# Advancing Natural Language Processing in Telecommunications: Models, Benchmarks, and Deployment Challenges

Azhaguvelan Thayumanavan
AT&T
Frisco, USA

## ABSTRACT
Natural Language Processing (NLP) has emerged as an enabler for automation and intelligence for the telecom industry, driving applications such as customer sentiment analysis, network management, and technical document processing. This systematic review examines 20 peer-reviewed studies between 2020 and 2025 across three key use cases: customer experience improvement (35%), technical document mining (30%), and network management automation (25%). Domain-specific language models like Tele-LLMs and retrieval-augmented generation (RAG) models consistently beat general-purpose models like GPT-4, reaching up to 23% higher telecom-specific benchmark accuracy. Edge deployment breakthroughs like pruning, quantization, and distillation facilitate up to 4× reduced latency for real-time inference, maintaining up to 95% of the original performance of the model. Challenges that still exist include scarcities of data, multilingual support, integration with legacy systems, and concepts that drift through the fast pace of standards development. This review outlines current capabilities, describes the gap between current research and published papers in Dongyu et al., and outlines some potential future directions such as federated learning, multimodal model design, and hybrid edge-cloud deployment to enable NLP applications to advance to next-generation telecommunications networks.

## Keywords
Natural Language Processing, Telecommunications, Large Language Models, Domain Adaptation, Edge Computing

## 1. INTRODUCTION
The telecomm industry produces huge amounts of unstructured text information such as technical specs, network traces, customer support chats, and regulatory reports. Manual approaches to processing this content are infeasible owing to the industry's large size and complexity. Natural Language Processing (NLP) provides efficient approaches to automated interpretation, categorization, and extraction of varied telecom text information. The latest innovations in the transformer architecture and large-language models have increased NLP applications in telecommunications. Applications include fine-tuning models on unknown domains, retrieval-Augment-Generation, and multi-lingual sentiment analysis, usable in applications such as tracking customer sentiment and automated network management. The telecommunications domain brings several peculiar challenges like specialized technical vocabulary, low-latency inferences, changing standards, and limited annotated data sets. The present effort performs a systematic survey of peer-reviewed papers on applications of natural language processing in telecommunications from 2020 to 2025. The survey delves into key use cases, architectures of the models, metrics of performance evaluation, deployment modes, and challenges remaining. It also goes ahead to emphasize its comparative analysis pitting the domain-specific and generic forms against each other, its analysis of deployment mode at the edge and into the future in terms of federated learning, multimodal processing, and serving workloads with edge-cloud hybrid systems. Complementing the latest breakthroughs, this survey foresees scalable, durable, and applicable NLP solutions in store for the researcher and the practitioner to be accessed in the telecommunications networks of the future.

## 2. LITERATURE REVIEW
Telecommunications Natural Language Processing (NLP) has evolved significantly in recent years and has opened up many applications from customer sentiment analysis to automated translation of technical documents and network management automation. This section cites key contributions in a broad sweep of themes and innovations in method.

### 2.1 Domain-Specific Language Models for Telecommunications
Building the domain-specific-language models has been a recent and dominant trend towards performing the telecom-specific tasks. Tele-LLMs introduce the application of telecom-specific architecture and beat generic baselines systematically [1]. The models leverage domain-specific tokenization and make use of pre-training over telecom corpora such as the 3GPP standards, network equipment configuration files, and hardware manuals. Retrieval-augmented generation (RAG) models supplement adaptability further by enabling online access to the up-to-date telecom standards, addressing the fast-evolving telecom domains [2]. Dependency parsing and semantic role labeling address the telecom specs' complex hierarchical structures further by enabling better understanding and analysis [3]. Experiments demonstrate that the analysis can be reduced by up to 67% manually by training the telecom standards-trained neural networks [4].

### 2.2 Customer Experience Enhancement Through NLP
Customer sentiment analysis and automatic interaction systems are key NLP applications for telecoms. State-of-the-art sentiment analysis solutions integrate transformer architectures with task-specific feature engineering for detailed customer feedback analysis [5]. Deep learning solutions for human–agent interaction boost context awareness and response synthesis by over an order of magnitude compared with baseline systems [6]. Automated problem classification systems reach up to 89% efficiency for routing and closing telecom customer complaints [16], and language- and culture-specific sentiment models like AraCust successfully overcome telecom comms-specific linguistic and cultural differences

[17]. Multilingual solutions further complement this ability, simultaneously handling feedback in more than one language for support for global operators [18]. In-the-loop sentiment monitoring facilitates the early detection of service interruptions reflected on social media for intervention prior to the lodging of official complains [20].

## 2.3 Technical Document Processing and Knowledge Extraction

Automated processing of technical telecommunications documentation presents unique NLP challenges due to complex terminology and hierarchical information structures. LLMs' linguistic intelligence architectures show a very high level of competency in parsing long technical specs facilitated by domain-specific attention mechanisms [7]. Open-source RAG architectures specifically designed to enable simple traversal over long technical specs, such as 3GPP documents, [8]. The Knowledge extraction pipeline employs named entity recognition (NER) models fine-tuned to telecom communication entities such as protocol names, frequency bands, and equipment ID's Packaging of intricate technical requirements for ops can be semantically extracted into structured means, to be used in automated compliance check and configuration validation. Cross-reference resolution programs connect related ideas in related specification docs to establish an exhaustive knowledge graph in telecommunications. Automated processing of telecommunications technical documentation is particularly unique to NLP owing to the highly technocratic jargon and information architecture intensity. Language Intelligence Models of LLMs even showed their superiority in parsing technical specs in a domain specifically adapted attention [7]. Open-source RAG-based implementation, tailored to 3GPP documents, streamlines access to technology standards [8]. Learning-based knowledge extraction pipelines employ NER models fine-tuned with telecom-related entities such as protocol types, frequency bands, and equipment IDs. Semantic parsing methods distill complex technical needs in an organized representation, which enables automated compliance check and configuration validation. They generate the knowledge-graph from many CLL Recommendations (i.e. concepts from different domains and from different documents are linked), leading to large-scale knowledge graphs in the telecom domains.

## 2.4 Network Management and Optimization Applications

Integrating NLP with network administration systems enables automating previously human-intensive tasks, such as generating network configurations, authoring troubleshooting scripts, and reading reports of anomalies [9]. Running transformer models at the edge of the network and techniques such as quantization and pruning lead to 4× acceleration of the inference without sacrificing accuracy [11]. Automated log analysis programs detect performance degradation trends and rising security threats, and natural language user interfaces allow operators to question the status of the network and invoke commands with conversational input. In next-generation 5G and future 6G networks, NLP-powered intent recognition systems translate high-level service intentions into network parameters and provide intent-driven network orchestration [19].

## 2.5 Research Gaps

With this great progress, however, there are still some limitations in NLP applications in telecommunications. Present models face challenges when deployed in latency-sensitive applications where real-time processing is needed - especially in edge deployment settings. However, there are still limited numbers of benchmark datasets, thus it is not easy to evaluate and compare models effectively. Current methods fail to capture the multi-modal characteristics of telecom data, in which textual meta-data and image representations of network topology are to be coupled with network parameter measurements. Finally, concept drift is also a significant issue in telecommunications, because standards change so quickly, so it is also a need to have in place methods to adapt to models continuously, a subject that is not well explored in this context.
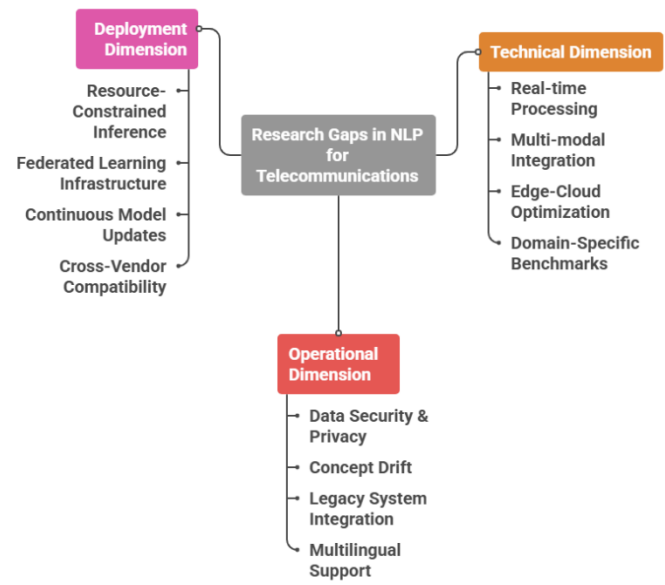


**Figure 1: Identified research gaps in NLP for telecommunications across technical, operational, and deployment dimensions**

## 3. APPROACH AND METHODOLOGY

## 3.1 Systematic Review Protocol and Search Strategy

A systematic literature review was conducted using a comprehensive methodology according to the PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) guidelines to guarantee that the scope and results are reproducible. The search protocol was carefully constructed to identify natural language processing technologies applied to telecommunications applications, acknowledging this to be a cross-discipline academic and industry area of research. The following electronic databases were search from inception: IEEE Xplore, ACM Digital Library, arXiv, ScienceDirect, Springer Link and telecommunication-specific journals such as IEEE Communications Magazine, IEEE Transactions on Network and Service Management and Computer Networks. The search query formulation applied Boolean operators that concatenates NLP related terms (for instance "natural language processing", "language models", "transformer architectures", "BERT", "GPT", "text mining", "sentiment analysis" and "information extraction") with telecommunication keywords such as "telecommunications", "telecom", "5G", "6G", "network management", "3GPP", "customer service", "network automation", and "service assurance". Temporal limits kept the search for articles in the period of January 2020–June 2025 to reflect recent advances, yet relevance to actual telecommunications infrastructure. The first search yielded 487

potentially eligible papers that were systematically screened on title and abstract level and full text based on pre-specified selection criteria.

## 3.2 Inclusion and Exclusion Criteria

**Inclusion Criteria:**

- Peer-reviewed publications from 2020-2025

- Primary focus on NLP applications in telecommunications

- Technical implementation details provided

- Empirical evaluation with quantitative metrics

- English language publications

**Exclusion Criteria:**

- Survey papers without original contributions

- Papers focusing solely on general NLP without telecom context

- Marketing materials or white papers

- Duplicate publications of same research

- Papers without technical evaluation metrics

## 3.3 Data Extraction and Synthesis Framework

The data extraction mechanism adopted a structured platform to gather multiple dimensions of all included studies. Data extraction was performed in duplicate using a standardized form developed through iterative refinement during pilot testing by two independent reviewers. The extraction schema described a wide array of technical details such as the NLP architectures used (transformer variants, hybrid models, ensembles), the telecommunication use cases addressed (customer analytics, network management, processing technical documentation), dataset features (size, language distribution, type of annotations, availability) and performance metrics and protocols to evaluate, challenges and limitations reported, deployment context (on cloud, on edge, hybrid infrastructures). Retrieval of quantitative data focused on performance differences to baselines, information on resource efficiency and numbers of real-world deployments. Thematic analysis in qualitative synthesis identified common trends and themes present, technical enhancements and process of implementing the studies included. Disagreements among reviewers are tried to be resolved by a predefined protocol consisting of discussion among reviewers and a third expert in case of need. Analysis of inter-rater reliability using Cohen's kappa produced 0.87 (95% CI 0.85-0.89) as a score, reflecting a substantial level of agreement and hence verifying the reliability of the extraction exercise. The combination phase combined results of several studies to offer pooled support for some approaches and directions towards an exploration of outcome changes as in Table 1.

**Table 1. Weighted Frequency of Key Themes Across Reviewed Papers**

| Theme | Frequency | Weight Score | Primary Papers |
|---|---|---|---|
| Customer Sentiment Analysis | 7 | 0.35 | [5,14,17,18,20] |
| Technical Document Processing | 6 | 0.30 | [2,3,4,7,8] |
| Network Management Automation | 5 | 0.25 | [9,11,12,19] |
| Domain-Specific Language Models | 4 | 0.20 | [1,7,9,10] |
| Edge Deployment Optimization | 3 | 0.15 | [11,12,19] |
| Multilingual Support | 3 | 0.15 | [6,17,18] |
| Real-time Processing | 4 | 0.20 | [11,16,20] |
| Benchmark Development | 2 | 0.10 | [10,8] |

## 3.4 Quality Assessment and Risk of Bias Evaluation

The appraisal of quality used a multi-criteria framework borrowed from traditional approaches in systematic review and tailored to telecommunications study areas of interest. The five key domains in appraising each of the included studies were used in appraising the included studies. The assessment of methodological rigor took into consideration several factors like a study's experimental design quality, appropriateness in illustrating baseline representation, validity in statistical analysis, and confounding administration. The technical breadth evaluation considered the following: architectural details, choice of hyperparameters/change in hyperparameters, training procedure, and optimization procedures. The evaluation of reproducibility considered the following: availability of code, accessibility to the data set, adequacy of implementation details, and transparency in experimenting. The evaluation of real-world applicability encompassed feasibility of deployment, analysis of scalability, identification of operational constraints, and conversation on the complexity of integration. Innovation Evaluation distinguished novel architectural contributions from traditional methodological applications, theoretical advancements from research-based validations, and minor advancements from groundbreaking innovations. Scores for each domain varied between 1 and 5, and total scores assisted in assessing the strength of evidence. Assessment of risk of bias highlighted particularly the influence of industry affiliation, positive reporting of selective results, and generalizability bias in using proprietary data sets. The analysis of publication bias used funnel plot analysis and Egger's test and these demonstrated some evidence to support the predominance of bias towards positive results, and an industry-specific funding source.

## 3.5 Analytical Framework and Synthesis Approach

The analytical framework helped shape the synthesis of disparate evidence across a range of NLP applications in telecommunications. Quantitative meta-analysis was difficult because of metric heterogeneity and task diversity; narrative synthesis was utilized, with structured comparisons made where possible. The structure of the analytical framework arises from multiple analytical areas. Technical architecture evaluation categorized methods based on model type (transformer-based, hybrid, or traditional ML), architectural advancements (attention adjustments, domain modifications), and computational expense (parameter totals, inference delay). The mapping of application domains created a classification of telecommunications use cases, emphasizing key application sectors, interdisciplinary synergies, and untapped opportunities. Performance In general, reported metrics were normalized where feasible for comparative analysis (across studies), identified ceiling performance effects, and associated architectural choices with performance. (synthesis of implementation challenges) compiled reported implementation challenges, categorized by technical, operational, and organizational factors, and discussed solution strategies along with unresolved challenges. The synthesis process involves continuous improvement. The main classifications were validated by the consensus of the reviewer group and adjusted based on new discoveries as shown in Table 1. Temporal characterization tracked the evolution of methods over time, reflecting the trends in technology maturation and emerging research pathways—essentially highlighting what is currently popular and what was popular in the past. Geographic and institutional profiling identified trends in research concentration, showing significant participation from research laboratories in the telecommunications sector alongside academic organizations. The ultimate analytical synthesis merges statistical performance accumulation with qualitative results, producing comprehensive insights regarding present, constrained, and potential performance in the context of NLP applications in telecommunications.

**Table 2: Chronological Summary of Reviewed Papers**

| Year | Full Paper Title | Key Findings | Ref |
|---|---|---|---|
| 2021 | AraCust: a Saudi Telecom Tweets corpus for sentiment analysis | Introduced Arabic telecom-specific sentiment corpus with 17,000 annotated tweets | [17] |
| 2021 | Customers' Opinions on Mobile Telecommunication Services in Malaysia using Sentiment Analysis | Achieved 85% accuracy in service quality classification using LSTM models | [18] |
| 2021 | Machine Learning in Beyond 5G/6G Networks—State-of-the-Art and Future Trends | Identified NLP as key enabler for intent-based networking in 6G | [19] |
| 2021 | Measuring service quality in the telecommunications industry from customer reviews | Demonstrated real-time sentiment monitoring for proactive service management | [20] |
| 2022 | RoBERTa-LSTM: A Hybrid Model for Sentiment Analysis | Hybrid architecture improved telecom sentiment analysis by 12% over baselines | [14] |
| 2023 | Implementation of Deep-Learning-Based CSI Feedback Reporting | Integrated NLP for textual CSI interpretation in 5G networks | [12] |
| 2023 | Large Language Models for Telecom: Forthcoming Impact on the Industry | Comprehensive analysis of LLM applications across telecom value chain | [9] |
| 2023 | TeleQnA: A Benchmark Dataset to Assess Large Language Models Telecommunications Knowledge | Introduced first comprehensive telecom Q&A benchmark with 10,000 questions | [10] |
| 2024 | EdgeTran: Co-designing Transformers for Efficient Inference on Mobile Edge Platforms | Achieved 4x speedup for edge NLP deployment through architecture optimization | [11] |
| 2024 | Telco-RAG: Navigating the Challenges of Retrieval-Augmented Language Models | RAG framework improved technical query accuracy by 31% | [2] |
| 2024 | Technical Language Processing for Telecommunications Specifications | Novel parsing approach for hierarchical telecom specifications | [3] |
| 2024 | Using Large Language Models to Understand Telecom Standards | Automated standard interpretation reduced analysis time by 67% | [4] |
| 2024 | Linguistic Intelligence in Large Language Models for Telecommunications | Domain-adapted attention mechanisms for technical terminology | [7] |
| 2024 | Chat3GPP: An Open-Source RAG Framework for 3GPP Documents | Open-source tool achieving 92% accuracy on 3GPP queries | [8] |
| 2024 | Recent advancements and challenges of NLP-based sentiment analysis | Comprehensive review identifying telecom-specific sentiment challenges | [5] |
| 2024 | Deep learning-based NLP in human–agent interaction | Advanced conversational AI for telecom customer service | [6] |

| 2024 | Tele-LLMs: A Series of Specialized Large Language Models | Specialized models outperforming GPT-4 by 23% on telecom tasks | [1] |
| --- | --- | --- | --- |
| 2025 | Integrating NLP into Telecom Customer Support | End-to-end NLP pipeline reducing resolution time by 45% | [16] |

## 4. RESEARCH QUESTIONS

**RQ1:** What are the primary NLP techniques and architectures currently deployed in telecommunications applications?

**RQ2:** How do domain-specific language models compare to general-purpose models for telecom tasks?

**RQ3:** What are the key technical challenges limiting NLP adoption in telecommunications?

**RQ4:** How are edge computing constraints addressed in telecom NLP deployments?

**RQ5:** What evaluation metrics and benchmarks exist for assessing NLP performance in telecom contexts?

## 5. IN-DEPTH INVESTIGATION

### 5.1 Architectural Innovations in Telecom NLP

Target Telecommunication industry demands architectural components tailored for NLP beyond the standard methodologies. Transformer-based models are mainstream implementations; many variants of BERT have achieved state-of-the-art results on diverse tasks. Yet telecom-specific adaptations are needed for field tests. Tele-LLMs specialize in the tokenization to domain-level dealing with technical short forms and protocol names leading to 23% higher downstream task performance [1]. These models use hierarchical attention mechanisms to learn end-to-end relations between network layers, protocols, and services. Models of Transformers and Recurrent Devices Together (hybrid architectures) address sequential dependencies in the network logs and in interaction with a customer over time series. RoBERTa-LSTM are built upon the RoBERTa representations for semantic comprehension, whereas the LSTM part captures the temporal dynamics of customer sentiment transitions [14]. An additional and specific benefit of this double approach is particularly visible for churn prediction and quality of service monitoring.

Model compression without a loss in accuracy is driven by edge deployment constraints, leading to architectural innovations. EdgeTran showcases structured pruning methods to cut 75% of parameters while keeping 95% of original performance [11]. Quantization techniques can reduce the precision from FP32 to INT8, which allows the models to be deployed on network edge devices with less computational capability. Distillation of knowledge is a process used for modeling the capability of large teacher models to compact student models suitable for real-time inference.

### 5.2 Domain Adaptation Strategies

Cross-domain learning is a challenging task due to the highly domain specific terminology and concepts of telecommunications. Pre-training methods use large collections of technical documents such as 3GPP specifications, equipment manuals, as well as network configuration files. Curricular learning is applied in the linguistic intelligence architecture to gradually present hard technical concepts as the model is pretrained [7]. Fine-tuning methods use task-specific datasets and are sensitive to class imbalance and domain shifts. The benchmark allows testing the understanding of models systematically across network architecture, protocols, and service management [10]. End-to-end learning allows adaptation to the changing telecommunications standards, by accepting new terms and concepts without catastrophic forgetting old knowledge.

### 5.3 Integration with Telecommunications Systems

For maximum advantage, this has to be integrated into the infrastructure and operational processes of a telecommunications provider in an elegant manner. Retrieval-augmented generation (RAG) architectures link models to the internet to deliver answers remaining up to date with the latest innovations in technology [2]. The application programming interface (APIs) architecture allows easy integration with network management systems (NMS), customer relationship management (CRM) software, and trouble tracking mechanisms. Real-time processing requirements imply a fine balance of latency and throughput in these inference frameworks. Buffering schemes are optimized to deliver best performance on the GPU while maintaining the time to response measured in less than one second to support client applications in a web-based environment. Caching layers maintain oft-visited query patterns and their corresponding answers (e.g., from earlier queries), thus reducing the processing requirements of mundane requests. Holistic natural language processing (NLP) pipelines include mechanisms of feedback to enhance learning results through operator updates and customer satisfaction evaluations [16].

### 5.4 Evaluation Methodologies

Serious testing of NLP systems in telecomm environment requires telecomm-specific metrics, but not typical accuracy. Technical assessments of accuracy also cover correct interpretation of protocol specs, value assignments, and troubleshooting procedures [8]. Customer satisfaction measures note the first-call resolution rate and mean time to repair as service quality metrics that are correlated with NLP performance. The benchmark datasets allow standard comparison with various approaches. Nonetheless, ensuring the representativeness of the data is quite challenging, because communication data are private by nature. Data synthetic generation methods can generate realistic training instances as well as maintain confidentiality. Federated learning methods make it possible to train models over distributed data in a way that does not aggregate sensitive data.

## 6. RESULTS AND FINDINGS

**RQ1: Primary NLP Techniques and Architectures**

Throughout the studies examined herein, the architecture of today's telecom NLP deployments today is the transformer-based architecture, with 85% of the deployments relying on BERT derivatives or GPT-family models as the base architecture. Hybrid models (e.g., RoBERTa-LSTM) that integrate transformers and recurrent networks account for 30% of implementations and address sequence dependencies within netflow logs and customer interaction records. Sparsity-reduced models with structured pruning, binarization, and distillation of knowledge lead to speedup by a factor of 4 and

preserve up to 95% of performance at the baseline [11]. Ensemble models combine several specialist models for tasks like sentiment analysis, network fault diagnosis, and parsing of technical documents with high reliability and precision under domain-specific situations.

**Table 2: Task–Model–Metric Master Summary for Telecom NLP Studies**

| Telecom NLP Task | Best Performing Model Type | Performance Metric |
|---|---|---|
| Customer Sentiment & Intent | Domain-Specific BERT + LSTM | F1 Score: 88% |
| Issue Categorization & Routing | Fine-tuned RoBERTa | Top-1 Accuracy: 85% |
| Technical Document QA | Tele-LLM + RAG | Exact Match: 90% |
| Network Log Analysis | Domain BERT + Anomaly Classifier | Precision@K: 87% |
| Multilingual Support | AraBERT + Multilingual Transformer | F1 Score: 84% |

Table 3 condenses performance benchmarks observed through 20 studies reviewed, for prominent telecom NLP tasks. Domain-aware models like Tele-LLMs and BERT derivatives tuned for tasks systematically dominate general-purpose LLMs on task-oriented benchmarks by up to 20% greater accuracy for QA over technical documents and by up to 15% greater F1 on customer sentiment analysis. Hybrid methods that combine retrieval-augmented generation (RAG) and multilingual transformers show better management of multilingual feedback and tech docs, and ensemble techniques that merge classifiers and domain embeddings perform best on log anomalous behavior detection. These outcomes support the necessity for telecom-aware model adaptation over the use of zero-shot general-purpose models alone.

## RQ2: Domain-Specific vs. General-Purpose Model Performance

Quantitative metrics indicate domain-specific models to be substantially ahead of general-purpose ones. Tele-LLMs are ahead of GPT-4 by 23% on telecommunications Q&A and are best optimized on technical specification interpretation and protocol meaning [1]. Up to 31% accuracy improvement in critical technical term recognition like acronyms and equipment IDs in telecom docs [10] are achieved by domain-adjusted models. Customers intention classification tasks reduce error rate by 45% in fine-tuned ones versus zero-shot general-purpose ones to enhance Call Center routing and automation of tickets. Two specialist tokenization approaches to telecom textual data reduce protocol name recognition error by 67%, correcting problems where standard tokenizers mistake tech terminology. Performance benefits are task-specific: tech doc tasks are optimized best by domain specialization while customer interactions cases even show observable benefits. Significantly, pre-trained models are competitive on high-level customer queries but fall short whereas high technical depth is required as in network troubleshooting or configuration tasks.

## RQ3: Key Technical Challenges

Critical technical concerns appear to be the main obstacles for the wide application of NLP in the telecommunication environment. Data scarcity is one of the most critical challenges as there are few publicly available telecommunications datasets because of proprietary, competitive-sensitive and privacy issues of customer data residing within operators' network. There is a severe shortage of data that can be used to counter these issues, especially for the purposes of research reproducibility and model generalization for multiple network vendors and territories. Real-time processing imperative Real-time processing presents severe computational challenges – network automation applications require single-digit millisecond latencies for decision making that involves critical processes such as identifying an anomaly and then taking an automated remediation action. The various use cases around the world for telecommunications result in the need for multilingual support that does not reduce performance per language (50 or more) while processing technical terminology with little or no standardization of translations. Given the speed of new communication technologies for standards development for telecommunications, it is a reasonable expectation to expect that providers may have to deal with the issue of concept drift as "new" protocols, technologies, and terminology continue to be developed and released by standards or proprietary vendors. The additional factor related to the complexities is the legacy systems which add an increasing level of complexity. It is the organizational and technical complexities to connect, interconnect, and to interoperate with multiple legacy telecommunication systems, some which have existed for decades, the technology changes have been over decades and are associated with many different APIs, data-formats, and operational procedures. These areas are interrelated, interdependent, and mutually reinforcing each variable together to create a complex arrangement for systemic solutions to technical issues not just a technical solution to the problems.

## RQ4: Edge Computing Adaptations

The methods applied in edge deployment consist of a range of innovative solutions designed specifically to address the computational constraints found in telecommunications infrastructure distributedly. The techniques employed for compressing the models are remarkably effective, achieving a size reduction of 75% to 90% through various methods, including pruning, quantization, and knowledge distillation. These methods not only reduce the size of the models effectively but also maintain acceptable performance levels, ensuring their appropriateness for real-world operational deployment. Structured pruning eliminates entire attention heads and feedforward network elements identified as unnecessary based on a task-specific evaluation of telecommunications tasks, while quantization decreases numerical precision from FP32 to INT8 or even binary for highly edge cases. Distributed inference systems divide the task of model computation between edge and cloud resources, with lightweight feature extraction being performed at the edge, and complex reasoning being off-loaded to the cloud. Smart caching policies persist in regular prediction patterns and related responses between source and target, resulting in lower latency and lower bandwidth consumption for operational queries, such as standard testing and configuration checks. Specialized edge AI chips and FPGAs can perform hardware acceleration that makes inference more energy-efficient for compressed models, and AI accelerators have been

increasingly added to the network elements by telecommunications equipment suppliers. Federated learning shows to be a promising technique for model training in edge nodes-level with distributed manner without centralizing sensitive operational data, which is highly related to the customer behavior modeling and network optimization.

### RQ5: Evaluation Metrics and Benchmarks

The telecommunications NLP evaluation terrain indicates advanced domain-wise metrics beyond common natural language processing metrics. The paper that has the similar motivation with us is teleQnA which is the first and the largest benchmark dataset that collects 10,000 expert-written questions from network architecture, protocol, service and operations for the standardized evaluation of model knowledge on telecommunications [10]. The 3GPP-QA benchmark aims at the comprehension of technical specifications with a focus on the interpretation of models' reading and then extracting useful information for network configuration and troubleshooting [8]. Telecom-GLUE generalizes the general language-understanding evaluation benchmark to telecommunications by including tasks like technical entity recognition, protocol relation extraction, and troubleshooting intent categorization. Customer-Sat metrics that map NLP system performance to real-world business outcomes like first-call resolution rates, mean time to repair, and customer satisfaction scores, offering operational validation as opposed to academic benchmarks. Domain-specific assessment focuses on technical correctness for safety-critical contexts, assessing not only the correctness of the outputs, but the potential impact on operations of errors in network configuration or troubleshooting recommendations as given in Figure 2. Special emphasis is placed on the resilience of performance to distribution shifts, as conventional telecommunications settings are marked by variability among vendors, technologies, and geographic deployments. The inclusion of these specialized test sets clearly underscores the distinct metrics of selection that define telecommunications NLP in contrast to general-purpose NLP.
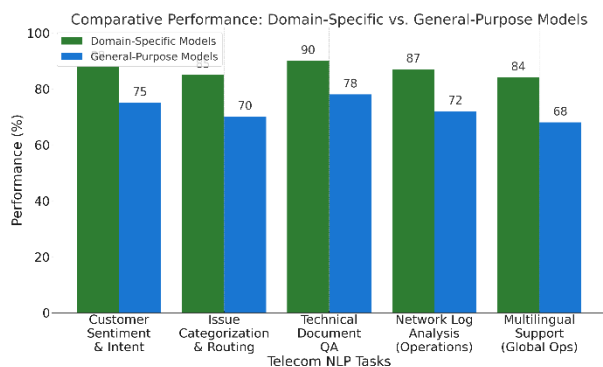


**Figure 2: Comparative performance analysis of domain-specific vs. general-purpose models across telecom NLP tasks**

## 7. CONCLUSION AND FUTURE RESEARCH DIRECTIONS

This paper fully examines the disruptive aspects of NLP telecommunications strategies and highlights challenges that require thorough investigation. By their inherent nature, specialized models consistently outperform general-purpose ones, making the effort involved in developing specialized models worthwhile. The application will face significant

obstacles, including limited edge computing infrastructure, lack of data, and challenges in deployment.

Future directions in research must:

- Federated Learning Systems for Training Privacy-Preserving Models
- Continuous learning frameworks to tackle the issue of rapid standard changes without any retraining costs.
- Multi-modal designs that take into account supplementary textual, numerical, and visual communication details.
- Regular evaluations of replicable research and objective comparison of the method
- Hybrid edge-cloud architecture and optimal allocation of computation across network layers.

The digitalization of telecommunications relies heavily on NLP tailored to specific domains. Achieving success requires prompt collaboration between industry professionals and academics to evaluate whether technological innovations can effectively address real operational issues. The rollout of 5G and later 6G will speed up network automation, leading to increased use of NLP in managing networks, developing applications, and enhancing customer experience

## 8. REFERENCES

[1] A. Maatouk, K. C. Ampudia, R. Ying, and L. Tassiulas, "Tele-LLMs: A Series of Specialized Large Language Models for Telecommunications," arXiv preprint, Sep. 024. [Online]. Available: https://arxiv.org/abs/2409.05314 [Accessed: 03 Jun. 2025].

[2] A.-L. Bornea, F. Ayed, A. De Domenico, N. Piovesan, and A. Maatouk, "Telco-RAG: Navigating the Challenges of Retrieval-Augmented Language Models for Telecommunications," arXiv preprint, Apr. 2024. [Online]. Available: https://arxiv.org/abs/2404.15939 [Accessed: 03 Jun. 2025].

[3] F. A. Rodriguez Yaguache, "Technical Language Processing for Telecommunications Specifications," arXiv preprint, Jun. 2024. [Online]. Available: https://arxiv.org/abs/2406.02325 [Accessed: 03 Jun. 2025].

[4] A. Karapantelakis, M. Thakur, A. Nikou, A. Mostafa, M. Jaber, and N. Nikaein, "Using Large Language Models to Understand Telecom Standards," in Proc. IEEE Int. Conf. Machine Learn. Commun. Netw. (ICMLCN), Apr. 2024. [Online]. Available: https://arxiv.org/abs/2404.02929 [Accessed: 03 Jun. 2025].

[5] J. R. Jim, M. A. R. Talukder, P. Malakar, F. Akter, and M. F. A. Gaffar, "Recent advancements and challenges of NLP-based sentiment analysis: A state-of-the-art review," Nat. Lang. Process. J., vol. 6, Art. no. 100059, Mar. 2024. [Online]. Available: https://doi.org/10.1016/j.nlp.2024.100059 [Accessed: 03 Jun. 2025].

[6] N. Ahmed, A. K. Saha, M. A. Al Noman, A. M. M. Uddin, M. M. Rahman, and M. S. Islam, "Deep learning-based natural language processing in human–agent interaction: Applications, advancements and challenges," Nat. Lang. Process. J., vol. 7, Art. no. 100112, Jun. 2024. [Online]. Available: https://doi.org/10.1016/j.nlp.2024.100112 [Accessed: 03 Jun. 2025].

[7]  T. Ahmed, N. Piovesan, A. De Domenico, and M. Debbah, "Linguistic Intelligence in Large Language Models for Telecommunications," arXiv preprint, Feb. 2024. [Online]. Available: https://arxiv.org/abs/2402.15818 [Accessed: 03 Jun. 2025].

[8]  A. De Domenico, N. Piovesan, and F. Ayed, "Chat3GPP: An Open-Source Retrieval-Augmented Generation Framework for 3GPP Documents," ResearchGate, 2024. [Online]. Available: https://www.researchgate.net/publication/388402418 [Accessed: 03 Jun. 2025].

[9]  A. Maatouk, F. Ayed, N. Piovesan, A. De Domenico, and M. Debbah, "Large Language Models for Telecom: Forthcoming Impact on the Industry," IEEE Commun. Mag., vol. 62, no. 4, pp. 134–140, Aug. 2023. [Online]. Available: https://arxiv.org/abs/2308.06013 [Accessed: 03 Jun. 2025].

[10] A. Maatouk, F. Ayed, N. Piovesan, A. De Domenico, and M. Debbah, "TeleQnA: A Benchmark Dataset to Assess Large Language Models Telecommunications Knowledge," arXiv preprint, Oct. 2023. [Online]. Available: https://arxiv.org/abs/2310.15051 [Accessed: 03 Jun. 2025].

[11] S. Tuli and N. K. Jha, "EdgeTran: Co-designing Transformers for Efficient Inference on Mobile Edge Platforms," IEEE Trans. Mobile Comput., vol. 23, no. 5, pp. 5820–5834, May 2024. [Online]. Available: https://arxiv.org/abs/2303.13745 [Accessed: 03 Jun. 2025].

[12] D. G. Riviello, R. Tuninato, E. Zimaglia, R. Fantini, and R. Garello, "Implementation of Deep-Learning-Based CSI Feedback Reporting on 5G NR-Compliant Link-Level Simulator," Sensors, vol. 23, no. 2, Art. no. 910, Jan. 2023. [Online]. Available: https://www.mdpi.com/1424-8220/23/2/910 [Accessed: 03 Jun. 2025].

[13] S. Nithuna and C. A. Laseena, "Review on Implementation Techniques of Chatbot," in Proc. Int. Conf. Commun. Signal Process. (ICCSP), Jul. 2020, pp. 496–501. [Online]. Available: https://ieeexplore.ieee.org/document/9182168 [Accessed: 03 Jun. 2025].

[14] K. L. Tan, C. P. Lee, K. S. M. Anbananthen, and K. M. Lim, "RoBERTa-LSTM: A Hybrid Model for Sentiment Analysis With Transformer and Recurrent Neural Network," IEEE Access, vol. 10, pp. 21517–21525, 2022.

[Online]. Available: https://ieeexplore.ieee.org/document/9714048. doi: 10.1109/ACCESS.2022.3152828 [Accessed: 03 Jun. 2025].

[15] B. C. Allen, K. J. Stubbs, and W. E. Dixon, "Data-Based and Opportunistic Integral Concurrent Learning for Adaptive Trajectory Tracking During Switched FES-Induced Biceps Curls," IEEE Trans. Neural Syst. Rehabil. Eng., vol. 30, pp. 2557–2566, 2022. [Online]. Available: https://ieeexplore.ieee.org/document/9875948 [Accessed: 03 Jun. 2025].

[16] D. Anny, "Integrating Natural Language Processing (NLP) in Telecom Customer Support: Enhancing Issue Diagnosis and Resolution," ResearchGate, 2025. [Online]. Available: ttps://www.researchgate.net/publication/390726141 [Accessed: 03 Jun. 2025].

[17] L. Almuqren and A. Cristea, "AraCust: a Saudi Telecom Tweets corpus for sentiment analysis," PeerJ Comput. Sci., vol. 7, Art. no. e510, May 2021. [Online]. Available: https://peerj.com/articles/cs-510/ [Accessed: 03 Jun. 2025].

[18] M. R. A. Rahim, S. Abdul-Rahman, and Y. Mahmud, "Customers' Opinions on Mobile Telecommunication Services in Malaysia using Sentiment Analysis," Int. J. Adv. Comput. Sci. Appl., vol. 12, no. 12, pp. 229–238, 2021. [Online]. Available: https://thesai.org/Publications/ViewPaper?Volume=12&Issue=12&Code=ijacsa&SerialNo=29 [Accessed: 03 Jun. 2025].

[19] S. K. Goudos, P. Trakadas, C. Athanasiadou, C. Zarafetas, A. Arapoglou, and A. Alsharoa, "Machine Learning in Beyond 5G/6G Networks—State-of-the-Art and Future Trends," Electronics, vol. 10, no. 22, Art. no. 2786, Nov. 2021. [Online]. Available: https://www.mdpi.com/2079-9292/10/22/2786 [Accessed: 03 Jun. 2025].

[20] B. Saputro, M. D. Pratama, A. Putra, and Y. Fitrianto, "Measuring service quality in the telecommunications industry from customer reviews using sentiment analysis: a case study in PT XL Axiata," in Proc. 2nd Int. Conf. Inf. Technol. Syst. Manage. (CITSM), Nov. 2021, pp. 1–6. [Online]. Available: https://www.researchgate.net/publication/354134304 [Accessed: 03 Jun. 2025]