# Pratim Bhosale

Developer Advocate

Backend Developer

Technical Writer

# Content

✓ Speech to Text transcription APIs are expensive

✓ What is Whisper & Whisper.cpp?

✓ What are Go bindings?!

✓ Transcription Usecases

✓ Demo

@BhosalePratim

# Speech to Text transcription APIs are expensive

# Pricing table

The prices in the table below apply to minutes of audio processed per month.

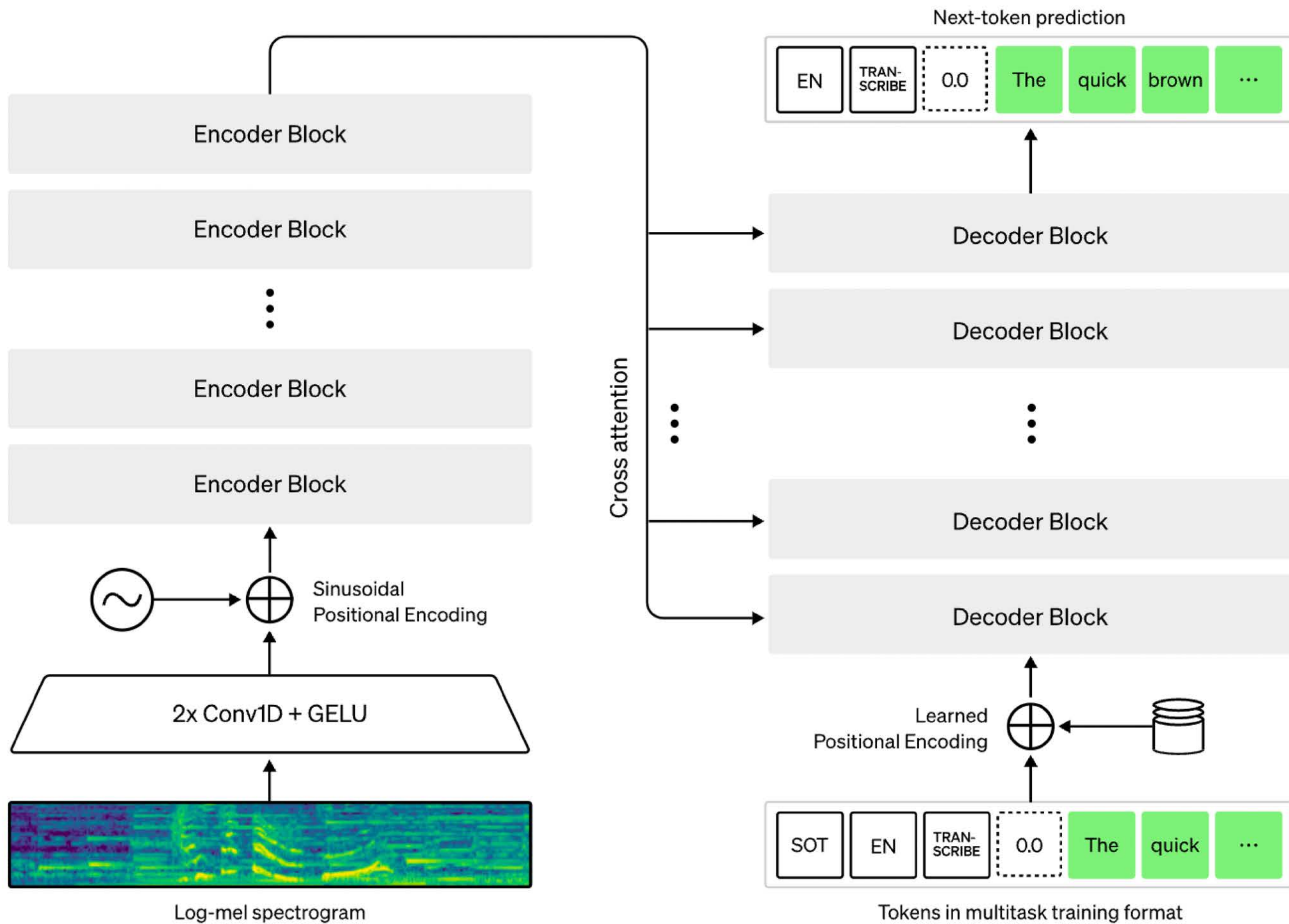| Category | Models | Pricing | |
|---|---|---|---|
| | | 0-60 Minutes/Month | Over 60 Minutes/Month |
| Speech Recognition (without data logging - default) | Standard[1] | Free | $0.024 / minute ** |
| | Medical[2] | Free | $0.078 / minute ** |
| Speech Recognition (with data logging opt-in) | Standard[1] | Free | $0.016 / minute ** |

**Google's Speech-to-Text API pricing**

| Tier | Volume (minutes/month) | Standard Batch Transcription ($/minute)* |
|---|---|---|
| T1 | First 250,000 minutes | $0.02400 |
| T2 | Next 750,000 minutes | $0.01500 |
| T3 | Next 4,000,000 minutes | $0.01020 |
| T4 | Over 5,000,000 minutes | $0.00780 |

**Amazon's speech to text api cost**

# What is Whisper and Whisper.cpp

## Whisper is the most underrated OpenAI model

Next-token prediction

| EN | TRAN-SCRIBE | 0.0 | The | quick | brown | ... |

Encoder Block

Encoder Block

Encoder Block

Encoder Block

Decoder Block

Decoder Block

Decoder Block

Decoder Block

Cross attention

Sinusoidal Positional Encoding

2x Conv1D + GELU

Log-mel spectrogram

Learned Positional Encoding

| SOT | EN | TRAN-SCRIBE | 0.0 | The | quick | ... |

Tokens in multitask training format

- A lightweight implementation of OpenAI's Whisper speech-to-text model
- Compatible with the Go stack, thanks to its Go bindings.
- It's a cost-effective alternative to Google, Amazon and IBM APIs.
- It's Open Source.
- Can be embedded.

# What are Go bindings?

- Go bindings are a way to interface and interact with a library or a package is written in another programming language (such as C or C++) from within Go code.

- They serve as a bridge between the Go code and the foreign language code, allowing the two to communicate and share data with each other.

- The whisper.cpp library takes care of the Go bindings.

# Getting started with Whisper.cpp

# Initializing the transcription model using the whisper.New()

```go
func transcribe(audioFilename string, modelPath string) error {
    // load whisper model
    model, err: = whisper.New(modelPath)
    if err ≠ nil {
        return fmt.Errorf("failed to load model: %w", err)
    }
    defer model.Close()

    log.Println("Successfully loaded the model")

    // Create processing context
    context, err: = model.NewContext()
    if err ≠ nil {
        return fmt.Errorf("failed to create context: %w", err)
    }
    fh, err: = os.Open(audioFilename)
    if err ≠ nil {
        return fmt.Errorf("failed to open audio file: %w", err)
    }
    defer fh.Close()

    if err ≠ nil {
        return fmt.Errorf("failed to open audio file: %w", err)
    }
```

**Decoding the WAV file and processing the context.**

```go
var data[] float32
    // Decode the WAV file - load the full buffer
dec: = wav.NewDecoder(fh)
if buf, err: = dec.FullPCMBuffer();
err ≠ nil {
    return err
} else if dec.SampleRate ≠ whisper.SampleRate {
    return fmt.Errorf("unsupported sample rate: %d", dec.SampleRate)
} else if dec.NumChans ≠ 1 {
        return fmt.Errorf("unsupported number of channels: %d", dec.NumChans)
    } else {
        data = buf.AsFloat32Buffer().Data
    }
    // Process the data
if err: = context.Process(data, nil);
err ≠ nil {
    return err
}

// Print out the results
for {
    segment, err: = context.NextSegment()
    if err == io.EOF {
        return nil
    } else if err ≠ nil {
        return err
    }
    fmt.Fprintf(w, "[%6s→%6s]", segment.Start.Truncate(time.Millisecond),
segment.End.Truncate(time.Millisecond))
    fmt.Fprintln(w, " ", segment.Text)
}
return nil
}
```

# Containerizing Whisper using Docker

```
# Install whisper
RUN git clone https://github.com/ggerganov/whisper.cpp.git &&\
    cd whisper.cpp && make &&\
    make libwhisper.so libwhisper.a &&\
    cp whisper.h /usr/local/include &&\
    cp ggml.h /usr/local/include &&\
    cp libwhisper.a /usr/local/lib &&\
    cp libwhisper.so /usr/local/lib &&\
    cd ..
```

# Usecases

- Transcribing meetings
- Audio Chatbots
- Automatic Translations
- Video Subtitles

# Demo

**Let's take a short video from youtube and convert it into text.**

Thank you!