

# Weekly Report 9

## Google Chirp Streaming Speech Recognition

Timothé Berland

December 3, 2025

## Introduction

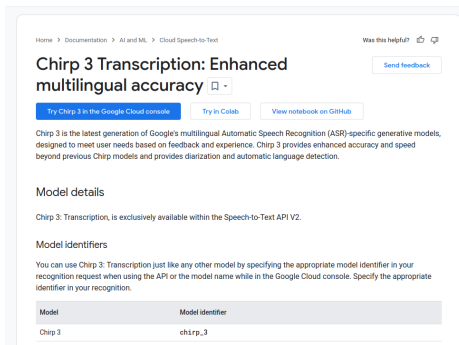
- Study Google Cloud Speech-to-Text V2 streaming API
- Explore the **Chirp model** for real-time transcription
- Integrate streaming transcription into previous Whisper-based system

## ① Overview of Chirp Streaming

## ② Streaming workflow

## ③ Model features & limitations

- Synchronous recognition
- Batch recognition
- **Streaming recognition**



The screenshot shows the Google Cloud documentation page for Chirp 3 Transcription. The page title is "Chirp 3 Transcription: Enhanced multilingual accuracy". It includes navigation links for "Try Chirp 3 in the Google Cloud console", "Try in Colab", and "View notebook on GitHub". The text describes Chirp 3 as the latest generation of Google's multilingual ASR-specific generative models, designed for enhanced accuracy and speed. It also mentions that Chirp 3 provides diarization and automatic language detection. Under the "Model details" section, it states that Chirp 3 Transcription is exclusively available within the Speech-to-Text API V2. The "Model identifiers" section explains that users can use Chirp 3 Transcription by specifying the appropriate model identifier in their recognition request. A table lists the model identifier for Chirp 3.

Model	Model identifier
Chirp 3	chirp_3

- Streaming mode allows **real-time partial and final transcripts**

① Overview of Chirp Streaming

② Streaming workflow

③ Model features & limitations

# Audio capture workflow

- Continuous audio capture from microphone
- Audio split into **small chunks** sent to API
- Same structure as previous Whisper Streamlit implementation

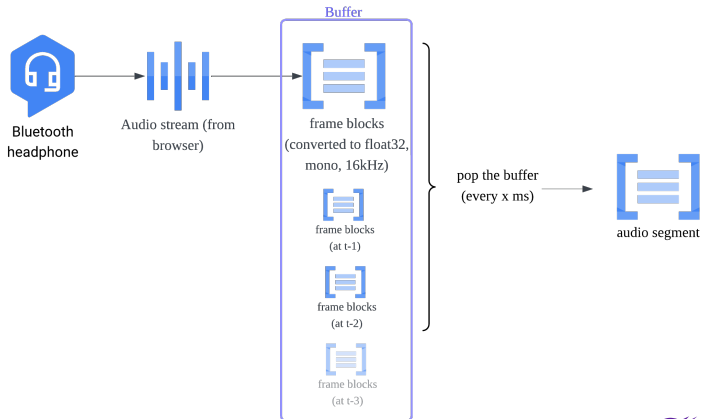


Figure: Real-time audio capture

# Streaming request workflow

- First message: recognizer name + streaming configuration
- Following messages: only audio bytes
- Chunk size: <15 KB per request
- Model returns:
  - **Interim results** (can change)
  - **Final results** (stable)

```
config_request = cloud_speech.StreamingRecognizeRequest(  
    recognizer=f"projects/{PROJECT_ID}/locations/{REGION}/recognizers/_",  
    streaming_config=streaming_config,  
)  
  
def request_generator():  
    yield config_request  
    yield from google_audio_generator(ctx)  
  
responses = client.streaming_recognize(requests=request_generator())
```

Figure: Chirp streaming inference

① Overview of Chirp Streaming

② Streaming workflow

③ Model features & limitations



- Punctuation, capitalisation, timestamps
- Language identification
- Stability score (0–1) for interim results
- Audio quality crucial: prefer clean PCM input

## Next steps

- Complete API implementation
- Optimize audio buffer
- Add optional translation (with a model like M2M100)