

A Comprehensive Survey on Abstractive Text Summarization with a Focus on Low Resource Languages and Challenges in Kannada Language

Puneeth R.^{1*}, J. Somasekar²

^{1,2}Department of Computer Science and Engineering, Faculty of Engineering and Technology, Jain (Deemed-to-be University), Bangalore, Karnataka, India

*Email: puneethvihaan@gmail.com¹

Abstract

The exponential growth of digital textual data has made automated text summarization a critical technology. Abstractive summarization, which generates novel human-like summaries using transformer-based models such as BERT, BART, and T5, has achieved strong results on high-resource languages, yet low-resource languages like Kannada remain severely underexplored. This survey reviews abstractive summarization methodologies from early neural architectures to advanced transformer frameworks, with special attention to knowledge-augmented models and multilingual pretraining. We critically examine challenges in Indic language summarization — including data scarcity, factual inconsistency, and evaluation limitations — and identify concrete research directions for Kannada: benchmark dataset creation, culturally informed modeling, cross-lingual transfer learning, and knowledge-aware summarization pipelines.

Keywords

Abstractive Text Summarization, Low-Resource Languages, Kannada NLP, Transformer Models, Multilingual Models, Data Scarcity.

Introduction

The unprecedented growth of digital content has created an environment of information overload across domains such as news media, scientific publications, and enterprise communication. Automated text summarization has emerged as a vital technology enabling efficient knowledge extraction [1]. Summarization approaches divide broadly into extractive methods, which select and concatenate key sentences, and abstractive methods, which interpret source semantics and generate novel fluent text. Abstractive approaches are closer to human-written summaries but require deeper semantic understanding [2].

Recent advances with Transformer architectures and Pre-trained Language Models (PLMs) notably BART and T5 have achieved state-of-the-art results, framing summarization as a

Submission: 13 March 2026; **Acceptance:** 16 June 2026; **Available online:** June 2026



Copyright: © 2026. All the authors listed in this paper. The distribution, reproduction, and any other usage of the content of this paper is permitted, with credit given to all the author(s) and copyright owner(s) in accordance to common academic practice. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license, as stated in the website: <https://creativecommons.org/licenses/by/4.0/>

sequence-to-sequence learning problem [3, 4]. However, these gains are disproportionately concentrated in English. Languages like Kannada, spoken by over 40 million people, face severe challenges: limited annotated corpora, few pretrained models, and insufficient linguistic tools [5]. Cross-lingual approaches often fail to capture Kannada's morphological nuances and syntax [6]. This survey addresses this gap by reviewing neural and transformer-based methods, analyzing the Indic language landscape, and proposing a research roadmap for Kannada summarization.

Evolution of Abstractive Text Summarization

• Task Formulation

The abstractive summarization task can be formally defined as generating a shortened sequence $S = v^1, v^2, \dots, v^m$ from a source document $D = w^1, w^2, \dots, w^n$, where $m \ll n$ and S is a fluent and informative condensation that is not merely a subset of D [4]. The goal is to learn a model parameterized by θ that maximizes the conditional probability of the summary given the equation 1:

$$P(S | D; \theta) = \prod_{t=1}^m P(v_t | v_{<t}, D; \theta) \quad (1)$$

where $v_{<t}$ denotes all previously generated tokens up to step t [7]. This formulation connects summarization to a conditional language modeling problem.

• From Rules to Neural Networks

Early systems relied on templates and fixed grammar rules, providing some control but generating unnatural summaries without deep semantic understanding especially problematic for agglutinative languages like Kannada. The introduction of sequence-to-sequence (seq2seq) neural architectures marked a paradigm shift: encoder-decoder models, augmented with attention mechanisms, learned summarization end-to-end from data, producing more accurate and coherent outputs [7]. However, these models required large labeled datasets, a bottleneck for low-resource languages.

• Transformer Revolution

Transformers replaced step-by-step recurrence with self-attention, enabling models to capture long-range dependencies in parallel. BART learns by reconstructing corrupted input, yielding deep understanding of sentence structure [3]. T5 frames every NLP task as text-to-text, enabling flexible fine-tuning across languages and domains [4]. More structured models like TP-Transformer incorporate tensor-product representations of grammatical roles, improving content selection and summary organization [8].

• Advanced Training Paradigms

To address training inefficiency and over-reliance on single reference summaries, curriculum learning introduces progressively complex examples, improving content selection [9]. Contrastive and multi-task learning enable models to distinguish high-quality from poor summaries, and to share knowledge across related tasks such as translation and paraphrasing [6]. The BRIO paradigm employs ranking-based fine-tuning against multiple candidate outputs, boosting performance in low-resource settings where reference diversity is limited [10].

Architectures and Strategies for Low-Resource Settings

The persistent challenges posed by low-resource languages such as limited annotated corpora, lack of pretrained models, and sparse linguistic tooling have catalyzed the development of specialized architectures and training paradigms aimed at maximizing learning efficiency under data scarcity.

- **Multilingual Pretrained Models**

For low-resource languages, multilingual pretrained models have become the primary strategy. mBART applies denoising autoencoding across many languages simultaneously, building a shared representation space that enables cross-lingual transfer [11]. mT5 extends the T5 framework to over 100 languages using a unified text-to-text format [12]. IndicBART and IndicT5 are specialized variants pre-trained on Indic language corpora, better capturing the grammar and vocabulary of languages like Kannada, Hindi, and Tamil [13].

- **Training Strategies**

Beyond architecture, specialized training strategies are critical. Curriculum learning starts with simpler documents and progressively introduces complex ones, stabilizing training on noisy low-resource data [9]. Contrastive learning teaches the model to distinguish accurate from inaccurate summaries, particularly valuable when annotated data is scarce [6]. Multi-task learning on related tasks summarization, translation, paraphrasing broadens the model’s language understanding and improves performance without requiring large task-specific datasets [2].

Table 1: Impact of Data Resources on Summarization Quality

Language	Resource Level	Summarization Pairs	% English Data	of	ROUGE-L Score	Hallucination Rate (%)
English	Very High	2,000,000+	100%		40–45	5–10
Hindi	Medium-High	~400,000	20%		30–35	12–18
Tamil	Low-Medium	~50,000	2.5%		20–25	20–28
Kannada	Very Low	< 20,000	< 1%		15–20	25–35+

Scores are illustrative based on common findings in low-resource NLP research.

Challenges in Low-Resource Language Summarization

One of the biggest challenges in building summarization systems for low-resource languages like Kannada is the lack of large, high-quality training datasets. Unlike English, which has many well-annotated corpora, Kannada has very few summarization datasets available. This shortage makes

it difficult to train neural models effectively and also affects how well we can evaluate their performance. Without enough data, models struggle to learn meaningful patterns, and researchers have limited tools to measure how good the summaries really are.

Another major difficulty comes from the language itself. Kannada, like many low-resource languages, has complex grammar. It uses agglutination, where words are formed by joining smaller parts, and has flexible word order, meaning the same sentence can be written in different ways. These features make it harder for models to understand and summarize text correctly. Most summarization models are designed for English and don't handle these unique language structures well, which leads to poor results when applied to Kannada.

Evaluating summarization systems in low-resource languages is also a challenge. Most researchers use ROUGE scores to measure summary quality, but these metrics were created for English and don't always work well for other languages. ROUGE mainly looks at word overlap and doesn't capture deeper meaning or cultural relevance. On top of that, tools for preprocessing and annotating Kannada text are still underdeveloped. This makes it hard to compare different models fairly or improve them over time [16].

Another issue is factual inconsistency, also known as hallucination. When models are trained on small datasets, they often generate summaries that sound correct but contain false or misleading information. This happens because the model hasn't seen enough examples to learn how to stay true to the original content. In low-resource settings like Kannada, where task-specific data is limited, this problem becomes even worse [10].

Finally, there are technical limitations. In many regions where low-resource languages are spoken, access to powerful computers and internet connections is limited. This makes it hard to train and run large neural models. To overcome this, researchers need to build efficient, lightweight models that can work well even with fewer resources. Designing such models is difficult but necessary to make summarization tools more accessible and practical in these areas.

Evaluation Metrics and Datasets

- **Evaluation Challenges**

Evaluating abstractive summarization systems for low-resource languages presents a distinct set of challenges that extend beyond those encountered in high-resource settings. These obstacles can be categorized into three primary areas: the limitations of automated metrics, the scarcity of reference materials, and the complexities of cultural and contextual appropriateness.

- **Limitations of Automated Metrics**

One of the major challenges in evaluating summarization systems for low-resource languages is the use of standard metrics like ROUGE. ROUGE works by comparing the words or phrases in a generated summary to those in a reference summary, mainly by counting overlaps. While this method is widely used and works reasonably well for English, it doesn't always reflect the true quality of summaries in other languages.

Languages like Kannada have different grammar rules, flexible word order, and rich word formations. Because of this, a summary might be accurate and well-written even if it doesn't match the reference summary word-for-word. ROUGE may give such summaries a low score simply because the words are arranged differently or use different forms, even though the meaning is

correct. This makes ROUGE less reliable for evaluating summaries in morphologically rich or syntactically flexible languages.

Moreover, ROUGE doesn't measure important aspects like fluency, coherence, or cultural appropriateness. It focuses only on surface-level matches, ignoring whether the summary reads naturally or captures the deeper meaning of the original text. As a result, relying only on ROUGE can lead to misleading evaluations, especially in low-resource language settings where linguistic diversity is high and reference summaries are limited [9].

- **Scarcity of High-Quality Reference Summaries**

In low-resource languages like Kannada, one of the biggest problems in evaluating summarization systems is the lack of good reference summaries. A reference summary is a human-written version of what a good summary should look like. These are used to compare and judge how well a machine-generated summary performs.

However, for many low-resource languages, there simply aren't enough of these high-quality reference summaries available. This makes it hard to train models properly and even harder to evaluate them. Without strong reference examples, it's difficult to tell whether a model is doing a good job or not. The few available references may not cover different topics, writing styles, or domains, which can lead to unfair or unreliable evaluation results.

This scarcity also limits progress in research. If models can't be tested against trusted benchmarks, it becomes challenging to improve them or compare different approaches. In short, without enough reference summaries, evaluation becomes weak, and the development of better summarization systems slows down.

- **Cultural and Contextual Appropriateness**

When evaluating summaries, especially in low-resource languages like Kannada, it's important to look beyond grammar and sentence structure. A summary might be linguistically correct but still feel out of place if it doesn't match the cultural tone or context of the original content. For example, certain expressions or references that work well in English might sound awkward, confusing, or even insensitive in Kannada.

This is because language carries more than just words it reflects local customs, values, and ways of thinking. A good summary should respect these cultural elements. It should use phrases and examples that make sense to native speakers and fit naturally within the social and cultural setting of the language. Without this awareness, even a technically accurate summary can fail to connect with its audience.

Evaluating cultural appropriateness is difficult to automate. Standard metrics like ROUGE don't measure whether a summary feels culturally respectful or contextually relevant. That's why human judgment is often needed to assess whether a summary truly fits the language and its users. This adds an extra layer of complexity to evaluation, but it's essential for building summarization systems that are not only accurate but also meaningful and inclusive.

- **Existing Datasets**

Scarcity of quality data hampers Kannada summarization model training and evaluation. The resources can be divided into three categories. Multilingual datasets (e.g. XL-Sum [11]) allow transfer learning but can include faulty or culturally inappropriate translations. Monolingual datasets are better for understanding, but Kannada datasets are few, small and topical. Synthetic

datasets, created via translation or using pretrained models, increase data quantity but contain mistakes, unnatural expressions and biases.

Summarization Research Across Indic Languages

Figure 1 shows significant differences in the summarization dataset sizes across Indic languages. Hindi leads the pack at more than 400,000 pairs of articles and summaries, thanks to significant resource creation and research attention. This is compared to Kannada, which has less than 20,000 pairs, which presents a significant challenge to building effective NLP and summarization systems. This contributes to a digital divide in AI. The disparity highlights the importance of targeted dataset generation, transfer learning, and community annotation efforts to promote balanced NLP advancement across Indic languages.

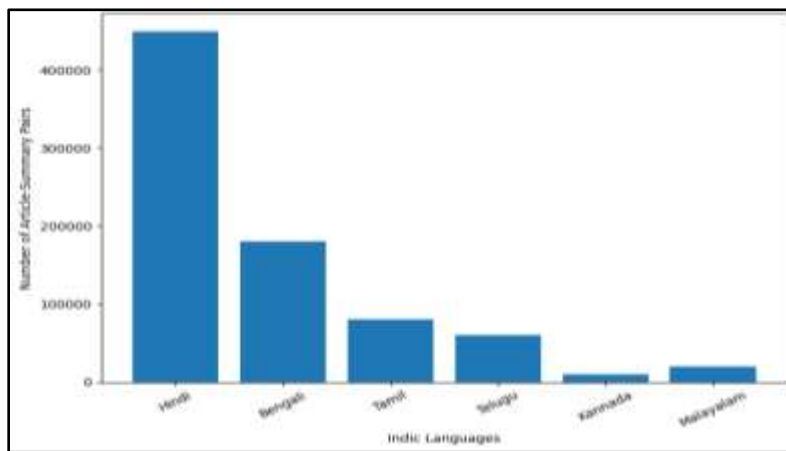


Figure 1. Dataset Availability for Indic Summarization

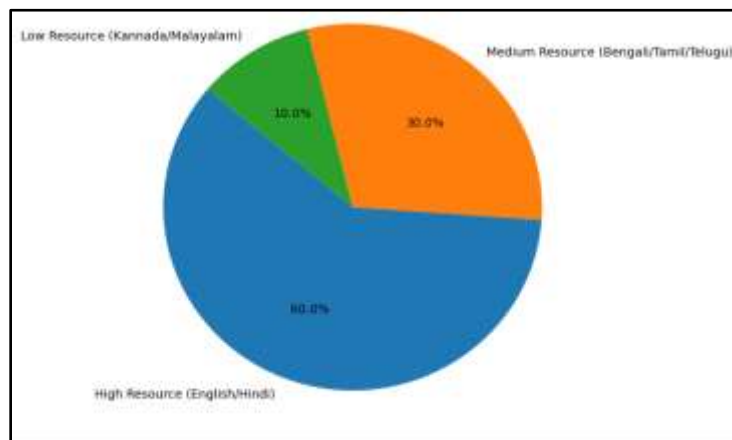


Figure 2. Summarization Resources by Language Type

The imbalance becomes clearer when languages are grouped into high-resource, medium-resource, and low-resource categories. Figure 2 demonstrates illustrates the imbalance in summarization resources when languages are grouped into three categories: high-resource, medium-resource, and

low-resource. High-resource languages mainly English and Hindi account for approximately 60% of the available datasets, tools, and model support. Medium-resource languages, such as Bengali and Tamil, make up around 30%, benefiting from moderate research attention and some publicly available corpora. In contrast, low-resource languages like Kannada and Malayalam represent only about 10% of the total resources. These languages suffer from limited annotated datasets, fewer pre-trained models, and minimal infrastructure for evaluation and deployment.

Beyond the differences seen between languages at a single point in time, the way summarization resources have grown over the years also reveals deep inequality. As shown in Figure 3, between 2010 and 2024, English has seen a massive increase in summarization datasets reaching millions of article-summary pairs. Hindi has also exhibit consistent growth, with new datasets being added gradually over time. These expansions reflect strong research interest, funding, and infrastructure support for high-resource languages.

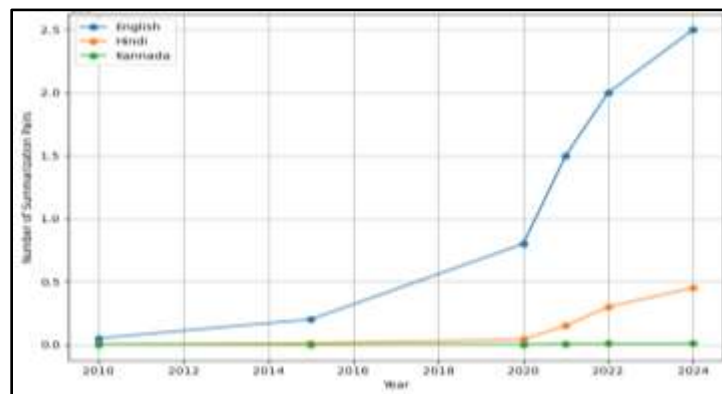


Figure 3. Growth of Summarization Datasets Over Time

Multilingual PLMs such as mT5 and IndicBART have become the primary tools for Indic NLP tasks [18, 19]. These models leverage cross-lingual transfer learning from high-resource languages. Still, they again suffer from the "curse of multilinguality," where performance per language is diluted compared to dedicated monolingual models [20]. Furthermore, their performance on Dravidian languages like Kannada typically lags behind that on Indo-Aryan languages like Hindi, due to linguistic differences and data imbalances within the pre-training corpora [21]. Here is a known problem and an major attempt to solve it. However, here is data proving the problem persists in a specific way. Consequently, our current best tools are fundamentally limited by this same data problem.

Progress in Indic NLP has been uneven. Hindi benefits from the most resources, including XL-Sum [17] and HindSum datasets. A key milestone was XL-Sum's release, enabling multilingual model evaluation across languages including Hindi, Bengali, and Gujarati. Multilingual PLMs such as mT5 [12] and IndicBART [13] have become primary tools, but they suffer from the "curse of multilinguality," where performance per language is diluted versus dedicated monolingual models. Dravidian languages like Kannada consistently lag behind Indo-Aryan languages due to greater linguistic distance and data imbalance in pretraining corpora [18].

Table 2: Comparative Overview of Text Summarization for Select Indic Languages

Language	Resource Level	Key Resources & Models	Primary Focus	Predominant Research Gaps
Hindi (Indo-Aryan)	Medium-High	XL-Sum [17], HindSum, mT5/mBART	Both paradigms	Domain adaptation; factual consistency metrics
Bengali	Medium	XL-Sum [17], BengaliBERT	Extractive	Large-scale abstractive datasets; evaluation beyond ROUGE
Tamil (Dravidian)	Low-Medium	Samanantar [19], TamilT5	Extractive	Abstractive benchmarks; Dravidian cross-lingual transfer
Telugu (Dravidian)	Low-Medium	Samanantar [19], TeluguBERT	Extractive	Curated datasets; complex syntax handling
Kannada (Dravidian)	Very Low	Samanantar [19], KannadaBERT [20], Kavi-NLP [21]	Almost exclusively Extractive [22]	No abstractive dataset; no evaluation benchmark; virtually no abstractive research
Malayalam (Dravidian)	Low	IndicCorp [23], pilot studies [24]	Extractive	Basic corpora creation; morphological tooling

• **Kannada: Challenges and Resources**

Kannada’s agglutinative morphology and syntactic complexity present distinct challenges. Words consist of multiple fused morphemes, producing a combinatorial explosion of surface forms that make tokenization and sequence modeling significantly harder [18]. Flexible word order, extensive case marking, and honorific forms further complicate seq2seq learning. The most acute challenge is severe data scarcity: Kannada lacks a large-scale dedicated summarization dataset, and available multilingual corpora such as IndicCorp [23] and Samanantar [19] are not annotated for summarization. While KannadaBERT [20] provides a text understanding baseline, no large-scale seq2seq model comparable to T5 or BART exists for Kannada. Research has been limited to extractive summarization [22] or pilot multilingual experiments [13]. Cultural context idiomatic expressions, locally relevant phrasing remains underexplored and is critical for producing summaries that resonate with native speakers.

Future Research Directions

Addressing Kannada's summarization gap requires a multi-pronged approach. The highest priority is the creation of a high-quality benchmark dataset. This can be achieved through manual curation, leveraging news headlines, or using strong multilingual teachers to generate silver-standard annotations [28]. Efficient learning paradigms must be explored, including advanced fine-tuning for multilingual models, parameter-efficient methods like LoRA, and knowledge-enhanced

frameworks that integrate culturally relevant sources [29, 12, 30]. Finally, developing standardized evaluation protocols beyond ROUGE, including human evaluation frameworks and Kannada-specific metrics for factual consistency, is essential for tracking meaningful progress [9].

To achieve these goals, focused research is needed in the following strategic areas:

- **Dataset Development and Curation**

The severe lack of data is the primary bottleneck in advancing Kannada summarization systems. Addressing this gap requires a multi-pronged strategy that balances scalability, linguistic diversity, and annotation quality. Future efforts must focus on:

- **Crowdsourced Annotation:** Developing sustainable frameworks and incentives for the crowdsourced creation of high-quality (document, summary) pairs.
- **Automatic Dataset Generation:** Innovating in techniques like zero-shot cross-lingual transfer and back-translation to automatically generate silver-standard training data from existing resources [28].
- **Cross-domain Datasets:** Creating datasets that span multiple domains (e.g., news, legal, healthcare) to improve model generalization and robustness for real-world applications.

Model Architecture and Learning Paradigms

Models need to be built or adapted to handle Kannada's unique language features effectively. The script includes complex characters that require careful preprocessing, and the grammar is rich, with long, detailed word forms. Kannada also allows flexible word order, which can confuse models trained on more rigid sentence structures. Without these adjustments, summarization systems may produce inaccurate or unnatural results. Designing language-aware models is key to improving performance and relevance.

- **Language-specific Adaptations:** Developing architectural components, such as specialized sub-word tokenizers or morphological encoders, explicitly tailored for agglutinative languages like Kannada.
- **Efficient Models:** Creating lightweight and efficient models (e.g., via pruning, quantization) suitable for deployment on devices in resource-constrained environments [29].
- **Continual Learning:** Developing frameworks that allow models to continuously learn and adapt from new data as it becomes available, without catastrophic forgetting of previous knowledge.

Transfer and Multi-Task Learning

Maximizing knowledge transfer from high-resource languages is crucial.

- **Optimal Source Language Selection:** Research into linguistic distance metrics to determine the optimal source languages (e.g., other Dravidian languages like Tamil) for transfer learning to Kannada.
- **Meta-Learning Approaches:** Developing meta-learning and few-shot learning frameworks that can quickly adapt to Kannada with minimal task-specific data.
- **Cross-Script Transfer:** Investigating novel methods for effective transfer learning across different writing systems and scripts, which remains a significant challenge.

➤ **Evaluation and Metrics**

Robust evaluation is needed to reliably measure progress.

- **Language-specific Metrics:** Moving beyond ROUGE by developing automated evaluation metrics that account for Kannada's specific linguistic properties, such as morphological richness and flexible word order.
- **Advanced Automated Evaluation:** Creating more sophisticated metrics that leverage language-specific embeddings and semantic representations to better correlate with human judgment.
- **Cultural Appropriateness Metrics:** Developing evaluation frameworks and metrics that assess the cultural and contextual appropriateness of generated summaries, ensuring they are relevant and accurate for native speakers.

Conclusion

This survey has outlined the remarkable progress in abstractive summarization driven by transformer models, while highlighting the significant disparity that exists for low-resource languages like Kannada. The challenges of data scarcity, model generalization, and factual consistency are profound but not insurmountable. The recent development of multilingual resources and advanced modeling techniques presents a unique opportunity to leapfrog these hurdles. Future work must be directed towards collaborative creation of foundational datasets, development of knowledge-aware models tailored for Kannada's linguistic and cultural context, and the establishment of robust evaluation benchmarks. Additionally, efforts must prioritize culturally sensitive design and inclusive annotation practices to ensure relevance and fairness. Strengthening community-driven research and open-access infrastructure will be key to sustaining long-term progress. Addressing these gaps is not only a technical necessity but a step toward linguistic equity. With sustained commitment, Kannada can become a model for inclusive NLP innovation across low-resource languages.

Acknowledgement

The authors would like to express their sincere gratitude to the Department of Computer Science and Engineering, Jain (Deemed-to-be University), Bengaluru, for providing the necessary academic environment, research facilities, and support to carry out this work. The authors are especially thankful to their research guide, faculty members, and staff for their continuous encouragement, valuable guidance, and constructive suggestions throughout the research. The authors also extend their appreciation to the anonymous reviewers and the journal editors for their insightful comments and valuable feedback, which helped improve the quality of this manuscript.

References

- A. Joshi et al., “The State and Fate of Linguistic Diversity and Inclusion in the NLP World,” in Proc. ACL, 2020, <https://doi.org/10.18653/v1/2020.acl-main.560>
- A. Kunchukuttan et al., “IndicBART: A Pre-trained Model for Natural Language Generation of Indic Languages,” <https://doi.org/10.48550/arXiv.2109.02903>
- A. Nenkova and K. McKeown, “Automatic Summarization,” Foundations and Trends in Information Retrieval, vol. 5, nos. 2–3, pp. 103–233, 2011, <https://doi.org/10.1561/15000000015>
- A. Ramesh and K. Shashirekha, “Kannada Text Summarization: A Study on Low-Resource Indic Language,” International Journal of Computer Applications, vol. 183, no. 19, pp. 20–26, 2021.
- A. Rizzello, “An Investigation on the Extractive Summarization of Kannada Text,” in Computational Intelligence and Applications. Singapore: Springer, 2023. DOI: https://doi.org/10.1007/978-981-99-1410-4_26
- C. Raffel et al., “Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer,” Journal of Machine Learning Research, vol. 21, no. 140, pp. 1–67, 2020. <https://doi.org/10.48550/arXiv.1910.10683>
- C.-Y. Lin, “ROUGE: A Package for Automatic Evaluation of Summaries,” in Proc. ACL Workshop on Text Summarization Branches Out, 2004. <https://aclanthology.org/W04-1013/>
- D. Kakwani et al., “IndicNLPsuite: Monolingual Corpora, Evaluation Benchmarks and Pre-trained Multilingual Language Models for Indian Languages,” in Findings of EMNLP, 2020, <https://doi.org/10.18653/v1/2020.findings-emnlp.445>
- D. S. Pankaj, “Challenges in Creating Text Summarization Models in Malayalam: A Study,” in Proc. International Conference on Innovative Computing and Cloud Computing, 2023. <https://doi.org/10.1109/ICCC57789.2023.10165363>
- E. J. Hu et al., “LoRA: Low-Rank Adaptation of Large Language Models,” in Proc. ICLR, 2022, <https://doi.org/10.48550/arXiv.2106.09685>
- G. Ramesh et al., “Samanantar: The Largest Publicly Available Parallel Corpora Collection for 11 Indic Languages,” <https://doi.org/10.48550/arXiv.2104.05596>
- H. P. Luhn, “The Automatic Creation of Literature Abstracts,” IBM Journal of Research and Development, vol. 2, no. 2, pp. 159–165, 1958, <https://doi.org/10.1147/RD.22.0159>
- I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to Sequence Learning with Neural Networks,” in Advances in Neural Information Processing Systems, vol. 27, 2014. <https://doi.org/10.48550/arXiv.1409.3215>
- J. Lin et al., “How to Train Your Dragon: Data Augmentation and Optimization for Low-Resource Summarization,” in Proc. COLING, 2022. <https://doi.org/10.48550/arXiv.2302.07452>
- Kannada BERT, Hugging Face Model Hub, 2021.
- L. Xue et al., “mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer,” in Proc. NAACL, 2021, <https://doi.org/10.18653/v1/2021.naacl-main.41>
- M. Fabbri et al., “SummEval: Re-evaluating Summarization Evaluation,” Transactions of the Association for Computational Linguistics, vol. 9, pp. 391–409, 2021, https://doi.org/10.1162/tacl_a_00373

- M. Lewis et al., “BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension,” in Proc. ACL, 2020, <https://doi.org/10.48550/arXiv.1910.13461>
- P. Dhakal and D. S. Baral, “Abstractive Summarization of Low-resourced Nepali Language using Multilingual Transformers,” <https://doi.org/10.48550/arXiv.2409.19566>
- P. K. S. Shivaraddi et al., “Kavi-Kannada Natural Language Processing System,” International Journal of Advanced Research in Science, Communication and Technology, vol. 2, no. 1, 2022, <https://doi.org/10.48175/IJARSCT-5837>
- P. Lewis et al., “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” in Advances in Neural Information Processing Systems, vol. 33, 2020, <https://doi.org/10.48550/arXiv.2005.11401>
- R. Kumar et al., “How Robust are Pre-trained Models to Domain Shift for Low-Resource Tasks? A Case Study in Summarization,” in Proc. Workshop on TextGraphs at ACL, 2023.
- S. Sotudeh et al., “Curriculum-Guided Abstractive Summarization,” <https://doi.org/10.48550/arXiv.2302.01342>
- T. Hasan et al., “XL-Sum: Large-Scale Multilingual Abstractive Summarization for 44 Languages,” in Findings of ACL-IJCNLP, 2021, <https://doi.org/10.18653/v1/2021.findings-acl.413>
- T. Liu et al., “Multilingual Denoising Pre-training for Neural Machine Translation,” Transactions of the Association for Computational Linguistics, vol. 8, pp. 726–742, 2020, https://doi.org/10.1162/tacl_a_00343
- W. Kryściński et al., “Evaluating the Factual Consistency of Abstractive Text Summarization,” in Proc. EMNLP, 2020, <https://doi.org/10.18653/v1/2020.emnlp-main.750>
- Y. Liu et al., “Abstractive Text Summarization Using the BRIO Training Paradigm,” <https://doi.org/10.48550/arXiv.2305.13696>