15.05.2025

# Listen who's whispering

## Speech recognition with APEX and Whisper

Dennis Amthor, APEX Connect 2025

code of change

Hyand

# Our key facts

**Hyand**

## Germany

- Brunswick
- Ratingen
- Hamburg
- Dortmund
- Cologne
- Frankfurt
- Munich
- Berlin

## Poland

- Warsaw

## Lithuania

- Vilnius
- Kaunas

## Romania

- Cluj-Napoca

## India

- Pune

**850+** Employees

**150+** Customers

**110+** million € turnover

Dennis Amthor

**Consultant APEX/JavaScript**

**Hyand**

# Idea/Motivation

Hyand

# Idea/Motivation

- Voice as data input inside APEX

  - Forms, meeting transcripts, notes

- Controlling APEX functions via voice

  - Searching, filtering

➢ Feasibility test

  - Complexity

  - Cost

  - Precision

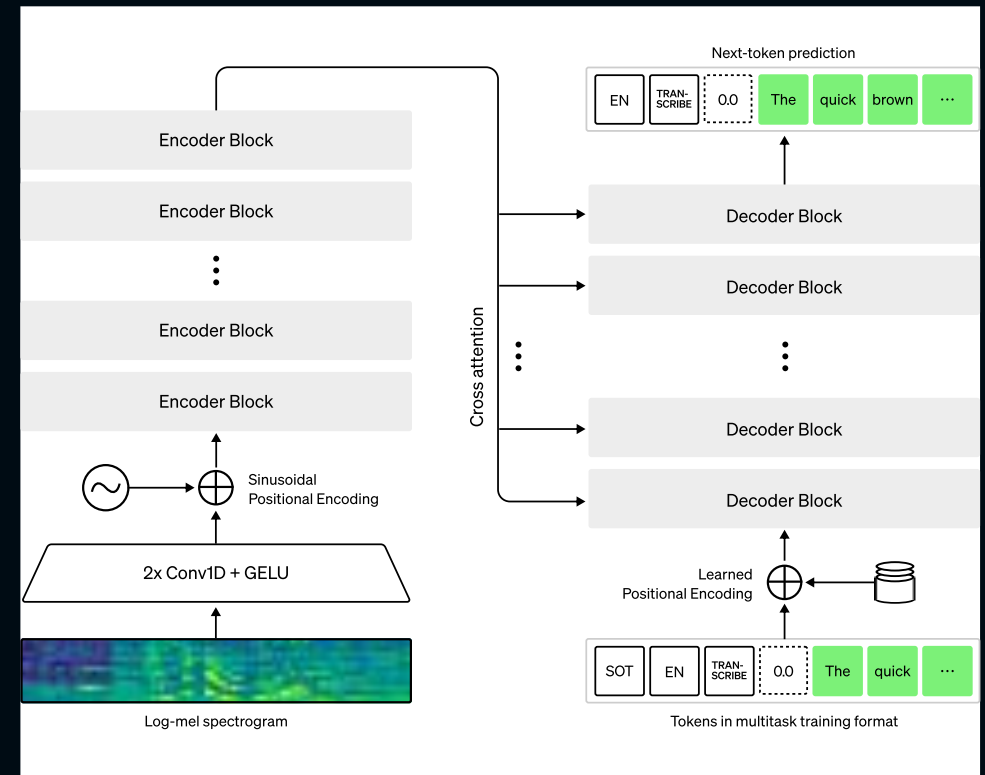  - Performance

  - Offline capability

Hyand

# What is Whisper?

Hyand

# What is Whisper?

- Automated Speech Recognition System (ASR) by OpenAI

- Trained with 680,000 hours of audio data (2/3 of which are English)

- Transcription and translation (in English)

- OpenAI: Data source larger and more "diverse" than on other models

  ➢ As a result, less precise overall, but better handling of accents and background noise

**Hyand**

# How it works

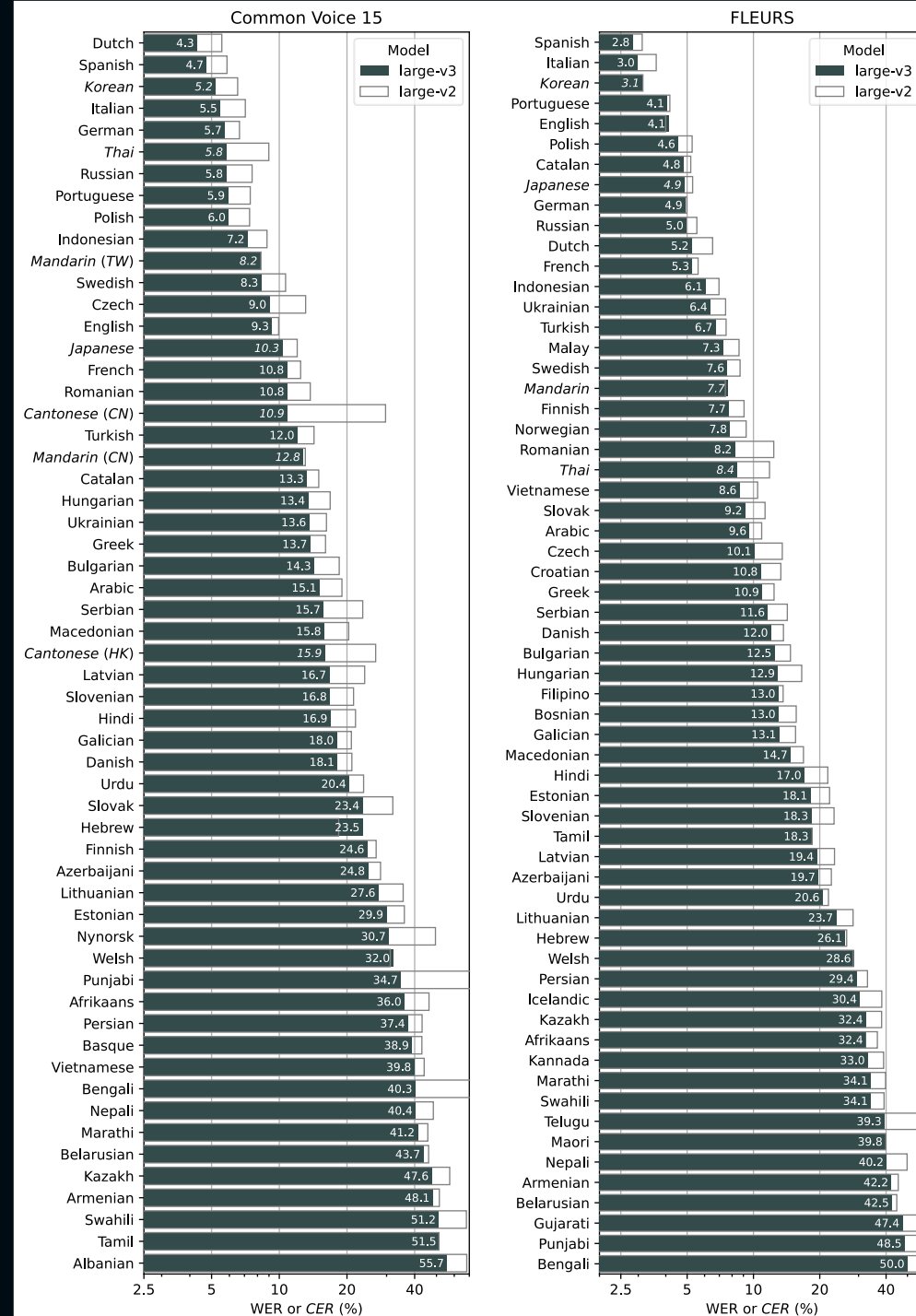- Chunking the audio input into blocks of 30 seconds

- Conversion to Mel-Spectogram

  - Representation of frequencies/time while taking human hearing ability into account

  - Higher resolution in a lower (=audible) frequency range

- Encoder-Decoder-Transformer-Model

  - Audio information as vectors

  - Self- and Cross-Attention

  - Different token types
    (Text, control, language, timestamps)

# Models

- 6 models available
  (tiny, base, small, medium, large, turbo)

- Differences in regard to

  - Performance

  - Storage

  - Error rate

- Alternative faster English-only variants available for tiny to medium

- Error rate also highly dependent on the language used

**Hyand**

Source: https://github.com/openai/whisper

## Common Voice 15

Model: large-v3, large-v2

| Language | WER or CER (%) |
|---|---|
| Dutch | 4.3 |
| Spanish | 4.7 |
| Korean | 5.2 |
| Italian | 5.5 |
| German | 5.7 |
| Thai | 5.8 |
| Russian | 5.8 |
| Portuguese | 5.9 |
| Polish | 6.0 |
| Indonesian | 7.2 |
| Mandarin (TW) | 8.2 |
| Swedish | 8.3 |
| Czech | 9.0 |
| English | 9.3 |
| Japanese | 10.3 |
| French | 10.8 |
| Romanian | 10.8 |
| Cantonese (CN) | 10.9 |
| Turkish | 12.0 |
| Mandarin (CN) | 12.8 |
| Catalan | 13.3 |
| Hungarian | 13.4 |
| Ukrainian | 13.6 |
| Greek | 13.7 |
| Bulgarian | 14.3 |
| Arabic | 15.1 |
| Serbian | 15.7 |
| Macedonian | 15.8 |
| Cantonese (HK) | 15.9 |
| Latvian | 16.7 |
| Slovenian | 16.8 |
| Hindi | 16.9 |
| Galician | 18.0 |
| Danish | 18.1 |
| Urdu | 20.4 |
| Slovak | 23.4 |
| Hebrew | 23.5 |
| Finnish | 24.6 |
| Azerbaijani | 24.8 |
| Lithuanian | 27.6 |
| Estonian | 29.9 |
| Nynorsk | 30.7 |
| Welsh | 32.0 |
| Punjabi | 34.7 |
| Afrikaans | 36.0 |
| Persian | 37.4 |
| Basque | 38.9 |
| Vietnamese | 39.8 |
| Bengali | 40.3 |
| Nepali | 40.4 |
| Marathi | 41.2 |
| Belarusian | 43.7 |
| Kazakh | 47.6 |
| Armenian | 48.1 |
| Swahili | 51.2 |
| Tamil | 51.5 |
| Albanian | 55.7 |

## FLEURS

Model: large-v3, large-v2

| Language | WER or CER (%) |
|---|---|
| Spanish | 2.8 |
| Italian | 3.0 |
| Korean | 3.1 |
| Portuguese | 4.1 |
| English | 4.1 |
| Polish | 4.6 |
| Catalan | 4.8 |
| Japanese | 4.9 |
| German | 4.9 |
| Russian | 5.0 |
| Dutch | 5.2 |
| French | 5.3 |
| Indonesian | 6.1 |
| Ukrainian | 6.4 |
| Turkish | 6.7 |
| Malay | 7.3 |
| Swedish | 7.6 |
| Mandarin | 7.7 |
| Finnish | 7.7 |
| Norwegian | 7.8 |
| Romanian | 8.2 |
| Thai | 8.4 |
| Vietnamese | 8.6 |
| Slovak | 9.2 |
| Arabic | 9.6 |
| Czech | 10.1 |
| Croatian | 10.8 |
| Greek | 10.9 |
| Serbian | 11.6 |
| Danish | 12.0 |
| Bulgarian | 12.5 |
| Hungarian | 12.9 |
| Filipino | 13.0 |
| Bosnian | 13.0 |
| Galician | 13.1 |
| Macedonian | 14.7 |
| Hindi | 17.0 |
| Estonian | 18.1 |
| Slovenian | 18.3 |
| Tamil | 18.3 |
| Latvian | 19.4 |
| Azerbaijani | 19.7 |
| Urdu | 20.6 |
| Lithuanian | 23.7 |
| Hebrew | 26.1 |
| Welsh | 28.6 |
| Persian | 29.4 |
| Icelandic | 30.4 |
| Kazakh | 32.4 |
| Afrikaans | 32.4 |
| Kannada | 33.0 |
| Marathi | 34.1 |
| Swahili | 34.1 |
| Telugu | 39.3 |
| Maori | 39.8 |
| Nepali | 40.2 |
| Armenian | 42.2 |
| Belarusian | 42.5 |
| Gujarati | 47.4 |
| Punjabi | 48.5 |
| Bengali | 50.0 |

# Usage

# Whisper

- Calls to OpenAI API
- Model: whisper-1 ($0.006/minute)
- Alternative: Local integration in Python using PyTorch
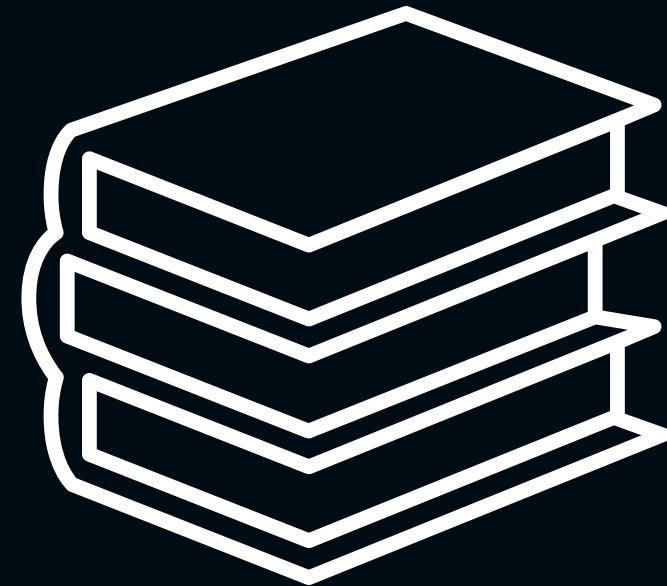
# Whisper Web

- Web integration using Transformers.js
- Pre-trained models
- Stored in the local cache
- Usage of own models possible after conversion

# Transformers
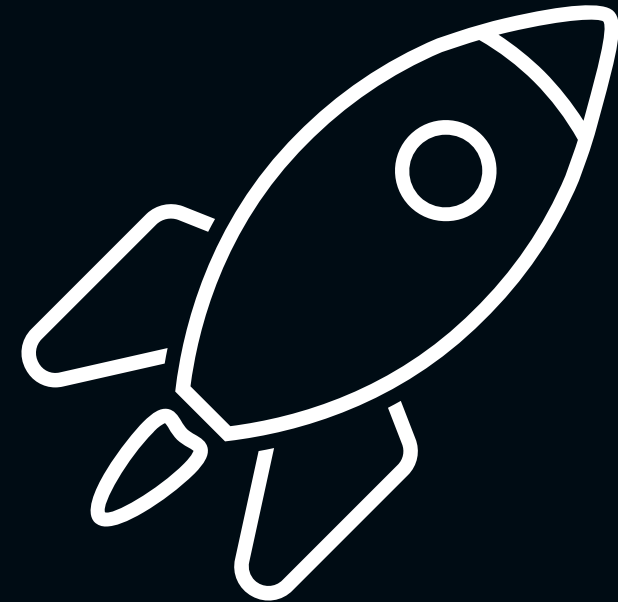
Huggingface.co

Hyand

# Transformers (Huggingface.co)

- Library of pre-trained models

  - Text, image, audio, video

  - Various NLP and multimedia tasks

- Models can be used directly and replaced in a modular way (no own training required)

- Unified API and JavaScript integration (Transformers.js)

- Usable directly in the browser (ONNX Runtime)

# Implementation

Hyand

# Implementation

- Calls towards Whisper Web/Whisper API using JavaScript (Webpack)

- Integration as APEX plugins

  - Region Plugin for testing/comparison (input methods, technologies)

  - DA plugin for direct use with buttons (examples)

- Whisper Web: Transformer.js with pre-trained model (Xenova/whisper-base)

- Whisper: OpenAI API: AI Service & apex_ai.generate

- Comparison with similar technologies

  - Web Speech API/Speech Recognition
  - (OCI Speech)

**Hyand**

# Demo

# Result

Hyand

# Result

| | |
|---|---|
| **S** | |
| **A** | |
| **B** | |
| **C** | |

Hyand

# Implementation effort

| | |
|---|---|
| **S** | Web Speech API |
| **A** | Whisper Web    Whisper API |
| **B** | |
| **C** | |

Hyand

# Configuration/Options

| | |
|---|---|
| **S** | |
| **A** | Whisper Web |
| **B** | Whisper API |
| **C** | Web Speech API |

Hyand

# Precision/Error rate

| S | Whisper API | |
|---|---|---|
| A | Web Speech API | |
| B | Whisper Web | |
| C | | |

Hyand

# Performance

| | |
|---|---|
| **S** | Web Speech API |
| **A** | Whisper API |
| **B** | |
| **C** | Whisper Web |

Hyand

# Cost

| S | Whisper Web | Web Speech API |
|---|---|---|
| A | Whisper API | |
| B | | |
| C | | |

# Privacy

| S | Whisper Web | |
|---|---|---|
| A | | |
| B | Whisper API | |
| C | Web Speech API | |

Hyand

# Result

| | Whisper Web (base model) | Whisper API | Web Speech API |
|---|:---:|:---:|:---:|
| Implementation effort | A | A | S |
| Configuration/Options | A | B | C |
| Precision/Error rate | B | S | A |
| Performance | C | A | A |
| Cost | S | A | S |
| Privacy | S | B | C |

Hyand

# Result

| | Whisper Web (base model) | Whisper API | Web Speech API |
|---|---|---|---|
| Implementation effort | A | A | OCI Speech |
| Configuration/Options | A | B | |
| Precision/Error rate | B | S | **?** |
| Performance | C | A | |
| Cost | S | A | |
| Privacy | S | B | C |

Hyand

# *Any questions?*

**Hyand**