# Increasing Usage of ChatGPT Mobile's Voice Input

**Team :** Product Team
**Contributors :** Pulkit kapoor
**Status :** review
**Launching on TBA**
**Resources :** milestone-1, milestone-2

**Problem Definition :** Many ChatGPT mobile users prefer typing instead of using the voice feature due to habit and past issues with misinterpretations and accuracy. This has led to low adoption and missed opportunities for engagement and growth in the mobile user base.

## What is the problem?

ChatGPT's voice system struggles to reliably interpret user intent during natural speech, leading to frequent inaccuracies and misinterpretations. Users also lack visibility into what the system understood or which previous voice inputs it is referencing, making it difficult to trust, verify, or correct responses. As a result, many users revert to typing and gradually stop using the voice feature altogether.

## Who is facing the problem?

Working professionals aged 21–34 who regularly use ChatGPT on mobile. This group represents one of the most active and engaged user segments for conversational and productivity use cases.

## What is the business value that will be unlocked by solving the problem?

Improving voice accuracy, transparency, and usability can increase the voice user base by at least 10%, drive higher engagement and session duration, improve retention, and strengthen monetization potential through premium conversions and feature adoption.

**How will the target users benefit if the problem is solved?**
Users will experience faster, more natural, and hands-free interactions, allowing them to multitask easily and express complex thoughts more comfortably on mobile.

**Why is it urgent to solve this problem now?**
India's voice assistant market is growing quickly as more people get comfortable using voice on their phones. Acting now will help ChatGPT attract early users, build strong usage habits, and stay ahead before competitors improve their voice features and gain more users.

**Goals**
The primary goal is to increase voice input usage on the ChatGPT mobile app by addressing accuracy issues, improving user trust, and reducing friction in adoption. The aim is to make voice the preferred mode of interaction for mobile users, particularly working professionals aged 21–34

**Functional Metrics**
**1) % of Users using Voice Input**
**Why it's important:** This directly measures adoption and indicates a successful behavioral shift from typing to speaking.

**2) Voice Feature Accuracy Rate:**
 Directly addresses the main pain point,misinterpretations and inaccuracies.

**3) Avg. Voice Input Duration per User:**
 Reflects engagement depth, longer sessions mean users find the feature more reliable and convenient.

**4) % of Avg. Time Spent on Voice Input:**
 Measures how integral voice input becomes to  overall usage, not just trial behavior.

**Non-Functional Metrics**

**1) Voice Response Latency :** Quick responses make conversations feel smooth and natural. If there is noticeable delay, users feel disconnected and are less likely to continue using the voice feature.

**2) Operational Reliability :** Voice interactions depend on real-time processing. High system availability ensures the feature works consistently when users choose to speak, building confidence and reducing frustration.

**3) Feature Discoverability:**Percentage of users who discover the voice feature during onboarding.

**Why These Metrics Matter**

- These metrics are important because they directly influence how often users engage with the voice feature, how long they continue using it, and whether they choose to upgrade to premium offerings.
- When voice accuracy and speed improve, users feel more confident and satisfied, which encourages repeat usage and builds habit. Increased adoption among working professionals drives deeper app engagement and strengthens the likelihood of premium conversion over time

**Validation of the Problem**

- 92% of the users knows chatgpt voice feature exist but only 14% are the daily users of the gpt voice feature
- 38.4% prefer using voice input while multitasking, 32.4% when tired of typing, 10.8% for quick questions, 13.5% have never used it, and a small percentage use it just to test the feature.
- The main reasons users do not use the GPT voice feature are because of habit followed by privacy concerns and a poor voice experience.

## Pain points

- A lack of trust caused by inconsistent voice accuracy leading to frustration and drop-off.

## Competitive insight

- Voice tools such as Google Assistant and Alexa/Siri have established strong standards for accuracy and reliability, shaping user expectations.

## Understanding the Target Audience

**Segment:** Working professionals aged 21–34

**Platform usage:** 52% use ChatGPT on mobile regularly

- **user Persona – Riya (24, Medical Intern)**
  - Uses ChatGPT regularly on her mobile for case discussions, quick medical research, and idea clarification.
  - Comfortable with typing and gradually stopped using the voice feature due to repeated inaccuracies in understanding her input.

**Student (Ajay, 21, Delhi):** Uses ChatGPT every day for assignments, translations, and quick help. Mostly types his questions on mobile between classes.

## Pain Points

- Voice input sometimes misunderstands what users say and gives incorrect responses.
- Users often need to repeat their questions, which interrupts the flow.
- Typing feels more dependable, even if it takes more time.
- Inconsistent accuracy reduces trust and causes users to stop using the voice feature altogether.

### Unmet Needs

● Voice recognition that accurately understands and responds correctly.
 ● A quick and smooth voice experience with little to no delay.
 ● Easy transition between typing and speaking without interruption.
 ● Hands-free use that makes multitasking more convenient.

## Solution 1:

 **Live Transcription with Edit Mode :** Real-time text display as users speak with instant correction capability, showing live transcription + tap-to-edit + confidence indicators for uncertain words.

Key Features:

- Words appear live on screen as user speaks
- Low-confidence words underlined in red
- Tap on any word to quickly edit it without having to restart the entire input.

## Solution 2:

**Voice Quality Improvements & Noise Handling :** Backend enhancements to improve voice recognition accuracy, specially tuned for Indian usage patterns, including diverse accents, background noise, and everyday mobile environments

Key feature:

● Enhanced recognition tuned for Indian English accents and speech patterns.
 ● Ability to understand and process incomplete or interrupted voice inputs.
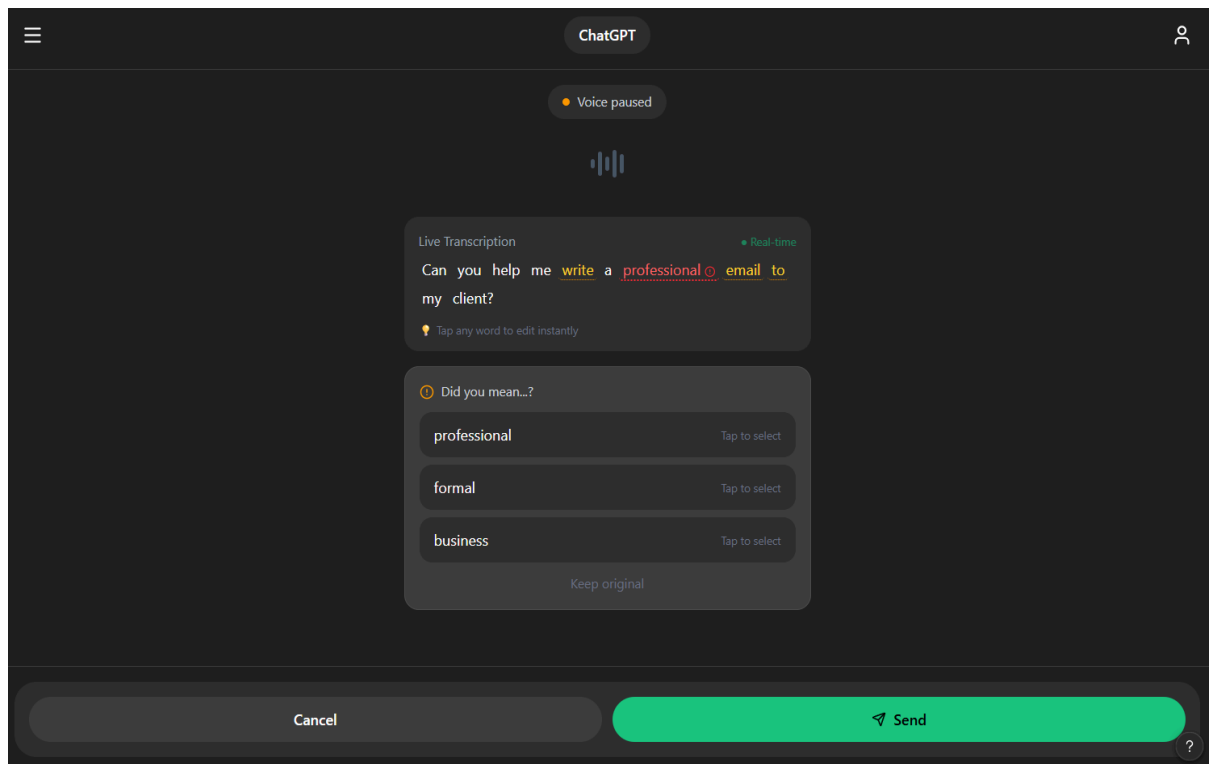 ● Intelligent retry suggestions like "Did you mean…?" with 2–3 options when the system is unsure.

.**Solution 3: Context-Aware Voice Mode:**The system adapts voice recognition based on user context (location, time, behavior patterns).
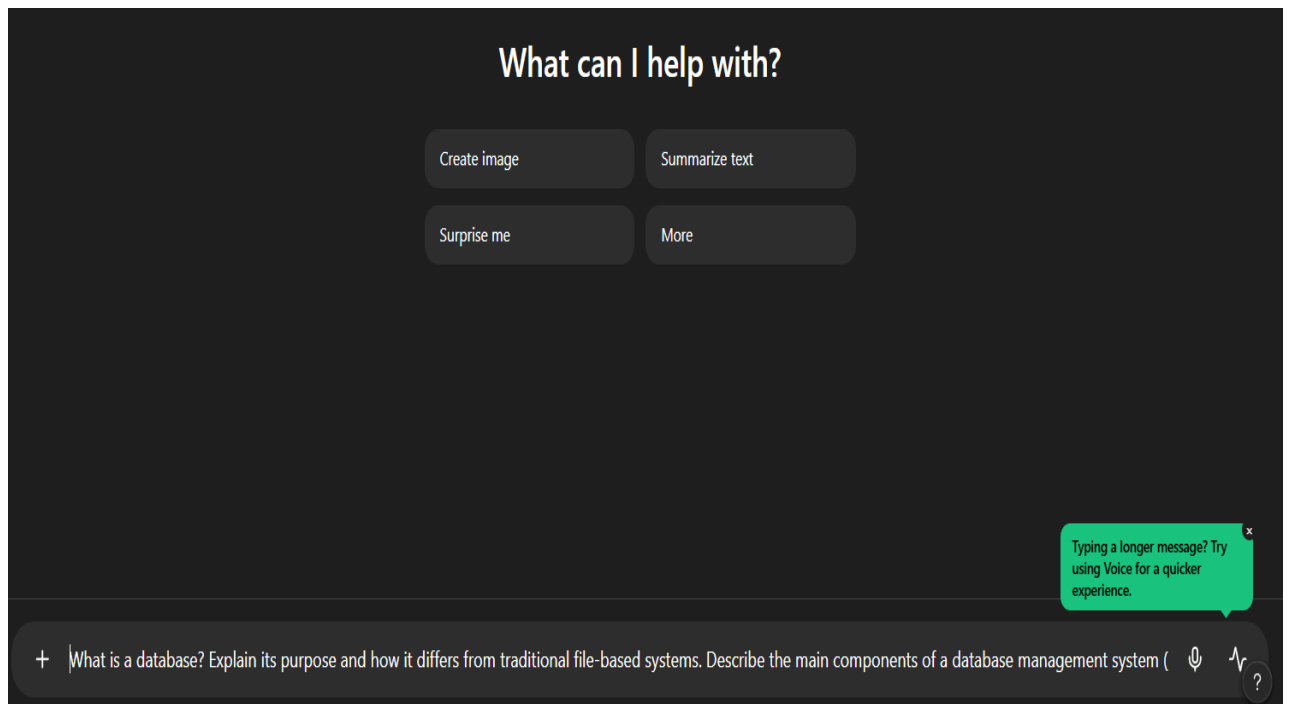
**Key Features:**
 ● Automatically switches to "High Accuracy Mode" in noisy environments.
 ● Detects medical/technical/business vocabulary based on usage history.
 ● Learns frequently used phrases and improves recognition over time.

**wireframe/prototype :**
**https://www.figma.com/make/8jv4j3zQxBMkI19RuIifgi/ChatGPT-Mobile-Voice-Interface?t=9ztIZxjpgSwwn9ZM-1**

**Launch Readiness :**

Key Milestones

• Design wireframes complete → Week 1-2

• Dev build with sticky mic & coach mark → Week 3-4

• QA + dogfooding → Week 5

• Soft launch in India (10% rollout) → Week 6

• General Availability: Week 8

**Internal Stakeholders**

PM: Defines success metics and monitors adoption.

Engineering: (Mobile Backend) Builds real-time subtitles, undo, pause.

Design: UX Designs subtitle UI, undo/pause controls, and color-coded confidence cues.

ML:  Speech Team Ensures subtitle accuracy, low-confidence detection.

QA :Tests across devices, network conditions, and languages

**Key Stakeholder Questions**

Engineering Question: How fast are subtitles shown?

Design / UX : Button placement in the use's field of view and range of motion.

ML : What confidence score threshold determines whether a word is marked as low confidence?

.

**Open Questions & Decisions Taken**
- What's the word limit for live transcript?
- Do we show confirmation for all voice inputs or only when accuracy is uncertain?

**Decisions Taken**
- There shouldn't be a hard word limit, but the UI should display only the most recent portion to maintain readability and performance
- We should show confirmation only when accuracy is uncertain, not after every voice input.

**Trade-offs**

- When confidence is low, prioritize accuracy by introducing a confirmation step instead of optimizing purely for conversational speed.