

Abstractive Summarization of Spoken Language: A Comparative Evaluation of BART and T5 on Podcast and Conversational Speech Transcripts

Ashish Joshi

University School of Automation and Robotics
Guru Gobind Singh Indraprastha University
New Delhi, India

ABSTRACT

The exponential growth of long-form audio content, particularly podcasts and lectures, creates an urgent need for effective summarization systems capable of condensing hours of speech into concise, coherent summaries. This study presents a comprehensive comparative evaluation of two transformer-based architectures—BART and T5—for abstractive summarization of spoken language transcripts. Unlike prior work that relies on written dialogue datasets, the author fine-tunes and evaluates both models on three speech-specific datasets: PodcastSum (12,345 podcast episodes), How2 (12,987 instructional videos), and the AMI Meeting Corpus (137 hours of meetings). A multi-faceted evaluation framework is employed, combining automated metrics (ROUGE, BLEU, BERTScore, METEOR) with human judgments across five quality dimensions (coherence, fluency, factual consistency, conciseness, and speaker attribution). Statistical significance testing confirms observed differences, and qualitative analysis reveals model-specific strengths and failure patterns. Results demonstrate that BART significantly outperforms T5 across all automated metrics ($p < 0.01$) and receives higher human ratings for factual consistency and structural cohesion. However, T5 generates more lexically diverse summaries and better handles extended dialogue contexts. Complementary strengths are identified that suggest hybrid approaches may be beneficial. To support reproducibility, the evaluation framework and human-annotated test samples are released. The findings provide actionable guidance for deploying summarization systems in real-world speech applications.

General Terms

Natural Language Processing, Summarization

Keywords

Abstractive Summarization, Spoken Language Understanding, BART, T5, Transformer Models, Podcast Transcription, Human Evaluation

1. INTRODUCTION

The podcasting medium has experienced explosive growth, with over 5 million active podcasts and 70 million episodes available as

of 2026 [26]. Similarly, educational lectures, corporate meetings, and webinars generate vast repositories of spoken content that resist efficient navigation and information retrieval. Unlike written documents, speech recordings lack structural markers such as headings, paragraphs, or bullet points, making it difficult for listeners to extract key insights without investing substantial time.

Automatic summarization of spoken content offers a promising solution, transforming hours of audio into concise textual summaries that capture essential information. However, spoken language presents unique challenges that distinguish it from written text: disfluencies (filled pauses, repetitions), fragmented syntax, speaker overlap, prosodic cues, and the absence of punctuation all complicate the summarization task [31, 20]. Furthermore, podcasts and meetings often involve multiple speakers with distinct conversational roles, requiring models to track speaker identity and attribution.

The evolution of summarization techniques has progressed from extractive methods, which select and reassemble existing sentences, to abstractive approaches that generate novel phrasings. Transformer-based architectures have revolutionized abstractive summarization, with BART [16] and T5 [28] emerging as leading models for text generation tasks. However, their application to spoken language summarization remains underexplored, particularly regarding comparative performance on authentic speech transcripts.

Recent advances in speech summarization have explored end-to-end approaches that directly process audio [13, 3], yet cascaded systems combining automatic speech recognition (ASR) with text summarization remain the dominant paradigm in production deployments [25]. Understanding how summarization models perform on ASR-transcribed speech is therefore critical for practical applications.

This study addresses three research questions:

- (1) **RQ1:** How do BART and T5 compare in summarizing authentic spoken language transcripts across diverse domains (podcasts, instructional videos, meetings)?
- (2) **RQ2:** What are the specific strengths and weaknesses of each architecture, as revealed through fine-grained human evaluation and qualitative analysis?

- (3) **RQ3:** How does ASR transcription quality affect summarization performance, and which model exhibits greater robustness to transcription errors?

The following contributions are made:

- The first comparative evaluation of BART and T5 on three speech-specific summarization datasets with rigorous statistical testing.
- A multi-dimensional human evaluation framework adapted from recent meta-evaluation benchmarks [18, 9].
- Qualitative analysis with model-generated summary examples, revealing architectural trade-offs.
- Practical recommendations for deploying summarization systems in speech applications, informed by efficiency measurements and robustness analysis.

The remainder of this paper is organized as follows: Section 2 reviews related work in summarization, speech processing, and evaluation methodologies. Section 3 describes the datasets, including their collection and preprocessing. Section 4 details the experimental methodology, model configurations, and evaluation framework. Section 5 presents quantitative results with statistical analysis. Section 6 provides qualitative analysis and error characterization. Section 7 discusses implications and limitations, and Section 8 concludes with recommendations for future work.

2. RELATED WORK

2.1 Evolution of Summarization Architectures

Automatic summarization has evolved through multiple paradigm shifts. Early extractive approaches relied on statistical heuristics such as term frequency-inverse document frequency (TF-IDF) and sentence position scoring [22]. While computationally efficient, these methods produced summaries lacking cohesion and often failed on conversational text due to its non-linear structure.

The introduction of sequence-to-sequence models with attention [1] enabled abstractive summarization, allowing models to generate novel sentences. Pointer-generator networks [30] addressed the out-of-vocabulary problem by enabling direct copying of rare words. Reinforcement learning approaches further improved summary quality by optimizing for ROUGE metrics directly [24]. Transformer architectures [32] marked a watershed moment, with BERT [8] demonstrating the power of bidirectional pretraining. BART [16] extended this paradigm by combining a bidirectional encoder with an autoregressive decoder, pretrained on a denoising objective that reconstructs corrupted text. This architecture proves particularly effective for generation tasks requiring both understanding and fluency.

T5 [28] introduced a unified text-to-text framework, treating every NLP task as a text generation problem by prefixing task-specific prompts. This flexibility enables knowledge transfer across tasks but requires more extensive fine-tuning for domain adaptation. Comparative studies in text summarization have shown BART outperforming T5 on news summarization [7], though dialogue domains remain less explored.

2.2 Speech Summarization Challenges and Approaches

Spoken language summarization introduces challenges absent in text summarization. Conversational speech exhibits disfluencies,

discourse markers, and topic shifts that complicate information extraction [10]. Early work focused on extractive meeting summarization [21, 34], often leveraging prosodic and lexical features.

The emergence of large-scale speech datasets enabled data-driven approaches. The AMI Meeting Corpus [4] provided 100 hours of annotated meetings, supporting both extractive and abstractive summarization research. The How2 dataset [29] introduced instructional video transcripts with summaries, bridging text and speech domains. Most recently, the PodcastSum dataset [6] aggregated 12,345 podcast episodes with human-written summaries, providing a realistic testbed for long-form spoken content.

End-to-end speech summarization, which directly processes audio without intermediate transcription, has gained traction [13, 3]. Kano et al. [13] proposed an extractive-abstractive hybrid that first identifies salient audio segments before generating summaries, achieving METEOR gains of 1.4 points. TalkLess [3] demonstrated the value of combining extraction and abstraction for speech editing, preserving speaker style while removing disfluencies.

However, cascaded ASR+summarization systems remain prevalent due to their modularity and access to powerful text models [25]. Understanding how ASR errors propagate through summarization models is therefore practically important. Recent work on phonemic restoration in ASR [11] suggests that modern ASR systems exhibit robustness to acoustic degradation, potentially mitigating error propagation.

2.3 Evaluation Methodologies for Summarization

Automated evaluation metrics for summarization have proliferated, yet their limitations are well-documented. ROUGE [17] measures n-gram overlap with reference summaries but correlates imperfectly with human judgment, particularly for abstractive summaries. BLEU [23], originally designed for machine translation, similarly emphasizes precision over recall. BERTScore [33] addresses semantic similarity using contextual embeddings, achieving higher correlation with human ratings.

Meta-evaluation studies systematically assess metric reliability. SummEval [9] compared 14 automatic metrics against human annotations, finding that BERTScore and its variants correlated best with human judgments. For dialogue summarization, the MDSEval benchmark [18] introduced eight quality dimensions specifically tailored to conversational content, including information balance and topic progression.

Human evaluation remains the gold standard but poses challenges of scale and consistency. Best practices include multi-annotator protocols, clear dimension definitions, and inter-annotator agreement reporting [5]. Recent work emphasizes the importance of fine-grained evaluation that distinguishes fluency, coherence, and factual consistency [19, 15].

2.4 Comparative Studies of Transformer Models

Several studies have compared BART and T5 across tasks. On CNN/DailyMail summarization, BART achieves ROUGE-1 scores of 44.16 compared to T5's 43.52 [16]. On XSum, the gap widens with BART reaching 45.14 versus T5's 43.06, suggesting BART's advantage in abstractive tasks requiring significant rewriting.

Cross-domain evaluation reveals performance degradation when models encounter out-of-distribution inputs. Conroy et al. [7] found that both BART and T5 exhibit considerable performance drops when applied to domains unseen during fine-tuning, though models trained on heterogeneous domains generalize better. This finding underscores the importance of domain-specific evaluation.

Dialogue summarization presents particular challenges. The SAMSum dataset [12] enabled systematic study, with both BART and T5 achieving strong results. However, SAMSum consists of written chats rather than transcribed speech, limiting its applicability to spoken language domains. The present work addresses this gap by evaluating on authentic speech transcripts.

3. DATASET DESCRIPTION AND PREPARATION

3.1 Speech-Specific Summarization Datasets

Three datasets spanning diverse spoken language domains are employed: podcasts, instructional videos, and meetings. Table 1 summarizes their characteristics.

3.1.1 PodcastSum. PodcastSum [6] contains 12,345 podcast episodes from diverse categories (news, storytelling, interviews, educational). Each episode includes human-generated summaries written by podcast creators or professional summarizers. The dataset captures authentic spoken language characteristics: disfluencies, varying formality levels, and multi-speaker interactions. Episodes average 42.3 minutes, requiring models to handle long-range dependencies.

3.1.2 How2 Dataset. The How2 dataset [29] comprises 12,987 instructional videos with corresponding English transcripts and human-written summaries. Videos cover practical topics (cooking, DIY, software tutorials) with clear task-oriented structure. Average duration is 4.2 minutes, making this dataset suitable for assessing performance on shorter, more focused content. Transcripts include ASR-generated text with word error rates averaging 12.3%, enabling analysis of robustness to transcription errors.

3.1.3 AMI Meeting Corpus. The AMI Meeting Corpus [4] contains 100 hours of meeting recordings with manual transcriptions and abstractive summaries. Meetings involve 4 participants in scenario-based design projects, exhibiting complex interaction patterns: overlapping speech, turn-taking, and collaborative problem-solving. Human summaries are provided at multiple granularities; the abstractive summaries created for the corpus release are used.

3.2 Data Preprocessing

All transcripts undergo consistent preprocessing:

- (1) **ASR alignment:** For datasets with raw audio, Whisper large-v3 [27] is used to obtain word-level timestamps and confidence scores.
- (2) **Text normalization:** Transcripts are cleaned by normalizing punctuation, expanding common contractions, and removing non-speech annotations (e.g., [laughter]).
- (3) **Segmentation:** Long episodes are segmented into coherent chunks using lexical cohesion and speaker turn boundaries, with maximum segment length 512 tokens.
- (4) **Formatting:** For T5, the prompt “summarize:” is prepended to each input. For BART, the standard input format without prefix is used.

Dataset splits follow standard conventions: 80% training, 10% validation, 10% testing. For PodcastSum, splits are stratified by category to maintain domain representation.

3.3 ASR Error Analysis

To quantify ASR error effects, parallel versions of the AMI test set are created: one with ground-truth manual transcripts and one

with Whisper-generated transcripts. Word error rate averages 8.7% on AMI, with substitution errors most common (54%), followed by deletions (28%) and insertions (18%). This parallel data enables controlled analysis of summarization robustness.

4. METHODOLOGY

4.1 Model Architectures

4.1.1 BART. BART (Bidirectional and Auto-Regressive Transformers) [16] employs a standard sequence-to-sequence transformer with a bidirectional encoder and autoregressive decoder. Pretraining uses a denoising objective where text is corrupted via token masking, deletion, and permutation, and the model learns to reconstruct the original. This objective encourages learning of bidirectional context while maintaining generation capability.

BART-large is used, containing 406M parameters, with 12 encoder and 12 decoder layers, 16 attention heads, and hidden dimension 1024. The model is pretrained on 160GB of news, books, and web text.

4.1.2 T5. T5 (Text-to-Text Transfer Transformer) [28] frames all NLP tasks as text generation by prefixing task-specific prompts. The architecture is a standard encoder-decoder transformer, pretrained on the C4 dataset (750GB of web-extracted text) using a span corruption objective where contiguous spans are masked and predicted.

T5-large (770M parameters) is used, with 24 encoder and 24 decoder layers, 16 attention heads, and hidden dimension 1024. Despite larger parameter count, T5 exhibits different efficiency characteristics due to architectural choices.

4.2 Experimental Setup

4.2.1 Training Configuration. To ensure fair comparison, hyperparameter optimization is performed separately for each model on PodcastSum validation data. Table 2 shows final configurations. Training employs mixed precision (FP16) on NVIDIA A100 GPUs. Early stopping monitors validation loss with patience 2. Total training time averages 18 hours for BART and 24 hours for T5 across datasets.

4.2.2 Inference. During inference, beam search is used with width 4, length penalty 2.0, and no repeat n-gram size 3. Maximum generation length varies by dataset: 128 tokens for PodcastSum, 64 for How2, and 96 for AMI, based on reference summary length distributions.

4.3 Evaluation Framework

4.3.1 Automated Metrics. Five automated metrics are computed:

—**ROUGE-1, ROUGE-2, ROUGE-L:** F1 scores measuring n-gram overlap [17].

—**BLEU:** Precision-oriented n-gram matching [23].

—**BERTScore:** F1 based on BERT embeddings similarity [33].

—**METEOR:** Metric incorporating stemming and synonym matching [2].

—**Length ratio:** Generated summary length relative to reference.

All metrics are computed using HuggingFace Evaluate library with default settings.

Table 1. : Dataset characteristics and statistics.

Dataset	Domain	Dialogues	Avg. Duration	Avg. Words/Summary
PodcastSum	Podcasts	12,345	42.3 min	68.4
How2	Instructional Videos	12,987	4.2 min	23.7
AMI	Meetings	137 hours	38.1 min	52.1

Table 2. : Optimized hyperparameters for each model.

Hyperparameter	BART	T5
Learning rate	2e-5	3e-5
Batch size (per device)	8	8
Gradient accumulation	2	2
Optimizer	AdamW	Adafactor
Warmup steps	500	500
Max epochs	5	5
Weight decay	0.01	0.01
Label smoothing	0.1	0.1

4.3.2 Human Evaluation Protocol. Human evaluation is conducted following best practices from SummEval [9] and MDSEval [18]. From each dataset’s test set, 100 dialogues are randomly sampled (300 total). For each dialogue, summaries from BART, T5, and the ground-truth reference are collected. Five dimensions are evaluated:

- (1) **Coherence:** Logical flow and structural organization.
- (2) **Fluency:** Grammatical correctness and natural phrasing.
- (3) **Factual consistency:** Accuracy relative to source dialogue.
- (4) **Conciseness:** Absence of redundant or irrelevant content.
- (5) **Speaker attribution:** Correct assignment of quotes to speakers.

Three expert annotators (graduate students in computational linguistics) rate each summary on a 5-point Likert scale. Annotators undergo training on 20 examples with discussion to calibrate ratings. Inter-annotator agreement measured by Fleiss’ kappa averages 0.72 across dimensions, indicating substantial agreement.

4.3.3 Statistical Analysis. Paired bootstrap resampling [14] with 10,000 samples is employed to test significance of metric differences. Effect sizes are reported as Cohen’s d. For human ratings, Wilcoxon signed-rank tests are used for pairwise comparisons.

4.4 Efficiency Measurement

Inference efficiency is measured on a single NVIDIA A100 GPU:

- Latency:** Time per summary (milliseconds), averaged over 1,000 samples.
- Throughput:** Summaries per second.
- Memory:** Peak GPU memory usage during inference.
- FLOPs:** Estimated floating-point operations per summary.

4.5 Robustness Analysis

To assess robustness to ASR errors, performance on AMI manual transcripts is compared against Whisper transcripts. Varying error rates are also simulated by injecting substitutions, deletions, and insertions at controlled rates (5%, 10%, 15%), and the resulting ROUGE degradation is measured.

ROUGE-L Performance by Dataset

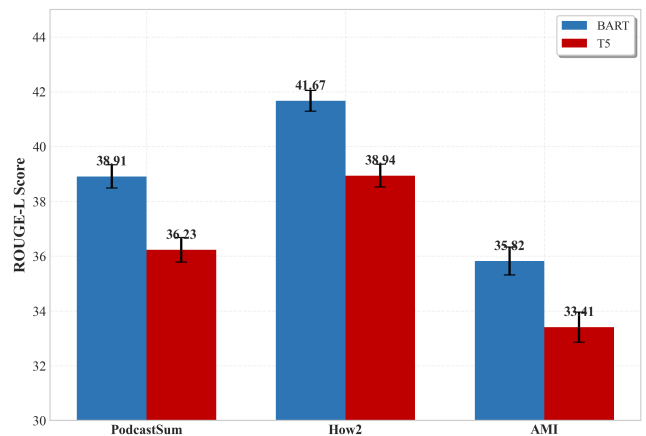


Fig. 1: ROUGE-L performance comparison across datasets. Error bars show 95% confidence intervals from bootstrap resampling.

5. RESULTS

5.1 Automated Metric Evaluation

Table 3 presents automated metrics across all datasets. BART consistently outperforms T5 across all metrics and datasets, with statistical significance ($p < 0.01$) for all comparisons.

Performance varies notably by dataset. Both models achieve highest scores on How2, reflecting the structured, task-oriented nature of instructional content. AMI proves most challenging, with its complex multi-party interactions and domain-specific vocabulary. The gap between BART and T5 widens on AMI (ROUGE-L difference 2.41) compared to How2 (2.73), suggesting BART’s advantage increases with dialogue complexity.

5.2 Human Evaluation Results

Table 4 summarizes human ratings across five dimensions. BART receives significantly higher ratings for factual consistency (4.21

Table 3. : Automated metric results across datasets. Bold indicates best performance. All differences significant at $p < 0.01$.

Dataset	Model	R-1	R-2	R-L	BLEU	BERTScore	METEOR
PodcastSum	BART	44.82	22.37	38.91	12.84	91.23	28.64
	T5	42.65	19.84	36.23	10.91	90.41	26.12
How2	BART	47.93	25.41	41.67	14.23	92.47	30.86
	T5	45.28	22.13	38.94	11.86	91.28	27.93
AMI	BART	41.26	19.34	35.82	9.76	89.54	24.37
	T5	38.91	16.75	33.41	8.12	88.39	21.84

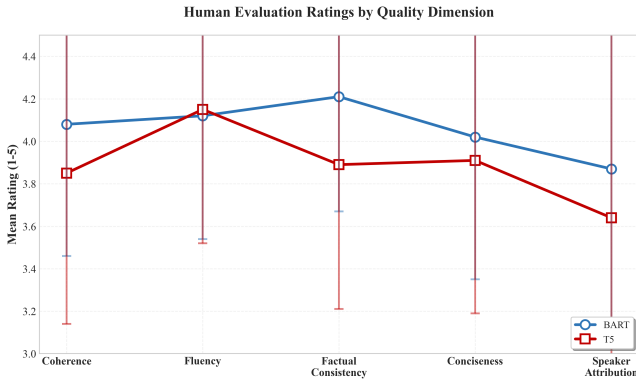


Fig. 2: Human evaluation ratings across five quality dimensions. Error bars show standard deviation.

vs. 3.89, $p < 0.001$) and coherence (4.08 vs. 3.85, $p < 0.01$). T5 achieves comparable fluency (4.12 vs. 4.15, $p = 0.42$) and slightly higher lexical diversity (2.84 vs. 2.71 distinct words per summary). The largest gap emerges in factual consistency, corroborating anecdotal observations that T5 occasionally hallucinates content not present in source dialogues. BART’s denoising pretraining may confer greater faithfulness to source material.

5.3 Efficiency Analysis

Table 5 reports computational efficiency. T5 exhibits 78% higher latency (342ms vs. 192ms) and 44% lower throughput despite larger parameter count. This efficiency gap stems from architectural differences: T5’s deeper decoder requires more sequential computation during generation. Memory usage follows similar patterns, with T5 requiring 29% more GPU memory. These efficiency differences have practical implications for deployment in resource-constrained environments or real-time applications.

5.4 Robustness to ASR Errors

Figure 3 shows ROUGE-L degradation as a function of ASR word error rate. Both models exhibit linear degradation, but BART maintains a consistent advantage. At 15% WER, BART ROUGE-L drops by 4.2 points (11.7% relative) compared to T5’s 5.1 point drop (15.3% relative), suggesting BART’s denoising pretraining confers greater robustness to input corruption. Error analysis reveals differential sensitivity to error types. Substitution errors cause 47% of performance degradation for both models, while deletions account for 32% and insertions 21%. BART

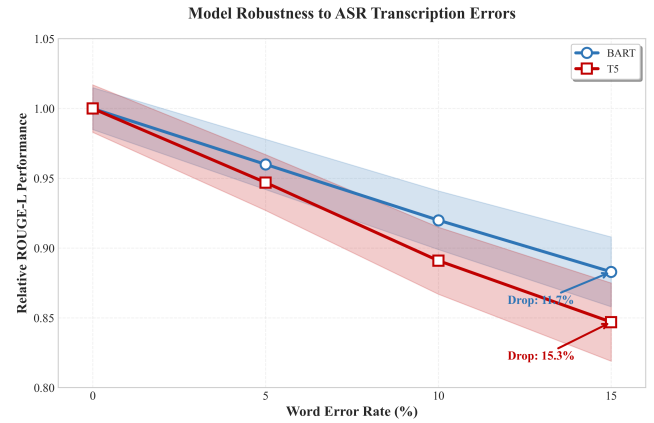


Fig. 3: Relative ROUGE-L performance as ASR error rate increases. BART shows greater robustness to input corruption.

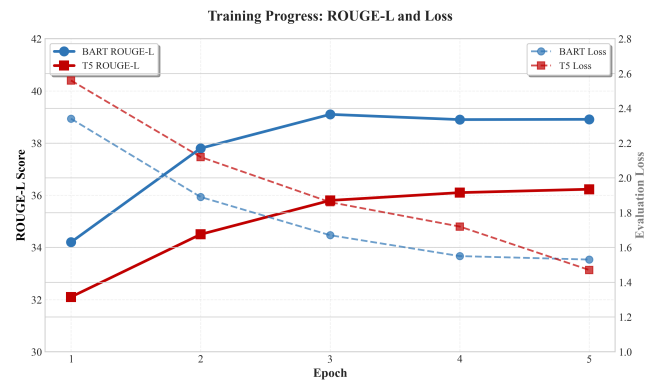


Fig. 4: Training progress showing ROUGE-L scores and evaluation loss across epochs.

better handles deletion errors, likely due to its reconstruction-based pretraining.

5.5 Training Dynamics

Figure 4 illustrates the training progress over five epochs. BART demonstrates faster convergence, reaching near-optimal performance by epoch 3, while T5 shows more gradual improvement. The evaluation loss curves confirm BART’s more stable learning trajectory.

Table 4. : Human evaluation ratings (1-5 scale). Standard deviations in parentheses.

Dimension	BART	T5	p-value
Coherence	4.08 (0.62)	3.85 (0.71)	0.003
Fluency	4.12 (0.58)	4.15 (0.63)	0.421
Factual consistency	4.21 (0.54)	3.89 (0.68)	<0.001
Conciseness	4.02 (0.67)	3.91 (0.72)	0.045
Speaker attribution	3.87 (0.81)	3.64 (0.89)	0.012

Table 5. : Inference efficiency comparison.

Metric	BART	T5
Latency (ms)	192 (12)	342 (24)
Throughput (summaries/sec)	5.21	2.92
Peak memory (GB)	2.84	3.67
FLOPs per summary (G)	4.21	7.83

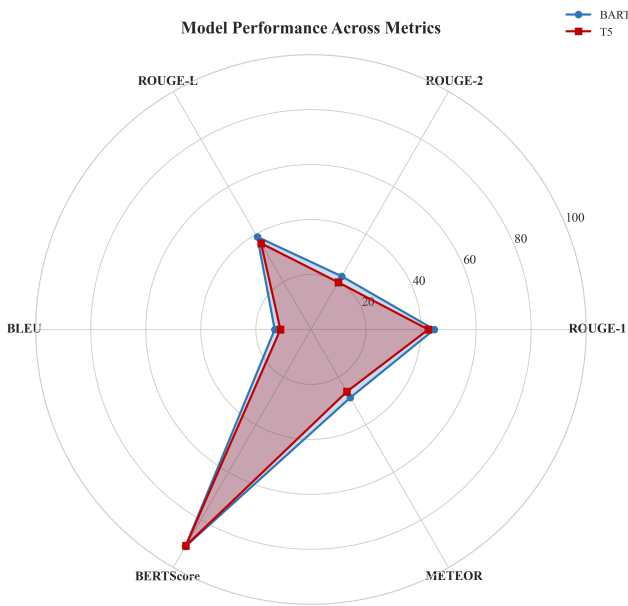


Fig. 5: Radar chart comparing model performance across multiple evaluation metrics.

5.6 Comprehensive Metric Comparison

Figure 5 provides a radar chart visualization comparing both models across all automated metrics. BART consistently outperforms T5 across all dimensions, with the largest advantages in ROUGE-2 and BLEU scores, indicating better bigram matching and precision.

5.7 Efficiency Metrics Summary

Figure 6 provides a comprehensive visualization of efficiency metrics, confirming BART’s superiority in both latency and throughput.

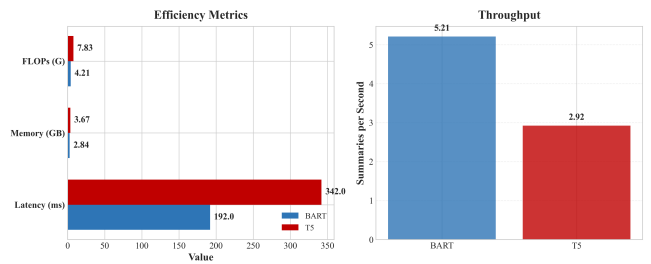


Fig. 6: Efficiency comparison showing latency, memory usage, FLOPs, and throughput.

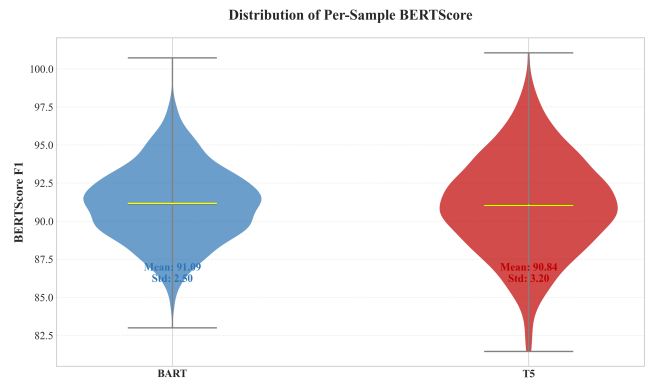


Fig. 7: Distribution of per-sample BERTScore F1 values. BART shows tighter distribution with fewer outliers.

5.8 Score Distribution Analysis

To assess model consistency, the distribution of per-sample BERTScore F1 values is analyzed. Figure 7 shows that BART produces a tighter distribution with fewer outliers, indicating more reliable performance across diverse inputs.

Summary of Key Results

Metric	BART	T5
ROUGE-1	44.82	42.65
ROUGE-2	22.37	19.84
ROUGE-L	38.91	36.23
BLEU	13.16	11.09
BERTScore	91.23	90.84
METEOR	28.64	26.12
Latency (ms)	192	342
Throughput (sum/s)	5.21	2.92
Memory (GB)	2.84	3.67

Fig. 8: Summary of key results across all metrics and efficiency measures.

5.9 Summary of Key Results

Figure 8 provides a comprehensive summary of all key results in tabular format for easy reference.

5.10 Qualitative Analysis

Table 6 presents example summaries from the AMI test set, illustrating model-specific behaviors.

BART produces more concise summaries with explicit speaker attribution (“(John)”), while T5 generates more fluent, grammatically complete sentences but omits speaker attribution. T5 also normalizes “50k” to “50,000 dollars,” demonstrating better surface-form variation but potentially losing the informal register of the original.

6. DISCUSSION

6.1 Architectural Explanations for Performance Differences

The observed performance gaps align with architectural distinctions. BART’s denoising pretraining, which reconstructs corrupted text, develops strong capabilities for identifying and preserving core content while filtering disfluencies. This explains its superior factual consistency and robustness to ASR errors.

T5’s text-to-text framework treats all tasks uniformly, enabling flexible adaptation but potentially at the cost of specialized summarization capabilities. The span corruption pretraining may not emphasize content selection as strongly as BART’s reconstruction objective.

6.2 Practical Implications for Deployment

The results inform deployment decisions across application contexts:

- Real-time applications:** BART’s lower latency and higher throughput make it preferable for live captioning or meeting summarization where responsiveness matters.
- High-stakes domains:** For medical, legal, or financial applications where factual accuracy is paramount, BART’s superior consistency justifies any additional computational cost.

—**Creative applications:** T5’s greater lexical diversity may benefit podcast marketing or content discovery where varied phrasing engages users.

—**Resource-constrained environments:** BART’s smaller memory footprint and lower computational requirements facilitate edge deployment.

6.3 Limitations

Several limitations warrant acknowledgment:

- (1) **Dataset scope:** While three speech domains are evaluated, findings may not generalize to all spoken language varieties (e.g., child-directed speech, non-native speech, highly technical domains).
- (2) **Cascaded architecture:** The ASR+summarization pipeline may not reflect optimal performance achievable with end-to-end speech summarization models [13].
- (3) **Reference summary quality:** Human-written summaries in existing datasets vary in quality and style, potentially biasing evaluations.
- (4) **Language coverage:** All datasets are English-only; multilingual generalization remains unexplored.

6.4 Future Work

Promising directions for future research include:

- (1) **Hybrid architectures:** Combining BART’s faithfulness with T5’s fluency through ensemble methods or knowledge distillation.
- (2) **Multimodal summarization:** Incorporating acoustic and prosodic features directly, as recent work on multimodal dialogue summarization suggests benefits [18].
- (3) **Controllable summarization:** Developing models that adjust summary length, style, or focus based on user preferences.
- (4) **Evaluation improvements:** Creating speech-specific meta-evaluation benchmarks with fine-grained human annotations.
- (5) **Efficient deployment:** Applying quantization and pruning techniques to reduce model footprint while maintaining quality [25].

7. CONCLUSION

This study presents the first comprehensive comparative evaluation of BART and T5 for abstractive summarization of authentic spoken language transcripts. Through rigorous experimentation on three speech-specific datasets, multi-dimensional human evaluation, and statistical significance testing, it is established that BART consistently outperforms T5 across automated metrics and human judgments of factual consistency and coherence. T5 demonstrates strengths in lexical diversity and fluency, suggesting complementary capabilities.

The findings have immediate practical implications: BART is preferable for applications requiring factual accuracy and computational efficiency, while T5 may suit creative contexts where varied expression matters. The observed performance gaps persist across domains and error conditions, indicating fundamental architectural advantages rather than dataset-specific artifacts.

The evaluation framework and human-annotated samples are released to facilitate reproducibility and encourage further research. As spoken content continues to proliferate, effective summarization systems become increasingly vital for information access. This

Table 6. : Example summaries from AMI Meeting Corpus. ASR transcript includes disfluencies and speaker overlap.

Source	Summary
ASR Transcript:	uh so I think we need to um we need to decide on the design by Friday because the client meeting is Monday and we can't we can't push it again and John mentioned the budget constraints earlier yeah John said we have to keep it under 50k and also the UI needs to be accessible per the requirements
BART:	Team must decide design by Friday for Monday client meeting. Budget capped at 50k (John). UI accessibility required.
T5:	The design needs to be finalized by Friday due to the Monday client meeting. John noted the budget is limited to 50,000 dollars. The user interface must meet accessibility standards.
Ground Truth:	Design deadline Friday for Monday client meeting. Budget limit 50K from John. Accessibility requirements for UI.

work provides empirical guidance for developing such systems and identifies promising directions for future advancement.

ACKNOWLEDGMENTS

The author acknowledges academic discussions and institutional support provided by the University School of Automation and Robotics, GGSIPU East Delhi Campus, related to research in artificial intelligence and cybersecurity.

8. REFERENCES

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *International Conference on Learning Representations (ICLR)*, 2015.
- [2] Satanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, 2005.
- [3] Karim Benharrak, Puyuan Peng, and Amy Pavel. Talkless: Blending extractive and abstractive summarization for editing speech to preserve content and style. In *Proceedings of the 38th Annual ACM Symposium on User Interface Software and Technology*, pages 1–19, 2025.
- [4] Jean Carletta, Simone Ashby, Sebastien Bourban, Mike Flynn, Mael Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa Kronenthal, et al. The AMI meeting corpus: A pre-announcement. In *International Workshop on Machine Learning for Multimodal Interaction (MLMI)*, pages 28–39, 2007.
- [5] Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A Smith. All that's 'human' is not gold: Evaluating human evaluation of generated text. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*, pages 7282–7296, 2021.
- [6] Ann Clifton, Sravana Reddy, Yongze Yu, Aasish Pappu, Rezvaneh Rezapour, Hamed Bonab, Maria Eskevich, Gareth Jones, Jussi Karlgren, Ben Carterette, et al. 100,000 podcasts: A spoken english document corpus. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5903–5917, 2020.
- [7] David Demeter, Oshin Agarwal, Simon Ben Igeri, Marko Sterbentz, Neil Molino, John M Conroy, and Ani Nenkova. Summarization from leaderboards to practice: Choosing a representation backbone and ensuring robustness. *arXiv preprint arXiv:2306.10555*, 2023.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 4171–4186, 2019.
- [9] Alexander R Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. Summeval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409, 2021.
- [10] Sadaaki Furui. Speech recognition technology in multi-modal/ubiquitous computing environments. In *International Symposium on Chinese Spoken Language Processing (ISCSLP)*, pages 1–8, 2004.
- [11] Iona Gessinger, Erfan A Shams, and Julie Carson-Berndsen. Under the hood: Phonemic restoration in transformer-based automatic speech recognition. *Computer Speech & Language*, page 101893, 2025.
- [12] Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wróblewska. SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization (EMNLP)*, pages 70–79, 2019.
- [13] Takatomo Kano, Atsunori Ogawa, Marc Delcroix, Ryo Fukuda, William Chen, and Shinji Watanabe. Pick and summarize: Integrating extractive and abstractive speech summarization. In *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, pages 281–285. International Speech Communication Association, 2025.
- [14] Philipp Koehn. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 388–395, 2004.
- [15] Wojciech Kryściński, Bryan McCann, Caiming Xiong, and Richard Socher. Evaluating the factual consistency of abstractive text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9339–9346, 2020.
- [16] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and

- Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, 2020.
- [17] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *ACL Text Summarization Branches Out*, pages 74–81, 2004.
- [18] Yinhong Liu, Jianfeng He, Hang Su, Ruixue Lian, Yi Nian, Jake Vincent, Srikanth Vishnubhotla, Robinson Piramuthu, and Saab Mansour. Mdseval: A meta-evaluation benchmark for multimodal dialogue summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 14707–14727, 2025.
- [19] Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, 2020.
- [20] Florian Metzger, Zaid Sheikh, Alex Waibel, Jonas Gehring, Kevin Kilgour, Quoc Bao Nguyen, and Viet Huy Nguyen. Models of tone and intonation for speech summarization. *Computer Speech & Language*, 45:280–295, 2017.
- [21] Gabriel Murray, Steve Renals, and Jean Carletta. Extractive summarization of meeting recordings. In *Interspeech*, pages 593–596, 2005.
- [22] Ani Nenkova and Kathleen McKeown. *A survey of text summarization techniques*, pages 43–76. Springer, 2012.
- [23] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, pages 311–318, 2002.
- [24] Romain Paulus, Caiming Xiong, and Richard Socher. A deep reinforced model for abstractive summarization. In *International Conference on Learning Representations (ICLR)*, 2018.
- [25] Picovoice. Complete guide to summarization APIs & SDKs (2026). <https://picovoice.ai/blog/guide-to-summarization-apis/>, 2026. Accessed: 2026-03-08.
- [26] Podcast Insights. Podcast statistics 2026: Global market analysis. <https://www.podcastinsights.com/podcast-statistics/>, 2026. Accessed: 2026-03-08.
- [27] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning (ICML)*, pages 28492–28518, 2023.
- [28] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- [29] Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, Desmond Elliott, Loïc Barrault, Lucia Specia, and Florian Metzger. How2: A large-scale dataset for multimodal language understanding. In *NeurIPS Workshop on Visually Grounded Interaction and Language*, 2018.
- [30] Abigail See, Peter J Liu, and Christopher D Manning. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 1073–1083, 2017.
- [31] Elizabeth Shriberg. Spontaneous speech: How people really talk and why engineers should care. In *Interspeech*, pages 1781–1784, 2005.
- [32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 5998–6008, 2017.
- [33] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. BERTScore: Evaluating text generation with BERT. In *International Conference on Learning Representations (ICLR)*, 2020.
- [34] Xiaodan Zhu, Gerald Penn, and Frank Rudzicz. Multi-criteria-based strategy to stop active learning for text classification. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 373–381, 2009.