

# Increase the Usage of voice input in the ChatGPT mobile application

**Team:** Product Manager

**Contributors:** Dhaval Patel

**Status:** Brainstorming/In Review/In Development

**Launching on:** To Be Decided

**Resources:** [Systems Thinking & Mapping Outcomes](#) and [User Research & Problem Framing](#)

**User Research(Survey & Interviews):** N = 69 Sample Size, [Survey Link](#)

## Top 3 Key Issues Identified from Milestone 02 Survey

Barrier	What It Includes	Data Says	User Impact	Biz. Impact
<b>Inaccuracy</b>	Accent issues, misunderstood words, wrong transcription	Most repeated complaint across responses	Trust erosion; users stop after 1–2 failures	High drop-off post-trial
<b>Background Noise</b>	Recording disruption, recalibration, noise interference	Frequently reported in home, office, commute contexts	Feels unreliable in real-world use	Limits contextual adoption
<b>Perceived Slowness</b>	“Typing feels faster”, delay in response	Very commonly selected reason	Users revert to typing habit	Low repeat usage

## Problem Definition

### 1. What is the Problem?

Voice input in ChatGPT mobile is underutilized among Indian young working professionals is low compared to the use of voice search on **YouTube/Siri/Alexa(37%), Google Voice Search(22%), What's App Voice Note(40%) & Other apps(1%)** due to:

- Accent recognition issues (Indian English, Hinglish) - **32%**
- Low discoverability - **24%**
- Voice Input takes longer time than Text input - **16%**
- Social discomfort in public spaces - **15%**
- Habit bias toward typing - **14%**
- Low awareness of real use cases - **4%**

**This leads to:**

- Missed engagement opportunities
- Poor accessibility in **Tier-2/3 markets**
- Competitive disadvantage in voice-first markets

### 2. Who is facing the Problem?

This issue is mainly faced by the following people :

- Working Professionals - **(47/69) - 68%**
- Having age group of 25-34 years - **(36/69) - 52%**
- Belong to the Tier2/3 cities - **(35/69) - 50%**
- Using the Mobile Phone to prompt several times a day - **(23/69) - 33%**
- Prefers to type mostly - **(46/69) - 67%**

Who are generally using ChatGPT for getting suggestions to work-related problem statements, creating a resume, writing code, learning English, reviewing a manually written document, and taking strategic business decisions.

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

Improving voice adoption in India unlocks growth in **engagement, retention, and subscription revenue**. With **400M+ smartphone** users in **Tier 2/3 cities**, capturing even **25% in the next 3 years** can significantly accelerate **Go and Pro subscriptions**.

Better discoverability and accent accuracy will:

- Increase voice adoption
- Improve premium conversions
- Drive longer, more frequent sessions
- Strengthen retention and lifetime value

Even a **+0.5 increase** in daily sessions per active user can generate millions of additional interactions at scale.

**Higher voice adoption → higher engagement → stronger retention → greater revenue.**

### 4. How will the target users benefit if the problem is solved?

Solving this problem will make ChatGPT faster, more accurate, and easier to use through voice.

**Users will:**

- Get more accurate, personalized, and to-the-point answers
- Solve queries faster without typing
- Capture ideas instantly with hands-free convenience
- Experience natural, engaging voice conversations
- Use ChatGPT as a daily voice companion for work, learning, translation, search, and multitasking

**Result:** Faster, hassle-free, and more reliable AI assistance that saves time and boosts productivity.

### 5. Why it is urgent to solve this problem now?

Voice adoption in India is growing rapidly, especially in Tier 2/3 cities where AI usage is still forming habits. This is a narrow window to shape behavior.

Competitors are investing heavily in voice-first AI. If ChatGPT doesn't act now, students and young professionals may build loyalty elsewhere.

**Acting early will:**

- Capture emerging voice habits
- Secure long-term user loyalty
- Strengthen competitive positioning
- Establish ChatGPT as the default voice AI in India

Delay risks losing a high-growth, voice-first market to faster-moving competitors.

# Goals

**Primary Objective:** Make voice input relevant, trusted, and habit-forming for learners and young professionals in India.

- **Drive Awareness & Adoption**

Increase discoverability and position voice input as a go-to tool for study, productivity, and brainstorming to boost activation and repeat usage.

- **Improve Trust & Accuracy**

Enhance recognition for Indian accents and Hinglish while strengthening user confidence through visible feedback and higher perceived accuracy.

- **Enable Natural, Contextual Conversations**

Improve context continuity and encourage deeper, more natural voice interactions that feel seamless and intelligent.

- **Differentiate the Voice Experience**

Introduce learner-focused, voice-first features that create a distinct and competitive advantage in the market.

**Success Outcome:** Higher voice adoption, longer sessions, stronger retention, and sustained usage habits.

## Functional Metrics (Quantifiable Targets)

1. **Voice Awareness Rate**

Increase user awareness of the voice feature from **50% → 80%**.

2. **Monthly Active Voice Users (MAVU)**

Grow active monthly voice users from **10% → 25% of DAUs**.

3. **Working Professional Voice Adoption**

Increase verified professional voice usage to **20%+ monthly penetration**.

4. **Voice Accuracy Satisfaction Score**

Improve perceived voice accuracy satisfaction by **+25%**.

5. **Context Retention Accuracy**

Achieve **90%+ session accuracy** in correctly referencing prior voice interactions.

6. **Accessibility Adoption Impact**

Drive a **+10% uplift** in adoption among Tier 2/3 and accessibility-driven users.

**Success Indicator:** Higher activation → higher repeat usage → longer sessions → measurable retention lift.

## Non-Functional Metrics (Performance & Experience Targets)

1. **Response Latency**

Maintain average voice response time at **< 1.5 seconds**.

2. **Error Correction Rate**

Reduce transcription correction rate by **30%**.

3. **Accessibility Navigation Success**

Achieve **≥ 90% task completion rate** via voice navigation.

4. **CSAT (Voice Experience)**

Increase voice-specific CSAT score to **≥ 4.5/5**.

## 5. NPS (Voice Feature – Professionals)

Improve Net Promoter Score for voice users to **+40 or higher**.

**Success Benchmark:** Fast, accurate, and reliable voice interactions that feel seamless, trustworthy, and recommendation-worthy.

## Non-Goals

### 1. No Core Model or Voice Stack Rebuild

No new speech-recognition ML models, backend infrastructure overhaul, or deep accuracy re-engineering in this phase.

### 2. No Subscription or Acquisition Focus

Conversion to Go/Pro plans and new user acquisition (installs, impressions, sign-ups) are not primary objectives.

### 3. No Hardware or Third-Party Integrations

No integrations with Alexa, Siri, smart speakers, external mics, or hotword activation (“Hey ChatGPT”).

### 4. No Expansion Beyond Target Segment

Scope is limited to Indian individual learners and professionals; institutional use cases, non-English language expansion, and global rollout are out of scope.

## Validation of the problem

- Survey Sample Size : 69 Users

Insights from the user research, [survey/ interviews](#)

#	Section	Key Data	Business Impact
1	Market	400M+ Tier 2/3 smartphone users	High-growth voice opportunity
2	Awareness	100% aware	Awareness ≠ usage
3	Trial	~65–70% tried once	Curiosity present
4	Regular Use	~20–25% active users	Large retention gap
5	Default Behavior	Mostly typing	Strong habit bias
6	Top Barriers	Accent issues, Noise, Typing faster, Social discomfort	Trust & UX friction
7	Accuracy Expectation	4–5/5 importance	Reliability critical
8	Experience Rating	Frequently 2–4/5	Expectation gap
9	Competitive Signal	Heavy voice use on WhatsApp, Google, Siri	ChatGPT not default voice tool
10	Growth Lever	+0.5 sessions/day uplift potential	Millions of additional interactions


# Competitive Insights

#	Platform	What Works	GPT Gap	Opportunity
1	WhatsApp	Hold-to-record	Low visibility	Gesture UI
2	Google Assistant	Always-on trigger	No contextual cues	Context nudges
3	Siri	Natural conversational flow	Low productivity focus	Task modes

## Understanding the target audience

Insights from the user research, [survey/ interviews](#)

### User Persona



**Tier-01 : Productivity Focussed Pro.**  
**Name :** Aditya Mehta  
**Role :** Product Manager, Bangalore  
**Age :** 27 | 4 YOE

**Key Traits:**

- Desk-based, attends frequent meetings
- Works under tight deadlines, managing multiple projects
- Commutes daily, often uses travel time to catch up on tasks

**Pain Points:**


- Writing emails, reports, meeting notes takes hours
- Capturing ideas during fast-paced meetings or commutes is difficult

**ChatGPT Voice Use Cases:**

- Dictate emails or reports while commuting
- Brainstorm and draft ideas hands-free
- Convert voice notes into actionable to-do lists

**Value:**

- **Boosts productivity & ROI** – saves time, fits multitasking workflows, even on the go



**Tier-02/3 : Accessibility Focussed Pro.**  
**Name :** Mitva Arora  
**Role :** Marketing Executive, Indore  
**Age :** 24 | 2 YOE

**Key Traits:**

- Frequently on mobile while moving between meetings or locations
- Struggles with typing, often switches between English and Hinglish

**Pain Points:**

- Typing long emails, notes, or content is cumbersome
- Ideas often get lost while on the go

**ChatGPT Voice Use Cases:**

- Quick dictation of notes, messages, or content in Hinglish
- Capture ideas or creative drafts hands-free
- Instant summaries /structured outputs without typing

**Value:**

- **Reduces friction & improves accessibility** – makes ChatGPT effortless for on-the-go users

[From the user research: survey conducted during milestone-02](#)

**Core Takeaway:** High awareness + high trial + low retention = Experience-driven adoption gap.

## Target Audience Overview

#	Category	Details	Insight
1	Primary Segment	Working Professionals (25–40 yrs)	High mobile usage, productivity-focused
2	Use Case	Work-related problem solving, brainstorming, decision support	Voice as productivity accelerator
3	Market Opportunity	5% adoption ≈ 5M users (India)	Large scalable impact

# Key Persona

#	Persona	Age	Voice Use Context	Core Need
1	Young Analyst	25–30	Brainstorming while multitasking	Speed & convenience
2	Mid-Level Manager	30–35	Structured inputs during meetings	Clarity & efficiency
3	Senior Specialist	35–40	In-depth validation & complex problem solving	Accuracy & depth

## Solution

### Key Barriers Summary

#	Barrier	What Users Said	Data Pattern
1	Accent Accuracy Gap	“Didn’t pick my accent.”“Misunderstood words.”	High repetitionTrust drop
2	Noise & Social Discomfort	“Background noise issue.”“Awkward speaking aloud.”	Context frictionOffice avoidance
3	Typing Habit Bias	“Typing feels faster.”“Took longer than typing.”	Post-trial dropHabit inertia

## Solution Directions and Evaluation (RICE Framework)

RICE Score = (Reach×Impact×Confidence) / Effort

#	Solution Direction	What It Solves	Description	RICE Breakdown
1	Accent-Adaptive Feedback System	Accuracy + Trust	Real-time transcript with low-confidence highlights and “Tap to Fix” correction before sending	Reach: 60% Impact: 3 Confidence: 3 Effort: 2 <b>RICE Score: 270</b>
2	Indian Accent Fine-Tuning (ASR)	Core Accuracy	Fine-tune ASR model on Indian English + Hinglish datasets; reduce word error rate	Reach: 80% Impact: 4 Confidence: 3 Effort: 4 <b>RICE Score: 240</b>
3	Hinglish Detection Layer	Mixed-Language Errors	Detect hybrid Hindi-English tokens using phonetic + contextual modeling	Reach: 55% Impact: 3 Confidence: 2 Effort: 3 <b>RICE Score: 110</b>

4	<b>Noise Stabilization Engine</b>	Environmental Accuracy	Adaptive noise suppression for office/public transport environments	Reach: 70% Impact: 3 Confidence: 2 Effort: 3 <b>RICE Score: 140</b>
5	<b>Quick Retry Voice Loop</b>	Trust Recovery	Instant “Re-speak last phrase” prompt without resetting session	Reach: 40% Impact: 2 Confidence: 3 Effort: 1 <b>RICE Score: 240</b>
6	<b>Personal Accent Memory Profile</b>	Long-Term Accuracy	System adapts to user-specific speech patterns over time	Reach: 50% Impact: 4 Confidence: 2 Effort: 4 <b>RICE Score: 100</b>

I am proposing a high-level solution focused on closing the **Accent Accuracy Gap**, built using a working-backwards approach:

- **UX Research** (accent + first-use testing)
- **Voice of Customer loops** (continuous correction signals)
- **Direct User Inputs** (correction-based learning)
- **Robust Accuracy Metrics** (word confidence tracking)
- **A/B Testing** (confidence UI vs no-confidence UI)
- **Competitive Benchmarking** (Google/WhatsApp voice patterns)

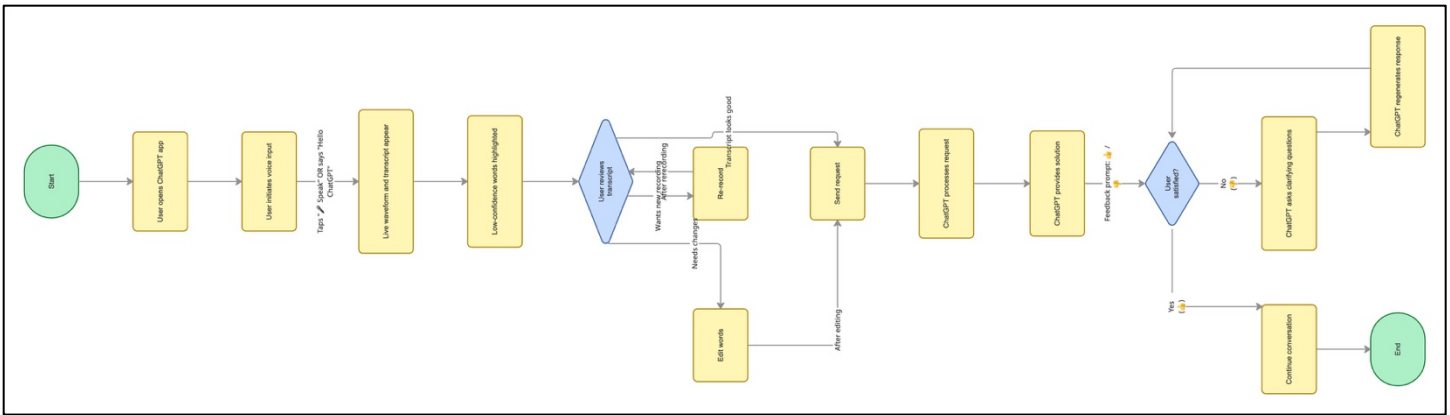
The goal is to improve both **real accuracy** and **perceived trust** in voice interactions.

## Proposed Key Features

- **Accent-Aware Mode** → Optimize recognition for Indian English + Hinglish users.
- **Live Confidence Transcript** → Show real-time transcript with low-confidence word highlights before sending.
- **Tap-to-Correct Layer** → Allow users to edit flagged words instantly before submission.
- **Noise-Stabilization Engine** → Adaptive background noise filtering for office/public environments.
- **Hold-to-Talk Interaction** → Familiar, WhatsApp-style voice gesture to reduce friction.
- **Persistent Context Memory** → Retain relevant session context to avoid repeated inputs.
- **Critical Questioning Mode** → Auto-trigger clarifying questions to improve solution quality.

## User Flow

- User starts voice session → live transcript appears with low-confidence words highlighted
- User edits or re-records before sending → ChatGPT processes and responds
- If response is unclear → ChatGPT asks clarifying questions and regenerates
- Continuous correction loop improves accuracy, trust, and conversation flow



*Made using : Miro*

## Key Logic (Backend & System Changes)

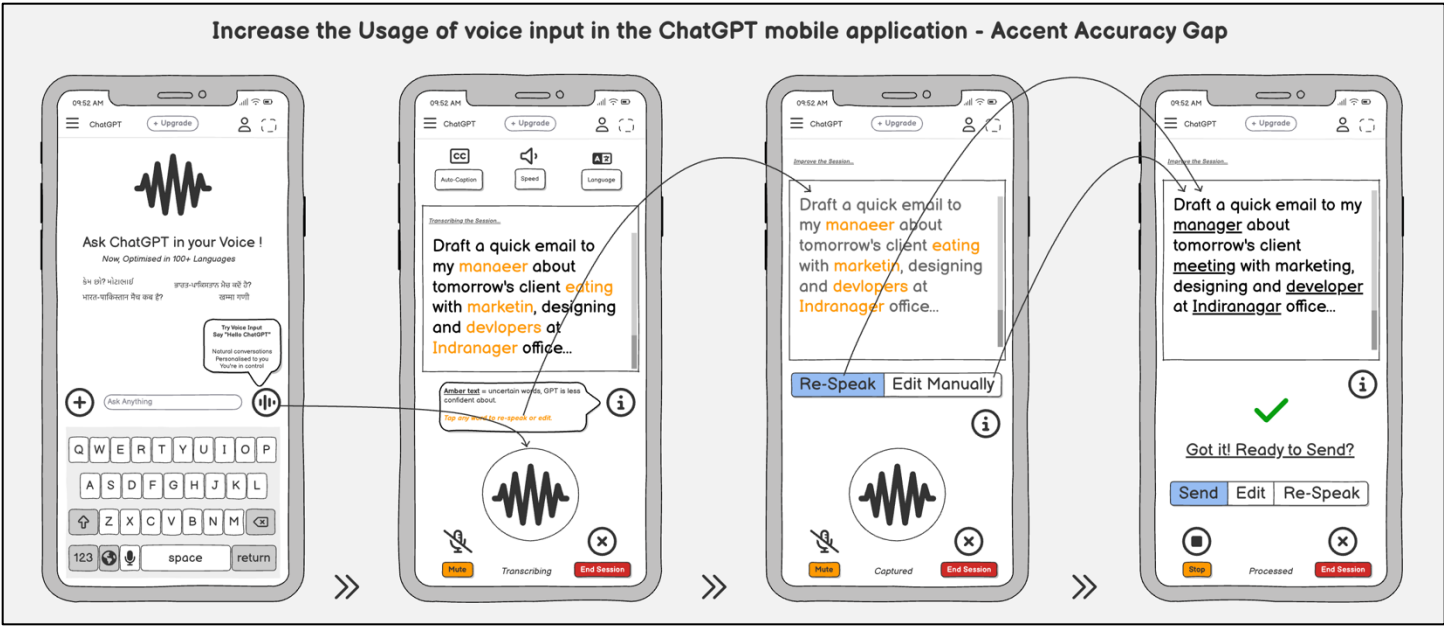
- **Accent Reduction**
  - Fine-tune ASR for Indian accents
  - User-level speech profiling
  - Correction-based adaptation
- **Confidence Visibility**
  - Word-level probability scoring
  - Low-confidence word highlighting
  - Smart alternative suggestions
- **Hinglish Handling**
  - Mixed-language detection
  - Hybrid phonetic + contextual matching
- **Noise Resilience**
  - Real-time voice isolation
  - Adaptive noise suppression
- **Trust Calibration**
  - Capture corrections + 🍌/🍌 signals
  - Continuous confidence recalibration

## Wireframing

The low-fidelity wireframes for illustrating the solution are as follows:

1. **Home Screen Prompt** – ChatGPT highlights the mic and encourages users to start a voice query
2. **Live Transcription** – User speaks and the system shows real-time transcript with low-confidence words marked
3. **Correction Options** – User can quickly re-speak or manually edit uncertain words before submission
4. **Review & Confirm** – User verifies corrected text and sends for final processing





Made using: [Balsamiq](#)

# Launch Readiness — Accent Accuracy Improvement

## Key Milestones & Timeline

Phase	Timeline	Milestone	Output
Phase 1	Week 1–3	Design Finalized	Accent-aware UX, Confidence UI, Edit Layer
Phase 2	Week 4–9	Development Complete	ASR fine-tuning, Hinglish detection, Confidence scoring
Phase 3	Week 10	Internal Dogfooding	Native Indian accent testing
Phase 4	Week 11	QA & Performance Testing	Latency < 1.5s, error rate validation
Phase 5	Week 12	Phased Rollout	1% → 10% → 50% → 100% traffic (IN regions)

## Launch Checklist

### Product

- Accent fine-tuned ASR deployed
- Confidence scoring threshold validated
- Hinglish detection stable
- Real-time transcription stable

### Engineering

- Latency < 1.5s
- Crash-free rate > 99%
- Noise suppression validated

### Data & Analytics

- Voice activation tracking live

- Correction rate tracking live
- Repeat voice usage metric active

Support & Ops

- Help Centre updated
- Support team trained on voice queries
- Escalation path defined

Internal Stakeholders

Team	Responsibility
Product	KPI ownership & rollout decision
ASR / ML	Accent fine-tuning
Mobile Eng	UI & integration
Backend Eng	Confidence scoring infrastructure
Data Science	Experiment analysis
QA	Stability & validation
Customer Support	User issue handling
Marketing	Launch communication
DevOps	Traffic ramp & monitoring

Experimentation Plan

Experiment	Hypothesis	Primary Metric
Confidence Highlight UI	Visible accuracy indicators increase user trust	Repeat Voice Usage ↑
Accent Badge (“Optimized for Indian Speech”)	Localized framing boosts activation	Voice Activation Rate ↑
Guided First Voice Session	Structured first-use reduces abandonment	Post-Trial Drop-off ↓
Hold-to-Talk vs Tap-to-Talk	Familiar interaction model increases usage	Sessions per User ↑

Success Metrics — Accent Accuracy Gap Resolution

KPI	Baseline	Target	Focus Area
Word Error Rate (Indian Accents)	—	↓ 30%	Core Accuracy
Misrecognition Rate (Flagged Words)	High	↓ 40%	Accuracy
Accent Accuracy CSAT	3.7 / 5	≥ 4.5 / 5	Trust

<b>First-Session Success Rate</b>	45%	$\geq 70\%$	First-Use Trust
<b>Repeat Voice Usage (30D)</b>	< 5%	$\geq 20\%$	Retention
<b>Post-Trial Drop-off</b>	High	$\downarrow 25\%$	Adoption
<b>Correction Success Rate (Amber Words)</b>	—	$\geq 85\%$	Perceived Accuracy
<b>Retry Loop Recovery Rate</b>	—	$\geq 75\%$	Trust Recovery
<b>Voice Session Duration</b>	0.6 min	$\geq 1.5$ min	Engagement

## Open Questions & Decisions Taken

Question	Decision
<b>Confidence UI in all sessions?</b>	Yes — default ON for trust building.
<b>Numeric confidence scores?</b>	No — color cues only.
<b>Threshold for highlighting?</b>	< 0.8 confidence score.
<b>Per-user or global learning?</b>	Start per-user; expand later.
<b>Mandatory accent calibration?</b>	No — optional.
<b>Fallback on low confidence?</b>	Trigger Quick Retry Loop.
<b>Voice default ON?</b>	No — keep opt-in.
<b>Support non-English in V1?</b>	No — English (Indian accent) only.

## Descoped (For MVP)

- Full offline voice processing
- Multi-language (Hindi, Tamil, etc.) expansion
- Hot word activation (“Hey ChatGPT”)
- Voice-to-voice natural speech synthesis upgrades
- Advanced global accent clustering
- Proactive smart voice nudges (awareness track handled separately)

## Trade-offs Made

- UI transparency > Deep retraining
- English-first > Multi-language breadth
- Per-user learning > Global complexity
- Accuracy > Slight latency increase
- Habit formation > Feature richness

### Important Links

1. User Research : [https://docs.google.com/forms/d/1-CElSTAowwrOpjXwq4Rw\\_I3fwy2InaG3D06XHYP5-GA/edit#responses](https://docs.google.com/forms/d/1-CElSTAowwrOpjXwq4Rw_I3fwy2InaG3D06XHYP5-GA/edit#responses)
2. Confluence : <https://dhavalpatelpm.atlassian.net/wiki/x/AgAE>
3. Miro : <https://miro.com/app/board/uXjVGcQnuWo=/>
4. Balsamiq : <https://balsamiq.cloud/sgyu8ns/pt28pme/r2278>