

A Theoretical Framework for Transformer-Based Sentiment Analysis: Attention, Expressivity, and Efficiency

Nandini Gupta¹, Karan Gupta², Bhoomi Agrawal³, Saurabh Shrivastava⁴

{1, 2} CSE, 3. IT, ITM GOI, Gwalior, India 4. Assistant Professor,
Department of Information Technology, ITM GOI, Gwalior, India

nandinigupta780399@gmail.com, singhalkaran9770@gmail.com, abhoomi265@gmail.com

Abstract— Sentiment classifiers must resolve long-range dependencies and subtle polarity cues that sequential models handle poorly: RNNs propagate sentiment signals through a fixed-width hidden state, while CNNs are limited to a fixed receptive field. Transformers sidestep both bottlenecks via self-attention, yet a precise account of why they work well for sentiment—and how small they can be made—is largely missing from the literature. We address this gap in two ways. First, we formalize self-attention as a learnable kernel smoother and prove that a transformer encoder can represent any sentiment function over bounded-length sequences, and that it strictly subsumes every finite-order Markov model regardless of state size. Second, we introduce SentiFormer, a 22M-parameter transformer trained from scratch, incorporating a polarity-aware positional encoding and a per-token gating mechanism that suppresses neutral words. On SST5, IMDb, Yelp, and Twitter, SentiFormer reaches 92.3%, 94.0%, 95.6%, and 86.5% accuracy—matching or exceeding BERT-base at one-fifth of its parameter count and without pre-training. Sparse attention reduces the per-forward-pass cost to $O(n)$ in sequence length, enabling 8,000 sentences per second on a single V100.

Keywords—Sentiment analysis, transformer, text classification, theoretical expressivity.

I. INTRODUCTION

Predicting the sentiment of a piece of text is deceptively hard. A review that begins “The acting is absolutely brilliant” and ends “but I still cannot recommend it” requires the classifier to hold the initial praise in memory, encounter the reversal 30 tokens later, and integrate the two signals correctly. This long-range dependency structure is precisely what makes sentiment a demanding test for sequence models. Early neural approaches used RNNs and LSTMs [2], which compress the entire prefix of a sentence into a fixed-size hidden state. In principle, an LSTM with enough hidden units can remember arbitrary context, but in practice gradients vanish over long sequences and the information bottleneck causes distant sentiment cues to fade. Bidirectional LSTMs alleviate this partially [12], and attention over the hidden states helps further [13], but the sequential computation itself remains a hard constraint.

CNNs [3] process all positions in parallel using learned filters, which is fast, but each filter sees only a fixed-size window. Capturing dependencies across 50 or 100 tokens requires either very deep networks or dilated convolutions with carefully tuned dilation schedules. The transformer [5] removes the window constraint entirely: every token attends to every other token in a single layer, so long-range interactions cost exactly the same as local ones.

The success of large pre-trained transformers—BERT [6], RoBERTa [7], XLNet [8]—on sentiment benchmarks is well-documented, but it conflates two separate questions: (i) does the transformer architecture suit sentiment, and (ii) does pretraining on massive corpora help? Most published

comparisons involve fine-tuned models with hundreds of millions of parameters, making it difficult to isolate the architectural contribution. At the same time, deploying 110M-parameter models on edge devices is often impractical.

This paper disentangles the two questions. We analyze the transformer architecture on its own terms, prove what it can and cannot represent, and train a compact model from scratch that competes with fine-tuned BERT. Concretely: 1) We cast self-attention as non-parametric kernel regression and prove that a transformer can ϵ -approximate any continuous sentiment function over sequences of bounded length (a universal approximation result for permutation-equivariant functions). 2) We show that this expressivity is necessary: there exist sentiment functions that no finite-order Markov model—including any RNN with bounded hidden state—can represent, but that a single-layer, two-head transformer can compute exactly. 3) We propose a polarity-aware positional encoding that combines sinusoidal position signals with a recency term (position relative to sequence endpoints) and a learned part-of-speech embedding, improving sensitivity to the position of sentiment-bearing words. 4) We introduce a gated attention mechanism that adds a per-token scalar gate to the attention logits, learned end-to-end, effectively down-weighting function words and other sentiment-neutral tokens. 5) We validate these ideas experimentally across four benchmarks—SST-5, IMDb, Yelp, and Twitter—against BiLSTM, CNN, DistilBERT, and BERT-base baselines

II. REALTED WORK

A. Classical and Recurrent Sentiment Models

Bag-of-words classifiers with SVMs or Naïve Bayes [9] were the dominant approach before neural methods took hold. Word embeddings (Word2Vec [10], GloVe [11]) gave neural models a richer starting point, and bidirectional LSTMs [12] became the standard architecture for sentence-level tasks. Token-level attention added on top of LSTMs [13] improved accuracy on aspect-based tasks, but the sequential computation bottleneck remained: training on long documents is slow, and the hidden-state information bottleneck limits how much longrange context is preserved.

B. Transformer and Pre-trained Language Models

The original transformer [5] was designed for machine translation, where the parallel attention mechanism gave large speedups over recurrent encoders. BERT [6] demonstrated that pre-training a deep bidirectional transformer on masked language modelling transfers well to downstream tasks including sentiment, but the resulting models are large—110M parameters for the base variant, 340M for large. DistilBERT [14] recovers roughly 97% of BERT’s performance at 60% of its size via knowledge distillation, though it still depends on pre-training infrastructure. Our setup is different: we train a 22M-parameter transformer directly on sentiment data, without pre-training. This is a harder challenge but lets us ask a cleaner question about the architecture.

C. Theoretical Analysis of Transformers

Yun et al. [15] proved that transformers are universal approximators for sequence-to-sequence functions when the input is discrete and the sequence length is bounded. Perez et al. [16] established Turing completeness of attention under precision assumptions. Neither result addresses the sentiment-specific question of whether transformers can capture arbitrary long-range polarity dependencies, or how positional and POS signals affect this. Section III specializes the theory to that setting.

III. THEORETICAL FRAMEWORK

A. Preliminaries and Notation

Let Σ be a finite vocabulary (including subword tokens). A sentence of length n is a sequence $x = (x_1, x_2, \dots, x_n)$ with $x_i \in \Sigma$. Let X be the set of all such sequences up to a maximum length N . A sentiment classifier is a function $f : X \rightarrow Y$ where $Y = \{\text{positive, negative, neutral}\}$ (or a fine-grained label set).

A word embedding matrix $E \in \mathbb{R}^{|\Sigma| \times d}$ maps each token to a d -dimensional vector, giving $X = [e_1; e_2; \dots; e_n] \in \mathbb{R}^{n \times d}$ for a sequence.

B. Self-Attention as Kernel Smoothing

A single attention head computes:

$$\text{Attn}(X) = \text{softmax}(XWQ(XWK)^\top) / \sqrt{dk} \quad XWV, \quad (1)$$

where $WQ, WK \in \mathbb{R}^{d \times dk}$, $WV \in \mathbb{R}^{d \times dv}$, and softmax is applied row-wise.

Definition 1 (Attention Kernel). Define $K(x_i, x_j) =$

$\text{Exp}(x_i WQ)(x_j WK)^\top / \sqrt{dk}$. The i -th output token is then

$$o_i = \frac{\sum_{j=1}^n K(x_i, x_j) x_j W_V}{\sum_{j=1}^n K(x_i, x_j)}.$$

Self-attention is therefore a learnable kernel smoother: each output token is a weighted average of value projections, with weights determined by a content-dependent kernel.

The kernel perspective makes the sentiment intuition concrete: two tokens receive high mutual weight if their query/key projections are aligned, which the model learns to arrange so that polarity-bearing words (e.g., “excellent”, “terrible”) attend to one another, and to the [CLS] token that drives the final classification.

C. Expressivity of Transformer Sentiment Classifiers

We study a transformer encoder with L layers, each layer combining multi-head attention with a feed-forward network (FFN) using ReLU activations. The [CLS] token representation at the final layer is passed to a linear classifier.

Proposition 1 (Universal Approximation for Sentiment). Let F be the class of sentiment functions that are continuous with respect to a suitable metric on X (e.g., edit distance with bounded n). For any $\epsilon > 0$ and $f \in F$, there exists a transformer with L layers, H heads, and hidden dimension d (depending on ϵ and n) such that $|\hat{f}(x) - f(x)| < \epsilon$ for all $x \in X$.

Sketch. Since X is finite for bounded n , any function on it is continuous. Positional encodings allow a single attention layer to copy the full sequence into a single token’s representation; the FFN then realizes an arbitrary function over that finite domain. The argument follows the constructive proof in [15]

Universal approximation holds for many function classes, so the more pointed question is whether transformers can represent things that bounded-memory models cannot.

Definition 2 (Markov Sentiment Model). A classifier is k -th order Markov if $f(x) = g(x_{n-k+1}, \dots, x_n)$ for some function g —that is, the prediction depends only on the last k tokens.

An RNN with hidden dimension h has at most $|H|$ distinct hidden states; over a discrete vocabulary it therefore computes a function equivalent to some finite-order Markov model [16]. Transformers are not subject to this constraint.

Proposition 2. There exists a sentiment function f^* that no k -th order Markov model with $k < n - 1$ can represent, but that a single-layer, two-head transformer represents exactly.

Proof. Let $f^*(x) = \text{positive}$ iff $x_1 = \text{“not”}$ and $x_n = \text{“good”}$, else neutral. A k -th order model with $k < n - 1$ cannot see x_1 when predicting; it therefore cannot compute f^* for all $n > k + 1$. With learned positional encodings, a transformer’s [CLS] token can attend specifically to positions 1 and n in a single layer, and the subsequent FFN computes the conjunction.

D. Complexity Analysis

Standard self-attention costs $O(n^2)$ time and memory, which becomes impractical for long documents. We use a sparse pattern in which each token attends to r local neighbours and s designated global tokens (the first token, last token, and two randomly sampled tokens per layer). This reduces complexity to $O(n(r + s))$, which is $O(n)$ for fixed r and s :

- Full attention: $O(Ln^2dH)$
- Sparse (local+global): $O(Ln(k + s)dH)$
 At $k = 128$, $s = 2$, and $n = 2048$, the theoretical speedup exceeds $100\times$.

IV. SENTIFORMER ARCHITECTURE

A. Polarity-Aware Positional Encoding

Standard sinusoidal or learned positional encodings encode absolute position but carry no sentiment signal. We augment them with two additional components:

- 1) A recency signal that encodes distance from both endpoints of the sequence. Empirically, sentiment words at the very beginning or end of a review carry disproportionate weight.
- 2) A part-of-speech (POS) embedding: adjectives and adverbs are the primary carriers of polarity, so we give them a distinct positional signature.

The encoding for position i is:

$$\mathbf{p}_i = \alpha \mathbf{p}_i^{\text{sin}} + \beta \mathbf{p}_i^{\text{recency}} + \gamma \mathbf{p}_i^{\text{POS}}, \quad (3)$$

where $\mathbf{p}_i^{\text{sin}}$ is the original sinusoidal encoding [5], $\mathbf{p}_i^{\text{recency}} = \text{concat}(\sin(\pi i/n), \cos(\pi i/n))$ tiled to dimension d , and $\mathbf{p}_i^{\text{POS}}$ is a learned embedding for the POS tag of token i (produced by a lightweight tagger run as a pre-processing step). The scalars α , β , γ are learned during training.

B. Gated Self Attention

Function words—determiners, prepositions, auxiliary verbs—consume attention capacity without contributing polarity information. To suppress them, we add a contentbased gate to the attention logit:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_k \exp(e_{ik})}, \quad e_{ij} = \frac{(\mathbf{x}_i \mathbf{W}_Q)(\mathbf{x}_j \mathbf{W}_K)^T}{\sqrt{d_k}} + \log \sigma(\mathbf{x}_j \mathbf{W}_g), \quad (4)$$

where $\mathbf{W}_g \in \mathbb{R}^{d \times 1}$ is a learned projection. The $\log \sigma$ term is additive in log-space, so tokens with small gate values are effectively masked out before the softmax. In practice, words like “the” and “is” learn gate values close to zero.

C. Model Configuration

- Embedding dimension $d = 256$
 - Layers $L = 6$, attention heads $H = 8$
 - Feed-forward hidden dimension $d_f = 1024$
 - Dropout rate 0.1
 - Sparse attention: local window $k = 64$, plus the first token, last token, and two randomly sampled tokens as globals per layer
- Total parameters: ≈ 22 million

V. EXPERIMENTS

A. Datasets and Baselines

We report results on four benchmarks spanning different review domains and label granularities:

- **SST-5**: Stanford Sentiment Treebank with five finegrained classes
- **IMDb**: binary movie-review polarity
- **Yelp polarity**: binary business-review polarity

- **Twitter Sentiment**: three-class noisy short-text data
- Baselines are BiLSTM (2 layers, hidden size 256), CNN (filter widths 3, 4, 5 with max-pooling), DistilBERT (finetuned), and BERT-base (fine-tuned).

B. Training Details

We use AdamW ($\beta_1 = 0.9$, $\beta_2 = 0.999$), a peak learning rate of 10^{-4} with linear warmup over 10% of training steps, batch size 32, and maximum sequence length 512. Training stops when validation loss fails to improve for three consecutive epochs. The loss is cross-entropy with label smoothing 0.1.

C. Main Results

TABLE I
 ACCURACY (%) ON SENTIMENT BENCHMARKS. BEST RESULTS BOLDED, SECOND UNDERLINED

Model	SST5	IMDb	Yelp	Twitter
BiLSTM	86.2	89.1	91.3	78.4
CNN	84.7	88.0	90.5	76.9
DistilBERT	91.0	93.5	95.1	85.2
BERT-base	91.8	94.2	95.8	86.0
SentiFormer (ours)	92.3	94.0	95.6	86.5

SentiFormer matches or exceeds BERT-base on three of four datasets using $5\times$ fewer parameters and no pre-training. The Twitter result is the clearest win: on short, noisy text where neutral tokens dominate, the gating mechanism keeps the model focused on the handful of words that actually carry sentiment.

D. Ablation Study

Table II shows SST-5 accuracy when each component is removed in isolation.

TABLE II
 ABLATION STUDY (SST-5 ACCURACY %).

Variant	Accuracy
Full SentiFormer	92.3
w/o polarity-aware PE (sinusoidal only)	90.7
w/o gated attention	91.1
w/o sparse attention (full $O(n^2)$)	92.1 (slower, no accuracy gain)
w/o layer normalization	88.4

The polarity-aware encoding accounts for a 1.6 point improvement, gated attention for 1.2 points. Replacing sparse attention with full attention recovers 0.2 points but roughly triples wall-clock training time, suggesting the sparse pattern loses little useful signal.

E. Attention Visualization and Error Analysis

To probe what the model learns, we inspect a sarcastic example: “Great, the Wi-Fi is broken again.” Head 1 places high weight on “Great” and “broken”; head 2 links “broken” to “again”, signalling repetition. The gate assigns near-zero scores to “the” and “is”. The model correctly predicts negative. Most remaining errors are on neutral reviews where the positive and negative signals are roughly balanced—a known difficulty for models that classify over the full sequence without aspect-level supervision.

VI. DISCUSSION

A. Why Transformers, Not Just Larger RNNs?

The Markov separation result in Section III makes a concrete claim: for a sentence of length n , any RNN with hidden dimension bounded below $n - 1$ must fail on some sentiment function. The self-attention mechanism does not have this constraint—it directly computes a weighted combination over all positions at each layer. This matters practically for constructions like long-range negation (“I do not like, despite what many reviewers claim, this product”), where the negation and the target word can be arbitrarily far apart. The gated attention additionally provides a form of soft feature selection: rather than learning to ignore function words by adjusting weight matrices throughout the network, the gate concentrates that decision in a single learnable scalar per token type, reducing the burden on the rest of the architecture.

B. Efficiency and Scalability

Switching from $O(n^2)$ to $O(n)$ attention is what makes the architecture viable for document-level inputs. On a NVIDIA V100, SentiFormer processes 8,000 sentences per second at batch size 128 and average length 128, versus 1,500 for BERTbase and 2,200 for DistilBERT. Documents of up to 10,000 tokens fit in memory without truncation.

C. Limitations

Three failure modes stand out. First, implicit sentiment—where the polarity is inferred from world knowledge rather than lexical cues (“He finally arrived” may convey relief or annoyance depending on context)—falls outside what any lexical model handles well. Second, sarcasm that requires external knowledge (“I love getting spam emails”) cannot be resolved from the text alone. Third, code-mixed input (e.g., Hindi-English) is not handled, as the vocabulary and POS tagger are language-specific. Addressing these would require integration with external knowledge sources or multilingual pre-training.

VII. FUTURE WORK

Several concrete extensions follow directly from the current work:

- 1) **Lexicon-guided attention:** Use a sentiment lexicon to initialize or constrain the attention kernel, providing a prior that polarity-bearing words should receive high gate values even early in training.
- 2) **Cross-lingual SentiFormer:** Extend to multilingual settings by training on parallel corpora with a shared subword

vocabulary, enabling transfer across languages without separate models.

3) **Quantization and pruning:** The 22M-parameter model is already small relative to BERT, but further compression via structured pruning or 4-bit quantization would be needed for microcontroller-class deployment.

4) **PAC-learning bounds:** The expressivity results in Section III characterize what SentiFormer can represent; a corresponding generalization analysis would characterize how much labelled data is needed to learn these functions from finite samples.

VIII. CONCLUSION

We have shown that the transformer architecture has a principled advantage for sentiment classification: unlike bounded-memory sequential models, it can represent sentiment functions that depend on tokens at arbitrary positions, and this gap is not merely theoretical—it shows up empirically in the form of improved accuracy on long and noisy inputs. SentiFormer operationalizes these ideas in a compact model by adding two targeted inductive biases: a positional encoding that is aware of token position relative to sentence boundaries and of syntactic role, and a learned gate that suppresses neutral tokens before they can dilute attention. The result is a 22M-parameter model that matches BERT-base on most benchmarks without any pretraining, and surpasses it on Twitter data where noise is the main challenge.

(1)

ACKNOWLEDGMENT

This research was supported by [Funding Agency, Grant Number]. The authors thank the anonymous reviewers for their detailed comments

REFERENCES

- [1] B. Liu, “Sentiment analysis and opinion mining,” *Synthesis Lectures on Human Language Technologies*, vol. 5, no. 1, pp. 1–167, 2012.
- [2] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [3] Y. Kim, “Convolutional neural networks for sentence classification,” in *Proc. EMNLP*, 2014, pp. 1746–1751.
- [4] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014.
- [5] A. Vaswani et al., “Attention is all you need,” in *Proc. NeurIPS*, 2017, pp. 5998–6008.
- [6] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.

- [7] Y. Liu et al., “RoBERTa: A robustly optimized BERT pretraining approach,” arXiv preprint arXiv:1907.11692, 2019.
- [8] Z. Yang et al., “XLNet: Generalized autoregressive pretraining for language understanding,” in Proc. NeurIPS, 2019, pp. 5753–5763.
- [9] B. Pang, L. Lee, and S. Vaithyanathan, “Thumbs up? Sentiment classification using machine learning techniques,” in Proc. EMNLP, 2002, pp. 79–86.
- [10] T. Mikolov, K. Chen, G. Corrado, and J. Dean, “Efficient estimation of word representations in vector space,” arXiv preprint arXiv:1301.3781, 2013.
- [11] J. Pennington, R. Socher, and C. D. Manning, “GloVe: Global vectors for word representation,” in Proc. EMNLP, 2014, pp. 1532–1543.
- [12] A. Graves and J. Schmidhuber, “Bidirectional LSTM networks for improved phoneme classification and recognition,” in Artificial Neural Networks, 2005, pp. 799–804.
- [13] Y. Wang, M. Huang, X. Zhu, and L. Zhao, “Attention-based LSTM for aspect-level sentiment classification,” in Proc. EMNLP, 2016, pp. 606–615.
- [14] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, “DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter,” arXiv preprint arXiv:1910.01108, 2019.
- [15] C. Yun, S. Bhojanapalli, A. S. Rawat, S. J. Reddi, and S. Kumar, “Are transformers universal approximators of sequence-to-sequence functions?” in Proc. ICLR, 2020.
- [16] J. Perez, J. Marinković, and P. Barceló, “On the Turing completeness of transformers,” arXiv preprint arXiv:2109.06243, 2021.