# Token-Wise Kernels (TWiKers) for Vicinity-Aware Attention in Transformers

Anonymous ACL submission

#### Abstract

Self-attention mechanisms in transformers enable tokens to interact across a sequence but lack an explicit inductive bias to capture local contextual dependencies, an inherent characteristic of human languages. We propose Token-Wise Kernels (TWiKers), a novel enhancement to transformers that learn token-specific convolutional kernels applied to the keys or values. Each token is assigned a small kernel, initialized to the "Central Dirac" (e.g., [0,1,0] for size=3), meaning the token "bears" the attention from all other tokens alone. During training, these kernels adapt, and greater deviation from the Central Dirac indicates stronger attention redistribution to neighboring tokens. This introduces the first transformer weights with direct semantic interpretability. Our experiments show that content words (e.g., nouns and verbs) retain self-focus, while function words (e.g., prepositions and conjunctions) shift attention toward their neighbors, aligning with their syntactic and semantic roles. We further apply TWiKers to distinguish literary genres, historical periods, and authors, demonstrating their effectiveness in capturing high-level stylistic patterns. Finally, by allowing them to vary with attention heads, we show the potential of TWiKers as a new inductive bias to enhance transformer training.

# 1 Introduction

004

011

012

014

018

023

035

040

042

043

Transformers have revolutionized natural language processing (NLP), powering large language models (LLMs) that achieve state-of-the-art performance across diverse tasks. Recent base models, such as DeepSeek-V3 (DeepSeek-AI et al., 2025), LLaMA-4 (Grattafiori et al., 2024), and Qwen-3 (Yang et al., 2025), have exhibited increasingly strong emergent abilities, fueling speculation that large language models may be approaching the threshold of artificial general intelligence (AGI).

One of the most remarkable aspects of transformers is the multi-head attention mechanism (Vaswani

et al., 2017), which not only offers scalability but also enhances interpretability. Deep embeddings facilitate distance-based comparisons, a fundamental principle behind retrieval-augmented generation (RAG) (Lewis et al., 2020)-a key ingredient of modern AI agents. Token (shallow) embeddings are also widely used for lexical analysis, including clustering (Cha et al., 2017; Zhang et al., 2023), visualization (Le and Lauw, 2017; Molino et al., 2019), and analogy reasoning (Zhu et al., 2018; Petersen and van der Plas, 2023). However, these embeddings lack inherent meaning on their own; their interpretability depends on distance measurements and comparisons. So far, no weights in transformers have been shown to encode direct semantic meaning at the parameter level.

044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

081

While one strength of transformers is their ability to capture long-range contextual dependencies, human languages exhibit strong vicinity reliance, particularly at the lexical level. For example, when reading "War and Peace", a human would naturally focus on "War" and "Peace" while ignoring "and", which carries less semantic weight. This selective attention to content words over function words is a fundamental characteristic of natural language, not unique to English but observed in most languages. Such locality has supported sliding-window attention, enabling models like Longformer (Beltagy et al., 2020) to achieve linear-time attention computation, along with its variations such as BigBird (Zaheer et al., 2020), Mamba (Gu and Dao, 2024), and LongLoRA (Chen et al., 2024). In computer vision, similar principles have been applied in models like Swin Transformer (Liu et al., 2021) and Neighborhood Attention Transformer (NAT) (Hassani et al., 2022). Another approach that exploits local dependencies is n-gram tokenization, which explicitly captures fixed-length word sequences (Mikolov et al., 2013b; Pennington et al., 2014; Bojanowski et al., 2017; Devlin et al., 2019). However, despite the prevalence of local dependencies in hu-



Figure 1: Overview of the TWiKer mechanism. After training,  $\omega$  deviating from the Central Dirac ([0,1,0]) indicates a shift in attention toward neighboring tokens. Here we omit TWiKers for keys and their variability across heads for simplicity.

# man languages, the transformer architecture lacks an explicit inductive bias to take advantage of this characteristic.

087

094

100

101

102

103

104

106

107

108

109

110

111

112

113

In this paper, we introduce **Token-Wise Kernels** (**TWiKers**), a novel enhancement to transformers that incorporates an inductive bias to reflect vicinity reliance while preserving transformer's global attention. We assign a small, trainable convolutional kernel to each token, enabling the model to learn how different tokens interact with their immediate neighbors through attention redistribution. In this way, TWiKers capture vicinity-aware semantic relationships, as illustrated in Figure 1.

The key novelties of TWiKers are as follows:

- 1. **Direct Semantic Meaning**: Unlike standard transformer weights, TWiKers learn interpretable patterns that align with syntactic and semantic roles of words. For example, content words (nouns, verbs) tend to retain self-focus, while function words (e.g., prepositions, conjunctions) emphasize their surroundings.
- 2. Automatic Lexical and Semantic Analysis: Since TWiKers encode token-specific contextual behavior, they can be directly analyzed to distinguish lexical categories, track historical language changes, and classify text styles without additional supervision.
- 3. Enhanced Training Efficiency: Given its semantic relevance, TWiKers provide a mean-

ingful inductive bias that may improve both pretraining and finetuning by helping transformers learn embeddings aligning better with human languages. 114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

154

155

156

157

158

159

160

162

We validate TWiKers through comprehensive experiments for English, demonstrating their alignment with linguistic principles and their effectiveness in real-world applications.

# 2 Related Work

# 2.1 Sliding-Window Attention

To address the quadratic complexity of full selfattention, the sliding-window methods confine attention to local regions. For example, Longformer (Beltagy et al., 2020) uses fixed-size local windows with select global tokens for linear complexity, while BigBird (Zaheer et al., 2020) integrates random and sparse global patterns to better approximate full attention. Recent methods like Mamba (Gu and Dao, 2024), LongLoRA (Chen et al., 2024), BASED (Arora et al., 2024), and CEPE (Yen et al., 2024) further optimize local attention. In computer vision, approaches such as Swin Transformer (Liu et al., 2021) and NAT (Hassani et al., 2022) similarly enhance efficiency by focusing attention on local regions.

Although sliding-window approaches resemble TWiKers in their emphasis on local context, their motivations and effects fundamentally differ. Sliding-window methods aim to improve efficiency by restricting attention to fixed-size windows, thereby compromising the transformer's global receptive field. In contrast, TWiKers explicitly encode local semantic interactions into tokenlevel parameters, enabling the model to capture local dependencies without sacrificing global attention. Nonetheless, both approaches are grounded in the vicinity-dominated nature of human languages.

# 2.2 N-Gram Tokenization

N-gram tokenization, also based on strong vicinity reliance, represents language as sequences of contiguous units. Traditional n-gram modelsoften enhanced by smoothing techniques such as Kneser-Ney (Kneser and Ney, 1995)–have demonstrated effectiveness in classical language modeling. Neural approaches further incorporate n-gram features: fastText (Bojanowski et al., 2017) enriches word embeddings with character-level ngrams, while BPE (Sennrich et al., 2016) and SentencePiece (Kudo and Richardson, 2018) construct

subword vocabularies based on frequent n-gram 163 patterns. Recent developments have extended the 164 power of n-gram modeling. N-Grammer (Thai 165 et al., 2020) augments transformers by integrating 166 latent n-gram representations directly into the architecture. Subsequent analytical work has employed 168 n-gram statistics to examine how language models 169 implicitly capture linguistic structures (Li et al., 170 2022), conceptually close to our methodology. The Infini-gram model (Liu et al., 2024) generalizes 172 n-gram methods to infinite-length sequences using 173 an advanced back-off mechanism. Again, n-gram 174 tokenizers highlight the strong local dependencies 175 in natural language, which modern subword tok-176 enizers under-exploit. This principle aligns with 177 our approach. However, TWiKers capture locality 178 through adaptive, semantically meaningful weights learned directly from data.

# 2.3 Token Embeddings in NLP Tasks

181

183

184

190

191

192

194

195

196

197

199

211

Token embeddings are shallow representations of tokens. While they are less effective than deep transformer embeddings for contextual understanding, they have proven valuable in lexical semantic studies. Foundational models such as LSA (Landauer and Dumais, 1997), word2vec (Mikolov et al., 2013b), GloVe (Pennington et al., 2014), and fastText (Bojanowski et al., 2017) laid the groundwork for applications including clustering (Hill et al., 2015; Vulić and Mrkšić, 2018), visualization (Mikolov et al., 2013a; Reif et al., 2019), and analogy reasoning (Mikolov et al., 2013b). Recent work has extended these embeddings to cognitive and psycholinguistic domains, where they are used to model human semantic memory, word associations, and lexical access (Günther et al., 2019; Nematzadeh et al., 2017; Chronis and Erk, 2020; Samir et al., 2020). However, existing token embeddings are largely derived from statistical cooccurrence and offer limited semantic interpretability via distance comparison. In contrast, TWiKers provide direct semantic interpretability, distinguishing lexical categories (e.g., content vs. function words) and enabling automatic, linguistically meaningful analysis without supervision.

#### 3 Methodology

#### 3.1 Token-Wise Kernels in Self-Attention

In a standard transformer architecture (Vaswani et al., 2017), the attention mechanism computes output representations using the scaled dot-product attention:

$$A = \operatorname{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^{\top}}{\sqrt{d}}\right)\boldsymbol{V},\qquad(1)$$

where  $Q, K, V \in \mathbb{R}^{L \times d}$ , with L being the sequence length and d the feature dimension (ignoring the multi-head dimension). It allows each token to attend to all others in the sequence simultaneously, capturing global dependencies.

To introduce an explicit inductive bias for vicinity awareness while preserving global dependencies, we associate each token in the vocabulary with two kernels of size n: a key kernel  $\omega^k \in \mathbb{R}^n$ and a value kernel  $\omega^v \in \mathbb{R}^n$ . These kernels modify the attention mechanism by convolving the keys and values with the kernels:

$$A = \operatorname{softmax}\left(\frac{\boldsymbol{Q}(\boldsymbol{\Omega}^{\mathsf{k}}\boldsymbol{K})^{\top}}{\sqrt{d}}\right)(\boldsymbol{\Omega}^{\mathsf{v}}\boldsymbol{V}), \quad (2)$$

where  $\Omega^k, \Omega^v \in \mathbb{R}^{L \times L}$  are banded matrices (with a fixed bandwidth of *n*) that assemble the pertoken kernels  $\omega_{ij}^k$  and  $\omega_{ij}^v$  (i = 1, 2, ..., L; j =1, 2, ..., n) in a sliding-window manner. For example, when L = 4 and n = 3:

$$\boldsymbol{\Omega}^{\mathbf{k}} = \begin{bmatrix} \boldsymbol{\omega}_{11}^{\mathbf{k}} & \boldsymbol{\omega}_{12}^{\mathbf{k}} & \boldsymbol{\omega}_{13}^{\mathbf{k}} & & \\ \boldsymbol{\omega}_{21}^{\mathbf{k}} & \boldsymbol{\omega}_{22}^{\mathbf{k}} & \boldsymbol{\omega}_{23}^{\mathbf{k}} & \\ & \boldsymbol{\omega}_{31}^{\mathbf{k}} & \boldsymbol{\omega}_{32}^{\mathbf{k}} & \boldsymbol{\omega}_{33}^{\mathbf{k}} \\ & & \boldsymbol{\omega}_{41}^{\mathbf{k}} & \boldsymbol{\omega}_{42}^{\mathbf{k}} \end{bmatrix} \begin{bmatrix} \boldsymbol{\omega}_{43}^{\mathbf{k}} & & \\ \boldsymbol{\omega}_{43}^{\mathbf{k}} & & \\ & & & \\ & & & & \\ &$$

Here,  $\omega_{11}^k$  and  $\omega_{43}^k$  are truncated at sequence boundaries to avoid padding. The value transformation matrix  $\Omega^v$  is assembled in the same manner. Equation (2) is provided here for clarity, while in practice, we adopt the standard fold-multiplicationunfold pipeline to preserve the  $\mathcal{O}(L)$  complexity of convolution.

Understanding the semantic significance of key convolution  $(\Omega^k K)$  and value convolution  $(\Omega^k V)$ is essential for interpreting the learned weights. Key convolution directly shifts attention weights by incorporating neighboring tokens' key representations, effectively redistributing attention to surrounding words. Value convolution, on the other hand, blends local context into the retrieved representations, allowing tokens to reflect semantic nuances from nearby tokens. Together, these mechanisms enhance the model's ability to encode syntactic relationships and contextual meaning, explicitly reinforcing the importance of local dependencies in natural language understanding. Notably, these 213

214

215

216

217 218 219

221

222 223 224

225

227

226

228 229 230

231

232

233

234

242

243

244

245

246

247

248

249

250

251

252

346

298

vicinity-aware behaviors are semantically meaningful *only because the kernels are token-specific*,
rather than position-based, distinguishing TWiKers
from any position-wise parameters, such as IA3 for
parameter-efficient finetuning (Liu et al., 2022).

#### 3.2 Enforcing Causality

261

262

263

264

267

269

271

272

276

277

278

279

282

283

287

290

291

295

296

In autoregressive language modeling, tokens must not attend to future tokens (Vaswani et al., 2017; Dai et al., 2019). However, Eq. (2) introduces information leakage, as TWiKers allow the *i*-th token to incorporate key representations from up to (n-1)/2 future tokens, violating causality.

To address this, we must restrict the range of key and value summation in Eq. (2). The attention weight,  $A = Q(\Omega^k K)^\top$ , is corrected to:

$$A_{ij} = \sum_{l=1}^{d} Q_{il} \sum_{m=1}^{\min(n,p+i-j)} \omega_{jm}^{k} K_{j-p+m,l} \quad (4)$$

where p = (n+1)/2. The upper bound of the inner summation is reduced from n (i.e., not considering leakage) to  $\min(n, p + i - j)$ , ensuring that query i only attends to past and present keys. In implementation, we only correct the affected main diagonal and the first p - 2 sub-diagonals in  $A_{ij}$ , ensuring an  $\mathcal{O}(L)$  complexity. The attention output is adjusted similarly by modifying the summation limits on value aggregation.

The corrections to the attention weights and output incur minimal overhead. However, they prevent TWiKer-based attention from being seamlessly compatible with KV caching and flash attention (Dao et al., 2022). While adapting TWiKers to these optimizations is feasible in principle, we omit such integration in our implementation, as TWiKers are applied only to the input layer. KV caching and flash attention remain fully applicable to all deeper layers.

#### 3.3 Enforcing Probabilistic TWiKers

To enhance the interpretability of TWiKers, we enforce them to be probabilistic distributions, ensuring that their values are non-negative and sum to one. For this purpose, we define the unconstrained trainable parameters  $\hat{\omega}^k$  and  $\hat{\omega}^v$ , which are transformed via a softmax function to compute the actual kernels used for convolution:

$$\omega_{ij}^{\mathbf{k},\mathbf{v}} = \frac{\exp\left(\hat{\omega}_{ij}^{\mathbf{k},\mathbf{v}}/\tau\right)}{\sum_{m=1}^{n}\exp\left(\hat{\omega}_{mj}^{\mathbf{k},\mathbf{v}}/\tau\right)}, \quad i = 1, 2, \dots, n,$$
(5)

where  $\tau$  is the softmax temperature, treated as a hyperparameter.

To ensure that TWiKers do not affect the model prior to training, we initialize the unconstrained kernels to a sharpened Central Dirac, such as [0, 10, 0]for n = 3. This initialization enforces self-focus at the beginning, allowing the model to learn meaningful vicinity-aware modifications during training.

# 4 Experiments

In this section, we finetune GPT-2 (Radford et al., 2019) for causal language modeling using various English texts. The corpora, summarized in Table 1, span diverse genres including poetry, novels, drama, translations, and scientific articles. While TWiKers are broadly applicable to other languages and newer base models, we focus on English and GPT-2 due to resource constraints (see Limitations). Detailed data declarations are provided in Appendix A, and full engineering details can be found in Appendix B.

Specifically, the following setups are applied to ensure fair comparison across corpora:

- 1. **Data sampling**: The original corpora vary in length. From each corpus, we sample 2200 segments, each containing complete sentences and capped at 1000 tokens. The first 2000 segments are used for training, and the remaining 200 for evaluation.
- 2. Two-stage finetuning: We finetune GPT-2 on each corpus independently. We observe that some corpora (e.g., HarryPotter) converge much faster than others (e.g., Shakespeare) when trained directly with TWiKers. This discrepancy likely arises because different corpora start at varying distances from the pretrained model's local minimum-modern texts tend to be closer, while older or translated texts are farther away. To improve comparability across corpora, we adopt a two-stage training setup: each corpus is first finetuned for 30 epochs without TWiKers, allowing the model to adapt to the corpus. Then, we activate TWiKers and continue training for another 30 epochs.
- 3. **TWiKer hyperparameters**: TWiKers applied to keys or values can both shift attention toward neighboring tokens. To enhance semantic interpretability, we do not activate TWiKers for keys and values at the same time.

| Corpus (Time Period)<br>Data Source              | Linguistic Characteristics   |
|--|--|
| Shakespeare (1590–1616)<br>Zahid (2021)          | <b>Shakespeare's plays</b> : Contains 17 plays. Richly poetic language marked by inverted syntax, metaphor, and rhetorical patterning. Distinct from the straightforward diction of modern texts.  |
| Victorian (1800–1900)<br>Chapman (2022)          | <b>British Poetry from the Victorian Era</b> : Contains 2216 transcribed poems.<br>Favor formal adjectives, refined noun phrases, and measured syntax, contrasting with modern poetry's free, experimental style.  |
| NewPoems (post 2000)<br>Poetry Foundation (2023) | <b>Contemporary Poetry</b> : Contains 5000 sampled poems. Free verse with irregular syntax, simplified phrasing, and playful imagery. Many are written for children, with simple and creative language.  |
| War&Peace (~1923)<br>McKay (2016)                | <b>English Translation of War and Peace</b> : Formal adjectives and adverbs combine with expansive, subordinate clause-rich noun phrases. Retains traces of Russian syntactic structure, such as frequent use of passive voice and expressive, multi-clause constructions.                       |
| RedChamber (~1979)<br>Internet Archive (2020)    | <b>English Translation of</b> <i>The Dream of the Red Chamber</i> : Employs nuanced adjectives and adverbs and balanced noun phrases to evoke a lyrical tone. Retains classical Chinese narrative style. In our clustering experiment, we include five versions translated by different authors. |
| Dickens (1836–1870)<br>McAdams (2020)            | <b>Novels by Charles Dickens</b> : Contains 15 selected books. Ornate prose with complex noun phrases, long compounds, and descriptive clauses. Language varies with character voice and social context.   |
| StKing (1980–2000)<br>Ajmain (2022)              | <b>Novels by Stephen King</b> : Contains 20 selected books. Direct, vivid language with active verbs, informal phrasing, and narrative clarity. Blends colloquial realism with psychological tension.  |
| HarryPotter (1997–2007)<br>Kapoor (2024)         | <i>Harry Potter</i> : Contains all seven books. Clear, child-friendly prose with simple sentence structures and vivid verbs. Language mixes fantasy world-building with British idiomatic expressions.   |
| Papers (post 2000)<br>Holbrook (2020)            | <b>Scientific Articles</b> : Contains 1000 sampled paragraphs. Dense, impersonal prose with nominalization, passive constructions, and terminology. Emphasis on clarity, structure, and formal consistency.  |

Table 1: Corpora used for experiments, spanning diverse genres, time periods, and writing styles.

Unless stated otherwise (e.g., in ablation studies), we apply value convolution only, as it demonstrates greater robustness. The kernel size is fixed at three and shared across all attention heads. The softmax temperature is set to 0.4. Learning rates are fixed at  $5 \times 10^{-5}$ for model weights and  $5 \times 10^{-3}$  for TWiKer parameters, the latter compensating for small gradients near the Central Dirac initialization.

## 4.1 Lexical Attention Patterns

347

349

351

355

356

358

364

367

TWiKers offer direct insight into the local attention behavior of individual words. This subsection analyzes results from the HarryPotter corpus to demonstrate this.

We begin by examining the learned TWiKer weights without any processing. As shown in Figure 2, content words, such as "Potter" and "gold", exhibit sharply peaked central weights, indicating that they bear attention primarily on themselves. In contrast, function words like "the" and "and" spread attention across neighboring positions, reflecting their syntactic role in structuring phrases rather than anchoring meaning. This difference aligns well with traditional linguistic distinctions between semantic and grammatical categories. 368

369

370

371

372

373

374

375

376

377

378

380

381

382

384

385

387

388

To quantify this behavior across broader lexical classes, we compute the average deviation of TWiKer weights from the Central Dirac for common parts of speech (PoS) tags. As shown in Figure 3, function words such as determiners and conjunctions exhibit greater deviations, while contentrich categories such as nouns and verbs remain closer to the central peak. This highlights the capacity of TWiKers to encode meaningful linguistic structure in an interpretable, unsupervised fashion.

**Lexical handedness** Another interesting property we observe is a directional asymmetry in the learned TWiKer weights. Among tokens whose central kernel weight is below 0.99, we categorize them as *left-handed* if the left value exceeds the right, and *right-handed* otherwise. In the HarryPotter corpus, we find that 9,570 tokens



Figure 2: Learned TWiKer kernels for selected tokens in HarryPotter. Each triplet of bars shows the kernel weights for left (L), center (C), and right (R) positions. Content words show dominant center weights, while function words spread their attention to adjacent tokens.



0.005 0.01 0.015 0.02 0.025 0.03 0.035 Mean L2 Distance from TWiKers to Central Dirac

Figure 3: Mean deviation of learned TWiKers from Central Dirac [0, 1, 0] across PoS tags in HarryPotter. Higher values indicate broader attention spread away from the token itself.

exhibit right-handedness, while only 84 are lefthanded–a striking imbalance. This pattern aligns with the well-established fact that English is a rightbranching language (Dryer, 1992; Du et al., 2020), where syntactic dependents such as complements and modifiers typically follow their heads. Function words (e.g., prepositions, subordinating conjunctions) often anticipate or introduce material to their right, naturally shifting attention forward in the sequence. This finding reinforces the idea that TWiKers internalize not only lexical category behavior but also broader structural tendencies in-

396

400

herent in natural language.

A key limitation in interpreting TWiKers arises from the use of subword tokenization in LLMs. To address this, we adopt two filtering measures. First, we exclude tokens that serve solely as suffixes– specifically, those not beginning with the "Ġ" (space) character. Second, for each token, we examine its PoS tags across the corpus and discard those that appear as affixes in more than 10% of their occurrences. This accounts for cases where a rare word begins with a common word, such as "upsurge" beginning with "up". We revisit this issue in more detail in Limitations. 401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

#### 4.2 Cross-Corpus Comparison

Figure 4 shows the overall TWiKer deviations from the Central Dirac across corpora. With deviations ranging from low to high, the corpora can be grouped into four categories: Academic, Poetic, Translations, and Novels. The results capture the strong influence of genre and stylistic conventions on attention patterns, with more structured or constrained texts yielding lower deviation, and freer, narrative-driven texts yielding higher deviation.

The Academic corpus (Papers) exhibits the lowest deviation, consistent with its rigid syntactic patterns and semantically dense constructions. Such writing minimizes contextual dependencies and maintains tight lexical focus.

The Poetic corpora follow, with their low deviation reflecting structured phrasing and rhythmic regularity. Notably, Victorian poetry shows lower



0.001 0.002 0.003 0.004 0.005 0.006 0.007 Mean L2 Distance from TWiKer to Central Dirac

Figure 4: Mean deviation of learned TWiKers from the Central Dirac [0, 1, 0] across nine corpora. Higher values suggest broader attention spread at the lexical level, often associated with more dynamic or loosely structured prose. Lower values indicate tighter, more self-contained word usage, reflecting semantically denser expression or a more formal tone.

deviation than NewPoems: the former adheres to metrical constraints and formal diction, while the latter–comprising free verse and children's poetry– permits flexible syntax and imaginative phrasing, increasing attention spread. Shakespeare occupies an intermediate position, reflecting its combination of poetic formality with syntactic inversion and dramatic rhythm.

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

Novels, both translated and native, display the highest deviations, reflecting their narrative characteristic and syntactic variety. However, translated works (War&Peace, originally written in 19th-century Russian, and RedChamber, from 18th-century vernacular Chinese) exhibit somewhat lower deviation than native English novels, likely reflecting the relative syntactic compactness of their source languages and regularization introduced during translation. Within novels, HarryPotter exhibits the highest deviation, reflecting its conversational style, flexible sentence structures, and its blend of fantasy and colloquial language aimed at younger audiences.

To further investigate the relationship between genre, style, and attention spread, we examine TWiKer deviations by PoS tags in Appendix C, focusing on three corpora that diverge notably from general English norms.



0.0005 0.0010 0.0015 0.0020 Mean L2 Distance from TWiKer to Central Dirac

Figure 5: Mean deviation of learned TWiKer's from Central Dirac [0, 1, 0] across five English translations of *The Dream of the Red Chamber*.

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

### 4.3 Clustering Translations

As a real-world application, we use TWiKers to cluster different English translations of *The Dream of the Red Chamber* (红楼梦), one of the most celebrated novels in Chinese literature. A cloud over the novel's history is the uncertainty of its authorship. It is established that Cao Xueqin (曹雪芹) wrote the first 80 chapters, whereas the authorship of the final 40 chapters–possibly by Gao E (高鹗)– remains debated. While we are unable to resolve this historical mystery using GPT-2, it inspires us to analyze five full English translations of the novel through the lens of TWiKers.

We compare five English versions of Dream of the Red Chamber. The earliest, by H. Bencraft Joly (Joly, 1893), covers Chapters 1-56 in formal, archaic Victorian prose. It was later extended to Chapter 80 by Florence and Isabel McHugh (McHugh and McHugh, 1958), based not on the Chinese original but on Franz Kuhn's German version, adding an extra interpretive layer. The widely circulated edition by Yang Hsien-yi and Gladys Yang (Yang and Yang, 1980), published in China, is clear and faithful, prioritizing literal accuracy and accessibility over literary embellishment. David Hawkes' acclaimed transla-



Figure 6: Clustering five English translations of *The Dream of the Red Chamber*. Each point represents one corpus ( $\sim$ 24 chapters), where the label shows the ground-truth (initial of the first translator's name and the starting chapter number; see Figure 5), and the marker shape indicates clustering results. We use a simple KMeans algorithm, starting from 100 different random states, and show the best results as above. Subfigures (a) and (b) are based on learned TWiKer weights, and (c), as a baseline, is based on PoS tag distributions.

tion (Hawkes and Minford, 1986) (Volumes I-III), completed by John Minford (Volumes IV-V), is widely accepted as the most literary version, with idiomatic prose and extensive cultural notes. Lastly, we include a machine-translated version by OpenAI's o3-mini, which is fluent and modern but may lack consistency in style between chapters.

We split the novel's 120 chapters into five segments, each containing  $\sim$ 24 chapters, and use them as individual corpora to train TWiKer-enhanced GPT-2. Figure 5 shows the mean deviation of TWiKers from the Central Dirac. Even this single scalar metric can loosely differentiate translators.

For finer-grained analysis, we compute the average TWiKer deviation across PoS tags in each corpus, and apply KMeans clustering. Figure 6a shows the results using all five translations. Clustering is nearly perfect: the only notable misplacement is McHugh (M57) being absorbed into the AI cluster, while KMeans separates the two Joly corpora (J1, J29) to satisfy the five-cluster constraint. When we exclude the AI translation, Figure 6b shows that all corpora are clustered correctly. As a baseline, we also cluster based on simple PoS tag distributions (Figure 6c). While PoS can reflect some stylistic distinctions, its granularity is insufficient for accurate clustering, mainly due to mixing between the human-translated versions.

# 4.4 Ablation Study

486

487

488

489

490

491

492

493

494

495

496

497

498

499

506

507

511

512

To assess the impact of architectural choices on TWiKer behavior, we conduct an ablation study focused on three factors: kernel size, whether TWiKers are applied to keys or values, and whether they vary across attention heads. We observe that overall lexical patterns–such as content words being self-focused and function words distributing attention–remain consistent across the different configurations. Specifically, increasing the kernel size broadens attention spread; applying TWiKers to keys introduces strong variation in deviation for certain PoS tags; and head-specific TWiKers smooth deviation patterns and improve training convergence. These results, detailed in Appendix D, support that TWiKers serve as an effective inductive bias with small parameter footprint. 516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

538

539

540

541

542

543

544

545

546

# 5 Conclusion

We have introduced TWiKers, a novel mechanism that equips transformers with token-specific convolutional kernels, providing a lightweight inductive bias toward vicinity reliance-an inherent property of human languages. Our experiments show that TWiKers capture meaningful lexical and syntactic behaviors without supervision: content words retain self-focus, while function words redistribute attention to neighboring tokens. This behavior generalizes across diverse corpora in English, reflecting both low-level linguistic regularities and highlevel stylistic variation. By offering the first transformer weights with direct semantic interpretability, TWiKers may open new directions for linguistic analysis and the development of efficient, interpretable neural weights for language modeling.

# 547 Limitations

548 Our study has two main limitations. First, tokens do not always correspond to words under modern 549 subword tokenization schemes. We address this 550 by excluding suffix tokens from our analysis and 551 consistently aligning tokens with complete words. 552 While this filtering reduces the statistical power of our results, word-token alignment holds for the majority of the text-reflecting a key design prin-555 ciple of subword tokenizers. For more linguistically demanding applications, it is possible to 557 558 pretrain models with larger, word-oriented vocabularies. Second, due to resource constraints, our experiments are conducted using GPT-2. Although an older model, GPT-2 retains the core architectural principles of causal decoder models, making 562 it appropriate for testing our hypotheses. As a 563 consequence, our analysis is limited to English. 564 Extending TWiKers to languages with diverse morphological and syntactic structures remains an important direction for future work. 567

## References

570

571

576

577

580

581

582

585

593

597

- Enan Ajmain. 2022. Stephen king books dataset. https://www.kaggle.com/datasets/ lujar1762/stephen-king-books. Kaggle dataset.
- Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, James Zou, Atri Rudra, and Christopher Re. 2024. Simple linear attention language models balance the recall-throughput tradeoff. In *ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models*.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL 2020), pages 1248–1261.
- Sister Magdalen Louise Blum. 1950. The imagery in the poetry of gerard manley hopkins. Ma thesis, University of New Mexico.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. In *Transactions of the Association for Computational Linguistics (ACL 2017)*, volume 5, pages 135–146.
- Miriam Cha, Youngjune Gwon, and H. T. Kung. 2017. Language modeling by clustering with word embeddings for text readability assessment. In Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, CIKM 2017, pages 2003–2006, New York, NY, USA. Association for Computing Machinery.

Alison Chapman. 2022. Digital victorian periodical poetry project (dvpp). https://dvpp.uvic.ca/. University of Victoria, last accessed April 24, 2022. 598

599

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

- Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2024. Longlora: Efficient fine-tuning of long-context large language models. In *The International Conference on Learning Representations (ICLR 2024).*
- George Chronis and Katrin Erk. 2020. When is a bishop not like a rook? when it's like a rabbi! multi-prototype bert embeddings for estimating semantic relationships. In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 227–244.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-XL: Attentive language models beyond a fixed-length context. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 2978–2988, Florence, Italy. Association for Computational Linguistics.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. *arXiv preprint arXiv:2205.14135*.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, and 181 others. 2025. Deepseek-v3 technical report. *Preprint*, arXiv:2412.19437.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*.
- Rodger Drew. 1996. Symbolism and Sources in the Painting and Poetry of Dante Gabriel Rossetti. Phd thesis, University of Glasgow.
- Matthew S Dryer. 1992. The greenbergian word order correlations. *Language*, 68(1):81–138.
- Wenyu Du, Zhouhan Lin, Yikang Shen, Timothy O'Donnell, Yoshua Bengio, and Yue Zhang. 2020. Exploiting syntactic structure for better language modeling: A syntactic distance approach. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6611– 6628.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. *Preprint*, arXiv:2407.21783.

756

757

758

759

760

761

Albert Gu and Tri Dao. 2024. Mamba: Linear-time sequence modeling with selective state spaces. In *First Conference on Language Modeling (COLM 2024)*.

658

662

664

670

671

672

673

674

675

676

677

678

679

685

687

691

697

703

704

707

- Fritz Günther, Luisa Rinaldi, and Marco Marelli. 2019. Vector-space models for the representation of word meanings: a survey. *Language, Cognition and Neuroscience*, 34(5):572–590.
- Anthony H. Harrison. 2004. Pre-raphaelite and ruskinian aesthetics. *The Victorian Web*.
- Ali Hassani, Steven Walton, Jiacheng Li, Shengjia Li, and Humphrey Shi. 2022. Neighborhood attention transformer. 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2023), pages 6185–6194.
- David Hawkes and John Minford. 1986. The Story of the Stone. Penguin Books. Five-volume edition; Hawkes translated Chapters 1–80, Minford translated 81–120.
- Felix Hill, Roi Reichart, and Anna Korhonen. 2015. Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4):665–695.
- Ryan Holbrook. 2020. Scientific papers dataset. https://www.kaggle.com/datasets/ ryanholbrook/scientific-papers. Kaggle dataset.
- John Dixon Hunt. 1968. *The Pre-Raphaelite Imagination, 1848-1900.* Routledge & Kegan Paul, London.
- Internet Archive. 2020. The dream of the red chamber (book i). https://archive.org/details/ the-dream-of-the-red-chamber-book-i\_ gutenberg-etext9603. Archive.org (Gutenberg text).
- Jacob Jewusiak. 2021. Tennyson's wrinkled feet: Ageing and the poetics of decay. *19: Interdisciplinary Studies in the Long Nineteenth Century*, 2021(32):1– 20.
- H. Bencraft Joly. 1893. *The Dream of the Red Chamber*. Kelly & Walsh Press, Hong Kong; Macao Commercial Printing Bureau. English translation of Chapters 1–56.
- Rupanshu Kapoor. 2024. Harry potter books dataset. https://www.kaggle.com/datasets/ rupanshukapoor/harry-potter-books. Kaggle dataset.
- Reinhard Kneser and Hermann Ney. 1995. Improved backing-off for m-gram language modeling. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 181–184. IEEE.
- Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In

Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations (EMNLP 2018), pages 66–71.

- Thomas K Landauer and Susan T Dumais. 1997. A solution to plato's problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*, 104(2):211.
- Tuan M. V. Le and Hady W. Lauw. 2017. Semantic visualization for short texts with word embeddings. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017*, pages 2074–2080.
- Patrick Lewis, Ethan Perez, Adam Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Kelvin Lu, Sebastian Riedel, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems (NeurIPS 2020)*, volume 33.
- Xuebo Li, Linyi Wu, Xinyang Li, and Cho-Jui Hsieh. 2022. Understanding transformers via n-gram statistics. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing* (*EMNLP 2022*), pages 9950–9960.
- Jiacheng Liu, Sewon Min, Luke Zettlemoyer, Yejin Choi, and Hannaneh Hajishirzi. 2024. Infini-gram: Scaling unbounded n-gram language models to a trillion tokens. In *First Conference on Language Modeling (COLM 2024)*.
- Peng Liu, Weizhu Xu, Yu Zhang, Xingang Lin, Xuezhi Ma, Jie Liu, and Steven C.H. Hoi. 2022. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS).*
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV 2021)*, pages 10012–10022.
- P. N. Madhusudana. 2022. Dramatic monologues: A study of robert browning's narrative techniques. *Journal of Emerging Technologies and Innovative Research (JETIR)*, 9(10):f47–f51.
- Josh McAdams. 2020. Jane austen and charles dickens collection. https:// www.kaggle.com/datasets/joshmcadams/ jane-austin-and-charles-dickens. Kaggle dataset.
- Florence McHugh and Isabel McHugh. 1958. *The Dream of the Red Chamber*. Pantheon Books, New York, NY. Translated from Franz Kuhn's German version.

Matt McKay. 2016. Text of war and peace. https: //github.com/mmcky/nyu-econ-370/blob/ master/notebooks/data/book-war-and-peace. txt. GitHub repository.

762

763

766

773

775

777

778

780

784

787

789

790

791

794

796

799

801

807

810

811

812

813

814

815

- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. In *Proceedings of the 1st International Conference on Learning Representations (ICLR 2013).*
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeffrey Dean. 2013b. Distributed representations of words and phrases and their compositionality. In Advances in Neural Information Processing Systems (NeurIPS 2013), volume 26.
- Nicole Miras. 2024. Art, myth, and literature: The pre-raphaelites. *The Crossroads Gazette*.
- Piero Molino, Yang Wang, and Jiawei Zhang. 2019. Parallax: Visualizing and understanding the semantics of embedding spaces via algebraic formulae. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, pages 165–180, Florence, Italy. Association for Computational Linguistics (ACL 2019).
- Aida Nematzadeh, Stephan C Meylan, and Thomas L Griffiths. 2017. Evaluating vector-space models of word representation, or, the unreasonable effectiveness of counting words near other words. *Proceedings of the 39th Annual Conference of the Cognitive Science Society*, pages 859–864.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP 2014), pages 1532–1543.
- Molly Petersen and Lonneke van der Plas. 2023. Can language models learn analogical reasoning? investigating training objectives and comparisons to human performance. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP 2023)*, pages 16414–16425, Singapore. Association for Computational Linguistics.
- Poetry Foundation. 2023. Poetryfoundation.org: Data summary. https://www.poetryfoundation.org. Public online archive.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and 1 others. 2019.
   Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Emily Reif, Ann Yuan, Martin Wattenberg, Fernanda B Viegas, Andy Coenen, Adam Pearce, and Been Kim. 2019. Visualizing and measuring the geometry of bert. In Advances in Neural Information Processing Systems (NeurIPS 2019), volume 32. Curran Associates, Inc.

Moussa Samir, Stephane Dufau, Gareth Gaskell, and Anastasia Ulicheva. 2020. Modeling semantic priming and lexical decision with distributed semantic spaces. In *Proceedings of the Cognitive Science Society*, volume 42, pages 1054–1060. 816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistic (ACL 2016)*, pages 1715–1725.
- Svetlana S. Takhtarova and Amelia Sh. Zubinova. 2018. The main characteristics of stephen king's idiostyle. *Vestnik Volgogradskogo gosudarstvennogo universiteta Serija 2 Jazykoznanije*, 17(3):139–147.
- Jeuti Talukdar. 2024. Narrative techniques in the novels of charles dickens: A comparative analysis. *International Journal of Creative Research Thoughts* (*IJCRT*), 12(2):154–160.
- Binh Thai, Yu Wu, Pranav Jain, Ankur P Ravula, and Mohit Iyyer. 2020. N-grammer: Augmenting transformers with latent n-grams. *arXiv preprint arXiv:2007.12766*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS 2017)*, volume 30.
- Ivan Vulić and Nikola Mrkšić. 2018. Specialising word vectors for lexical entailment. *Proceedings* of NAACL, pages 1134–1145.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. Qwen3 technical report. *Preprint*, arXiv:2505.09388.
- Xianyi Yang and Gladys Yang. 1980. *A Dream of Red Mansions*. Foreign Languages Press, Beijing. Threevolume English translation of the full novel.
- Howard Yen, Tianyu Gao, and Danqi Chen. 2024. Longcontext language modeling with parallel context encoding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2588–2610, Bangkok, Thailand. Association for Computational Linguistics.
- Manzil Zaheer, Gokhan Guruganesh, Souvik Dubey, James Ainslie, Claudio Alberti, Santiago Ontanon, Paul Pham, Abhishek Ravula, Qifan Wang, and Li Yang. 2020. Big bird: Transformers for longer sequences. In *Advances in Neural Information Processing Systems (NeurIPS 2020)*, volume 33, pages 17283–17303.
- Asim Zahid. 2021. Shakespeare plays dataset. https://www.kaggle.com/datasets/ asimzahid/shakespeare-plays. Kaggle dataset.

- 872 873
- 875
- 87

- 878 879
- 880 881
- 88

890

900

901

902

903

905

908

909

910

911

913

914

915 916

917

918

919

920

Yuwei Zhang, Zihan Wang, and Jingbo Shang. 2023. ClusterLLM: Large language models as a guide for text clustering. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 13903–13920, Singapore. Association for Computational Linguistics (ACL 2023).

Xunjie Zhu, Tingfeng Li, and Gerard de Melo. 2018. Exploring semantic properties of sentence embeddings. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 632–637, Melbourne, Australia. Association for Computational Linguistics (ACL 2018).

# A Data Declaration

We reviewed all nine corpora (see Table 1 for detailed descriptions of texts, time periods, and sources) to ensure that no personally identifying or offensive content was present. All materials are drawn from published literary works (public domain or widely distributed) and sampled scientific abstracts; we did not include private correspondence or unpublished personal data. Automated scripts scanned for full name patterns, email style strings, and offensive keywords; any hits were manually inspected and where necessary redacted. In the case of scientific articles, we also removed author bylines, institutional affiliations, and acknowledgments to protect anonymity.

All our data are in English. For originally non-English works (War and Peace and The Dream of the Red Chamber), we use their English translations; we also note multiple translator variants and demographic context (e.g. British vs. Russian vs. Chinese authors) in Table 1. For each corpus we record the number of works (e.g. 17 Shakespeare plays, 2 216 Victorian poems, 5 000 contemporary poems, 1 000 scientific article paragraphs, etc.), the source citation, and the predominant linguistic phenomena (e.g. inverted syntax and metaphor in Shakespeare, nominalization and passive constructions in scientific prose).

Across all corpora we processed approximately 1.2 million tokens. Each corpus was split at the document level into 80% train, and 20% test sets stratified by author and genre to preserve stylistic diversity. Detailed token counts per split (and per PoS tag) are provided in the supplementary Jupyter notebook, alongside document counts and PoS tag distributions.

# **B** Engineering Details

We trained GPT-2 Base (117 M parameters) using a single NVIDIA V100 (40 GB) GPU. Total compute per corpus averaged under one GPU-hour (including both forward and backward passes), with all experiments running on the same V100 instance.

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

All main experiments reported in Section 4.1 used fixed hyperparameters: a learning rate of  $5 \times 10^{-5}$  for the Transformer weights and  $5 \times 10^{-3}$ for the TWiKer kernel parameters; a batch size of 6 for both training and evaluation; a TWiKer kernel size of 3 applied to the value projections in the attention mechanism;  $2 \times 30$  training epochs; and a softmax temperature of 0.4 for normalizing TWiKers. Hyperparameter sweeps and ablation studies are discussed separately in Appendix D.

TWiKers are implemented through local modifications to Huggingface's transformers library. For data processing and analysis, we use SpaCy's en\_core\_web\_sm model for part-of-speech tagging and NLTK's default rule-based tokenizer for sentence segmentation.

All results are reproducible via one-click experiment scripts and plotting utilities included in our released codebase and dataset package.

# C TWiKers across PoS Tags in Corpora

In this Appendix, we present continued results for Section 4.2. Figure 7 shows the mean deviation of learned TWiKers from the Central Dirac kernel across nine corpora, broken down by PoS tags. These results reveal consistent trends in attention spread across lexical categories, while also highlighting stylistic variation among genres and time periods.

**Charles Dickens** and **Victorian Poetry** are the only corpora in which *prepositions* likely exhibit greater deviation than *determiners*. In Dickens, this may reflect a narrative style that tends to emphasize spatial density and rhythmic layering (Talukdar, 2024). For example, in *Great Expectations*:

"In a corner of the forge, the fire was burning brightly, and Joe was at his bellows, energetically puffing away."

Here, *prepositions* such as "in", "of", and "at" likely function as structural anchors, distributing descriptive weight across the sentence. This style could be seen as aligning with Victorian literary aesthetics, where detailed spatial descriptions and atmospheric depth were common. TWiKers can



Figure 7: Mean deviation of learned TWiKer's from Central Dirac [0, 1, 0] across PoS tags in nine corpora.

learn to spread attention accordingly, capturing the rhetorical centrality of prepositional phrases in Dickens's prose.

970 971

972

973

975

976

977

981

982

984

985

**Victorian poetry**, though also showing elevated prepositional deviation, appears to follow a different stylistic rationale. Many literary scholars have noted that poets like Alfred Tennyson, Gerard Manley Hopkins, and Dante Gabriel Rossetti often favor determiner-noun imagery over clause-based narrative progression (Jewusiak, 2021; Blum, 1950; Drew, 1996). This stylistic choice likely reflects an emphasis on visual immediacy and symbolic precision, where *prepositions* often serve dual roles: indicating location and reinforcing prosodic balance. For instance, in Tennyson's *Tithonus*:

"The woods decay, the woods decay and

# Or Hopkins's The Windhover:

thing!"

987

988

990 991

992

993

994

995

996

997

998

999

989

Such usage suggests that *prepositions* and *determiners* function not merely as grammatical elements but as imagistic anchors. In contrast to narrative poets like Robert Browning, who rely heavily on *conjunctions* for logical progression ("And then she smiled..."), these poets emphasize stasis, vision, and repetition (Madhusudana, 2022). This static and visual emphasis connects closely with contemporary Victorian movements, such as the Pre-Raphaelite focus on symbolic and detailed vi-

"The achieve of, the mastery of the

sual imagery (Harrison, 2004; Hunt, 1968; Miras, 2024).

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1013

1014

1015

1016

1018

1019

1020

1021

1022

1024

1025

1027

1029

1030

1031

1032

1034

1036

1038

1039

1041

1042

1043

1044

1045

1046

1047

1049

Additionally, **Victorian** poetry is the only corpus in which *determiners* likely deviate more than *conjunctions*. This could be attributed to the design of our dataset, which includes poets who often prioritize determiner-led imagery over logical connectives. For example, in Tennyson's *The Lady of Shalott*:

> "The mirror crack'd from side to side; 'The curse is come upon me,' cried The Lady of Shalott."

Each instance of "the" may function as a visual or symbolic anchor–"mirror", "curse", "lady"– while *conjunctions* are comparatively minimized. This focus on determiner-led imagery is not universal among Victorian poets; for example, Browning and Christina Rossetti are known fo their reliance on clause-driven narrative progression. Our corpus likely foregrounds poets with a more determinercentric style.

**Stephen King** presents a third, striking divergence: his is the only corpus where *interjections* appear to show the highest TWiKer deviation. This may be due to his focus on emotional immediacy, especially in horror and psychological suspense, where interjections often serve as narrative turning points (Takhtarova and Zubinova, 2018). From *The Green Mile*:

"We each owe a death, there are no exceptions, I know that, but sometimes, oh God, the Green Mile is so long."

And in Carrie:

"No. Oh dear God, please no. (please let it be a happy ending)"

These utterances do not carry strict syntactic function, but they likely help regulate pacing, convey fear, and anchor character perspective. TWiKers may capture this by assigning wider attention to such tokens, reflecting their dependence on surrounding discourse rather than immediate syntactic neighbors.

Taken together, these stylistically grounded deviations could support a key claim: TWiKers do not merely encode syntactic proximity-they can internalize genre conventions, authorial style, and literary tradition. The model's attention behavior highly resonates with deep patterns in English literary history, offering an interpretable bridge between data-driven learning and humanistic reading.



Figure 8: Training loss curves for ablation variants of TWiKer on the HarryPotter corpus. Final evaluation loss and accuracy are shown in the legend.

## **D** Ablation Study

In this section, we examine how various architectural choices influence the behavior of TWiKers, using the HarryPotter corpus. The reference configuration uses a kernel size of 3, with TWiKers applied to the values in the attention mechanism, shared across all attention heads. This setup underpins the results presented in Section 4.1 and Section 4.2. 1050

1051

1052

1053

1054

1055

1058

1059

1060

1061

1063

1064

1065

1067

1069

1070

We consider three ablation variants, each modifying a single factor while keeping all others fixed:

- **Kernel size = 5**: Increases the TWiKer kernel width, allowing tokens to incorporate a broader local context.
- **TWiKer on Keys**: Applies TWiKers to the keys instead of the values, shifting the locality bias from the value aggregation to the query-side matching process.
- Head Variant: Assigns a separate TWiKer to each attention head within the input layer, enabling head-specific attention patterns.

Figure 8 shows the training loss curves under 1071 each configuration. As TWiKers introduce only a 1072 small number of additional parameters, they do not substantially affect optimization dynamics on 1074 their own. However, when allowed to vary by head 1075 (Head Variant), we observe slight improvements 1076 in both convergence rate and final evaluation ac-1077 curacy. This suggests that TWiKers can serve as a 1078 lightweight and semantically grounded inductive 1079 bias in language modeling. Nevertheless, as noted 1080 in Limitations, all results are based on GPT-2. We do not claim general efficiency or scalability of 1082

TWiKers at larger model scales, and leave this for future investigation. 1084

1083

1103

Figure 9 shows the mean deviation of learned 1085 TWiKer kernels from the Central Dirac across PoS 1086 tags under different configurations. Across all these 1087 variants, the overall pattern holds: function words 1088 (e.g., determiners, conjunctions) tend to shift at-1089 tention to neighbors, while content words (e.g., 1090 nouns, verbs) retain self-focus. Increasing the ker-1091 nel size to five leads to broader deviation, espe-1092 cially for function words. Subordinate conjunc-1093 tions show an outstanding relative increase in de-1094 viation when TWiKers are applied to keys, likely 1095 because their clause-linking function interacts more 1096 strongly with the query-side of attention. Allow-1097 ing variation across heads (Head Variant) results 1098 in smoother distance distributions across PoS cat-1099 egories, suggesting a regularizing effect from dis-1100 tributing the locality pattern across multiple atten-1101 tion paths. 1102

#### Ε Use of AI Assistants

We used ChatGPT-40 and DeepSeek R1 to help 1104 write Python code and improve sentences. No part 1105 of the code or paper was generated by AI without 1106 human guidance and verification. 1107



Figure 9: Mean deviation of learned TWiKers from the Central Dirac [0, 1, 0] across PoS tags for different architectural configurations. *Reference*: kernel size = 3, TWiKer on values, head-invariant.