# Whisper

It is a toolkit can do:

## English transcription

🔊 Ask not what your country can do for …

● Ask not what your country can do for …

## Any-to-English speech translation

🔊 El rápido zorro marrón salta sobre …

● The quick brown fox jumps over …

## Non-English transcription

🔊 언덕 위에 올라 내려다보면 너무나 넓고 넓은 …

● 언덕 위에 올라 내려다보면 너무나 넓고 넓은

## No speech

🔊 (background music playing)

● ∅

# Overview

## Context

### Paper Title

"Robust Speech Recognition by Large Scale Weak Supervision"

### Key Feature

- Generalization capabilities
- Transfer setting in zero-shot
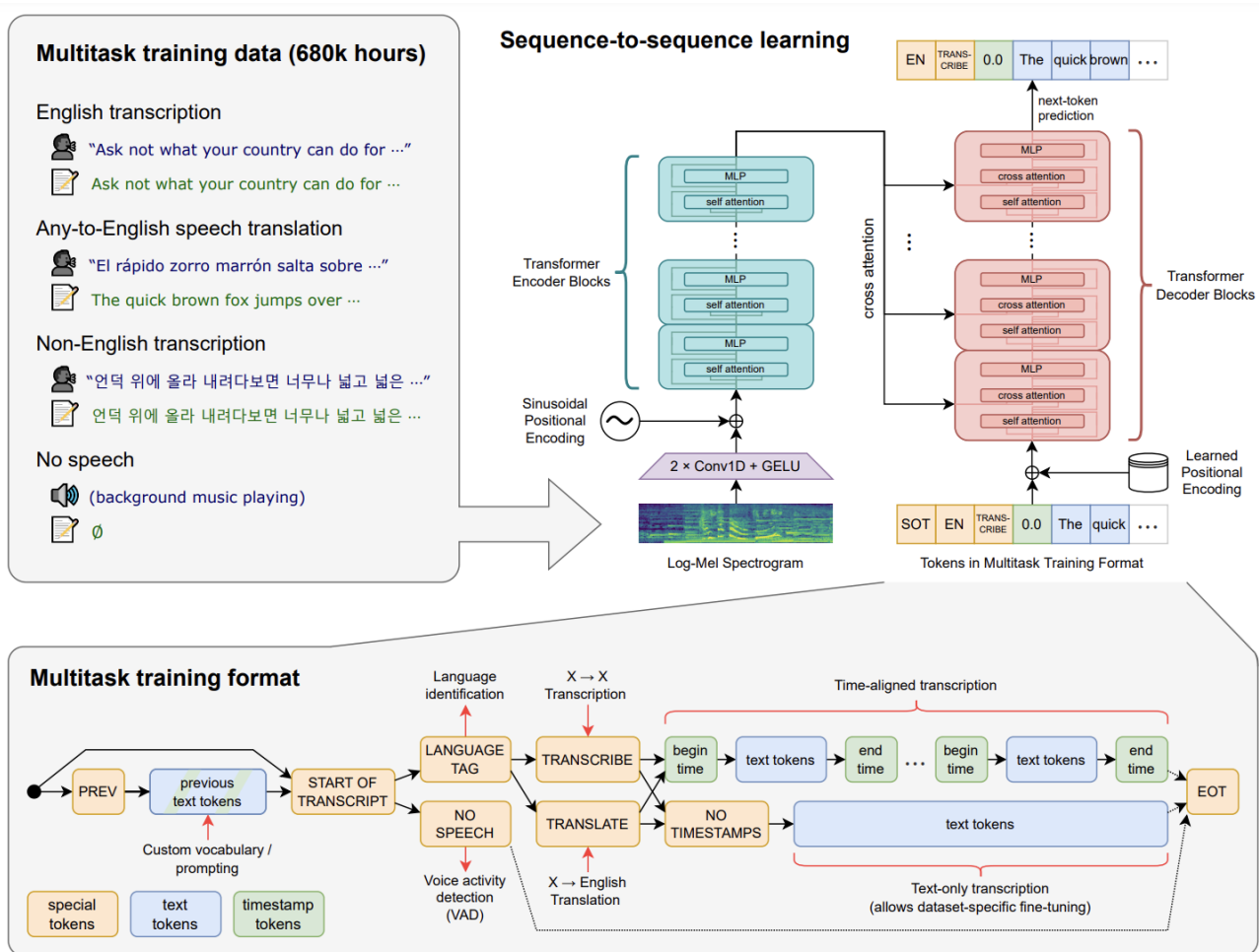- No need for fine-tuning

## Problem

- Unsupervised models: need fine-tuning to be truly effective

- Supervised models: robust but not enough

## Approach

- Scaling weakly supervised speech recognition the next order of magnitude to 680,000 hours of labeled audio data

- Without any need for fine-tuning
- The dataset is a global affair

# Architecture

## Overview of Architecture

**Multitask training data (680k hours)**

**English transcription**
"Ask not what your country can do for ⋯"
Ask not what your country can do for ⋯

**Any-to-English speech translation**
"El rápido zorro marrón salta sobre ⋯"
The quick brown fox jumps over ⋯

**Non-English transcription**
"언덕 위에 올라 내려다보면 너무나 넓고 넓은 ⋯"
언덕 위에 올라 내려다보면 너무나 넓고 넓은 ⋯

**No speech**
(background music playing)
∅

**Sequence-to-sequence learning**

| EN | TRANS-CRIBE | 0.0 | The | quick | brown | ⋯ |

next-token prediction

MLP
cross attention
self attention

MLP
cross attention
self attention

MLP
cross attention
self attention

Transformer Decoder Blocks

cross attention

Transformer Encoder Blocks

MLP
self attention

MLP
self attention

MLP
self attention

Sinusoidal Positional Encoding

2 × Conv1D + GELU

Log-Mel Spectrogram

Learned Positional Encoding

| SOT | EN | TRANS-CRIBE | 0.0 | The | quick | ⋯ |

Tokens in Multitask Training Format

**Multitask training format**

special tokens | text tokens | timestamp tokens

PREV → previous text tokens → START OF TRANSCRIPT
Custom vocabulary / prompting

Language identification → LANGUAGE TAG
NO SPEECH
Voice activity detection (VAD)

X → X Transcription → TRANSCRIBE
TRANSLATE
X → English Translation

begin time | text tokens | end time | ⋯ | begin time | text tokens | end time
Time-aligned transcription

NO TIMESTAMPS

text tokens
Text-only transcription (allows dataset-specific fine-tuning)

EOT

## Discussion Question #1

How does Whisper handle multitasking?

## Discussion Question #2

Why does Whisper use basic Transformers for its model?

## Model Variants

| Model | Layers | Width | Heads | Parameters |
|---|---|---|---|---|
| Tiny | 4 | 384 | 6 | 39M |
| Base | 6 | 512 | 8 | 74M |
| Small | 12 | 768 | 12 | 244M |
| Medium | 24 | 1024 | 16 | 769M |
| Large | 32 | 1280 | 20 | 1550M |

## Performance Metrics

1. Remove any phrases between matching brackets ( [, ] ).

2. Remove any phrases between matching parentheses ( (, ) ).

3. Remove any of the following words: hmm, mm, mhm, mmm, uh, um

4. Remove whitespace characters that comes before an apostrophe ʼ

5. Convert standard or informal contracted forms of English into the original form.

6. Remove commas (, ) between digits

7. Remove periods (. ) not followed by numbers

8. Remove symbols as well as diacritics from the text, where symbols are the ch starting with M, S, or P, except period, percent, and currency symbols that may b

9. Detect any numeric expressions of numbers and currencies and replace with a fo thousand dollars" → "$10000".

10. Convert British spellings into American spellings.

11. Remove remaining symbols that are not part of any numeric expressions.

12. Replace any successive whitespace characters with a space.

**Pseudocode**

**Algorithm 1** Whisper

1: **Class**: Whisper
2: **Parent Class**: nn.Module
3: **procedure** INITIALIZATION($dims$ (ModelDimensions))
4:     Initialize $dims$
5:     Create $encoder$ using AudioEncoder with:
6:         $n\_mels$
7:         $n\_audio\_ctx$
8:         $n\_audio\_state$
9:         $n\_audio\_head$
10:         $n\_audio\_layer$
11:     Create $decoder$ using TextDecoder with:
12:         $n\_vocab$
13:         $n\_text\_ctx$
14:         $n\_text\_state$
15:         $n\_text\_head$
16:         $n\_text\_layer$
17:     Define $all\_heads$ as a tensor of zeros of shape ($n\_text\_layer$, $n\_text\_head$).
18:     Set the last half of $all\_heads$ to True.
19:     Register $alignment\_heads$ as a buffer with the sparse version of $all\_heads$. Make it persistent.
20: **end procedure**

**Algorithm 2** AudioEncoder

1: **Class**: AudioEncoder
2: **Parent Class**: nn.Module
3: **procedure** INITIALIZATION($n\_mels, n\_ctx, n\_state, n\_head, n\_layer$)
4:     Initialize parent class
5:     Create $conv1$ with parameters $n\_mels, n\_state$
6:     Create $conv2$ with parameters $n\_state, n\_state$
7:     Register buffer $positional\_embedding$ using sinusoids function
8:     Define $blocks$ as a list of ResidualAttentionBlock with length $n\_layer$
9:     Define $ln\_post$ as LayerNorm with $n\_state$
10: **end procedure**
11: **procedure** FORWARD($x$ (Tensor))                    ▷ $x$ is the mel spectrogram of the audio
12:     $x \leftarrow$ Apply GELU activation after passing through $conv1$
13:     $x \leftarrow$ Apply GELU activation after passing through $conv2$
14:     Permute dimensions of $x$
15:     Assert shape of $x$ matches $positional\_embedding$
16:     Add $positional\_embedding$ to $x$
17:     **for** each block in $blocks$ **do**
18:         $x \leftarrow$ Apply block to $x$
19:     **end for**
20:     $x \leftarrow$ Apply $ln\_post$ to $x$
21:     **return** $x$
22: **end procedure**
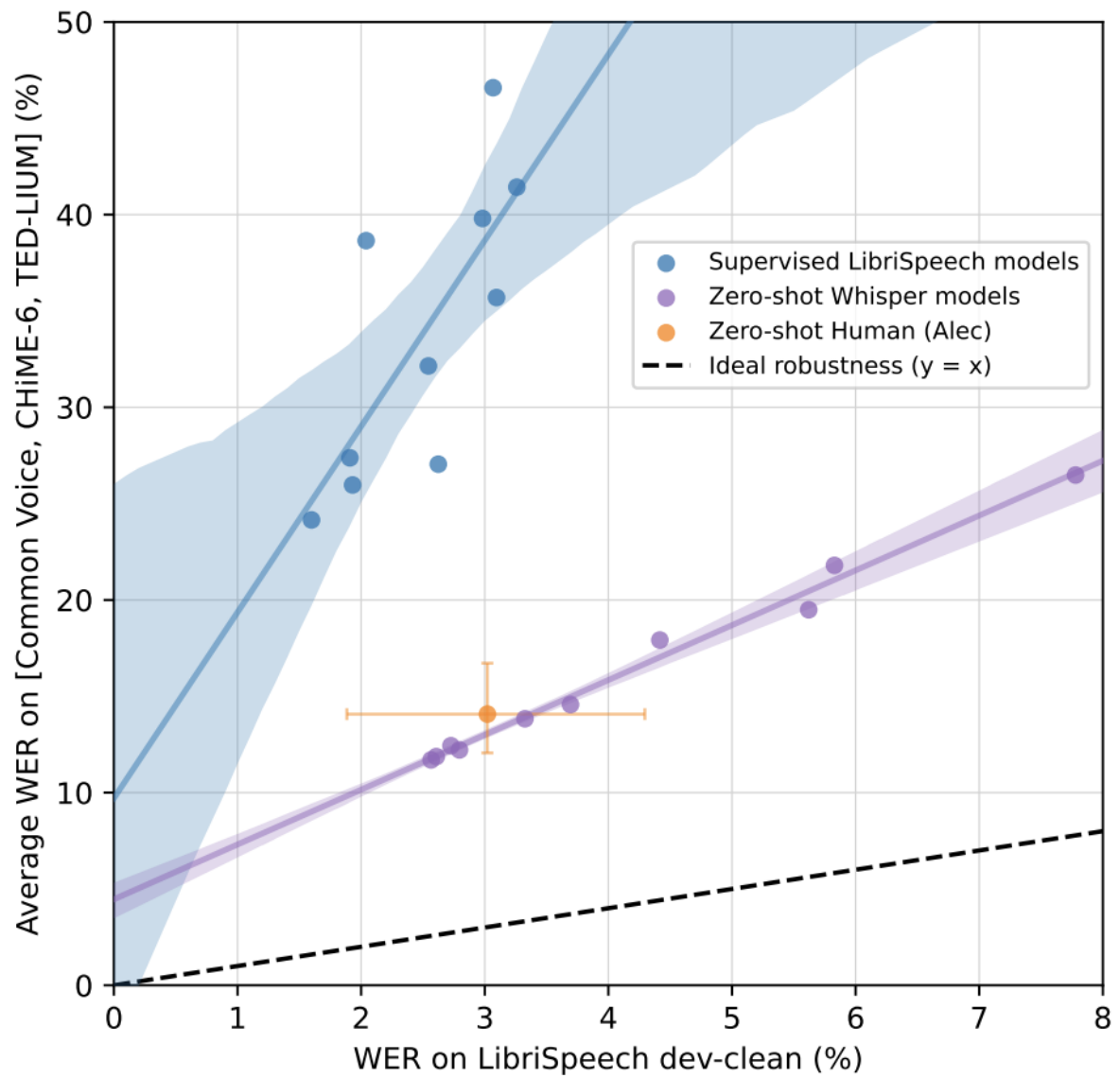
**Algorithm 3** TextDecoder

1: **Class**: TextDecoder
2: **Parent Class**: nn.Module
3: **procedure** INITIALIZATION($n\_vocab, n\_ctx, n\_state, n\_head, n\_layer$)
4:     Initialize parent class
5:     Create $token\_embedding$ with parameters $n\_vocab, n\_state$
6:     Initialize $positional\_embedding$ with size $n\_ctx \times n\_state$
7:     Define $blocks$ as a list of ResidualAttentionBlock with cross attention, length $n\_layer$
8:     Define $ln$ as LayerNorm with $n\_state$
9:     Create a mask $mask$ with size $n\_ctx \times n\_ctx$ and set upper triangle to $-\infty$
10:     Register buffer $mask$
11: **end procedure**
12: **procedure** FORWARD($x$ (Tensor), $xa$ (Tensor), $kv\_cache$ (Optional[dict])) ▷ $x$ is the text tokens, $xa$ is the encoded audio features
13:     **if** $kv\_cache$ exists **then**
14:         $offset \leftarrow$ shape of the first value in $kv\_cache$
15:     **else**
16:         $offset \leftarrow 0$
17:     **end if**
18:     Update $x$ with $token\_embedding$ and $positional\_embedding$ based on offset
19:     Convert $x$ to the same dtype as $xa$
20:     **for** each block in $blocks$ **do**
21:         $x \leftarrow$ Apply block to $x$, $xa$ with $mask$ and $kv\_cache$
22:     **end for**
23:     $x \leftarrow$ Apply $ln$ to $x$
24:     Calculate $logits$ using $x$ and the transpose of $token\_embedding$ weight
25:     Convert $logits$ to float type
26:     **return** $logits$
27: **end procedure**

```
mask
tensor([[0., -inf, -inf,  ..., -inf, -inf, -inf],
        [0., 0., -inf,  ..., -inf, -inf, -inf],
        [0., 0., 0.,  ..., -inf, -inf, -inf],
        ...,
        [0., 0., 0.,  ..., 0., -inf, -inf],
        [0., 0., 0.,  ..., 0., 0., -inf],
        [0., 0., 0.,  ..., 0., 0., 0.]])
```
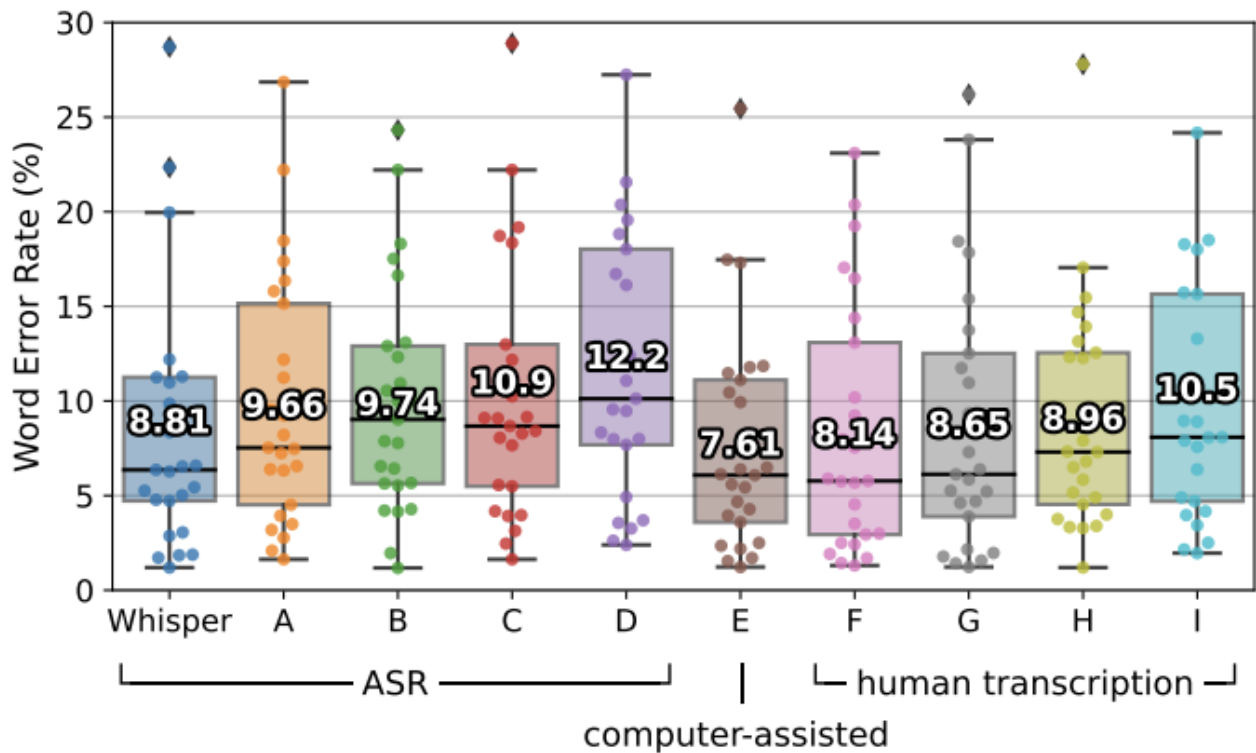
## Empirical Results

### Close the gap to human robustness

- Lack of Robustness in Supervised Models

**As good as human being?**

Close to performance of professional human transcribers!



## Code demonstration

https://colab.research.google.com/drive/1M8zNZ24lGcf05j-u53y73D-OhOv6z0I0?usp=drive_link#scrollTo=j9UgVYrod4SB

## Critical Analysis

### Were there any errors?

- The predictions may include texts that are not actually spoken in the audio input

- Lower accuracy on low-resource and/or low-discoverability languages or languages

- Prone to generating repetitive texts

## Broader Impact

We hope Whisper's high accuracy and ease of use will allow developers to add voice interfaces to a much wider set of applications.
However,
it also raises dual-use concerns.

## Seek for more?

### Research Index of OpenAI

| Mar 14, 2023 | GPT-4 | Read paper ↗ |
| Sep 21, 2022 | Introducing Whisper | Read paper ↗ |
| Apr 13, 2022 | Hierarchical text-conditional image generation with CLIP latents | Read paper ↗ |
| Jan 27, 2022 | Aligning language models to follow instructions | Read paper ↗ |
| Sep 23, 2021 | Summarizing books with human feedback | |
| Jul 7, 2021 | Evaluating large language models trained on code | Read paper ↗ |
| Mar 4, 2021 | Multimodal neurons in artificial neural networks | Read paper ↗ |
| Jan 5, 2021 | DALL·E: Creating images from text | |
| Jan 5, 2021 | CLIP: Connecting text and images | Read paper ↗ |
| Sep 4, 2020 | Learning to summarize with human feedback | Read paper ↗ |
| Jun 17, 2020 | Image GPT | Read paper ↗ |
| May 28, 2020 | Language models are few-shot learners | Read paper ↗ |
| Apr 30, 2020 | Jukebox | Read paper ↗ |
| Oct 15, 2019 | Solving Rubik's Cube with a robot hand | Read paper ↗ |
| Sep 17, 2019 | Emergent tool use from multi-agent interaction | Read paper ↗ |
| Apr 25, 2019 | MuseNet | |

Link: https://openai.com/research?contentTypes=milestone

# Repository



# Resource Links

1. Paper: https://cdn.openai.com/papers/whisper.pdf
2. Code: https://github.com/openai/whisper
3. Model Card: https://github.com/openai/whisper/blob/main/model-card.md
4. Introducing Whisper: https://openai.com/research/whisper
5. Hugging Face Community https://huggingface.co/openai/whisper-large

# Citation For paper

- Chan et al. Simply mix all available speech recognition data to train one large neural network.
- Galvez et al. The people's speech: A large-scale diverse english speech recognition dataset for commercial usage.
- Chen et al. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio.
- Baevski et al. wav2vec 2.0: A framework for self-supervised learning of speech representations.
- Baevski er al. Unsupervised speech recognition. Advances in Neural Information Processing Systems
- Zhang et al. BigSSL: Exploring the frontier of large-scale semi-supervised learning for automatic speech recognition.