# The Multicultural Medical Assistant: Can LLMs Improve Medical ASR Errors Across Borders?

**Ayo Adedeji** [1]   **Mardhiyah Sanni** [2]   **Emmanuel Ayodele** [2]   **Sarita Joshi** [1]   **Tobi Olatunji** [2]

[1] Google Cloud, 1600 Amphitheatre Parkway, Mountain View, CA 94043, USA
[2] Intron Health, San Francisco, CA, USA and Lagos, Nigeria

## Abstract

The global adoption of Large Language Models (LLMs) in healthcare shows promise to enhance clinical workflows and improve patient outcomes. However, Automatic Speech Recognition (ASR) errors in critical medical terms remain a significant challenge. These errors can compromise patient care and safety if not detected. This study investigates the prevalence and impact of ASR errors in medical transcription in Nigeria, the United Kingdom, and the United States. By evaluating raw and LLM-corrected transcriptions of accented English in these regions, we assess the potential and limitations of LLMs to address challenges related to accents and medical terminology in ASR. Our findings highlight significant disparities in ASR accuracy across regions and identify specific conditions under which LLM corrections are most effective.

## 1 Introduction

In recent years, medical Automatic Speech Recognition (ASR) systems have become integral to healthcare, revolutionizing processes such as physician-dictated notes, telemedicine, and doctor-patient conversations (Johnson et al., 2014; Zhang et al., 2023). By easing the administrative burden on healthcare providers, these systems allow them to focus more on patient care and less on documentation (Saxena et al., 2018).

Despite their contributions to efficiency in modern healthcare recordkeeping, significant challenges remain. Achieving high accuracy across various medical terminologies and diverse demographic accents continues to be a formidable task (Mani et al., 2020; DiChristofano et al., 2023). ASR systems often struggle with the precise recognition of specialized medical terminology, including drug names and diagnoses (Hodgson and Coiera, 2015). This limitation can lead to errors that undermine the quality and reliability of medical records.

Recent advances in Large Language Models (LLMs) have emerged as a promising direction for addressing these challenges. Their ability to understand and process human language nuances positions them as potential tools to improve the accuracy of medical transcription. Our study evaluates how effectively LLMs can improve medical transcription accuracy across healthcare settings in Nigeria, the United Kingdom, and the United States.

Our **contributions** are:

1. The first large-scale evaluation of both ASR performance and LLM-based corrections across healthcare settings in Nigeria, the United Kingdom, and the United States, analyzing 191 medical conversations spanning multiple specialties.

2. Our evaluation framework and metrics, released for reproducibility and future research in cross-regional medical ASR.

## 2 Related Work

### 2.1 ASR in Medical Settings

Recent research has focused on improving the accuracy of ASR to drive its acceptance in clinical practice. Although ASR systems have shown promise in dictating medical reports, they often struggle with the nuanced context and speaker diarization inherent in patient-clinician conversations. The high word error rate (WER) in these scenarios indicates poor performance in contextual understanding and diarization (Park et al., 2023; Tran et al., 2023). The wide variability in accents between healthcare providers and patients exacerbates these issues, leading to possible misinterpretations of critical medical information (Afonja et al., 2024; Zaporowski, 2024).

Addressing these challenges requires extensive training data, yet privacy concerns and regulatory

restrictions have led to a scarcity of publicly available medical conversation datasets (Korfiatis et al., 2022; Le-Duc, 2024). This limited data availability has restricted researchers' ability to develop robust and adaptable ASR systems for global healthcare environments.

## 2.2 Error Correction Approaches

Recent advances in error correction methods have shown promising results. Leng et al. (Leng et al., 2021) introduced FastCorrect 2, a non-autoregressive model that takes advantage of multiple ASR candidates to improve correction accuracy through a novel alignment algorithm. Boros et al. (Boros et al., 2024) evaluated fourteen foundation LLMs for post-transcription correction, providing valuable information on the capabilities and limitations of LLMs in this domain. Complementing these approaches, Radhakrishnan et al. (Radhakrishnan et al., 2023) developed a cross-modal fusion technique that combines acoustic information with external linguistic representations, demonstrating significant improvements in WER reduction.

## 2.3 LLMs in Medical Transcription

Initial experiments with GPT-4 to create structured documentation from clinical conversations revealed consistent challenges in information preservation, often introducing errors through omission or addition of content (Kernberg et al., 2024). However, more promising results have emerged from LLM-enhanced ASR systems, which have demonstrated improved speaker diarization and reduced WER. (Adedeji et al., 2024; Wang et al., 2024).

## 2.4 Cross-Regional ASR Studies

While several studies have examined ASR performance in specific healthcare contexts, comprehensive cross-regional evaluations remain limited. Existing research has focused mainly on single-region or single-accent scenarios, leaving a significant gap in our understanding of the performance of ASR systems in global healthcare settings (DiChristofano et al., 2023). Our work addresses this gap by providing a systematic evaluation of both ASR performance and LLM-based corrections across Nigeria, the United Kingdom, and the United States, offering insight into the practical applicability of ASR and LLM-correction approaches in global healthcare settings.

| Dataset | Region | Num. Conv. | Avg. Turns |
|---|---|---|---|
| Intron Health Teleconsultations | Africa | 25 | 99 |
| PriMock57 | UK / Europe | 57 | 92 |
| Fareez Medical Interviews | United States | 109 | 112 |

Table 1: Overview of the three medical conversation datasets used in this study, showing the geographic distribution, number of conversations, and average number of turns per conversation.

## 3 Materials

### 3.1 Nigerian Dataset

For our Nigerian dataset, we used a collection of simulated medical conversations provided by Intron Health (Olatunji et al., 2023). This dataset comprises 25 doctor-patient consultations that capture approximately 4 hours of spoken dialogue.

These consultations span multiple specialties, with obstetrics and gynecology being the most prevalent (5 cases), followed by infectious diseases (4), gastroenterology (2), cardiology (2), and endocrinology (2). Additional specialties include neurology (2), orthopedics (2), general surgery (2), and single cases in pulmonology, otolaryngology, family medicine, and hematology, providing comprehensive coverage across medical disciplines.

Each consultation in this dataset is rich in content, averaging 99 conversational turns and 1,437 spoken words. To ensure demographic representation, the dataset features a mix of 4 female and 7 male speakers, all aged between 25 and 35 years. Notably, both the patient and the doctor roles are portrayed by Nigerian medical professionals.

### 3.2 United Kingdom Dataset

For our United Kingdom dataset, we used PriMock57, a collection of simulated medical conversations developed by Babylon Health (Korfiatis et al., 2022). This dataset comprises 57 doctor-patient consultations that capture approximately 9 hours of spoken dialogue.

These consultations span multiple conditions, with cardiovascular problems being the most prevalent (11 cases), followed by gastrointestinal disorders (8), respiratory conditions (8), migraines (6), and other infections (8). Additional categories include fever (4), dermatological conditions (4), anaphylactic reactions (3), mental health concerns

(3), and physical injuries (2), providing comprehensive coverage of primary care presentations.

To ensure demographic diversity, the dataset includes 7 clinicians and 57 actors portraying patients, with a balanced gender distribution and an age range predominantly between 25 to 45 years. The linguistic landscape of the UK's healthcare system is reflected in the variety of accents: clinicians speak mainly British English with some Indian accents, while patient accents span British English (47.4%), various European (31.6%), and other English and non-English accents (21%).

### 3.3 United States Dataset

For our United States dataset, we used a subset of simulated medical conversations developed by Fareez et al. (Fareez et al., 2022). This subset comprises 109 doctor-patient consultations that capture approximately 22 hours of spoken dialogue, strategically sampled from a larger pool of 272 conversations.

To ensure specialty diversity, we included all non-respiratory cases and randomly sampled from the respiratory cases, which formed the majority of the original dataset. These consultations span multiple conditions, with respiratory cases being the most prevalent (51 cases), followed by musculoskeletal conditions (46), gastrointestinal (6), cardiac (5), and dermatologic (1).

The dataset features a balanced gender distribution between doctors (57% male, 43% female) and patients (55% male, 45% female). Both doctor and patient roles were portrayed by senior medical students and resident doctors, leveraging their clinical experience to simulate realistic patient-doctor interactions.

## 4 Approach

Our approach combines ASR systems and LLMs to transcribe, diarize, correct, and analyze medical conversations. We employ six ASR systems, chosen for their capabilities in handling medical terminology and diverse accents, and three LLMs, selected for their reasoning capabilities. For a detailed description of these models, see Appendices A.1 and A.2.

To identify and categorize medical concepts within our transcriptions, we utilized Google Cloud's Healthcare Natural Language API (Google Cloud, 2024b). This API allows for the precise extraction and categorization of medical entities,
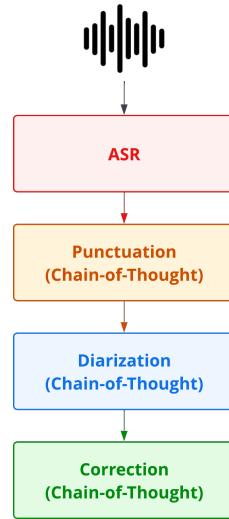


Figure 1: The steps of our Chain-of-Thought (CoT) prompting pipeline for medical conversation processing.

which is crucial for our analysis. More details about this API can be found in Appendix A.3.

### 4.1 Pipeline

Our pipeline for processing and analyzing medical conversations builds upon the methodology established in our previous work (Adedeji et al., 2024), with key enhancements to address cross-regional accents and variations in English. Figure 1 illustrates the main stages of our approach.

The pipeline consists of the following stages:

1. **Transcription:** Raw audio from our datasets is processed by each ASR system to generate initial transcriptions.

2. **Punctuation Enhancement:** We instruct the LLM with a Chain-of-Thought (CoT) prompt to improve punctuation to normalize orthographic variations and provide structured outputs for subsequent steps. A representative prompt can be found in Appendix A.4.

3. **Diarization:** We instruct the LLM with a CoT prompt to perform speaker diarization from scratch, without relying on any existing speaker labels. The LLM analyzes conversation patterns to assign each word and sentence to distinct speakers. A representative prompt can be found in Appendix A.5.

4. **Correction:** We instruct the LLM with a CoT prompt to improve transcription accuracy while preserving medical context. A representative prompt can be found in Appendix A.6.

5. **Error Analysis and Entity Recognition:** We

use Google Cloud's Healthcare Natural Language API for medical entity identification and categorization and the Jiwer Python library (Jiwer, 2024) for WER computation. This allows us to calculate both standard WER for general transcription accuracy and Medical Concept Word Error Rate (MC-WER) for evaluating medical concept transcription accuracy.

## 4.2 Evaluation Metrics

To quantify the effectiveness of our pipeline, we employ metrics that capture three key aspects of performance: transcription accuracy, diarization performance, and preservation of medical concepts.

### 4.2.1 Transcription Accuracy

We used WER as our primary metric to assess overall transcription accuracy. WER is calculated as:

$$WER = \frac{S + D + I}{N} \qquad (1)$$

where $S$ is the number of substituted words, $D$ is the number of deleted words, $I$ is the number of inserted words, and $N$ is the total number of words in the reference transcript.

### 4.2.2 Diarization Performance

We adapted a word-level error rate metric, common in multi-party conversation analysis (Yu et al., 2022; Shafey et al., 2019), to jointly evaluate speaker diarization and transcription accuracy. While traditional diarization error rate (DER) measures speaker attribution accuracy as:

$$DER = \frac{\text{Incorrectly attributed words}}{\text{Total words}} \times 100 \qquad (2)$$

our analysis uses WER computed against aligned hypothesis and ground truth speaker segments, effectively capturing both speaker attribution and transcription errors. We computed this speaker-level WER separately for doctor and patient speech segments, where a higher WER reflects both poor speaker attribution and transcription errors. This approach eliminates the need for timestamp alignment while providing granular insights into system performance across speaker roles.

### 4.2.3 Medical Concept Accuracy

We used MC-WER to specifically evaluate the transcription accuracy of medical terminology. While standard WER treats all words independently, MC-WER operates on complete medical concepts identified through Healthcare NLP annotation. MC-WER is calculated as:

$$MC\text{-}WER = \frac{S_m + D_m + I_m}{M} \qquad (3)$$

where $S_m$ is the number of medical concept substitutions (e.g., "diuretics" for "Dioralyte" or "high tension" for "hypertension"), $D_m$ is the number of medical concept deletions, $I_m$ is the number of medical concept insertions, and $M$ is the total number of medical concepts in the reference transcript. Medical concepts are aligned between hypothesis and reference transcripts as complete units, regardless of whether they comprise single or multiple words. We calculated MC-WER for both lemmatized and non-lemmatized versions of medical concepts:

- **Lemmatized MC-WER:** Captures errors in the base forms of medical terms (e.g., treating "antibiotics" and "antibiotic" as equivalent).

- **Non-lemmatized MC-WER:** Preserves morphological variations (e.g., distinguishing between singular/plural forms), which can be clinically significant.

The results presented in this paper use the non-lemmatized version to maintain the integrity of important medical distinctions.

## 5 Experiments

We conducted all Whisper Large V3 and NVIDIA Canary-1B inference experiments using a single NVIDIA T4 GPU. For both models, the transcriptions were processed in 20 second speech slices to optimize memory usage while maintaining transcription accuracy. For the remaining ASR systems, we utilized their publicly available enterprise APIs with default configurations optimized for medical speech. Similarly, all LLM operations were performed through public-facing APIs, with temperature values optimized for each model within the 0 to 0.1 range based on preliminary testing. The LLM correction and diarization steps were applied in 10-line segments. No additional hyperparameter tuning or specific resource configurations were needed for these tasks.
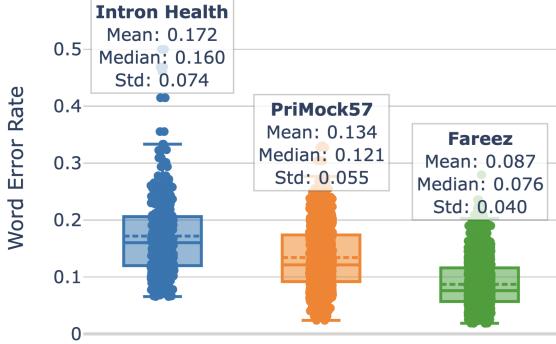
Figure 2: Distribution of WER across baseline ASR transcriptions for all speakers in each dataset.
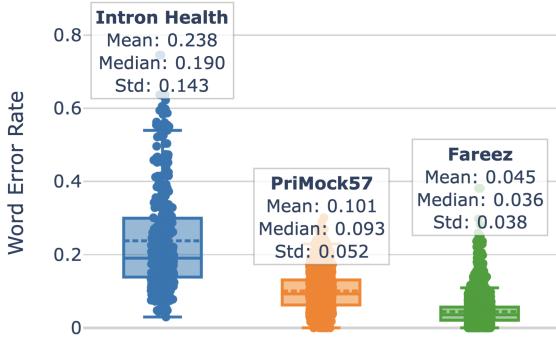


Figure 3: Distribution of MC-WER across baseline ASR transcriptions for all speakers in each dataset.
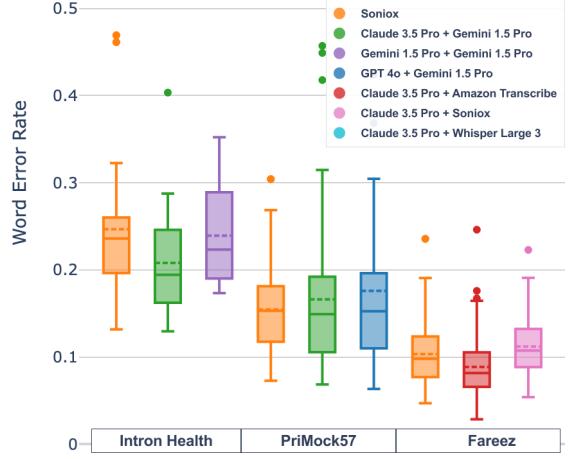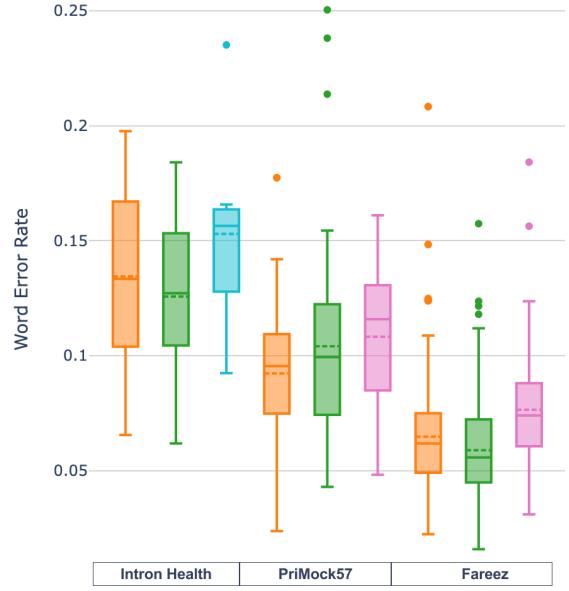
# 6 Results

## 6.1 Cross-Dataset WER Analysis

Our analysis revealed systematic differences in transcription accuracy across the datasets (Figures 2 and 3). For general transcription accuracy (WER), the Nigerian-accented Intron Health dataset showed the highest error rate (mean: 0.172), followed by the UK-accented PriMock57 (mean: 0.134), and the US-accented Fareez dataset (mean: 0.087). This represents a 4.7 percentage point difference between PriMock57 and Fareez and an 8.5 percentage point difference between Intron Health and Fareez. Medical concept transcription accuracy (MC-WER) showed larger disparities: Intron Health (mean: 0.238), PriMock57 (mean: 0.101), and Fareez (mean: 0.045). This represents a 5.6 percentage point difference between PriMock57 and Fareez and a 19.3 percentage point difference between Intron Health and Fareez.

## 6.2 LLM Diarization Performance

Across the three datasets, LLM-based diarization (where an LLM performs de novo speaker labeling of an ASR-generated transcript) demonstrated varying degrees of improvement when compared



(a) Patient Speech WER Distribution



(b) Doctor Speech WER Distribution

Figure 4: WER Distributions for patient and doctor speech for the baseline ASR system, Soniox (orange), and top performing LLM and ASR pairs.

to Soniox, our baseline ASR system chosen for its native diarization capabilities. For patient speech (Figure 4a), the Intron Health dataset showed the most marked improvement, with Claude 3.5 Sonnet + Gemini 1.5 Pro (mean: 0.208) outperforming Soniox (mean: 0.247) by 3.9 percentage points.

Multiple LLM-ASR pairs achieved comparable performance in PriMock57, with the best pair yielding mean WER of 0.15, a modest 0.5 percentage point improvement over Soniox (mean: 0.155). The Fareez dataset maintained the lowest overall error rates, with the top LLM-ASR pair achieving mean WER of 0.089 versus Soniox's 0.1034, an improvement of 1.44 percentage points.
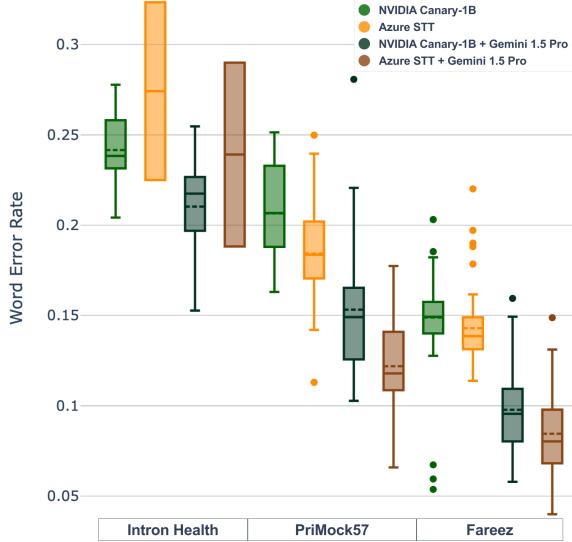
Figure 5: Comparison of WER before and after correction for lower-performing ASR systems across datasets. The graph shows the WER distribution of NVIDIA Canary-1B (green), Azure STT (orange), and their LLM-corrected versions using Gemini 1.5 Pro (olive for NVIDIA correction, brown for Azure STT correction).

Doctor speech (Figure 4b) demonstrated consistently lower error rates than patient speech across all datasets. For doctor speech, the top performing LLM and ASR pair (Claude 3.5 Sonnet + Gemini 1.5 Pro) achieved mean WERs of 0.127, 0.104, and 0.059 for Intron Health, PriMock57, and Fareez datasets, respectively, compared to Soniox's baseline mean WERs of 0.135, 0.095, and 0.064. The performance differences between LLM-based diarization and Soniox varied by dataset, with small improvements in Intron Health and Fareez, but slightly worse performance in PriMock57 for doctor speech.

### 6.3 LLM Correction Performance

When evaluating LLM corrections with Gemini 1.5 Pro as our representative model, we observed universal improvements across lower-performing ASR systems (Figure 5). Correction results for other LLMs are available in Appendices A.7, A.8, and A.9. In the Intron Health dataset, NVIDIA Canary-1B's mean WER improved from 0.241 to 0.21, and Azure STT's from 0.274 to 0.239. On PriMock57, corrections reduced Canary-1B's mean WER from 0.207 to 0.153 and Azure STT's from 0.184 to 0.122. The Fareez dataset showed the strongest relative improvements, with Canary-1B's mean WER decreasing from 0.149 to 0.098 and Azure STT's from 0.142 to 0.085.
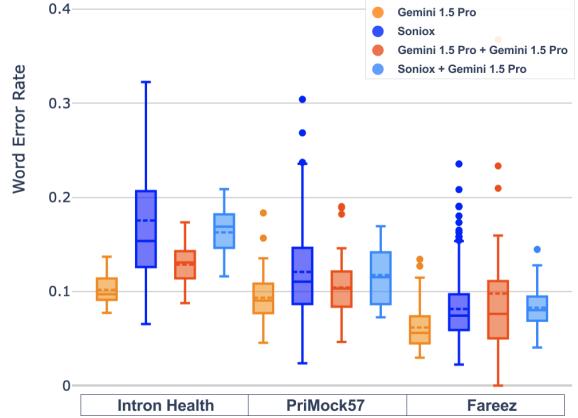


Figure 6: Comparison of WER before and after correction for higher performing ASR systems across datasets. The graph shows the WER distribution of Gemini 1.5 Pro (orange), Soniox ASR (dark blue), and their LLM-corrected versions using Gemini 1.5 Pro (red for Gemini correction, light blue for Soniox correction).

For higher performing systems such as Gemini 1.5 Pro and Soniox, corrections yielded minimal improvements or negative effects (Figure 6). Testing on Intron Health showed Gemini 1.5 Pro's mean WER increased from 0.102 to 0.129, while Soniox's decreased from 0.175 to 0.163. The pattern continued with PriMock57, where Gemini 1.5 Pro's mean WER rose from 0.093 to 0.104 while Soniox's slightly improved from 0.12 to 0.117. On Fareez, both systems showed reduced performance, with Gemini 1.5 Pro's mean WER increasing from 0.062 to 0.098 and Soniox's from 0.081 to 0.083.

### 6.4 Error Types and Correction Patterns

Analysis of correction patterns, using NVIDIA Canary-1B with Gemini 1.5 Pro as a representative model pairing, reveals that effectiveness varies significantly between different types of medical concepts and error categories (Figure 7). In terms of substitutions, the most substantial improvements were observed in medical problems (conditions and diseases; -63 substitutions), medications (therapeutic drugs and preparations; -44 substitutions), and laboratory data (bodily sample test results; -26 substitutions). In terms of insertions, LLM correction was particularly effective in reducing medication route errors (medication administration locations; -54 insertions) and medical problems (-42 insertions). However, we observed some categories where LLM correction introduced additional errors, notably in medications (+12 insertions), substance abuse terminology (psychoactive substance use; +10 insertions), and biological measurement results

| | Deleted | Inserted | Substituted |
|---|---|---|---|
| Substance Abuse | 6 | 10 | 3 |
| Severity | -8 | 2 | 2 |
| Procedure | -15 | 1 | -5 |
| Problem | -21 | -42 | -63 |
| Medication Unit | 0 | 0 | 0 |
| Medication Strength | 0 | 0 | 0 |
| Medication Status | 2 | 0 | 0 |
| Medication Route | 0 | -54 | 0 |
| Medication Frequency | 0 | 0 | 0 |
| Medication Form | 0 | 4 | 0 |
| Medication | 0 | 12 | -44 |
| Medical Device | 2 | -2 | -2 |
| Laboratory Result | -4 | -4 | 0 |
| Laboratory Data | -10 | 0 | -26 |
| Body Measurement Value | 4 | 2 | -2 |
| Body Measurement Unit | 0 | 0 | -2 |
| Body Measurement Result | -2 | 6 | 4 |
| Body Measurement | 2 | 6 | 2 |
| Body Function Result | 0 | -2 | 0 |
| Body Function | -10 | 2 | 0 |
| Anatomical | -6 | 0 | -14 |

Figure 7: Distribution of error differences by medical concept category, comparing pre-correction (NVIDIA Canary-1B ASR baseline) with post-correction (with Gemini 1.5 Pro) on the Intron Health dataset. Negative values indicate fewer errors after correction, while positive values indicate more errors after correction.

(basic vital signs and measurements; +6 insertions). For deletions, correction showed improvements in medical problems (-21 deletions), procedures (diagnostics or treatments; -15 deletions), laboratory data (-10 deletions), and body functions (-10 deletions). The correction process also introduced omissions, particularly for substance abuse terminology (+6 deletions) and body measurement values (+4 deletions).

Analysis of character-level error patterns (Figure 8) reveals that most medical concept substitution errors involved minor variations of less than 5 characters, comprising 68% in Intron Health, 73% in PriMock57, and 79% in Fareez. LLM correction successfully resolved approximately two thirds of these low character difference errors in all datasets (67%, 64%, and 67%, respectively), demonstrating particular effectiveness in addressing orthographic variations in medical terminology.

# 7 Discussion

## 7.1 Impact of Accents on Medical ASR

Our results demonstrate significant accent-based disparities in ASR performance across medical contexts. The systematic differences between datasets likely reflect the ASR systems' training data composition, with performance declining for non-US
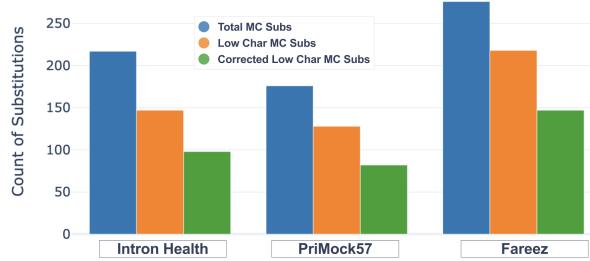


Figure 8: Comparison of total medical concept (MC) substitutions (blue), low character-difference MC substitutions (<5 characters) (orange), and corrected low character-difference MC substitutions (green) across datasets with NVIDIA Canary-1B ASR with Gemini 1.5 Pro correction.

accents. This training bias becomes particularly pronounced for non-Western accents, as shown by the substantially higher error rates in Nigerian-accented speech. Interestingly, while accented speech might be expected to disproportionately impact medical terminology recognition, our MC-WER analysis reveals a more nuanced pattern.

## 7.2 Comparing WER and MC-WER

Our analysis revealed an unexpected pattern: MC-WER values in the PriMock57 and Fareez datasets were consistently lower than their corresponding WER values. This apparent contradiction of medical terms being easier to transcribe requires careful interpretation. MC-WER and WER operate on fundamentally different sampling distributions: while WER considers all words independently, MC-WER focuses exclusively on a subset of medically relevant terms. This means that direct numerical comparisons between WER and MC-WER should be interpreted with caution.

Consider how MC-WER handles compound errors: when "hypertension" becomes "high tension", traditional WER counts both substitution and insertion errors, while MC-WER records only a single substitution. More critically, in a substitution where "amoxicillin" becomes "ampicillin," MC-WER treats this as a simple substitution, despite the significant clinical implications of confusing these antibiotics, which differ in their spectrum of activity and clinical indications. This suggests that while MC-WER effectively captures the preservation of medical terms as linguistic units, it may not reflect the clinical impact of transcription errors.

These findings highlight the need for additional metrics that can bridge the gap between statistical error rates and clinical significance.

## 7.3 Comparing LLM and ASR Diarization

Our results indicate that LLMs consistently diarize at or above the performance of ASR systems in all datasets tested. When examining speaker-specific error patterns, we observed consistently lower error rates in doctor speech compared to patient speech across all datasets. This pattern aligns with the linguistic characteristics of medical conversations: doctors typically follow more standardized patterns of medical discourse, while patients exhibit greater variability in accent and expression. This disparity is most pronounced in the Intron Health dataset, where both doctors and patients speak Nigerian-accented English, but patients' more varied and colloquial expression patterns led to higher error rates. A similar pattern appears in the PriMock57 dataset, where patient speech exhibits various European and non-English accents, while physicians maintain more consistent British English medical discourse, highlighting the ongoing challenges ASR systems face with accent diversity.

## 7.4 Differential Impact of LLM Correction

The effectiveness of LLM correction showed a distinct relationship with baseline ASR performance. Lower-performing ASR systems like Azure STT and NVIDIA Canary-1B saw substantial WER reductions with LLM correction, while high-performing ASR systems like Gemini 1.5 Pro and Soniox showed minimal gains or degradation. These results suggest that high-performing ASR systems may be approaching the theoretical performance limits on these datasets. Interestingly, while Gemini 1.5 Pro's attempts to correct its own transcriptions consistently led to degraded performance, Soniox showed modest improvements when corrected by Gemini 1.5 Pro. This contrast suggests that combining models with complementary strengths—potentially arising from different training focuses, data sources, or optimization objectives—can lead to better overall performance, aligning with ensemble and mixture-of-experts principles in machine learning.

The success of LLM correction was strongly correlated with the type of transcription error. LLMs excelled at correcting minor character-level variations (e.g., "fexifenadine" to "fexofenadine"), but struggled with semantically plausible alternatives (e.g., "relieved" to "really eased") and more distant medical term substitutions (e.g., "metformin" to "warfarin"). These findings suggest that integrated ASR-LLM systems might benefit from specialized pathways for handling semantic, distant medical term substitutions, and orthographic errors.

## 8 Conclusion

This study presents a large-scale evaluation of ASR performance and LLM-based corrections in healthcare settings in Nigeria, the United Kingdom, and the United States. Our findings demonstrate that LLMs diarize better than ASR systems in challenging accent scenarios, while their effectiveness in error correction varies with baseline ASR performance. Our analysis of accents and medical terminology highlights the limitations of current metrics such as WER and MC-WER, emphasizing the need for metrics that better reflect the clinical impact of transcription errors.

By evaluating off-the-shelf ASR systems and LLMs without specialized fine-tuning, our findings demonstrate both the potential and current limitations of these tools in improving medical transcription accuracy in global settings. This is particularly relevant for resource-limited environments where access to customized solutions may be impractical. While LLMs show promise in improving medical transcription accuracy, their effectiveness varies by accent and context. Future work should focus on developing more robust solutions that better handle accents and preserve medical terms, particularly for accents underrepresented in current ASR training data.

## 9 Limitations

Several limitations should be acknowledged in our study. Our focus on English-language medical conversations excludes local languages and code-switching scenarios common in multilingual healthcare settings. The simulated nature of our datasets, while allowing controlled comparison, may not fully capture real-world clinical environments with background noise and varied acoustic conditions.

Our evaluation of LLM correction focused on text-based corrections without taking advantage of acoustic features of the original audio, potentially overlooking valuable phonetic cues for accent-specific error correction. Additionally, while MC-WER quantifies errors in medical terminology, it does not capture the clinical significance of transcription errors, treating all substitutions equally regardless of their potential impact on patient care.

# References

Ayo Adedeji, Sarita Joshi, and Brendan Doohan. 2024. The sound of healthcare: Improving medical transcription asr accuracy with large language models. *Preprint*, arXiv:2402.07658.

Tejumade Afonja, Tobi Olatunji, Sewade Ogun, Naome A. Etori, Abraham Owodunni, and Moshood Yekini. 2024. Performant asr models for medical entities in accented speech. *Preprint*, arXiv:2406.12387.

Anthropic. 2024. Claude 3.5 sonnet. https://www.anthropic.com/news/claude-3-5-sonnet.

AWS. 2024. Amazon transcribe medical. https://aws.amazon.com/transcribe/medical/.

Azure. 2024. Azure ai speech. https://azure.microsoft.com/en-us/products/ai-services/ai-speech.

Emanuela Boros, Maud Ehrmann, Matteo Romanello, Sven Najem-Meyer, and Frédéric Kaplan. 2024. Post-correction of historical text transcripts with large language models: An exploratory study. In *Proceedings of the 8th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature (LaTeCH-CLfL 2024)*, pages 133–159, St. Julians, Malta. Association for Computational Linguistics.

Alex DiChristofano, Henry Shuster, Shefali Chandra, and Neal Patwari. 2023. Global performance disparities between english-language accents in automatic speech recognition. *ArXiv*, abs/2208.01157.

Hugging Face. 2024. Nvidia nemo canary 1b. https://huggingface.co/nvidia/canary-1b.

Faiha Fareez, Tishya Parikh, Christopher Wavell, Saba Shahab, Meghan Chevalier, Scott Good, Isabella De Blasi, Rafik Rhouma, Christopher McMahon, Jean-Paul Lam, Thomas Lo, and Christopher W. Smith. 2022. A dataset of simulated patient-physician medical interviews with a focus on respiratory cases. *Scientific Data*, 9(313).

Google Cloud. 2024a. Gemini API support on Vertex AI. https://cloud.google.com/blog/products/ai-machine-learning/gemini-support-on-vertex-ai. Accessed: 2024-08-23.

Google Cloud. 2024b. Healthcare natural language API. https://cloud.google.com/healthcare-api/docs/concepts/nlp. Accessed: 2024-08-23.

Tobias Hodgson and Enrico Coiera. 2015. Risks and benefits of speech recognition for clinical documentation: a systematic review. *Journal of the American Medical Informatics Association*, 23(e1):e169–e179.

Jiwer. 2024. Evaluation metrics for Automatic Speech Recognition. https://github.com/jitsi/jiwer.

M. Johnson, S. Lapkin, V. Long, P. Sanchez, H. Suominen, J. Basilakis, and L. Dawson. 2014. A systematic review of speech recognition technology in health care. *BMC Med Inform Decis Mak*, 14(94):1–18.

Annessa Kernberg, Jeffrey A Gold, and Vishnu Mohan. 2024. Using chatgpt-4 to create structured medical notes from audio recordings of physician-patient encounters: Comparative study. *Journal of Medical Internet Research*, 26:e54419.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Large language models are zero-shot reasoners. *Preprint*, arXiv:2205.11916.

Alex Korfiatis, Francesco Moramarco, Radmila Sarac, and Aleksandar Savkov. 2022. PriMock57: A dataset of primary care mock consultations. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 588–598, Dublin, Ireland. Association for Computational Linguistics.

Khai Le-Duc. 2024. VietMed: A dataset and benchmark for automatic speech recognition of Vietnamese in the medical domain. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 17365–17370, Torino, Italia. ELRA and ICCL.

Yichong Leng, Xu Tan, Rui Wang, Linchen Zhu, Jin Xu, Wenjie Liu, Linquan Liu, Xiang-Yang Li, Tao Qin, Edward Lin, and Tie-Yan Liu. 2021. FastCorrect 2: Fast error correction on multiple candidates for automatic speech recognition. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4328–4337, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Anirudh Mani, Shruti Palaskar, and Sandeep Konam. 2020. Towards understanding ASR error correction for medical conversations. In *Proceedings of the First Workshop on Natural Language Processing for Medical Conversations*, pages 7–11, Online. Association for Computational Linguistics.

Tobi Olatunji, Tejumade Afonja, Aditya Yadavalli, Chris Chinenye Emezue, Sahib Singh, Bonaventure FP Dossou, Joanne Osuchukwu, Salomey Osei, Atnafu Lambebo Tonja, Naome Etori, et al. 2023. Afrispeech-200: Pan-african accented speech dataset for clinical and general domain asr. *Transactions of the Association for Computational Linguistics*, 11:1669–1685.

OpenAI. 2024. Gpt-4. https://openai.com/research/gpt-4.

OpenAI. 2024. Whisper large v3. https://huggingface.co/openai/whisper-large-v3. Accessed: 2024-06-14.

Tae Jin Park, Kunal Dhawan, Nithin Rao Koluguri, and Jagadeesh Balam. 2023. Enhancing speaker diarization with large language models: A contextual beam search approach. *ArXiv*, abs/2309.05248.

Srijith Radhakrishnan, Chao-Han Yang, Sumeer Khan, Rohit Kumar, Narsis Kiani, David Gomez-Cabrero, and Jesper Tegnér. 2023. Whispering LLaMA: A cross-modal generative error correction framework for speech recognition. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10007–10016, Singapore. Association for Computational Linguistics.

K. Saxena, R. Diamond, R. F. Conant, T. H. Mitchell, I. G. Gallopyn, and K. E. Yakimow. 2018. Provider adoption of speech recognition and its impact on satisfaction, documentation quality, efficiency, and cost in an inpatient ehr. *AMIA Jt Summits Transl Sci Proc*, pages 186–195. Epub 2018 May 18.

Laurent El Shafey, Hagen Soltau, and Izhak Shafran. 2019. Joint speech recognition and speaker diarization via sequence transduction. *Preprint*, arXiv:1907.05337.

Soniox. 2024. Soniox. https://soniox.com/.

B. D. Tran, R. Mangu, M. Tai-Seale, J. E. Lafata, and K. Zheng. 2023. Automatic Speech Recognition Performance for Digital Scribes: A Performance Comparison Between General-Purpose and Specialized Models Tuned for Patient-Clinician Conversations. In *AMIA Annu Symp Proc. 2022*, pages 1072–1080.

Quan Wang, Yiling Huang, Guanlong Zhao, Evan Clark, Wei Xia, and Hank Liao. 2024. Diarizationlm: Speaker diarization post-processing with large language models. *Preprint*, arXiv:2401.03506.

Fan Yu, Zhihao Du, Shiliang Zhang, Yuxiao Lin, and Lei Xie. 2022. A comparative study on speaker-attributed automatic speech recognition in multi-party meetings. *Preprint*, arXiv:2203.16834.

S. Zaporowski. 2024. The impact of foreign accents on the performance of whisper family models using medical speech in polish. In *Harnessing Opportunities: Reshaping ISD in the post-COVID-19 and Generative AI Era (ISD2024 Proceedings)*, Gdańsk, Poland. University of Gdańsk.

J. Zhang, J. Wu, Y. Qiu, A. Song, W. Li, X. Li, and Y. Liu. 2023. Intelligent speech technologies for transcription, disease diagnosis, and medical equipment interactive control in smart hospitals: A review. *Comput Biol Med*, 153:106517. Epub 2023 Jan 5.

## A  Appendix

### A.1  Speech Recognition Models

In our analysis, we employed six ASR systems, selected for their advanced capabilities in processing complex medical language and their potential to handle diverse accents and linguistic variations. These systems are:

1. Google's Gemini-1.5-Pro (Version 001): Google's latest advancement in multimodal AI as of July 2024, capable of processing both textual and audio inputs simultaneously (Google Cloud, 2024a).

2. Microsoft Azure's Speech-to-Text: a commercially available service with support for multiple languages and specialized medical vocabulary (Azure, 2024).

3. OpenAI's Whisper 3: the latest iteration of OpenAI's open-source speech recognition model, known for its multilingual capabilities (OpenAI, 2024).

4. NVIDIA NeMo Canary-1B: a multi-lingual, multi-tasking model supporting ASR in English, German, French, and Spanish, with 1 billion parameters (Face, 2024).

5. Amazon Transcribe Medical: a commercially available service specialized for healthcare. We used the conversation model with specialty set to Primary Care (AWS, 2024).

6. Soniox: a commercially available service. We used the en_v2 model, which offers high accuracy across various accents (Soniox, 2024).

### A.2  Large Language Models

To investigate the potential of Large Language Models (LLMs) in improving medical ASR outputs, we selected three models:

1. Google's Gemini-1.5-Pro (Version 001): Google's latest advancement in multimodal AI as of July 2024, capable of processing both textual and audio inputs simultaneously (Google Cloud, 2024a).

2. Anthropic's Claude 3.5 Sonnet: Known for its strong performance in reasoning tasks, it is well suited to analyze and refine medical transcriptions (Anthropic, 2024).

3. OpenAI's GPT-4o: A leading model in natural language processing, offering advanced capabilities in understanding context and nuance in medical conversations (OpenAI, 2024).

These models were chosen based on their demonstrated capabilities in understanding and generating human-like text, their ability to process and reason over complex information, and their potential for zero-shot and few-shot learning in specialized domains (Kojima et al., 2023).

## A.3 Healthcare Natural Language API

To accurately identify and categorize medical concepts within our transcriptions, we employed Google Cloud's Healthcare Natural Language API. This choice builds upon our previous work (Adedeji et al., 2024), where we demonstrated the API's effectiveness in parsing and structuring unstructured medical text. The Healthcare Natural Language API offers several key functionalities crucial to our study:

1. Extraction of medical concepts, including diseases, medications, and procedures.

2. Categorization of these concepts into over thirty distinct entities.

3. Analysis of functional features such as temporal relationships and certainty assessments.

A significant advantage of this API is its ability to map extracted concepts to standard medical vocabularies like RxNorm, ICD-10, and SNOMED CT. This mapping is essential for maintaining consistency and accuracy in our medical concept identification across multi-regional datasets.

## A.4 Example Punctuation Prompt

For Punctuation, the following prompt template was utilized:

> You are a helpful speech-to-text transcription assistant. Your task is to correct punctuation in medical dialogue transcripts, ensuring accurate reflection of natural pauses and speaker transitions. Here's how to approach the task step by step:
>
> 1. Contextual Reading: Analyze each sentence for natural pauses, transitions, and speaker changes.
> 2. Sentence Splitting: Split run-on sentences into separate statements when there's a change in speaker.
> 3. Question Identification: Mark questions

appropriately with question marks, paying attention to intonation and structure.
> 4. Run-On Correction: Divide run-on sentences into properly punctuated independent clauses.
> 5. Speaker Transition: Use appropriate punctuation to separate clauses spoken by different speakers.
> 6. Capitalization: Begin each new statement or speaker transition with a capital letter.
> 7. Response Delineation: End each response with appropriate terminal punctuation.
>
> 8. Examples #1-10:
> {examples}
>
> Transcript:
> {transcript_segments}
>
> Expected output structure as a JSON array:
>
> ```
> [
>   {
>     "sentence": "[punctuated sentence]"
>   },
>   {
>     "sentence": "[punctuated sentence]"
>   }
> ]
> ```

## A.5 Example Diarization Prompt

For Diarization, the following prompt template was utilized:

> You are a helpful speech-to-text transcription assistant. Your current task is to diarize a conversation with no speaker labels. You will use your advanced understanding of medical terminology, dialogue structure, and conversational context to diarize the text accurately. Here's how to approach the task step by step:
> 1. Contextual Reading: Read each sentence thoroughly, absorbing its content, tone, sentiment, and vocabulary.
> 2. Sentence Splitting: Actively split sentences into separate statements when there's a change in speaker. Look for cues like pauses, speech direction changes, thought conclusions, questions, and answers.
> 3. Reasoning: Consider whether the lan-

guage is technical (suggesting a medical professional) or expresses personal experiences/emotions (suggesting the patient).
4. Look-Around Strategy: Analyze the five sentences before and after the current one to understand the conversation flow. Questions may be followed by answers, and concerns by reassurance.
5. Label with Justification: Assign a label 'Doctor' or 'Patient' to each sentence, providing a brief justification based on your analysis. Ensure each justification pertains to only one person.
6. Consistent Attribution: Maintain a thorough approach throughout the transcript, treating each sentence with equal attention to detail.
7. Extremely Granular Attribution: Break down the conversation into the smallest parts (question, answer, utterance) for clarity. Each clause should be precisely attributed to either the doctor or the patient, with no overlap in speaker identity.
8. Examples #1-10:
{examples}

Transcript:
{transcript_segments}

Expected output structure as a JSON array:

```
[
 {
    "sentence": "[reference]",
   "justification": "[justification]",
    "speaker": "Patient or Doctor"
 },
 {
    "sentence": "[reference]",
   "justification": "[justification]",
    "speaker": "Patient or Doctor"
 }
]
```

Here's how to approach the task step by step:

1. Contextual Reading: Analyze each sentence for potential transcription errors, considering medical terminology and context.
2. Medical Term Verification: Pay special attention to medical terms, medication names, and procedures that may have been misinterpreted.
3. Accent Consideration: Account for diverse accents and potential misinterpretations in the transcription.
4. Homophone Analysis: Identify and correct instances where similar-sounding words may have been confused.
5. Semantic Preservation: Ensure corrections maintain the original meaning while improving accuracy.
6. Contextual Coherence: Verify that corrected terms align with the medical context of the conversation.
7. Examples #1-10:
{examples}

Transcript:
{transcript_segments}

Expected output structure as a JSON array:

```
[
 {
  "reference_sentence": "[original]",
  "rationale": "[brief explanation]",
  "corrected_sentence": "[corrected]",
   "speaker": "Patient or Doctor"
 },
 {
  "reference_sentence": "[original]",
  "rationale": "[brief explanation]",
  "corrected_sentence": "[corrected]",
   "speaker": "Patient or Doctor"
 }
]
```

## A.6 Example Correction Prompt

For Correction, the following prompt template was utilized:

You are a helpful speech-to-text transcription assistant. Your task is to review and correct transcription errors in medical dialogues, focusing on accuracy and medical context.

## A.7 Intron Health WER Table

| LLM | STT | Method | WER |
|---|---|---|---|
| – | Gemini 1.5 Pro | ASR | 10.18% ± 1.81 |
| Gemini 1.5 Pro | Gemini 1.5 Pro | Corrected | 12.86% ± 2.09 |
| – | Amazon Transcribe Medical | ASR | 14.21% ± 2.89 |
| Claude 3.5 Sonnet | Gemini 1.5 Pro | Diarized | 15.32% ± 6.14 |
| GPT-4o | Soniox | Corrected | 16.28% ± 2.72 |
| Gemini 1.5 Pro | Soniox | Corrected | 16.33% ± 3.90 |
| – | Soniox | ASR | 17.54% ± 7.89 |
| Gemini 1.5 Pro | Amazon Transcribe Medical | Corrected | 18.24% ± 3.30 |
| Claude 3.5 Sonnet | Whisper Large 3 | Corrected | 18.34% ± 2.90 |
| Gemini 1.5 Pro | Whisper Large 3 | Corrected | 18.35% ± 4.24 |
| GPT-4o | Amazon Transcribe Medical | Corrected | 19.07% ± 3.24 |
| – | Whisper Large 3 | ASR | 19.33% ± 2.58 |
| Gemini 1.5 Pro | NVIDIA Canary 1b | Corrected | 21.03% ± 2.60 |
| GPT-4o | NVIDIA Canary 1b | Corrected | 21.22% ± 2.75 |
| Gemini 1.5 Pro | Azure STT | Corrected | 23.90% ± 2.13 |
| GPT-4o | Azure STT | Corrected | 24.11% ± 0.00 |
| – | NVIDIA Canary 1b | ASR | 24.16% ± 2.07 |
| – | Azure STT | ASR | 27.42% ± 5.69 |

Table 2: Comparison of WER ASR baselines and their top two LLM corrections in the Intron Health dataset.

## A.8 PriMock57 WER Table

| LLM | STT | Method | WER |
|---|---|---|---|
| – | Gemini 1.5 Pro | ASR | 9.34% ± 2.69 |
| – | Amazon Transcribe Medical | ASR | 9.98% ± 2.13 |
| Gemini 1.5 Pro | Gemini 1.5 Pro | Corrected | 10.43% ± 3.10 |
| Gemini 1.5 Pro | Soniox | Corrected | 11.74% ± 3.17 |
| – | Soniox | ASR | 12.09% ± 4.95 |
| Gemini 1.5 Pro | Amazon Transcribe Medical | Corrected | 12.10% ± 2.96 |
| Gemini 1.5 Pro | Azure STT | Corrected | 12.19% ± 2.58 |
| GPT-4o | Gemini 1.5 Pro | Corrected | 12.43% ± 4.88 |
| GPT-4o | Soniox | Corrected | 12.49% ± 3.76 |
| GPT-4o | Amazon Transcribe Medical | Corrected | 12.73% ± 2.48 |
| GPT-4o | Whisper Large 3 | Corrected | 13.60% ± 2.16 |
| Claude 3.5 Sonnet | Whisper Large 3 | Corrected | 13.74% ± 2.95 |
| GPT-4o | Azure STT | Corrected | 14.56% ± 3.05 |
| Gemini 1.5 Pro | NVIDIA Canary 1b | Corrected | 15.32% ± 4.23 |
| – | Whisper Large 3 | ASR | 16.76% ± 2.52 |
| Claude 3.5 Sonnet | NVIDIA Canary 1b | Corrected | 17.35% ± 4.55 |
| – | Azure STT | ASR | 18.43% ± 2.67 |
| – | NVIDIA Canary 1b | ASR | 20.65% ± 2.61 |

Table 3: Comparison of WER ASR baselines and their top two LLM corrections in the PriMock57 dataset.

## A.9 Fareez WER Table

| LLM | STT | Method | WER |
|---|---|---|---|
| – | Gemini 1.5 Pro | ASR | 6.19% ± 2.30 |
| – | Amazon Transcribe Medical | ASR | 6.33% ± 1.89 |
| Gemini 1.5 Pro | Amazon Transcribe Medical | Corrected | 7.42% ± 2.43 |
| – | Soniox | ASR | 8.14% ± 3.33 |
| Gemini 1.5 Pro | Soniox | Corrected | 8.28% ± 1.98 |
| Gemini 1.5 Pro | Azure STT | Corrected | 8.45% ± 2.15 |
| Claude 3.5 Sonnet | Gemini 1.5 Pro | Diarized | 8.73% ± 10.72 |
| Claude 3.5 Sonnet | Amazon Transcribe Medical | Diarized | 9.05% ± 3.33 |
| Claude 3.5 Sonnet | Soniox | Corrected | 9.15% ± 2.23 |
| Gemini 1.5 Pro | NVIDIA Canary 1b | Corrected | 9.78% ± 2.19 |
| Gemini 1.5 Pro | Gemini 1.5 Pro | Corrected | 9.79% ± 9.11 |
| Claude 3.5 Sonnet | Whisper Large 3 | Diarized | 10.34% ± 5.37 |
| Claude 3.5 Sonnet | Azure STT | Diarized | 11.08% ± 6.28 |
| Claude 3.5 Sonnet | NVIDIA Canary 1b | Diarized | 11.49% ± 5.33 |
| Claude 3.5 Sonnet | Whisper Large 3 | Corrected | 11.95% ± 2.47 |
| – | Azure STT | ASR | 14.29% ± 1.98 |
| – | Whisper Large 3 | ASR | 14.30% ± 2.79 |
| – | NVIDIA Canary 1b | ASR | 14.89% ± 2.64 |

Table 4: Comparison of WER ASR baselines and their top two LLM corrections in the Fareez dataset.

## A.10 Intron Health MC-WER Table

| LLM | STT | Method | WER |
|---|---|---|---|
| – | Gemini 1.5 Pro | ASR | 10.69% ± 3.72 |
| Gemini 1.5 Pro | Gemini 1.5 Pro | Corrected | 12.54% ± 4.38 |
| GPT-4o | Soniox | Diarized | 16.83% ± 5.78 |
| Claude 3.5 Sonnet | Soniox | Diarized | 17.14% ± 6.79 |
| Claude 3.5 Sonnet | Gemini 1.5 Pro | Corrected | 17.45% ± 6.27 |
| – | Soniox | ASR | 18.12% ± 7.60 |
| – | Whisper Large 3 | ASR | 19.46% ± 6.74 |
| Claude 3.5 Sonnet | Whisper Large 3 | Corrected | 19.95% ± 7.73 |
| GPT-4o | Whisper Large 3 | Corrected | 20.34% ± 6.89 |
| – | Amazon Transcribe Medical | ASR | 20.89% ± 8.19 |
| Gemini 1.5 Pro | Amazon Transcribe Medical | Corrected | 21.16% ± 8.58 |
| GPT-4o | Amazon Transcribe Medical | Corrected | 21.66% ± 7.11 |
| Claude 3.5 Sonnet | NVIDIA Canary 1b | Corrected | 25.52% ± 6.67 |
| Gemini 1.5 Pro | NVIDIA Canary 1b | Corrected | 25.76% ± 6.69 |
| – | NVIDIA Canary 1b | ASR | 30.16% ± 8.62 |
| GPT-4o | Azure STT | Corrected | 46.17% ± 11.86 |
| Gemini 1.5 Pro | Azure STT | Corrected | 47.82% ± 12.27 |
| – | Azure STT | ASR | 48.67% ± 12.60 |

Table 5: Comparison of MC-WER ASR baselines and their top two LLM corrections in the Intron Health dataset.

## A.11 PriMock57 MC-WER Table

| LLM | STT | Method | WER |
|---|---|---|---|
| – | Gemini 1.5 Pro | ASR | 7.38% ± 4.17 |
| Gemini 1.5 Pro | Gemini 1.5 Pro | Corrected | 8.13% ± 4.60 |
| Gemini 1.5 Pro | Soniox | Corrected | 8.35% ± 4.25 |
| – | Soniox | ASR | 8.38% ± 4.30 |
| Claude 3.5 Sonnet | Soniox | Diarized | 8.84% ± 4.34 |
| GPT-4o | Gemini 1.5 Pro | Corrected | 9.26% ± 5.29 |
| – | Amazon Transcribe Medical | ASR | 9.47% ± 4.23 |
| – | Whisper Large 3 | ASR | 9.57% ± 4.51 |
| Gemini 1.5 Pro | Amazon Transcribe Medical | Corrected | 10.40% ± 4.93 |
| Claude 3.5 Sonnet | Whisper Large 3 | Corrected | 10.63% ± 5.24 |
| GPT-4o | Whisper Large 3 | Corrected | 10.82% ± 5.49 |
| Gemini 1.5 Pro | Azure STT | Corrected | 11.62% ± 4.90 |
| – | Azure STT | ASR | 11.76% ± 5.06 |
| GPT-4o | Amazon Transcribe Medical | Corrected | 11.82% ± 5.58 |
| Gemini 1.5 Pro | NVIDIA Canary 1b | Corrected | 13.25% ± 5.97 |
| – | NVIDIA Canary 1b | ASR | 13.26% ± 5.84 |
| Claude 3.5 Sonnet | Azure STT | Corrected | 13.68% ± 6.60 |
| Claude 3.5 Sonnet | NVIDIA Canary 1b | Corrected | 15.22% ± 6.98 |

Table 6: Comparison of MC-WER ASR baselines and their top two LLM corrections in the PriMock57 dataset.

## A.12 Fareez MC-WER Table

| LLM | STT | Method | WER |
|---|---|---|---|
| – | Soniox | ASR | 3.28% ± 3.02 |
| Claude 3.5 Sonnet | Soniox | Diarized | 3.47% ± 3.20 |
| – | Gemini 1.5 Pro | ASR | 3.58% ± 3.47 |
| Gemini 1.5 Pro | Soniox | Corrected | 3.81% ± 3.58 |
| – | Azure STT | ASR | 4.22% ± 3.01 |
| Claude 3.5 Sonnet | Gemini 1.5 Pro | Corrected | 4.37% ± 3.40 |
| – | Amazon Transcribe Medical | ASR | 4.41% ± 3.45 |
| Gemini 1.5 Pro | Amazon Transcribe Medical | Corrected | 4.58% ± 4.11 |
| – | Whisper Large 3 | ASR | 5.00% ± 3.63 |
| Gemini 1.5 Pro | Azure STT | Corrected | 5.07% ± 3.49 |
| Claude 3.5 Sonnet | Amazon Transcribe Medical | Corrected | 5.16% ± 3.85 |
| Claude 3.5 Sonnet | Azure STT | Corrected | 6.53% ± 3.96 |
| Gemini 1.5 Pro | NVIDIA Canary 1b | Corrected | 6.84% ± 4.30 |
| Gemini 1.5 Pro | Whisper Large 3 | Corrected | 6.94% ± 4.30 |
| Claude 3.5 Sonnet | Gemini 1.5 Pro | Diarized | 7.07% ± 9.35 |
| Claude 3.5 Sonnet | Whisper Large 3 | Corrected | 7.25% ± 4.62 |
| – | NVIDIA Canary 1b | ASR | 7.35% ± 4.05 |
| Claude 3.5 Sonnet | NVIDIA Canary 1b | Corrected | 7.62% ± 4.25 |

Table 7: Comparison of MC-WER ASR baselines and their top two LLM corrections in the Fareez dataset.