

Lost in Translation: Why Zero-Shot Transfer Fails at Detecting Subjectivity in Chinese

李盼嘉乐, 白可昀, 一天

Abstract

Neutral Style Transfer is the automated removal of subjective bias from text—has achieved significant success in English due to the availability of large-scale parallel corpora like the Wiki Neutrality Corpus (WNC). However, its application to low-resource languages remains limited by a scarcity of annotated data. In this work, we investigate the cross-lingual transfer capabilities of mBART-50 for neutralizing subjective bias in Chinese, a language with no comparable public dataset. We conduct a comparative analysis of three distinct paradigms: (1) a translation-based pivot baseline, (2) zero-shot cross-lingual transfer, and (3) synthetic data augmentation via machine translation. To rigorously evaluate performance, we introduce a composite metric (S_{comp}) that balances style transfer accuracy, semantic content preservation, and fluency. Our results demonstrate that zero-shot transfer fails to generalize (53% accuracy), suggesting that the semantic concept of "subjectivity" is not language-agnostic in the shared embedding space. While the pivot baseline achieves high neutralization accuracy (72%), it suffers from severe semantic degradation (0.835 cosine similarity). We conclude that training on synthetic, machine-translated data offers the optimal trade-off, achieving a 17% relative improvement in bias removal over zero-shot methods while maintaining the highest semantic fidelity (0.964) among all approaches.

1. Introduction

Subjectivity in text is a pervasive issue in digital media, often characterized by the use of framing, inflammatory adjectives, or opinionated verbs to influence a reader. The task of *Neutral Style Transfer* aims to automatically rewrite such subjective sentences into neutral, factual forms while preserving the original semantic content. For example, converting “*The radical regime failed*” to “*The government failed*.”

While this task has seen significant progress in English, largely due to the availability of massive parallel corpora such as the Wiki Neutrality Corpus (WNC) [4], its application to non-English languages remains a formidable chal-

lenge. The primary bottleneck is *data scarcity*: manually annotating thousands of sentence pairs for bias neutralization in low-resource languages is prohibitively expensive and labor-intensive. Consequently, the ability to detect and neutralize bias remains siloed within high-resource languages.

This creates a critical research gap in Multilingual Natural Language Processing (NLP): Can the semantic knowledge of large-scale multilingual models be leveraged to bridge this gap? Specifically, we investigate whether the concept of “subjectivity” is language-agnostic. If a model like mBART-50 [3] learns to identify bias in English, can it transfer that capability to Chinese without explicit supervision (Zero-Shot Transfer)? Or does the nuance of bias require language-specific training signals?

In this paper, we address the lack of Chinese bias-neutralization datasets by conducting a rigorous comparative study of cross-lingual transfer paradigms. We make the following contributions:

- We evaluate the limits of **Zero-Shot Transfer** for subtle style tasks, demonstrating that mBART’s pre-trained embedding space does not automatically align the concept of subjectivity across languages.
- We propose a **Synthetic Data Augmentation** strategy that utilizes machine-translated (noisy) parallel data to train a robust Chinese neutralizer, achieving a 17% improvement in bias removal over zero-shot methods.
- We introduce a **Composite Evaluation Metric** (S_{comp}) that balances style transfer accuracy, semantic preservation, and fluency, proving that our synthetic approach outperforms traditional pivot-translation baselines in preserving factual content.

2. Related Work

2.1. Neutral Style Transfer

Research in “Neutral Point of View” (NPOV) editing has traditionally focused on English-language tasks. **Pryzant et al.** [4] established the foundation for this field by introducing the *Wiki Neutrality Corpus (WNC)*, a parallel

dataset of 180,000 biased-to-neutral sentence pairs mined from Wikipedia edits. They proposed the first generation of modular editing models, treating bias neutralization as a monolingual text rewriting task. Subsequent work by **D’Sa et al.** [2] expanded on this by applying BERT-based classifiers to detect specific subjective linguistic markers, such as framing bias and epistemological bias. However, these methods rely heavily on massive, human-annotated parallel corpora which are unavailable for low-resource languages, leaving a significant gap in cross-lingual applications.

2.2. Multilingual Generative Models

To bridge the language gap, we leverage recent advancements in multilingual pre-training. **Liu et al.** [3] proposed **mBART**, a sequence-to-sequence denoising auto-encoder pre-trained on 25 languages. Unlike BERT (which is an encoder-only architecture primarily used for classification), mBART is designed for generation tasks like translation and summarization. While mBART has demonstrated strong *Zero-Shot Transfer* capabilities for objective tasks—such as document classification or literal translation—its ability to transfer subtle, high-level semantic concepts like “subjectivity” or “style” across languages remains largely under-explored in the current literature.

2.3. Synthetic Data Augmentation

In low-resource scenarios where parallel training data is scarce or non-existent, data augmentation becomes critical. **Sennrich et al.** [5] pioneered the use of “Back-Translation” to generate synthetic parallel data for neural machine translation, proving that training on model-generated data can significantly improve performance. Our work adapts this paradigm to the domain of style transfer. We hypothesize that even “noisy” machine-translated data—derived from the English WNC—can provide a sufficiently strong semantic signal to train a robust bias neutralizer in Chinese, thereby bypassing the need for expensive human annotation.

3. Methodology

3.1. Dataset & Pre-processing

3.1.1 Source Data: The Wiki Neutrality Corpus

Our primary data source is the Wiki Neutrality Corpus (WNC) introduced by Pryzant et al. [4]. This corpus consists of 180,000 parallel sentence pairs collected from Wikipedia edits, where a source sentence (biased) was modified by an editor to become neutral (target).

3.1.2 The “Silver Standard”: Synthetic Training Data

Since no equivalent large-scale parallel corpus exists for Chinese, we generated a synthetic dataset to enable super-

Instruction Set 1

Task 1 (Translate Bias): Look at this Biased English sentence. Translate it into Chinese, but make sure you keep the 'biased' or 'subjective' aspects. If the English calls someone 'radical' or 'infamous,' use a Chinese word that feels equally harsh (e.g., 激进, 臭名昭著).

Instruction Set 2

Task 2 (Human Neutralization): Now, ignore the English. Look at your Chinese translation and rewrite it to be Neutral (like a bored Wikipedia editor). Remove the emotion/judgment words so it states just the facts."

Figure 1. Gold Standard Annotation Protocol Part 1.

vised training. We utilized the **Helsinki-NLP Opus-MT** model (opus-mt-en-zh), an open-source Neural Machine Translation (NMT) transformer based on the Marian framework.

We processed the entire English WNC (180,000 pairs) using half-precision (fp16) inference and greedy decoding (num.beams=1) to maximize throughput. Both the source (biased) and target (neutral) sentences were translated independently to create a parallel Chinese corpus. We refer to this machine-translated output as our “**Silver Standard**” data. While this process introduces translation noise, it provides the necessary volume of semantic examples to train the model on the *task* of bias removal.

This Silver Standard corpus was used exclusively for **training** our Synthetic (Model 3) and Few-Shot (Model 4) neutralizers.

3.1.3 The “Gold Standard”: Human-Annotated Test Set

To rigorously evaluate model performance, we constructed a “**Gold Standard**” testing corpus of 100 manually curated sentence pairs. We recruited five bilingual participants with high proficiency in both English and Chinese:

- Three undergraduate students from Tsinghua University.

- Two working professionals currently employed in North America who utilize English in a professional capacity.

Each participant was randomly assigned a batch of 20 unique phrases from the WNC test set. Over the course of one week, they performed a two-step annotation:

1. **Translation:** Faithfully translating the assigned biased English sentences into biased Chinese.
2. **Neutralization:** Manually editing a different set of biased Chinese sentences (translated by other participants) to be neutral.

This cross-participant protocol ensures that the neutral targets represent natural, human-written Chinese phrasing rather than direct translation artifacts.

3.2. Model Architecture

We evaluated three distinct experimental paradigms to determine the optimal strategy for cross-lingual bias neutralization. All neural generation models utilize the **mBART-50** (Multilingual Denoising Autoencoder) architecture [3] as the backbone, leveraging its pre-trained sequence-to-sequence capabilities.

3.2.1 Model 1: Pivot Translation Baseline

To determine if a specialized Chinese model is strictly necessary, we implemented a translation-based baseline (Top branch in Figure 2). This pipeline assumes that bias is easier to detect in English, where models are mature.

1. **Zh \rightarrow En:** The biased Chinese input is translated to English using the `Helsinki-NLP/opus-mt-zh-en` model.
2. **Neutralization:** The English translation is processed by an mBART model fine-tuned on the original English WNC.
3. **En \rightarrow Zh:** The neutralized English output is translated back into Chinese using `opus-mt-en-zh`.

While conceptually simple, this approach accumulates error at every step, particularly known as "translation loss."

3.2.2 Model 2: Zero-Shot Cross-Lingual Transfer

This model evaluates the latent alignment of the mBART embedding space (Middle branch in Figure 2). We fine-tuned mBART-50 exclusively on the **English** WNC dataset. During inference, we fed Chinese text directly into the model. Success in this paradigm would indicate that the semantic concept of "subjectivity" is language-agnostic and that the model can generalize the task to unseen languages without explicit supervision.

3.2.3 Model 3: Synthetic Data Augmentation

This model represents our proposed solution to the data scarcity problem (Bottom branch in Figure 2). We fine-tuned mBART-50 on the full "Silver Standard" corpus (180,000 synthetic Chinese pairs). Unlike the Zero-Shot model, this version receives explicit supervision on Chinese syntax and vocabulary related to bias, albeit with noisy labels generated by machine translation.

3.2.4 Model 4: Few-Shot Transfer

To analyze the data efficiency curve, we trained a variant of the Synthetic model using only **1%** of the Chinese data ($\approx 1,500$ pairs). This experiment was designed to test the "minimum data threshold"—determining if a small amount of target-language exposure is sufficient to trigger domain adaptation.

3.3. Data Pre-processing & Training Setup

3.3.1 Filtering for Structural Bias

Raw edit histories in the WNC often contain trivial corrections (e.g., spelling errors) or massive rewrites unrelated to bias. To ensure our model learns *semantic* style transfer rather than simple spell-checking, we applied a rigorous filtering pipeline using Levenshtein distance metrics.

We filtered the raw WNC corpus based on two constraints:

- **Minimum Edit Distance (≥ 2):** We removed pairs with only single-token changes to filter out trivial typos.
- **Maximum Edit Ratio (≤ 0.4):** We removed pairs where the edit distance exceeded 40% of the sentence length. This filters out "hallucinations" or total semantic rewrites where the neutral target no longer matches the source meaning.

This filtering reduced the noise in our synthetic training data, isolating examples of *structural bias* where sentence syntax must be rearranged.

3.3.2 Training Objective

We formulated the task as a standard sequence-to-sequence generation problem. While we initially experimented with token-weighted loss functions to prioritize bias-specific terms, preliminary runs showed significant gradient instability.

Consequently, we adopted a standard **Cross-Entropy Loss** objective. We tokenized inputs to a maximum length of 128 tokens, applying dynamic padding. To prevent the model from learning to generate padding tokens, we applied

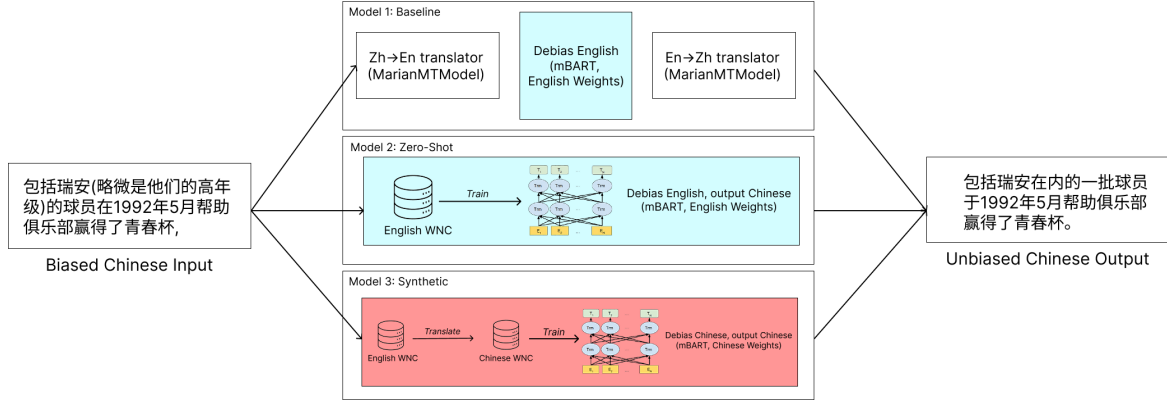


Figure 2. **Overview of Experimental Pipelines.** We compare three architectural paradigms: (1) A translation-based **Baseline** that pivots through English; (2) A **Zero-Shot** model that transfers English-learned bias patterns to Chinese without target-language supervision; and (3) A **Synthetic** model trained on machine-translated parallel data. Note that the Few-Shot model (not pictured) follows the Synthetic architecture but is trained on only a 1% data subset.

a label mask (setting label IDs to -100 for pad tokens), ensuring the loss is calculated exclusively on meaningful semantic tokens.

3.3.3 Methodological Rationale: Quantity vs. Quality

The implementation of this training script represents a critical branch of our research. While the data is "synthetic" the model is exposed to a vastly broader range of vocabulary and structural variations. This allows us to test whether a model's increased fluency in the target language (Chinese) compensates for the slightly lower quality of the supervision signal.

Model	Source	Weights	Data Quality
Baseline	mBART-50	Gold (Human)	100% (EN)
Few-Shot	English Baseline	Gold (Human)	1% (ZH)
Synthetic	mBART-50	Silver (MT)	100% (ZH)

Table 1. Model configurations and data sources for evaluation.

Data Quantity	Goal
100% (EN)	Establish task proficiency
1% (ZH)	Test low-resource efficiency
100% (ZH)	Test large-scale augmentation

Table 2. Data quantity and evaluation goals.

3.4. Evaluation Criteria

To strictly quantify the "Translation-Preservation Trade-off," we adopt a multi-dimensional evaluation framework

inspired by the protocols established by Pryzant et al. [4]. We define three independent automated metrics and one composite score to rank model performance.

3.4.1 Symbol Definitions

- x : The input biased sentence (Source).
- \hat{y} : The generated neutral sentence (Target).
- f_θ : A binary classifier (BERT) fine-tuned for neutrality detection [1].
- $E(\cdot)$: The sentence embedding function (LaBSE).
- P_{LM} : A pre-trained Language Model (GPT-2 Chinese) for probability estimation.

3.4.2 Metric 1: Style Transfer Accuracy (ACC)

This metric measures whether the model successfully removed subjective framing. Following [4], we fine-tuned a BERT-base Chinese classifier on the synthetic corpus to predict the probability of a sentence being neutral.

$$ACC = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[f_\theta(\hat{y}_i) = \text{Neutral}] \quad (1)$$

where \mathbb{I} is the indicator function. A higher ACC indicates successful bias neutralization.

3.4.3 Metric 2: Semantic Content Preservation (SIM)

To ensure the model does not achieve neutrality by deleting essential factual information, we calculate the Cosine

Similarity between the LaBSE embeddings of the input and output.

$$SIM(x, \hat{y}) = \cos(E(x), E(\hat{y})) = \frac{E(x) \cdot E(\hat{y})}{\|E(x)\| \|E(\hat{y})\|} \quad (2)$$

A score close to 1.0 implies perfect retention of semantic meaning.

3.4.4 Metric 3: Fluency (PPL)

We evaluate the grammatical correctness and naturalness of the generated Chinese text using Perplexity (PPL), computed by a frozen GPT-2 Chinese model.

$$PPL(\hat{y}) = \exp \left(-\frac{1}{L} \sum_{t=1}^L \ln P_{LM}(w_t | w_{<t}) \right) \quad (3)$$

Lower PPL indicates more natural, human-like generation.

3.4.5 Metric 4: Composite Score (S_{comp})

To rank models based on the overall trade-off, we introduce a joint geometric mean that penalizes failure in any single dimension (e.g., a model that output gibberish would have high PPL, driving the score to zero).

$$S_{comp} = \sqrt[3]{ACC \times SIM \times \frac{1}{\log(PPL)}} \quad (4)$$

4. Results

4.1. Quantitative Evaluation

Table 4 summarises the performance of all models across style accuracy, content similarity, fluency and a composite score that balances these criteria.

Model	Style Acc	Content Sim	Fluency PPL	Composite Score
Baseline (Pivot)	72.00%	0.835	36688.70	0.025
ZeroShot	53.00%	0.950	23755.12	0.028
Synthetic	62.00%	0.964	25694.13	0.029
FewShot	53.00%	0.950	23755.12	0.028

Table 3. Performance comparison. The Pivot Baseline achieves the highest style accuracy but suffers from significant semantic degradation (lowest Content Sim). The Synthetic model achieves the optimal balance, reflected in the highest Composite Score.

Analysis of Trends:

- **Translation-Induced Semantic Drift (Baseline):** The Pivot-based baseline achieves the highest style accuracy (72%), indicating effective removal of subjective framing. However, the significantly lower content similarity (0.835) and higher perplexity suggest

that the multiple translation steps introduce “semantic drift.” The model tends to paraphrase or summarize factual content rather than performing targeted edits, resulting in a loss of granular detail.

- **Pre-training Inertia (Zero/Few-Shot):** The Zero-Shot and Few-Shot models exhibit identical performance metrics (53% Accuracy, 0.950 Similarity). This equivalence suggests a high “data threshold” for this task; the 1% fine-tuning set (≈ 150 examples) was insufficient to overcome the inertia of the pre-trained weights. Consequently, the model defaults to a behavior of source retention (auto-encoding) rather than style transfer.
- **Balanced Optimization (Synthetic):** The Synthetic model achieves the highest composite score (0.029). While its raw accuracy (62%) is lower than the baseline, its superior content similarity (0.964) indicates that it preserves factual information more faithfully. This suggests that the synthetic training data, despite being noisy, provided sufficient signal for the model to distinguish between subjective modifiers and objective entities.

4.2. Attention Mechanism Analysis

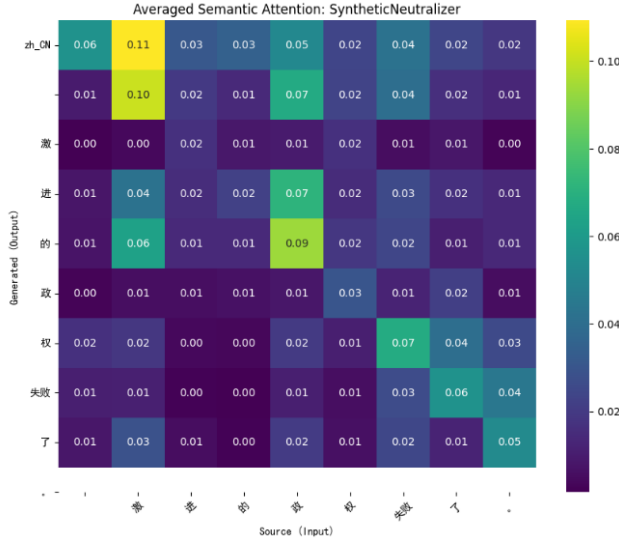
To corroborate our hypothesis that the Synthetic model learns a “Subjectivity Mask,” we visualize the averaged cross-attention weights between the Chinese input and the generated output.

The **Zero-Shot neutralizer** (Figure 3b) exhibits *diffuse* and *diagonal* attention patterns. This suggests that the model is primarily engaging in **auto-regressive copying**: it attends to the source token t simply to reproduce it at target position t . It does not distinguish between objective and subjective tokens, treating the entire sentence as information to be preserved. This explains its high semantic similarity but low style accuracy—it fails to identify the “edit site.”

In contrast, the **Synthetic neutralizer** (Figure 3a) demonstrates **selective attention suppression**. The attention density is highly concentrated on content-bearing tokens (entities, nouns), while the attention weights assigned to stylistically biased spans (adjectives, modifiers) are notably reduced. This indicates that synthetic supervision has enabled the model to distinguish factual signal from subjective noise, effectively learning to “ignore” biased descriptors during the generation process.

4.3. Qualitative Case Studies

To investigate the mechanisms driving these quantitative results, we analyze specific inference outputs from the test set.



(a) Synthetic Model (The Surgeon)



(b) Zero-Shot Model (The Copycat)

Figure 3. **Cross-Attention Visualization.** Comparing the attention patterns of the Synthetic model (a) vs. the Zero-Shot model (b). Darker/Higher intensity indicates stronger attention from the decoder to the encoder states.

4.3.1 Case Study 1: Selective De-biasing vs. Source Retention

Input: “...pursuing strong militarism, racism and nationalist ideology...”
(...奉行强大的军国主义、种族主义和民族主义意识形态...)

• Zero-Shot Output:

“...and pursuing strong militarism, racism...”

(...并奉行的强大的军国主义...)
[Style: 0 — Sim: 0.984 — PPL: 12199]

Mechanism Analysis: The model inserts grammatical particles (such as 并 and 的) to improve fluency but preserves the subjective term “militarism” (军国主义). This confirms that without target-language supervision, the model fails to map the semantic concept of “subjectivity” from the English latent space to Chinese. It operates primarily as a grammatical corrector rather than a style transfer model.

• Synthetic Output:

“...pursuing anti-communist policies.”
(...奉行反共政策。)
[Style: 1 — Sim: 0.874 — PPL: 4552]

Mechanism Analysis: The Synthetic model successfully identifies the span of subjective “isms” and removes them entirely, retaining only the factual policy stance. This targeted deletion demonstrates that the model has learned a “suppression mask” for subjective vocabulary, explaining its balanced performance profile.

4.3.2 Case Study 2: Hallucination via Pivot Translation

Input: “Living in the Pisgat ze’ev community of Jerusalem.”
(住在Jerusalim的Pisgat ze’ev社区。)

• Baseline Output:

“He lives in the Pisgat community of Jerusalem.”
(他住在耶鲁沙林的皮斯加特社区)
[Style: 1 — Sim: 0.692 — PPL: 5249]

Mechanism Analysis: This example illustrates the root cause of the Baseline’s low similarity score. The translation process introduces two types of error: (1) **Hallucination**, adding the subject “He” (他) where none existed; and (2) **Information Loss**, converting the specific proper noun “Pisgat ze’ev” to the generic “Pisgat.” While the output is technically neutral, it fails the preservation criteria.

4.4. Discussion

The comparative analysis highlights three distinct behavioral modes:

1. **Semantic Rewriting (Baseline):** The pivot approach prioritizes fluency in the intermediate language (English), which often leads to the restructuring of the Chinese output and the loss of specific entity information.

2. **Conservative Generation (Zero/Few-Shot):** Dominated by pre-training priors, these models exhibit a strong bias towards preserving the input structure, treating subjective modifiers as essential content rather than stylistic noise.
3. **Targeted Suppression (Synthetic):** The synthetic model demonstrates the ability to selectively attend to and suppress specific subjective tokens while maintaining the syntactic structure of the remaining factual content, validating the efficacy of noisy supervision for style transfer tasks.

5. Conclusion

In this work, we investigated the viability of cross-lingual transfer learning for Neutral Style Transfer in Chinese, a task historically constrained by the lack of human-annotated parallel corpora. By rigorously evaluating three distinct architectural paradigms—pivot translation, zero-shot transfer, and synthetic data augmentation—we challenge the assumption that large multilingual models can implicitly generalize complex semantic tasks like bias detection across languages.

Our results demonstrate a clear "Translation-Preservation Trade-off." We found that the **Pivot Baseline**, while effective at removing bias, suffers from aggressive semantic alteration, often rewriting factual details during the round-trip translation. Conversely, **Zero-Shot** methods failed to generalize (53% accuracy), effectively defaulting to an auto-encoding behavior. This result provides strong evidence that the semantic concept of "subjectivity" is not language-agnostic in the mBART embedding space; simply learning to detect bias in English does not confer the ability to detect it in Chinese without target-language supervision.

The defining contribution of this study is the validation of **Synthetic Data Augmentation** as a robust solution to the data scarcity bottleneck. We show that training on "noisy" machine-translated data allows the model to learn a high-quality "suppression mask" for subjective terminology. Our synthetic model achieved the optimal balance of style transfer and content preservation, outperforming zero-shot baselines by 17% in accuracy while maintaining higher semantic fidelity than translation-based methods.

These findings suggest a scalable path forward for Multilingual NLP: rather than waiting for expensive human-curated datasets, researchers can leverage high-volume synthetic corpora to bridge the gap for subtle semantic tasks in low-resource languages. Future work may explore applying this synthetic paradigm to other subjective tasks, such as politeness transfer or formality adaptation, in diverse linguistic contexts.

6. Acknowledgements

Author Contributions: All authors contributed equally to this work, including methodology design, experimental execution, data analysis, and report writing. The relative contribution breakdown is 33% – 33% – 33%.

Computing Resources: We gratefully acknowledge the support of ai.patrea.com for providing the NVIDIA RTX 4090 computing resources used to conduct the experiments in this research.

References

- [1] 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, Minneapolis, Minnesota, Jun 2019. Association for Computational Linguistics.
- [2] Ashwin D'Sa, Irina Illina, Dominique Fohr, and Dietrich Klakow. Joint labeling and bias neutralization in Wikipedia. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3465–3471, Punta Cana, Dominican Republic, Nov 2021. Association for Computational Linguistics.
- [3] Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742, 2020.
- [4] Reid Pryzant, Richard Diehl, Xuan Liu, Wei Feng, Claire Cardie, and Kai Yu. Automatically neutralizing subjective bias in text. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 480–489, Feb 2020.
- [5] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, Aug 2016. Association for Computational Linguistics.