaCy. Morphological features contribute 19 dimensions (one-hot encoded attributes such as Number, Tense, Person, Mood, etc.) also from SpaCy. Dependency edges are represented by 70-dimensional one-hot vectors encoding the dependency relation type (nsubj, dobj, amod, etc.) provided by DiaParser.

This combination provides both semantic representations (BERT embeddings) and explicit syntactic annotations (POS tags, morphology, dependency relations), allowing the model to leverage both learned distributional semantics and linguistic structure.

Model Architectures

We compare two architectural paradigms - **matching** (with cross-graph attention) and **embedding** (independent processing) - for both tree-based and transformer-based models.

Tree Matching Network (TreeMatchingNet)

The Tree Matching Network adapts Graph Matching Networks to dependency trees. The architecture consists of four main components:

1. TreeEncoder: Maps input features to initial node and edge representations:

$$h_i^{(0)} = \text{MLP}_{\text{node}}(x_i), \quad \forall i \in V \quad (1)$$

$$e_{ij} = \text{MLP}_{\text{edge}}(x_{ij}), \quad \forall (i, j) \in E \quad (2)$$

The node encoder MLP_{node} transforms input node features to hidden node states. The edge encoder MLP_{edge} transforms edge features to hidden edge states.

2. Graph Propagation: We employ T propagation layers. Following the GMN architecture [11], we use *shared* propagation parameters across all T layers (contrast with *unshared* variants where each layer has independent parameters). Each layer performs message passing along tree edges followed by cross-graph attention:

$$m_{j \rightarrow i} = f_{\text{message}}(h_i^{(t)}, h_j^{(t)}, e_{ij}) \quad (3)$$

$$\mu_{j \rightarrow i} = f_{\text{match}}(h_i^{(t)}, h_j^{(t)}) \quad [j \text{ from other graph}] \quad (4)$$

$$h_i^{(t+1)} = f_{\text{node}} \left(h_i^{(t)}, \sum_j m_{j \rightarrow i}, \sum_{j'} \mu_{j' \rightarrow i} \right) \quad (5)$$

Here, f_{message} is an MLP that concatenates source node, target node, and edge features. The function f_{match} implements cross-graph attention (detailed below). The function f_{node} is a GRU cell that updates node states based on both within-graph messages and cross-graph matching signals.

3. Cross-Graph Attention: The cross-graph attention mechanism computes how well each node in one graph matches nodes in the other graph:

$$a_{j \rightarrow i} = \frac{\exp(\text{sim}(h_i^{(t)}, h_j^{(t)}))}{\sum_{j'} \exp(\text{sim}(h_i^{(t)}, h_{j'}^{(t)}))} \quad (6)$$

$$\mu_{j \rightarrow i} = a_{j \rightarrow i} \cdot (h_i^{(t)} - h_j^{(t)}) \quad (7)$$

$$\sum_j \mu_{j \rightarrow i} = h_i^{(t)} - \sum_j a_{j \rightarrow i} h_j^{(t)} \quad (8)$$

This formulation has a useful property: when two graphs match perfectly and attention weights concentrate on the correct matches, $\sum_j \mu_{j \rightarrow i} \rightarrow 0$, causing the cross-graph communication to vanish. Conversely, differences between graphs are captured in the matching vectors and amplified through propagation.

4. Graph Aggregation: After propagation, we aggregate node representations to produce graph-level embeddings:

Algorithm 1 Tree Matching Network Forward Pass

Require: Trees T_A and T_B with node and edge features

Ensure: Graph embeddings $e_A, e_B \in \mathbb{R}^d$ where d is graph representation dimension

```

1:  $N_A, E_A \leftarrow \text{TreeEncoder}(T_A)$ 
2:  $N_B, E_B \leftarrow \text{TreeEncoder}(T_B)$ 
3: for layer  $l = 1$  to  $T$  do
4:    $N_A, E_A \leftarrow \text{Propagation}_l(N_A, E_A)$  ▷
     Within-graph messages
5:    $N_B, E_B \leftarrow \text{Propagation}_l(N_B, E_B)$ 
6:    $N_A, N_B \leftarrow \text{CrossAttention}_l(N_A, N_B)$  ▷
     Cross-graph matching
7: end for
8:  $e_A \leftarrow \text{Aggregator}(N_A)$ 
9:  $e_B \leftarrow \text{Aggregator}(N_B)$ 
10: return  $e_A, e_B$ 

```

$$h_{\text{graph}} = \text{MLP}_G \left(\sum_{i \in V} \sigma(\text{MLP}_{\text{gate}}(h_i^{(T)})) \odot \text{MLP}(h_i^{(T)}) \right) \quad (9)$$

The gated weighted sum allows the model to learn which nodes are most important for the graph-level representation. The final MLP_G projects to the graph representation dimension.

Configuration: Our specific TreeMatchingNet instantiation uses the hyperparameters shown in Table 2. This configuration yields approximately 36M parameters total.

Table 2: TreeMatchingNet configuration.

Component	Value
Node features	804
Edge features	70
Node state dim	1536
Edge state dim	768
Prop. layers (T)	5 (shared)
Graph rep. dim	2048
Total params	~36M

Tree Embedding Network (TreeEmbeddingNet)

The Tree Embedding Network is architecturally identical to TreeMatchingNet with one difference: it omits the cross-graph attention mechanism (line

6 in Algorithm 1). Each tree is processed completely independently, and comparison happens only at the final embedding level via cosine similarity.

This architecture allows us to isolate the contribution of cross-graph attention versus structural graph propagation. If TreeEmbeddingNet performs comparably to TreeMatchingNet, it would suggest that the structural bias alone is sufficient. Conversely, a performance gap would indicate that cross-graph attention provides additional benefits.

Configuration: TreeEmbeddingNet uses identical hyperparameters to TreeMatchingNet (Table 2), yielding approximately 36M parameters for fair comparison.

BERT Matching Network (BertMatchingNet)

To provide a fair transformer baseline, we develop BertMatchingNet with GMN-style cross-attention. The architecture consists of a standard BERT encoder with cross-attention layers inserted after each transformer layer. The cross-attention mechanism uses the same formulation as TreeMatchingNet (Equations 8-10), enabling information exchange between sentence A and sentence B representations during encoding.

Rationale for custom BERT: Using pre-trained BERT-base would provide an unfair advantage (3.3B words pretraining vs our 7M sentences). Training from scratch on identical data as TMN ensures a controlled comparison. We train a custom WordPiece tokenizer on the same pretraining data.

Configuration: Our BertMatchingNet configuration is shown in Table 3. We use fewer layers and smaller hidden dimensions than standard BERT-base to roughly match TreeMatchingNet’s parameter count, yielding approximately 41M parameters.

Table 3: BertMatchingNet configuration.

Component	Value
Vocabulary	5K
Hidden size	1024
Layers	4
Attention heads	16
Intermediate	4096
Max seq. len.	128
Total params	~41M

BERT Embedding Network (BertEmbeddingNet)

BertEmbeddingNet is standard BERT without cross-attention modifications. It serves two purposes: (1) baseline BERT performance without architectural modifications, and (2) testing whether cross-attention helps or hurts BERT. It processes sentences independently and compares their [CLS] token representations via cosine similarity.

Configuration: BertEmbeddingNet uses the same hyperparameters as BertMatchingNet (Table 3), yielding approximately 41M parameters.

Training Protocol

All models undergo identical three-phase training to ensure fair comparison.

Phase 1: Contrastive Pretraining

Objective: Learn general semantic similarity patterns from large-scale data.

Datasets: We use the combined WikiQS and AmazonQA corpus for contrastive pretraining (see Table 1 for complete statistics and sampling details).

Loss: Standard InfoNCE with temperature parameter τ :

$$\mathcal{L} = -\log \frac{\sum_{p \in P_i} e^{\text{sim}(z_i, z_p)/\tau}}{\sum_{p \in P_i} e^{\text{sim}(z_i, z_p)/\tau} + \sum_{n \in N_i} e^{\text{sim}(z_i, z_n)/\tau}} \quad (10)$$

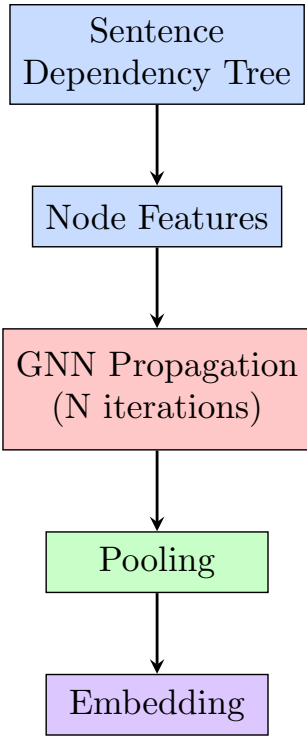
where P_i is the set of positive examples for anchor i , N_i is the set of negative examples, and we use $\tau = 0.05$ for all experiments.

Hyperparameters: Batch size 256, max batches per epoch 600, learning rate 10^{-6} (Adam), max epochs 50, patience 999 (no early stopping).

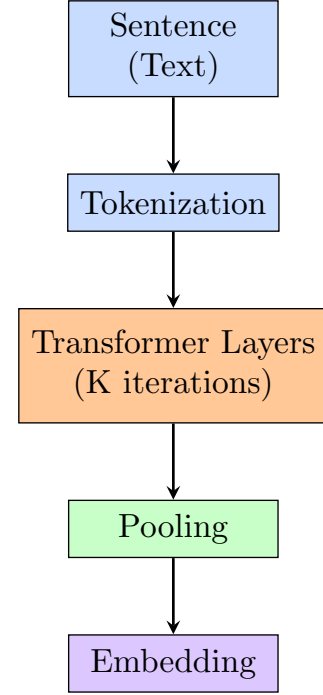
Phase 2: Primary Training (Multi-Objective Contrastive)

Objective: Adapt to task-specific similarity structure using a novel multi-objective InfoNCE formulation.

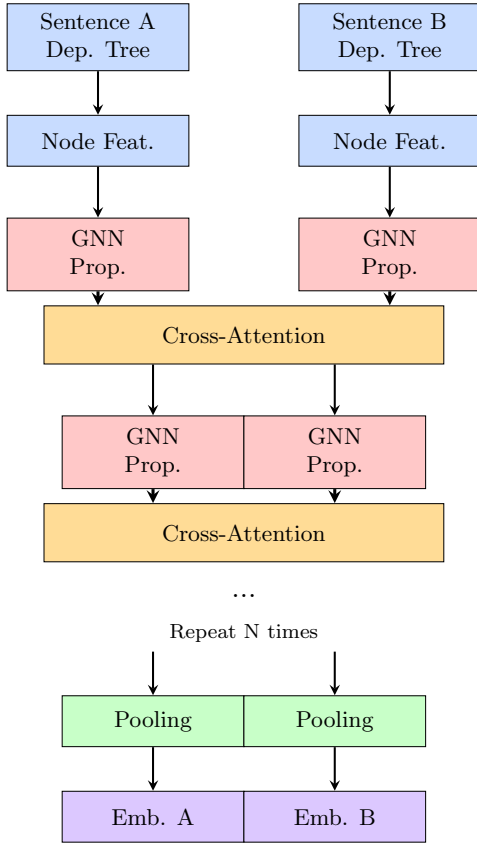
Traditional InfoNCE maximizes positive similarity and minimizes negative similarity. For 3-class NLI, we extend this to handle three distinct relationship types:



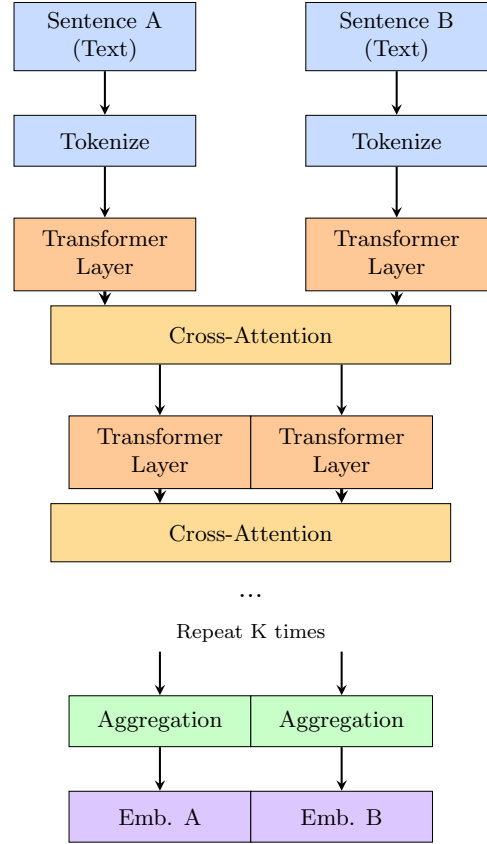
(a) Tree Embedding Network



(b) BERT Embedding Network



(c) Tree Matching Network



(d) BERT Matching Network

Figure 1: Complete architecture comparison. Top row (a-b): Embedding models process sentences independently, comparing only at the final embedding level. Bottom row (c-d): Matching models use cross-attention during encoding, enabling direct interaction between sentence representations. Left column (a,c): Tree-based models use graph propagation on dependency trees. Right column (b,d): BERT-based models use transformer encoders.

$$\begin{aligned} \text{sim}_{\text{pos}} &= \text{cosine}(e_A, e_B) \\ &\quad (\text{high for entailment}) \end{aligned} \quad (11)$$

$$\begin{aligned} \text{dist}_{\text{cos}} &= -\text{sim}_{\text{pos}} \\ &\quad (\text{high for contradiction}) \end{aligned} \quad (12)$$

$$\begin{aligned} \text{sim}_{\text{mid}} &= 1 - |\text{sim}_{\text{pos}}| \\ &\quad (\text{high for neutral}) \end{aligned} \quad (13)$$

The overall loss is a weighted combination:

$$\begin{aligned} \mathcal{L} &= w_{\text{pos}} \cdot \mathcal{L}_{\text{InfoNCE}}(\text{sim}_{\text{pos}}) \\ &\quad + w_{\text{dist}} \cdot \mathcal{L}_{\text{InfoNCE}}(\text{dist}_{\text{cos}}) \\ &\quad + w_{\text{mid}} \cdot \mathcal{L}_{\text{InfoNCE}}(\text{sim}_{\text{mid}}) \end{aligned} \quad (14)$$

The weights w_{pos} , w_{dist} , and w_{mid} reflect class importance: entailment receives the highest weight as the strongest positive signal, contradiction receives moderate weight as indicated by cosine distance, and neutral receives the lowest weight as it represents an ambiguous middle category. For SNLI (3-class), we use $w_{\text{pos}} = 0.55$, $w_{\text{dist}} = 0.30$, $w_{\text{mid}} = 0.15$. For SemEval (2-class similarity), we use $w_{\text{pos}} = 0.65$, $w_{\text{dist}} = 0.35$, $w_{\text{mid}} = 0.0$.

Hyperparameters: Batch size 256, max batches per epoch 600, learning rate 10^{-6} (Adam), max epochs 100, patience 999 (no early stopping).

Phase 3: Fine-Tuning (Direct Supervised Learning)

Objective: Maximize classification accuracy through direct supervision.

Method: Threshold-based classification from continuous similarity scores:

$$\text{label}(s) = \begin{cases} \text{Contradiction} & \text{if } s < \theta_{\text{low}} \\ \text{Neutral} & \text{if } \theta_{\text{low}} \leq s \leq \theta_{\text{high}} \\ \text{Entailment} & \text{if } s > \theta_{\text{high}} \end{cases} \quad (15)$$

where we use $\theta_{\text{low}} = -0.33$ and $\theta_{\text{high}} = 0.33$ to divide the similarity space into three approximately equal regions. These threshold values were fixed rather than tuned.

Loss: Cross-entropy over predicted class logits.

Hyperparameters: Batch size 256, max batches per epoch 600, learning rate 5×10^{-7} (Adam), max epochs 100, patience 999 (no early stopping).

Randomized Pairing Strategy for Matching Models

During pretraining of matching models (TreeMatchingNet and BertMatchingNet), we employ a randomized pairing strategy within each batch to prevent the model from exploiting the pairing structure itself rather than learning meaningful semantic relationships. Specifically, while each batch contains items with known positive relationships, we randomize which items are actually paired together during the forward pass through the cross-attention mechanism.

For example, if item A has positive match B in the batch, the model might process A paired with item C (a negative) during the forward pass. The InfoNCE loss then encourages A to be similar to its actual positive B (present in the batch) rather than to C (the item it was paired with). This forces matching models to learn robust representations that identify relationships between items regardless of the pairing order, rather than simply learning to output high similarity whenever two items are fed through cross-attention together.

TreeMatchingNet handled this randomization robustly, successfully learning to distinguish semantic relationships independent of the forward-pass pairing. BertMatchingNet, however, struggled significantly under this regime, suggesting that the cross-attention mechanism failed to learn meaningful similarity patterns when pairs were randomized.

Notably, embedding models (TreeEmbeddingNet and BertEmbeddingNet) do not face this challenge, as they process all items independently during the forward pass. Pairing only occurs at the loss computation level, where the model must identify which embeddings should be close or far apart. This architectural difference—processing independently versus processing pairs with cross-attention—may partially explain why the performance gap between TMN and BERT is smaller for embedding models.

A key diagnostic metric provides insight into this overfitting behavior: the standard deviation of embedding norms across the batch. In TreeMatchingNet (successful contrastive learning), this standard deviation increases $33.8\times$ during training ($0.012 \rightarrow 0.390$), reflecting the model learning to position embeddings at varied distances based on semantic content. In TreeEmbeddingNet (overfitting), the same metric increases only $5.7\times$ ($0.011 \rightarrow 0.063$), suggesting the model adopts a strat-

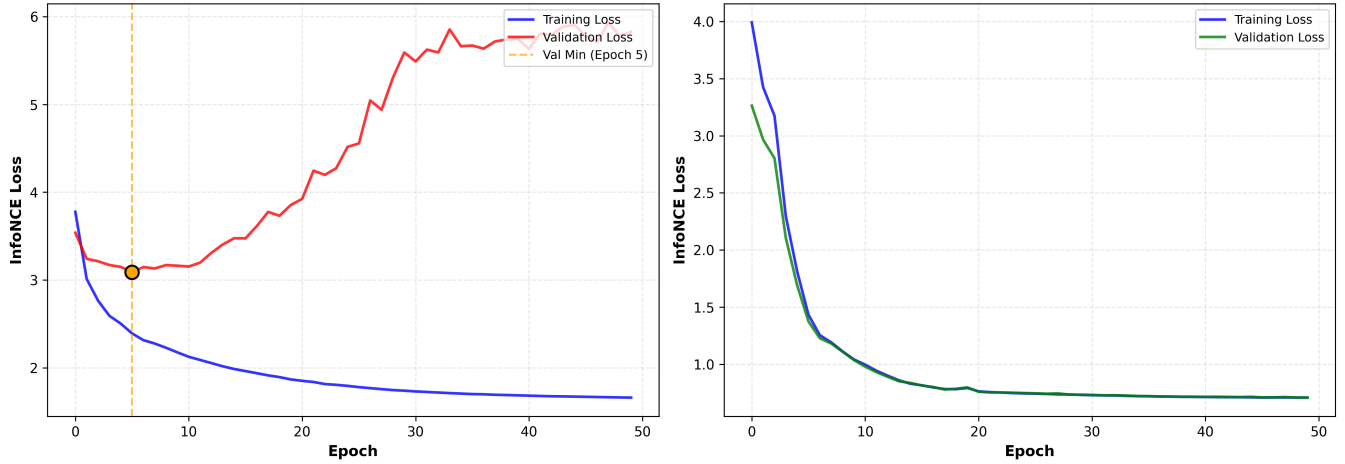


Figure 2: Training dynamics during contrastive pretraining phase. (Left) TreeEmbeddingNet shows classic overfitting: validation loss increases after an early minimum despite continued training loss decrease. (Right) TreeMatchingNet with curriculum learning shows stable convergence with both training and validation loss decreasing monotonically. Cross-graph attention appears to provide implicit regularization through comparative learning.

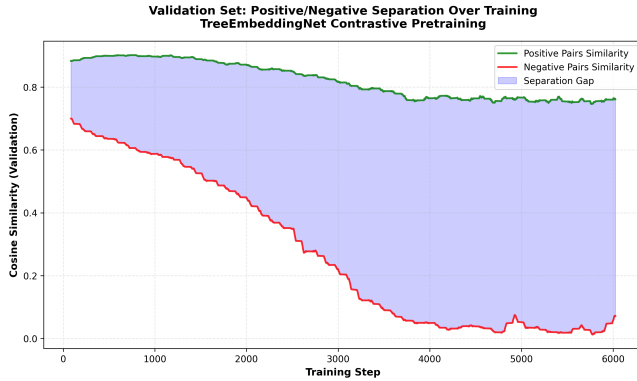


Figure 3: Positive and negative pair separation during contrastive pretraining for TreeEmbeddingNet. Mean similarity scores show successful learning: positive pairs (should be similar) increase from ~ 0.015 to ~ 0.25 , while negative pairs (should be dissimilar) remain near zero. Despite this clear separation, validation loss still degrades (Figure 2), suggesting overfitting to dataset artifacts rather than semantic relationships.

egy of pushing difficult examples toward middling similarity values rather than learning true semantic distinctions. This behavior may reflect insufficient model capacity for the contrastive task without cross-attention, leading to a suboptimal equilibrium where moderate similarity with many examples yields lower loss than attempting precise semantic positioning.

Experiments and Results

Experimental Setup

Component	Configuration
Models Evaluated	TreeMatchingNet (36M params) TreeMatchingNet Small (60K params) TreeEmbeddingNet (36M params) BertMatchingNet (41M params) BertEmbeddingNet (41M params)
Hardware	NVIDIA RTX 3090 (24GB VRAM)
Training Time	10-14 days per model (all phases)
Evaluation	SNLI test set (9,824 pairs)

Main Results: Matching Models

Table 4 presents the primary comparison between TreeMatchingNet and BertMatchingNet.

TreeMatchingNet achieves 60.70% test accuracy compared to BertMatchingNet’s 35.38%, represent-

Table 4: SNLI Test Set Performance (Matching Models)

Model	Parameters	Accuracy
TreeMatchingNet	36M	60.70%
BertMatchingNet	41M	35.38%
Random Baseline	—	33.33%
SOTA (RoBERTa-large)	355M	~90%

Table 5: TreeMatchingNet Confusion Matrix (60.70% Accuracy)

True	Predicted			Total
	C	N	E	
C	2009	654	564	3227
N	1017	1397	795	3209
E	400	484	2640	3524

ing a $1.7\times$ performance advantage with comparable parameter counts (36M vs 41M). While both models fall short of state-of-the-art performance (RoBERTa-large achieves approximately 90% with 355M parameters and extensive pretraining on 160GB of text), the comparison reveals that structural inductive biases provide substantial advantages at moderate parameter scales when training data is limited.

Confusion Matrix Analysis

TreeMatchingNet Performance

Table 5 shows the confusion matrix for TreeMatchingNet. The model achieves balanced performance across all three classes with a strong diagonal pattern indicating meaningful class distinctions.

Per-class metrics (Table 6) reveal varying performance across classes. The model achieves strongest performance on Entailment (74.91% recall, 66.02% precision), adequate performance on Contradiction (62.26% recall, 58.64% precision), and weaker performance on Neutral (43.53% recall, 55.11% precision).

The Neutral class shows the weakest performance, with 31.7% of neutral examples misclassified as Contradiction and 24.8% misclassified as Entailment. This pattern is consistent with the inherent ambiguity of the Neutral category, which is defined by the absence of a clear relationship rather than the presence of one.

Table 6: TreeMatchingNet Per-Class Metrics

Class	Precision	Recall	F1	Support
Contradiction	58.64%	62.26%	60.39%	3227
Neutral	55.11%	43.53%	48.64%	3209
Entailment	66.02%	74.91%	70.18%	3524

Table 7: BertMatchingNet Confusion Matrix (35.38% Accuracy)

True	Predicted			Total
	C	N	E	
C	0	0	3227	3227
N	0	0	3209	3209
E	0	0	3524	3524

BertMatchingNet Performance

Table 7 shows that BertMatchingNet predicts every single test example as Entailment, achieving 35.38% accuracy (the proportion of entailment examples in the test set).

This complete failure to learn Contradiction and Neutral classes results in 0% precision, recall, and F1-score for these categories, while Entailment achieves 35.38% precision and 100% recall by default. The model performs worse than random guessing (which would achieve 33.33% per class on balanced data) for two of the three classes.

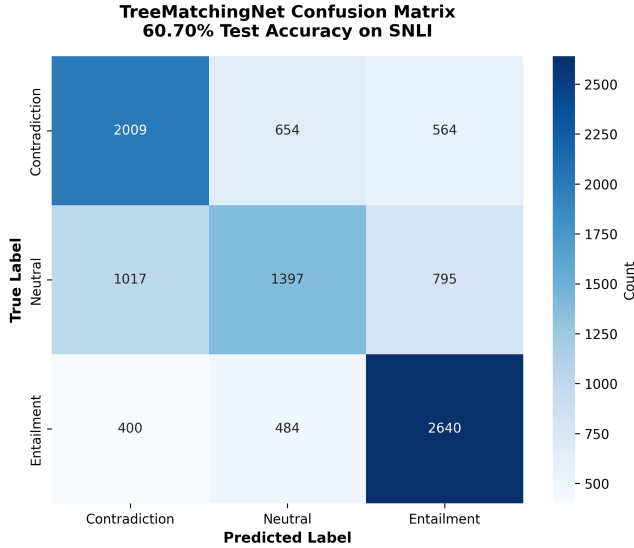
Scaling Plateau Analysis

Table 10 shows that increasing TreeMatchingNet parameters from 60K to 36M (a $600\times$ increase) yields minimal performance improvement.

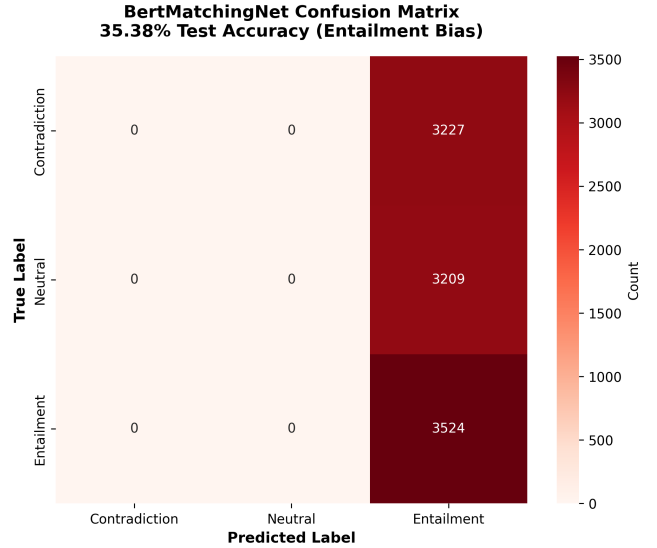
This scaling plateau indicates an architectural bottleneck preventing effective utilization of additional parameters. In typical deep learning scenarios, increasing parameters by $600\times$ with sufficient training data yields substantial improvements. The near-zero improvement here suggests that additional capacity cannot be effectively utilized by the current architecture.

Analysis of the architecture components suggests the aggregation layer as the likely bottleneck. The current implementation uses gated weighted sum pooling:

$$h_{\text{graph}} = \text{MLP}_G \left(\sum_{i \in V} \sigma(\text{MLP}_{\text{gate}}(h_i^{(T)})) \odot \text{MLP}(h_i^{(T)}) \right) \quad (16)$$



(a) TMN: Balanced predictions across all classes



(b) BERT: All predictions assigned to entailment

Figure 4: Confusion matrix visualization. (a) TreeMatchingNet shows balanced classification with a strong diagonal pattern, successfully learning all three classes. (b) BertMatchingNet exhibits complete failure, predicting every example as entailment regardless of true label. This stark contrast illustrates the difference in learning dynamics between structure-based and sequence-based approaches under identical training conditions.

Table 8: TreeMatchingNet Scaling Behavior

Model Variant	Parameters	Accuracy
TMN (Small)	60K	~60%
TMN (Large)	36M	60.70%

This formulation collapses an $N \times D$ matrix of node representations into a single D -dimensional vector, potentially destroying the rich structural information learned during propagation. A small model with limited propagation capacity produces simple node features that this aggregation can adequately summarize. A large model with rich propagation produces complex node features, but the same simple aggregation cannot effectively utilize this additional information.

Embedding Model Results

To isolate the contribution of cross-graph attention from structural inductive biases, we trained TreeEmbeddingNet and BertEmbeddingNet. These models process trees or sequences independently, with comparison occurring only at the embedding level via cosine similarity. Matching models use cross-attention during the forward pass,

Table 9: SNLI Validation Set: Matching vs Embedding Models

Model	Params	Accuracy
TreeMatching	36M	60.70%
BertMatching	41M	35.38%
TreeEmbedding	36M	~51%
BertEmbedding	41M	~40%

while embedding models compute representations independently.

Table 9 presents preliminary results from validation curves. TreeEmbeddingNet achieves approximately 51% accuracy while BertEmbeddingNet reaches approximately 40%. The performance gap between tree-based and BERT-based models remains substantial (~11 percentage points) but is noticeably smaller than the matching model gap (25.3 percentage points).

TreeEmbeddingNet outperforms BertEmbeddingNet by ~19 points, confirming that the structural inductive bias of dependency trees provides benefits independent of cross-attention mechanisms. The 19-point gap for embedding models is smaller than the 25-point gap for matching mod-

els, suggesting that cross-attention amplifies the advantage of structural representations—trees benefit more from cross-graph matching than sequences do.

BertEmbeddingNet ($\sim 40\%$) substantially outperforms BertMatchingNet (35.38%), indicating that the cross-attention mechanism interfered with BERT’s ability to learn meaningful representations. This aligns with BERT’s original design for embedding-style processing (masked language modeling, sentence pair classification without explicit cross-attention). TreeMatchingNet (60.70%) shows preforms better than TreeEmbeddingNet ($\sim 51\%$), while BertMatchingNet shows severe degradation from BertEmbeddingNet. This asymmetry suggests fundamental architectural compatibility differences.

The narrower performance gap for embedding models aligns with architectural design principles. The matching architecture with randomized pairing (Section 3.2.4) imposed a challenging learning regime that BERT’s architecture was not designed to handle. Tree-based GNNs, by contrast, naturally handle both matching and embedding paradigms due to their flexible message-passing framework.

Table 10: Parameter scaling comparison between toy and full-sized models on SNLI test set.

Model	Parameters	Accuracy
TreeMatching (toy)	$\sim 60\text{K}$	$\sim 60\%$
TreeMatching (full)	36M	60.70%
BertMatching (full)	41M	35.38%

TreeMatchingNet exhibits approximately the same performance between the toy model ($\sim 60\text{K}$ parameters) and the full-sized model (36M parameters), both achieving $\sim 60\%$ accuracy despite a $600\times$ increase in parameters. This plateau suggests an architectural bottleneck in the aggregation mechanism that prevents effective parameter utilization, regardless of the capacity added through additional propagation layers or hidden dimensions.

Analysis and Discussion

Why Does TreeMatchingNet Outperform BERT?

We propose three complementary hypotheses to explain TreeMatchingNet’s $1.7\times$ performance advan-

tage.

Structural Inductive Bias

Dependency trees explicitly encode syntactic relationships through labeled edges. Each edge represents a known linguistic relationship type (nsubj, dobj, amod, etc.), providing the model with prior structural information. In contrast, BERT processes token sequences and must implicitly learn which token pairs correspond to syntactically meaningful relationships.

This structural bias aligns with the task requirements. NLI fundamentally involves understanding how sentence components relate to each other. Dependency trees provide these relationships explicitly, while transformers must discover them through attention mechanisms operating on positional encodings and learned representations.

Efficient Information Flow

Quantitative comparison reveals a substantial difference in connectivity patterns. For a typical 50-word sentence, TMN creates approximately 50 edges (one per dependency relation), while BERT’s self-attention considers all 2,500 token pairs. Critically, every TMN edge represents a syntactically meaningful connection, whereas most BERT token pairs (e.g., "the" attending to "park" in "the dog runs in the park") lack direct syntactic relationships.

This focused connectivity enables TMN to concentrate computational resources on linguistically relevant paths through the graph. BERT must learn to identify relevant connections among a much larger set of potential relationships, effectively solving a harder learning problem with the same amount of training data.

Pre-Encoded Feature Quality

TMN node features combine BERT embeddings (768-dim contextual representations), POS tags (17-dim syntactic categories), and morphological features (19-dim linguistic annotations) for a total of 804 dimensions. This rich initial representation provides both semantic content and explicit syntactic information.

While BERT models also use rich embeddings, they do not have direct access to POS and mor-

phological features during training. TMN’s explicit encoding of these features may reduce the learning burden by providing structural information that BERT must extract implicitly.

Why Does BERT Matching Underperform?

BertMatchingNet’s complete failure (100% entailment predictions) requires careful analysis. The embedding model results provide crucial insights into this failure mode.

Cross-Attention Interference

Our BertMatchingNet employs unlearned cross-attention (matching TreeMatchingNet’s mechanism for fairness). However, BERT’s effectiveness relies on learned attention patterns developed during pretraining. Inserting fixed, non-adaptive cross-attention mechanisms disrupts these learned patterns, preventing effective information flow.

The BertEmbeddingNet results confirm this hypothesis: standard BERT without cross-attention achieves $\sim 40\%$ accuracy compared to BertMatchingNet’s 35.38%, demonstrating that cross-attention mechanisms interfere with BERT’s architecture. While 40% still underperforms TreeEmbeddingNet’s 51%, it represents a substantial 13% relative improvement, indicating that BERT can learn meaningful patterns when not constrained by the matching architecture.

Training from Scratch

Our experimental design requires training BERT from scratch on 7M sentences (100M words) to match TMN’s training data. Standard BERT-base pretraining uses 3.3B words from Books and Wikipedia. While our model has similar parameter count (41M vs 110M for BERT-base), the substantially reduced pretraining data may prevent adequate semantic understanding.

This undertraining could explain the entailment bias: the model learns a superficial pattern (predict the most common class) rather than deep semantic relationships. However, we cannot use pretrained BERT for fair comparison, as it would introduce a massive advantage in training data.

Randomized Pairing Challenge

As described in Section 3.2.4, matching models employ randomized pairing during pretraining: items are not always paired with their positive match during the forward pass, forcing the model to learn robust similarity representations rather than exploiting pairing structure. TreeMatchingNet handled this challenging regime successfully, while BertMatchingNet struggled significantly.

The embedding model results illuminate why this was so challenging for BERT: BertEmbeddingNet achieves $\sim 40\%$ without randomized pairing, while BertMatchingNet achieves only 35.38% with randomized pairing. This 5-point degradation (13% relative) indicates that BERT’s architecture cannot effectively learn from randomized pairs processed through cross-attention. The transformer’s attention mechanisms, optimized for learning from sequential structure, fail when forced to learn from mismatched pairs where the cross-attention target is deliberately chosen to be negative.

In contrast, TreeMatchingNet (60.70%) shows no degradation from TreeEmbeddingNet ($\sim 51\%$), instead outperforming it, indicating that graph-based message passing provides robustness to pairing randomization that transformers lack.

The Scaling Plateau Problem

The observed scaling plateau (60K params \approx 36M params \approx 60% accuracy) points to the aggregation mechanism as the primary source of information loss. The current pooling aggregation collapses N node representations (each 1536-dimensional) into a single 2048-dimensional graph embedding through gated weighted sum pooling. This lossy compression represents the greatest potential for representation collapse in the architecture. While the GNN propagation layers learn rich, node-level distinctions about syntactic relationships and semantic content, the aggregation step treats all nodes similarly through a simple pooling operation, potentially destroying the structural information encoded during propagation.

This observation motivates our proposed approach: using GNN-processed nodes as pre-encoded token embeddings for transformer-based aggregation. Rather than collapsing node representations through pooling, multi-headed self-attention over the graph nodes could preserve node-level distinc-

tions while enabling the model to learn which nodes are most relevant for the final representation. The GNN would handle structural encoding along dependency edges—a task it performs efficiently—while the transformer would handle aggregation over the enriched nodes, potentially achieving deeper semantic understanding with reduced training requirements compared to learning both structure and semantics from token sequences alone.

The Neutral Class Challenge

Both TMN models exhibit weaker performance on the Neutral class (43.53% recall) compared to Contradiction (62.26%) and Entailment (74.91%). This pattern reflects the inherent difficulty of the Neutral category.

Entailment and Contradiction are defined by the presence of specific logical relationships: P logically implies H (entailment) or P logically excludes H (contradiction). Neutral is defined by absence: neither entailment nor contradiction holds. In our continuous similarity space, entailment corresponds to high similarity (clear target), contradiction to low similarity (clear target), and neutral to mid-range similarity (less distinctive target).

Training signal analysis reveals the challenge:

$$\text{sim}_{\text{mid}} = 1 - |\text{sim}_{\text{pos}}| \quad (17)$$

This formulation encourages moderate absolute similarity values, but "moderate" is inherently less distinctive than "very high" or "very low". The model must learn to produce specific mid-range values, a more nuanced task than maximizing or minimizing similarity.

Additionally, many neutral examples require world knowledge or deep semantic inference. For instance, "A person is outdoors" and "A person is at a park" are neutral because the first statement is consistent with but does not entail the second. This reasoning is more subtle than recognizing direct entailment or clear contradiction.

Future Work

Transformer-Based Aggregation

The identified scaling plateau motivates exploring transformer-based aggregation as an alternative to

simple pooling. Such an architecture may leverage the benefits of tree-based semantic structures without running into the plateau, gaining efficiency in terms of the needed size of multi-headed self-attention. Treating GNN-enriched nodes as token encodings in a BERT-style model, with positional encoding that leverages the tree structure to aggregate the nodes into a single embedding, is worthy of exploration.

Ablation studies are necessary follow-ups to identify which architectural components contribute most to performance. Of particular interest is examining whether the contrastive pretraining stage can be skipped without substantial performance degradation, which would significantly reduce training time.

Several promising directions merit investigation: evaluation on additional datasets to assess generalization, exploration of larger-scale models with increased computational resources, and the proposed transformer-based aggregation architecture. These routes of inquiry may clarify the sources of tree-based advantages and identify pathways toward state-of-the-art performance with improved parameter efficiency.

Conclusion

We investigated whether graph neural networks operating on dependency parse trees can provide more parameter-efficient sentence embeddings than transformer-based models for natural language inference. We adapted Graph Matching Networks to linguistic dependency trees, creating Tree Matching Networks (TMN), and conducted comprehensive controlled comparison against BERT baselines with matched parameters ($\sim 36\text{M}$) and identical three-phase training protocols.

Key findings: (1) TreeMatchingNet achieves 60.70% SNLI accuracy vs BertMatchingNet's 35.38%, demonstrating a clear performance advantage; (2) embedding model experiments show the structural advantage persists (TreeEmbedding $\sim 51\%$ vs BertEmbedding $\sim 40\%$), with narrower gap indicating cross-attention amplifies structural benefits; (3) BertEmbeddingNet additionally outperforms BertMatchingNet, confirming that cross-attention interferes with BERT's architecture; (4) increasing TMN parameters $600\times$ (60K to 36M) yields minimal improvement, revealing a scal-

ing bottleneck consistent with prior literature on structure-based NLP approaches, but which may be overcome by introducing scalable attention-based architectures into the aggregation step of TMN.

Implications: Explicit structural representations significantly outperform sequence-based transformers at moderate scales, with benefits persisting across both matching and embedding paradigms. Inductive biases aligned with linguistic structure provide substantial learning efficiency gains. Cross-attention mechanisms amplify structural advantages for tree-based models but interfere with transformer architectures, particularly under randomized pairing regimes. Simple pooling aggregation creates bottlenecks preventing effective parameter scaling; attention-based aggregation may address this limitation.

Future directions: Attention-based aggregation promises to break through the performance plateau. Additional ablations will identify optimal configurations and further clarify the sources of tree-based advantages.

This work demonstrates that leveraging linguistic structure through graph-based neural architectures provides a promising path toward parameter-efficient natural language understanding. The dramatic performance advantage over transformers at moderate scales validates the structural bias hypothesis and opens new research directions in efficient NLP.

Code Availability

Code and data processing pipeline are publicly available:

- Tree Matching Networks: <https://github.com/jlunder00/Tree-Matching-Networks>
- Data Processing (TMN_DataGen): https://github.com/jlunder00/TMN_DataGen

Both repositories are works in progress and include documentation, configuration files, and training scripts.

Acknowledgments

This work was conducted as part of a Master’s thesis at Eastern Washington University. We thank our advisors and colleagues for their support and feedback throughout this research.

References

- [1] Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, 2015.
- [2] Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. Signature verification using a "siamese" time delay neural network. In *Advances in Neural Information Processing Systems*, volume 6, 1994.
- [3] Razvan Bunescu and Raymond Mooney. A shortest path dependency kernel for relation extraction. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 724–731, 2005.
- [4] Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively, with application to face verification. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, volume 1, pages 539–546. IEEE, 2005.
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.
- [6] Yuan Ding and Martha Palmer. A dependency-based word grouping approach to statistical machine translation. In *Proceedings of the Second Workshop on Building and Using Parallel Texts*, pages 1–8, 2005.
- [7] Daniel Gildea and Daniel Jurafsky. Automatic labeling of semantic roles. *Computational Linguistics*, 28(3):245–288, 2002.
- [8] Adi Haviv, Roi Reichart, and Roy Schwartz. A transformer with stack attention. In *International Conference on Learning Representations (ICLR)*, 2024.

- [9] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- [10] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*, 2020.
- [11] Yujia Li, Chenjie Gu, Thomas Dullien, Oriol Vinyals, and Pushmeet Kohli. Graph matching networks for learning the similarity of graph structured objects. In *International Conference on Machine Learning*, pages 3835–3845. PMLR, 2019.
- [12] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural networks. In *International Conference on Learning Representations*, 2016.
- [13] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [14] Devendra Sachan, Yuhao Zhang, Peng Qi, and William L Hamilton. Syntax-infused transformer and bert models for machine translation and natural language understanding. In *Proceedings of the 2021 Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 1425–1438, 2021.
- [15] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [16] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
- [17] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, 2013.
- [18] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In *Advances in Neural Information Processing Systems*, volume 27, 2014.
- [19] Kai Sheng Tai, Richard Socher, and Christopher D Manning. Improved semantic representations from tree-structured long short-term memory networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1556–1566, 2015.
- [20] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- [21] Sinong Wang, Han Fang, Madian Khabsa, Hanzhi Mao, and Hao Ma. Entailment as few-shot learner. In *arXiv preprint arXiv:2104.14690*, 2021.
- [22] Sergey Zagoruyko and Nikos Komodakis. Learning to compare image patches via convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4353–4361, 2015.
- [23] Yida Zhao, Chao Lou, and Kewei Tu. Dependency transformer grammars: Integrating dependency structures into transformer language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1518–1532, 2024.