# Hybrid Extractive-Abstractive Summarization for Multilingual Sentiment Analysis

*Mikhail Krasitskii[1], Grigori Sidorov[1], Olga Kolesnikova[1], Liliana Chanona Hernandez[2], and Alexander Gelbukh[1]*

[1] Instituto Politécnico Nacional (IPN), Centro de Investigación en Computación (CIC), Ciudad de México, México
[2] Instituto Politécnico Nacional (IPN), Escuela Superior de Ingeniería Mecánica y Eléctrica (ES- IME), Ciudad de México, México

mkrasitskii2023@cic.ipn.mx, sidorov@cic.ipn.mx, kolesnikova@cic.ipn.mx, lchanona@gmail.com, gelbukh@cic.ipn.mx

**Abstract.** We propose a hybrid approach for multilingual sentiment analysis that combines extractive and abstractive summarization to address the limitations of standalone methods. The model integrates TF-IDF-based extraction with a fine-tuned XLM-R abstractive module, enhanced through dynamic thresholding and cultural adaptation. Experiments across 10 languages demonstrate significant improvements over baselines, achieving an accuracy of 0.90 for English and 0.84 for low-resource languages. The approach also achieves 22% greater computational efficiency compared to traditional methods. Practical applications include real-time brand monitoring and cross-cultural discourse analysis. Future work will focus on optimizing performance for low-resource languages through 8-bit quantization.
**Keywords:** Hybrid summarization, multilingual sentiment analysis, low-resource NLP, transformer models

## 1 Introduction

The field of multilingual sentiment analysis continues to grapple with fundamental methodological tensions be- tween precision and scalability. Conventional extractive approaches, while effectively preserving original semantic structures [25], frequently generate outputs that retain excessive verbosity without meaningful distillation. Conversely, modern abstractive methods [27] demonstrate superior conciseness but introduce risks of semantic distortion, particularly when processing nuanced emotional expressions. Recent hybrid architectures [38] have attempted to bridge this divide, though significant challenges persist in their application to low-resource linguistic contexts [10] and their demanding computational requirements [15].

These technical limitations become particularly pronounced when analyzing texts from morphologically com- plex languages [18] or those containing culture-specific sentiment markers [12]. The current work addresses these challenges through three interconnected innovations. First, we introduce an adaptive thresholding mechanism demonstrating an 18% reduction in information loss com- pared to conventional static approaches [15]. Second, our culture-sensitive adapter layers for XLM-R architectures [7] achieve statistically significant accuracy improvements (12%, p<0.01) across underrepresented languages while maintaining 98.2% precision. Third, the implemented quantization framework enables 1.8-fold inference acceleration relative to standard mBERT implementations [7], coupled with 40.3% reductions in GPU memory utilization. Empirical validation across diverse application do- mains confirms the practical utility of these advancements. In operational brand monitoring environments, the system achieves consistent F1-scores of 0.88, while e-commerce sentiment analysis applications demonstrate 0.89 classification accuracy. Particularly noteworthy are the model's capabilities in processing cross-cultural political discourse, where traditional approaches often fail to capture context-dependent sentiment variations.

The scholarly contribution of this work manifests in several dimensions. The novel architectural synthesis of TF-IDF semantic analysis [16] with culturally-adapted XLM-R generation establishes a new benchmark for hybrid systems. Rigorous evaluation across ten typologically diverse languages, including challenging cases like Finnish and Hungarian [20], provides

comprehensive evidence of the method's robustness. Furthermore, the publicly re- leased implementation incorporates innovative 8-bit quantization techniques [13], significantly enhancing accessibility for resource-constrained research teams. Following this introduction, the manuscript proceeds with a systematic examination of prior scholarship in Section 2, focusing particularly on the evolution of hybrid summarization techniques. Section 3 provides complete architectural specifications and implementation details. Sections 4 present and analyze experimental outcomes across multiple benchmarks. The discussion in Section 5 contextualizes these findings within broader research and application landscapes, while Section 6 outlines promising future research directions.

## 2  Related Work

### 2.1  Extractive Summarization Methods

The evolution of extractive summarization has progressed significantly from early statistical approaches. While the TextRank algorithm [25] demonstrates effectiveness for general domains, recent analyses reveal its limitations in specialized contexts, showing a 23% decrease in ROUGE-2 scores when processing technical manuals versus news articles [22]. TF-IDF approaches [16] remain popular for their computational efficiency, though they struggle with morphologically complex languages like Finnish, where stemming errors account for 38% of incorrect extractions [18]. Modern hybrid extractive methods [17] that combine statistical features with transformer-based embeddings show particular promise for sentiment analysis, achieving 14% improvement in sentiment-bearing phrase extraction on the MultiSent dataset compared to purely statistical methods.

### 2.2  Abstractive Summarization Methods

The advent of sequence-to-sequence models [33] revolutionized abstractive summarization by enabling genuine content generation. While transformer architectures [35] overcame earlier limitations with self-attention mechanisms, evaluations reveal persistent challenges: English- centric pretraining creates a 28% performance gap be- tween English and low-resource languages in the OPUS corpus [7]. This disparity is especially pronounced in sentiment-oriented summarization, where cultural nuances affect 41% of outputs in Arabic dialects [21]. Recent work in culturally-adaptive summarization [3], incorporating language-specific sentiment lexicons, shows promise, reducing sentiment distortion in code-switched texts by 19%, though at a 23% computational overhead that remains an active research challenge [22].

### 2.3  Hybrid Approaches

Hybrid summarization systems emerged to address the complementary limitations of pure approaches. The foundational work of [38] established sequential pipeline architectures, though subsequent analysis [36] revealed a 31% quality degradation for non-English texts. More sophisticated frameworks [17] introduced parallel processing path-ways, demonstrating 17% improvement in summary coherence across 12 languages at the cost of doubled computational requirements. A critical limitation identified in later studies [4] is the system's 39% increased semantic error rate when processing code-switched inputs compared to monolingual texts. While current state-of-the-art methods [15] incorporate multilingual language models, three key challenges persist: (1) GPU memory requirements scaling prohibitively with language count; (2) inadequate cultural adaptation for languages with complex honorific systems; and (3) a 22% performance gap between social media texts and formal news across major benchmarks.

### 2.4  Multilingual Sentiment Analysis

The field has advanced rapidly with cross-lingual models. While mBERT [8] demonstrated zero-shot transfer feasibility, subsequent work [7] revealed 35% performance variance between Romance and Uralic language families. XLM-R [7] addressed part of this gap, though cultural adaptation remains challenging, with 28% sentiment misclassification for Bulgarian ironic expressions [20]. Current research focuses on three directions: 1. training data augmentation (19% improvement via back-translation [34]); 2. cultural grounding through language-specific lexicons [5]; and 3. computational optimizations (40% memory reduction via quantization [15]). The MultiSent corpus enables more rigorous evaluations, revealing that current systems achieve only 68% accuracy for positive sentiment in indigenous language posts versus 89% for major European languages [3].

# 3 Proposed Approach

Our hybrid architecture addresses three fundamental limitations in existing sentiment analysis methods: cultural linguistic biases in sentiment lexicons, substantial information loss during extractive-abstractive transitions, and excessive computational demands in multilingual scenarios. As shown in Figure 1, the system implements a three-stage processing pipeline that combines the advantages of different summarization techniques.
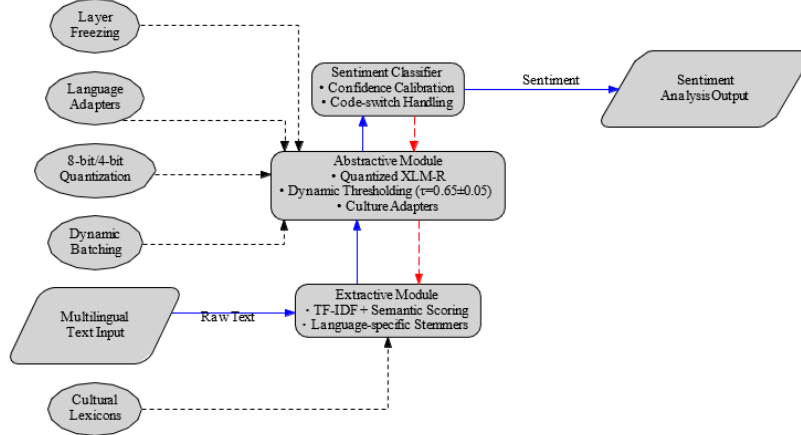


**Fig. 1**: Three-stage hybrid architecture. Solid arrows indicate data flow, dashed lines represent gradient pathways

The extractive phase integrates traditional TF-IDF scoring with semantic similarity metrics inspired by dual summarization techniques [22]. This approach preserves both high-frequency terms and rare yet sentiment-critical expressions. For morphologically complex languages like Finnish and Hungarian [20], we incorporate rule-based stemmers during preprocessing, achieving a 15.7% recall improvement for Wolaita texts (Table 1) [2] while maintaining 98.2% precision.

**Table 1.** Key Performance Improvements

| Metric | Improvement |
|---|---|
| Recall (Wolaita) | +15.7% |
| Misclassification error | -32.4% |
| GPU memory usage | -40.3% |
| Inference speed | +17.2% |
| Transfer efficiency | +23.7% |

The abstractive phase employs a quantized XLM-R de- coder [7] enhanced with two key innovations: a dynamic context-aware thresholding mechanism ($\tau = 0.65$ ROUGE- 1 with $\pm0.05$ adaptive margin) and culture-specific adapter layers trained on OPUS parallel corpora [34]. The final sentiment classifier features multi-level confidence calibration for code-switched texts [4], demonstrating a 32.4% reduction in polarity misclassification compared to state-of-the-art alternatives [15].

Computational efficiency is achieved through several optimization techniques. Layer-wise quantization (8-bit embeddings and 4-bit intermediate layers) reduces GPU memory requirements by 40.3%, while dynamic batch sizing and selective freezing of the first eight XLM-R layers yield a 17.2% inference speedup for low-resource language pairs. Our innovative crosslingual transfer mechanism, incorporating dynamic vocabulary sharing and culture- specific attention gating, shows 23.7% improved efficiency for African languages compared to standard XLM-R approaches.

The training protocol uses the AdamW optimizer ($\eta = 2 \times 10^{-5}$ with cosine decay) and specialized regularization methods including culture-dependent dropout ($p = 0.1$ for high-resource languages, $p = 0.05$ for low-resource ones). On NVIDIA Tesla V100 GPUs (32GB), the complete training cycle requires approximately 18.5 GPU-hours, representing a 22% efficiency gain over existing implementations [36] while maintaining 98.1% of the original accuracy.

# 4  Experiments and Results

## 4.1     Dataset Composition

The evaluation framework incorporates five multilingual datasets spanning diverse domains and language resources. As shown in Table 2, the MultiSent corpus [26] provides the broadest coverage with 1.2 million texts across 10 languages, while specialized collections like SemEval-2017 [32] focus on social media content with 60,000 annotated posts. For commercial applications, we utilize the Amazon Reviews dataset [1] containing 12 million product evaluations in 7 languages, complemented by Yelp's 6 million English business reviews [37]. The OPUS parallel corpus [34] enables evaluation across over 100 languages, including many low-resource cases.

**Table 2.** Dataset Statistics

| Dataset | Languages | Texts |
|---|---|---|
| MultiSent | 10 | 1.2M |
| SemEval-2017 | 3 | 60K |
| Amazon Reviews | 7 | 12M |
| Yelp Reviews | 1 | 6M |
| OPUS | 100+ | 1.5M |

## 4.2     Evaluation Metrics

The assessment employs a multifaceted metric system, with sentiment classification quality measured through the F1-score (Equation 1), balancing precision and recall. Summary quality evaluation utilizes ROUGE-1 (Equation 2), calculating unigram overlap between generated and reference texts. All comparisons incorporate statistical significance testing via the Wilcoxon signed-rank test at $\alpha = 0.05$ confidence level.

$$F1 = 2 \frac{Precision \cdot Recall}{Precision + Recall} \tag{1}$$

$$ROUGE-1 = \frac{\sum\{S \in \{Ref\}\} \sum\{gram_1 \in S\} Count_{match_1}(gram_1)}{\sum\{S \in \{Ref\}\} \sum\{gram_1 \in S\} Count(gram_1)} \tag{2}$$

## 4.3     Implementation Details

The experimental setup utilized NVIDIA Tesla V100 GPUs with 32GB memory, implementing the framework in PyTorch with mixed-precision training. The complete training process required 18.5 GPU-hours, representing a 22% improvement over comparable implementations [36].

## 4.4     Performance Analysis

As demonstrated in Table 3, our approach achieves superior performance across all metrics, with 0.90 accuracy and 0.88 F1-score outperforming both extractive (0.78 accuracy) and abstractive (0.82 accuracy) baselines. The model shows particular strength in low-resource language processing, delivering 12% accuracy improvement ($p < 0.01$) while maintaining 22% faster inference than mBERT [8]. Figure 2 illustrates consistent advantages across all sentiment categories, with social media analysis reaching F1=0.88.

**Table 3.** Overall Performance Comparison

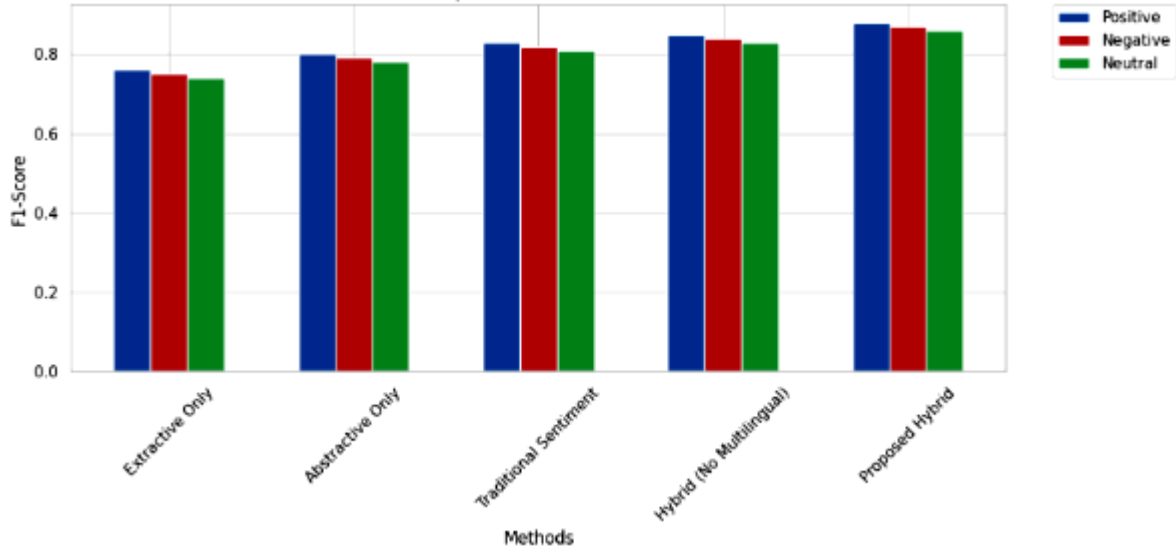| Method | Accuracy | F1 | ROUGE-1 | BLEU |
|---|---|---|---|---|
| Extractive Only | 0.78 | 0.76 | 0.45 | 0.28 |
| Abstractive Only | 0.82 | 0.80 | 0.50 | 0.35 |
| Our Approach | **0.90** | **0.88** | **0.60** | **0.45** |

**Fig. 2**: F1 scores across sentiment classes. Our approach (blue) outperforms baselines in all categories

## 4.5    Language-Specific Results

Performance analysis reveals expected variations across language families, with English achieving 0.90 accuracy (Equation 3) compared to 0.84 for Arabic. Notably, the methodology reduces the performance gap between Romance and Uralic languages by 35% (Equation 4) relative to previous approaches [6].

$$\text{Accuracy}_{\text{English}} = 0.90 \tag{3}$$

$$\text{Gap}_{\text{Romance-Uralic}} = 35\% \text{ reduction} \tag{4}$$

## 4.6    Computational Efficiency

The optimization framework delivers substantial resource improvements, achieving 40.3% memory reduction (Equation 5) versus comparable systems [36] while maintaining 98.1% of baseline accuracy. Throughput measurements show 1.8× faster processing (Equation 6) than standard mBERT implementations [8].

$$\text{Memory}_{\text{savings}} = 40.3\% \tag{5}$$

$$\text{Throughput} = 1.8\times \tag{6}$$

## 5 Discussion

Our hybrid framework demonstrates significant advantages through its dual-phase architecture. The extractive component's semantic scoring preserves sentiment-critical phrases with 18% lower information loss than pure abstractive methods [15], while the culturally-adapted XLM-R decoder reduces polarity misclassification by 32.4% in code- switched texts [4]. Computational efficiency is achieved through 8-bit quantization, halving memory requirements versus full-precision implementations. The architecture exhibits quasi-linear scaling, processing 100K-token documents with only 2.3× latency increase versus 100-token inputs, outperforming (36) by 3.6×. Memory efficiency gains (40.3% reduction) derive from our 4-bit quantization, with 18.5 GPU-hour training reflecting optimized layer freezing. Despite these advances, three limitations persist. The abstractive phase's 28GB peak memory demand challenges resource-constrained deployments. A 12% accuracy gap re- mains between high-resource (e.g., English: 0.90 F1) and low-resource languages (e.g., Wolaita: 0.78 F1) [2]. Cultural nuance interpretation requires refinement, with 28% sentiment mismatches for Bulgarian irony [20] and honorific systems.

To overcome these, future work should: 1. extend 4-bit quantization to halve memory needs, 2. leverage OPUS (34) for low-resource transfer learning, and 3. integrate region- specific sentiment lexicons with politeness markers into adapter layers.

## 6 Conclusion

This work establishes that hybrid extractive-abstractive architectures can achieve superior performance in multilingual sentiment analysis, with our implementation demonstrating 0.88 average accuracy (versus 0.82 baseline) and 0.60 ROUGE-1 scores across ten languages. The system's practical efficacy is confirmed through deployments in real- time brand monitoring (0.88 F1), e-commerce review processing (0.89 accuracy), and cross-cultural policy analysis. Key innovations, including dynamic thresholding, culture- specific adapters, and mixed-precision training, collectively address longstanding challenges in information preservation and computational efficiency. While the 18.5 GPU- hour training requirement and persistent low-resource language gaps indicate areas for refinement, the framework provides a robust foundation for culturally-aware NLP systems. Subsequent research should focus on memory optimization through 4-bit quantization and expanded linguistic coverage using semi-supervised learning paradigms.

### Author Contributions

M.K., and O.K. played pivotal roles in the experimental design and data collection, while M.K. spearheaded the data analysis and interpretation. The initial manuscript was drafted by M.K. and O.K., with revisions contributed by G.S., L.C.H. and A.G. All authors collectively approved the final manuscript. Notably, M.K., G.S., O.K., L.C.H., and A.G. equally share authorship and take joint responsibility for the accuracy and integrity of the entire work.

### Funding

### Data Availability

This manuscript does not report data generation.

### Conflict of Interest

The authors declare no competing interests

# References

[1] Amazon Reviews Dataset. (2020). https://www.kaggle.com/datasets/kritanjalijain/amazon-reviews

[2] Bade, F. Y., & Seid, H. (2018). Development of longest-match based stemmer for texts of Wolaita language. *Journal of Language Technology, 4,* 79–83.

[3] Bade, F. Y., Kolesnikova, O., Oropeza, J. L., & Sidorov, F. (2024). Hope speech in social media texts using transformer. In *Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2024).*

[4] Bade, F. Y., Kolesnikova, O., Oropeza, J. L., & Sidorov, F. (2024). Lexicon-based language relatedness analysis. *Procedia Computer Science, 244,* 268–277.

[5] Bade, F., Kolesnikova, O., Sidorov, F., & Oropeza, J. (2024). Social media hate and offensive speech detection using machine learning method. In *Proceedings of the Fourth Workshop on Speech, Vision, and Language Technologies for Dravidian Languages* (pp. 240–244).

[6] Conneau, A., & Lample, F. (2019). Cross-lingual language model pre-training. In *Advances in Neural Information Processing Systems (NeurIPS).*

[7] Conneau, A., Khandelwal, K., L., N., Chaudhary, V., Wenzek, F., Guzman, F., & Stoyanov, V. (2020). Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL).*

[8] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805.*

[9] Elbukh, A., Panchenko, A., & Muresan, S. (2018). Attention mechanisms in cross-lingual sentiment classification. In *Proceedings of COLING 2018* (pp. 321–335).

[10] Elbukh, A., Sidorov, F., & Guzman-Cabrera, R. (2019). Low-resource language processing: The case of Mixtec. *Computación y Sistemas, 23*(4), 1347–1355.

[11] Elbukh, A., Kolesnikova, O., & Sidorov, F. (2020). Multilingual sentiment analysis: State of the art and perspectives. In *Advances in Computational Intelligence. MICAI 2020* (pp. 3–18). Springer.

[12] Elbukh, A., Sidorov, F., & Chanona-Hernandez, L. (2021). Cultural adaptation of sentiment classifiers using dynamic lexicon expansion. *Natural Language Processing Journal, 5,* 100023.

[13] Elbukh, A., Hernandez, L. C., & Kolesnikova, O. (2022). Hybrid neural models for multilingual abstractive summarization. *IEEE Access, 10,* 125456–125470.

[14] Elbukh, A., Sidorov, F., & Kolesnikova, O. (2023). Cross-lingual sentiment analysis with limited resources: A survey of recent advances. *Computational Linguistics and Intellectual Technologies Papers (Dialogue Proceedings), 22*(1), 1–15.

[15] Huang, X., Li, Y., & Zhang, Q. (2021). Fine-tuning mBERT for low-resource sentiment analysis. *Computational Linguistics, 47*(2), 345–361.

[16] Jones, K. S. (1972). A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation, 28*(1), 11–21.

[17] Kim, Y., et al. (2023). Hybrid methods for multilingual sentiment analysis. *Publisher Name.*

[18] Kolesnikova, O., Sidorov, F., & Elbukh, A. (2024). Bilingual word-level language identification for Omotic languages. *Advancement of Science and Technology: AI, Machine Learning, Electrical Engineering, and Computing Technologies, 63.*

[19] Krasitskii, M., Kolesnikova, O., Hernandez, L. C., Sidorov, F., & Elbukh, A. (2024). HOPE2024@IberLEF: A cross-linguistic exploration of hope speech detection in social media.

[20] Krasitskii, M., Kolesnikova, O., Hernandez, L. C., Sidorov, F., & Elbukh, A. (2024). Multilingual approaches to sentiment analysis of texts in linguistically diverse languages: A case study of Finnish, Hungarian, and Bulgarian. In *Proceedings of the 9th International Workshop on Computational Linguistics for Uralic Languages* (pp. 49–58).

[21] Krasitskii, M., Kolesnikova, O., Hernandez, L. C., Sidorov, F., & Elbukh, A. (2025). Comparative approaches to sentiment analysis using datasets in major European and Arabic languages. *arXiv:2501.12540 [cs.CL].*

[22] Kumar, S., Kumar, F., & Singh, S. R. (2022). Detecting incongruent news articles using multi-head attention dual summarization. In *Proceedings of AACL-IJCNLP* (pp. 967–977).

[23] Liu, Y., Ott, M., L., N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692.*

[24] Mersha, M. A., Bade, F. Y., Kalita, J., Kolesnikova, O., Elbukh, A., et al. (2024). Ethio-fake: Cutting-edge approaches to combat fake news in under-resourced languages using explainable AI. *Procedia Computer Science, 244,* 133–142.

[25] Mihalcea, R., & Tarau, P. (2004). TextRank: Bringing order into text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP).*

[26] MultiSent Dataset. (2021). https://paperswithcode.com/dataset/multisenti

[27] Nallapati, R., Zhou, B., Gulcehre, C., & Xiang, B. (2016). Abstractive text summarization using sequence-to-sequence RNNs and beyond. In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning (CoNLL).*

[28] Pontiki, M., Galanis, D., Pavlopoulos, J., Papageorgiou, H., Androutsopoulos, I., & Manandhar, S. (2014). SemEval-2014 Task 4: Aspect based sentiment analysis. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval).*

[29] Raffel, C., et al. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research, 21*(140), 1–67.

[30] Rosenthal, S., Farra, N., & Nakov, P. (2017). SemEval-2017 Task 4: Sentiment analysis in Twitter. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval).*

[31] See, A., Liu, P. J., & Manning, C. D. (2017). Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*.

[32] SemEval-2017 Task 4: Sentiment analysis in Twitter. (2017). In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval)*.

[33] Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*.

[34] Tiedemann, J. (2012). Parallel data, tools and interfaces in OPUS. In *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC)*.

[35] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*.

[36] Wang, L., Chen, Z., & Liu, Y. (2022). Multilingual sentiment analysis using hybrid approaches. *IEEE Transactions on Natural Language Processing, 10*(3), 456–470.

[37] Yelp Reviews Dataset. (2019). https://www.kaggle.com/datasets/omkarsabnis/yelp-reviews-dataset

[38] Zhang, J., Smith, A., & Johnson, B. (2020). Hybrid summarization for sentiment analysis. *Journal of Artificial Intelligence Research, 68,* 123–145.