# Speech technology
## Trends, limitations & future

Vivian van Oijen
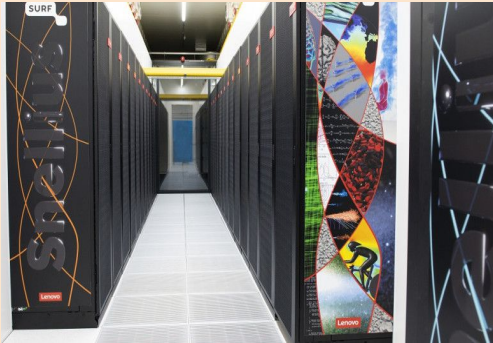
SURF

# Hello!

SURF

# Hello!



Machine learning
@SURF

# Hello!





Machine learning
@SURF

# Hello!



Machine learning @SURF





Speech technology

# Khanmigo is your always-available writing coach.

Can you help me solve this?

I need a grading rubric

**For teachers**

Knock something off your to-do list in minutes.

Sign up for free

**For learners**

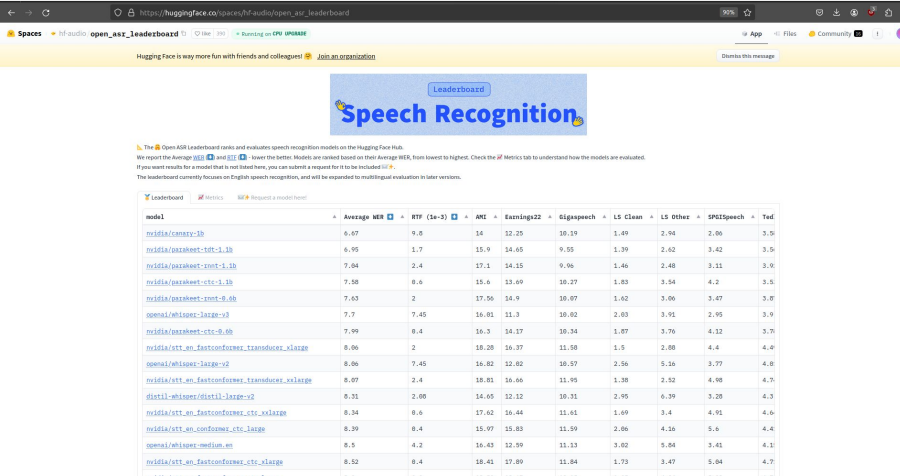Prep for tests without breaking a sweat.

Get Khanmigo

**For parents**

On-demand help makes homework time a breeze.

Get Khanmigo

SURF

# State of the Art

**SURF**

# State of the Art

Hugging Face Leaderboard



https://huggingface.co/spaces/hf-audio/open_asr_leaderboard

# State of the Art

Study case 1: Teams

# State of the Art

## Study case 2: Canary & Parakeet

# State of the Art

## Study case 2: Canary & Parakeet



🤗 **Hugging Face**    🔍 Search models, datasets, users...

Hugging Face is way more fun with friends and colleagues! 🤗 __Join an organization__

nvidia / **canary-1b** 🗐    ♡ like  200



🤗 **Hugging Face**    🔍 Search models, datasets, users...

Hugging Face is way more fun with friends and colleagues! 🤗 __Join an organization__

nvidia / **parakeet-rnnt-1.1b** 🗐    ♡ like  99

**SURF**

# State of the Art

## Study case 2: Canary & Parakeet



## How to use Parakeet-TDT

To run speech recognition with Parakeet-TDT, you'll need to install NVIDIA NeMo. It can be installed as a pip package, as shown below. Cython and PyTorch (2.0 and above) should be installed before trying to install NeMo.

```
pip install nemo_toolkit['asr']
```

Once NeMo is installed, you can use Parakeet-TDT to recognize your audio files as follows:

```
import nemo.collections.asr as nemo_asr
asr_model = nemo_asr.models.ASRModel.from_pretrained(model_name="nvidia/parakeet-tdt-1.1b")
transcript = asr_model.transcribe(["some_audio_file.wav"])
```

SURF

# State of the Art

Study case 3: Whisper

# Paper

**Robust Speech Recognition via Large-Scale Weak Supervision**

Alec Radford [* 1]   Jong Wook Kim [* 1]   Tao Xu [1]   Greg Brockman [1]   Christine McLeavey [1]   Ilya Sutskever [1]

## Abstract

We study the capabilities of speech processing systems trained simply to predict large amounts of transcripts of audio on the internet. When scaled to 680,000 hours of multilingual and multitask supervision, the resulting models generalize well to standard benchmarks and are often competitive with prior fully supervised results but in a zero-shot transfer setting without the need for any fine-tuning. When compared to humans, the models approach their accuracy and robustness. We are releasing models and inference code to serve as a foundation for further work on robust speech processing.

methods are exceedingly adept at finding patterns within a training dataset which boost performance on held-out data from the same dataset. However, some of these patterns are brittle and spurious and don't generalize to other datasets and distributions. In a particularly disturbing example, Radford et al. (2021) documented a 9.2% increase in object classification accuracy when fine-tuning a computer vision model on the ImageNet dataset (Russakovsky et al., 2015) without observing any improvement in average accuracy when classifying the same objects on seven other natural image datasets. A model that achieves "superhuman" performance when trained on a dataset can still make many basic errors when evaluated on another, possibly precisely because it is exploiting those dataset-specific quirks that humans are oblivious to (Geirhos et al., 2020).

https://arxiv.org/abs/2212.04356

# Paper

# Paper

# Paper

# Paper

**Robust Speech Recognition via Large-Scale Weak Supervision**

Alec Radford [* 1]  Jong Wook Kim [* 1]  Tao Xu [1]  Greg Brockman [1]  Christine McLeavey [1]  Ilya Sutskever [1]
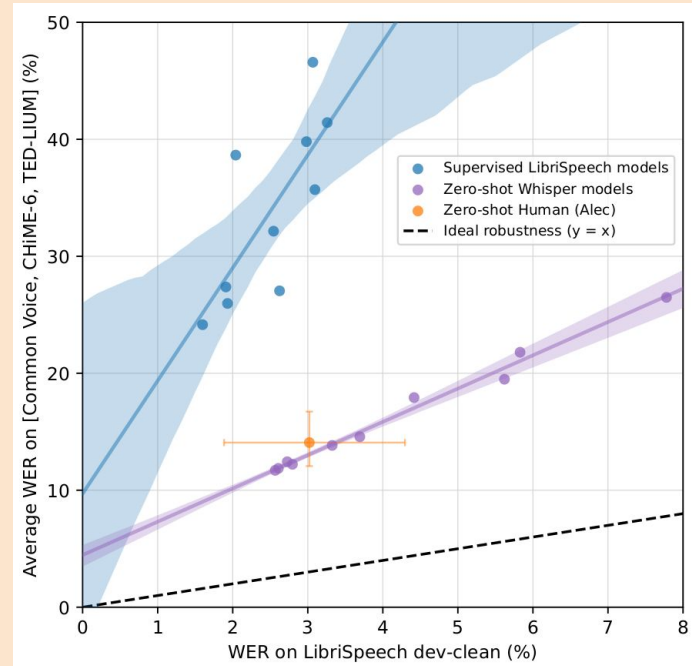
## Abstract

We study the capabilities of speech processing systems trained simply to predict large amounts of transcripts of audio on the internet. When scaled to 680,000 hours of multilingual and multitask supervision, the resulting models generalize well to standard benchmarks and are often competitive with prior fully supervised results but in a zero-shot transfer setting without the need for any fine-tuning. When compared to humans, the models approach their accuracy and robustness. We are releasing models and inference code to serve as a foundation for further work on robust speech processing.
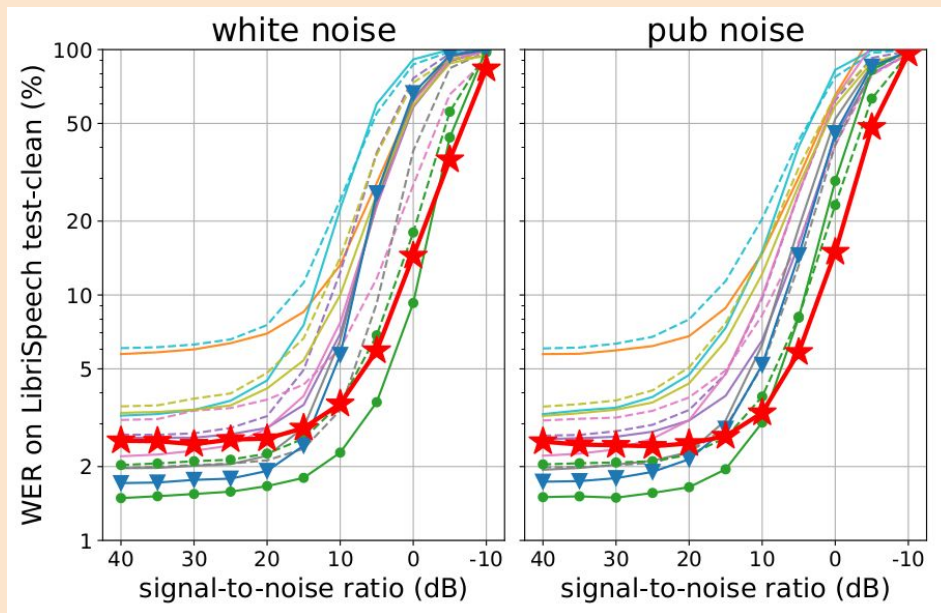
methods are exceedingly adept at finding patterns within a training dataset which boost performance on held-out data from the same dataset. However, some of these patterns are brittle and spurious and don't generalize to other datasets and distributions. In a particularly disturbing example, Radford et al. (2021) documented a 9.2% increase in object classification accuracy when fine-tuning a computer vision model on the ImageNet dataset (Russakovsky et al., 2015) without observing any improvement in average accuracy when classifying the same objects on seven other natural image datasets. A model that achieves "superhuman" performance when trained on a dataset can still make many basic errors when evaluated on another, possibly precisely because it is exploiting those dataset-specific quirks that humans are oblivious to (Geirhos et al., 2020).
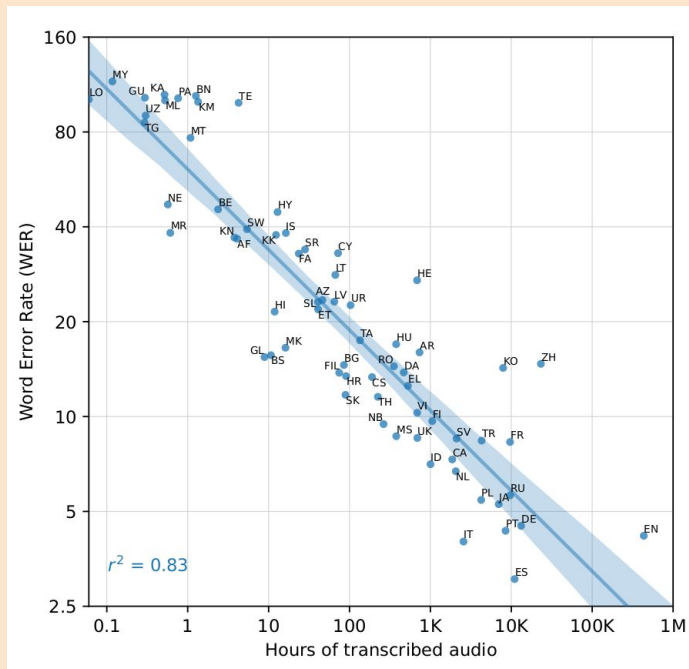
https://arxiv.org/abs/2212.04356

# Paper

# "Demo"

## Robust Speech Recognition via Large-Scale Weak Supervision

Alec Radford [* 1] Jong Wook Kim [* 1] Tao Xu [1] Greg Brockman [1] Christine McLeavey [1] Ilya Sutskever [1]

### Abstract

We study the capabilities of speech processing systems trained simply to predict large amounts of transcripts of audio on the internet. When scaled to 680,000 hours of multilingual and multitask supervision, the resulting models generalize well to standard benchmarks and are often competitive with prior fully supervised results but in a zero-shot transfer setting without the need for any fine-tuning. When compared to humans, the models approach their accuracy and robustness. We are releasing models and inference code to serve as a foundation for further work on robust speech processing.

methods are exceedingly adept at finding patterns within a training dataset which boost performance on held-out data from the same dataset. However, some of these patterns are brittle and spurious and don't generalize to other datasets and distributions. In a particularly disturbing example, Radford et al. (2021) documented a 9.2% increase in object classification accuracy when fine-tuning a computer vision model on the ImageNet dataset (Russakovsky et al., 2015) without observing any improvement in average accuracy when classifying the same objects on seven other natural image datasets. A model that achieves "superhuman" performance when trained on a dataset can still make many basic errors when evaluated on another, possibly precisely because it is exploiting those dataset-specific quirks that humans are oblivious to (Geirhos et al., 2020).

https://arxiv.org/abs/2212.04356

# The other way around: text to speech!







Introducing speech-to-text, text-to-speech, and more for 1,100+ languages

May 22, 2023 · 7 minute read

Meta AI

SURF

# The other way around: text to speech!



```
[viviano@gcn59 Bark]$ python -m bark --text "Hello, my name is Suno and I like pizza [laughs]" --output_filename "bark_output3.wav"
/gpfs/home2/viviano/SpeechTech/Bark/venv/lib/python3.11/site-packages/torch/nn/utils/weight_norm.py:28: UserWarning: torch.nn.utils.weight_norm
 is deprecated in favor of torch.nn.utils.parametrizations.weight_norm.
  warnings.warn("torch.nn.utils.weight_norm is deprecated in favor of torch.nn.utils.parametrizations.weight_norm.")
100%|                                                 | 316/316 [00:03<00:00, 82.49it/s]
100%|                                                 | 16/16 [00:08<00:00,  1.90it/s]
Done! Output audio file is saved at: './bark_output3.wav'
```

**Deepgram**

**Nova**
Unmatched performance and value

Our next-gen model surpasses all competitors in speed, accuracy, and cost. Compared to the nearest competitor, Nova is 22% more accurate, more than 20 times faster, and over 3x cheaper.

**Whisper**
Improvements you can't miss

Our fully managed Whisper APIs are faster, more reliable, and cheaper than OpenAI's. Includes built-in diarization, word-level timestamps, and an 80x higher file size limit.

**Custom**
Boost performance using your data

Custom trained speech models give accuracy a noticeable boost, especially on unique customer jargon. High throughput models are also available to meet enterprise scalability requirements.

**SPEECHMATICS**

# Great AI transcription: The bedrock of value

Accurate transcription. Valuable on its own. And the foundation for so much more.

Speak to sales

**IIElevenLabs**
# Generative Voice AI

**Amberscript**

# Transform your audio and video to text and subtitles

Our cutting-edge generative AI, paired with top-tier language professionals, collaboratively deliver highly accurate solutions tailored to your business needs.

Request a quote    Try It free

Some commercial alternatives…

SURF

# Limitations

SURF

# Limitations

**Services**

**Open models**

SURF

# Limitations

**Services**

- Little transparency

- Requires sharing your data

**Open models**

SURF

# Limitations

**Services**

- Little transparency

- Requires sharing your data

**Open models**

- Sometimes less quality

- Less user friendly

SURF

# Limitations



Dutch Open Speech Recognition Benchmark

Results of Dutch ASR models, collected by the community

View on GitHub

https://opensource-spraakherkenning-nl.github.io/ASR_NL_results/

SURF

# Future

SURF

# Future

How do we overcome current limitations
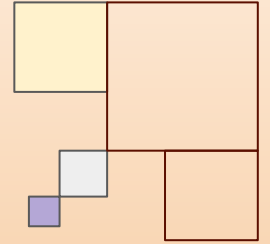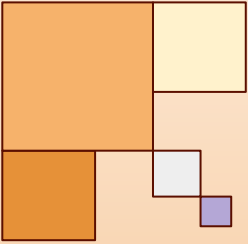
1. **Quality**

   We need data

2. **Infrastructure**

   Trusted party that provides API and UI

SURF

# Thank you!

Vivian van Oijen

SURF