

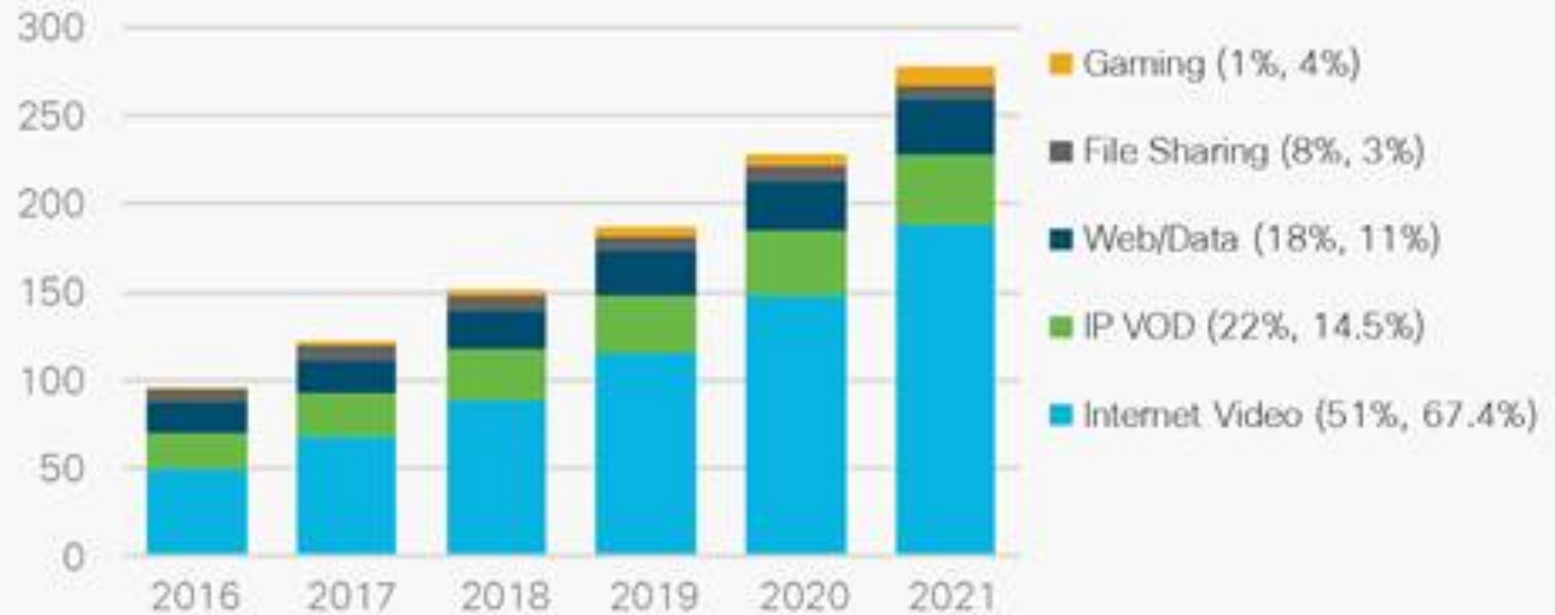
Speak search: ASR-free search for speech corpus

Divya Pitta, Fang-Yi Chiu

Internet Application Traffic Forecasting

24% CAGR
2016-2021

Exabytes
per month



Figures (n) refer to 2016, 2021 traffic shares.

Source: Cisco VNI Global IP Traffic Forecast, 2016-2021.

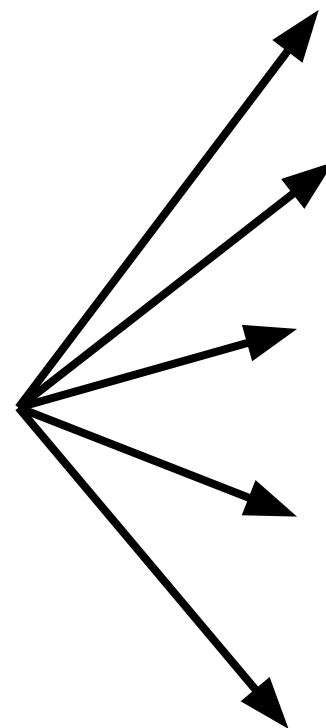
Outline

- Task Scope
- Dataset Setting
- Baseline: ASR
- Speech Encoding Result
- Future Work

Task Scope

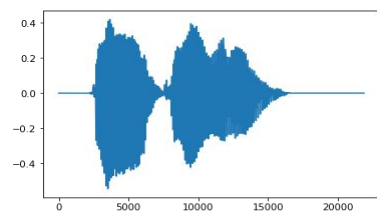


Source: Amazon Alexa



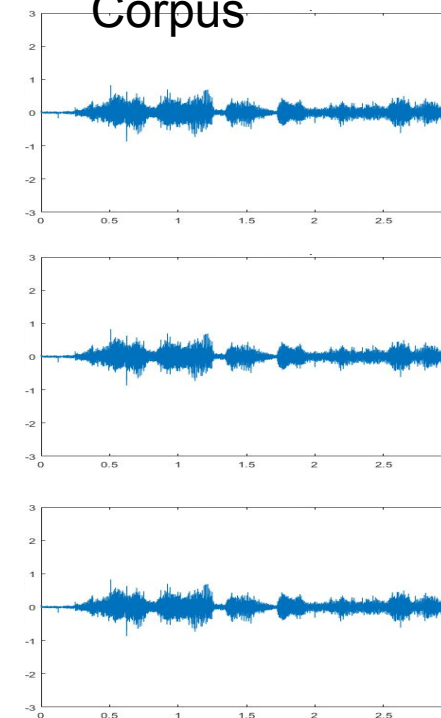
Task Scope

Query
Speech

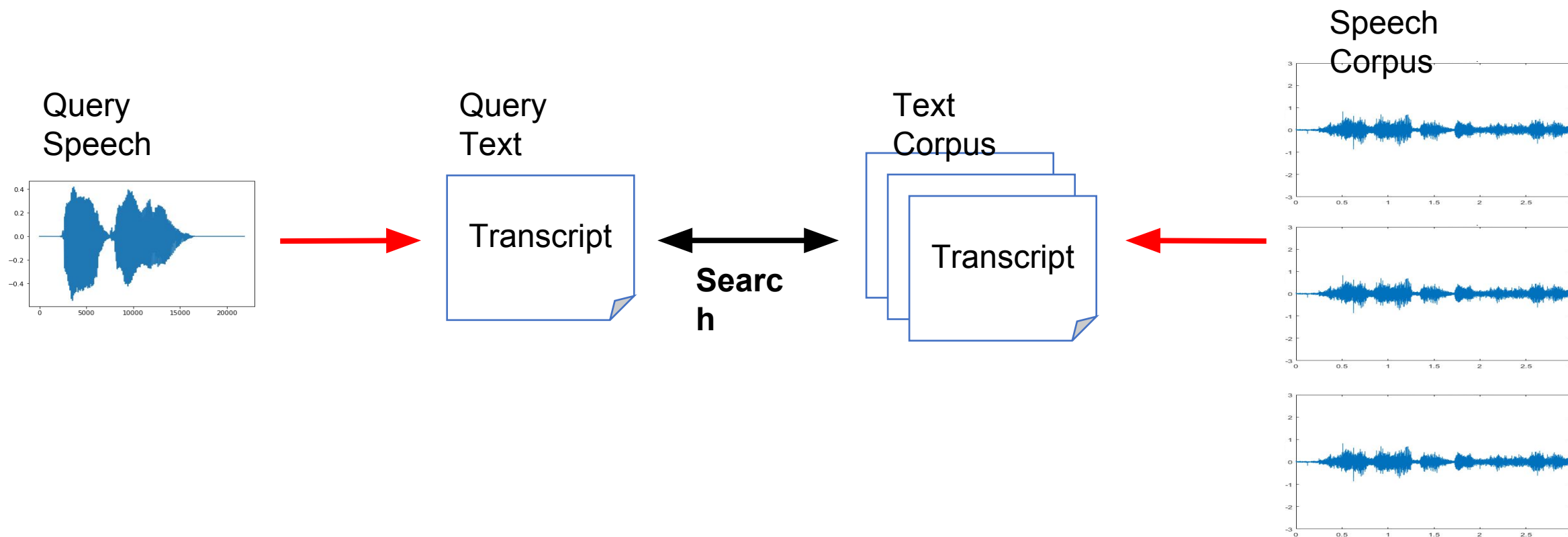


Search

Speech
Corpus



Current Solutions



Current Solutions

Hansel

->

cancel

stevie

->

TV

Hansel

<=

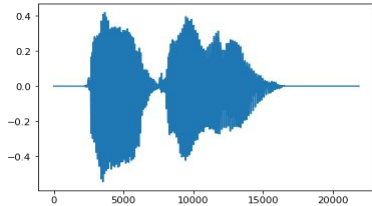
Hansel

stevie

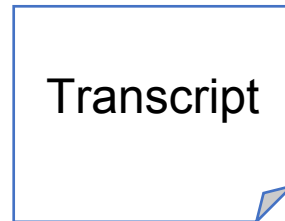
<=

stevie

Query
Speech

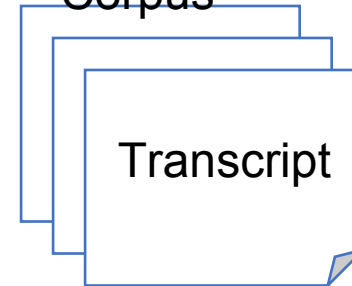


Query
Text

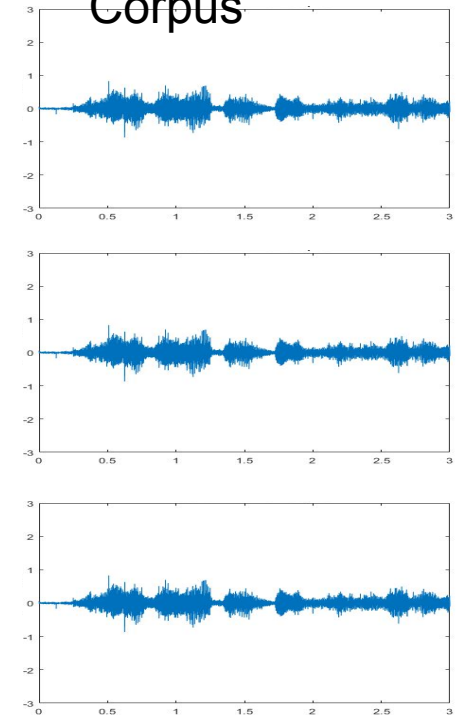


Search

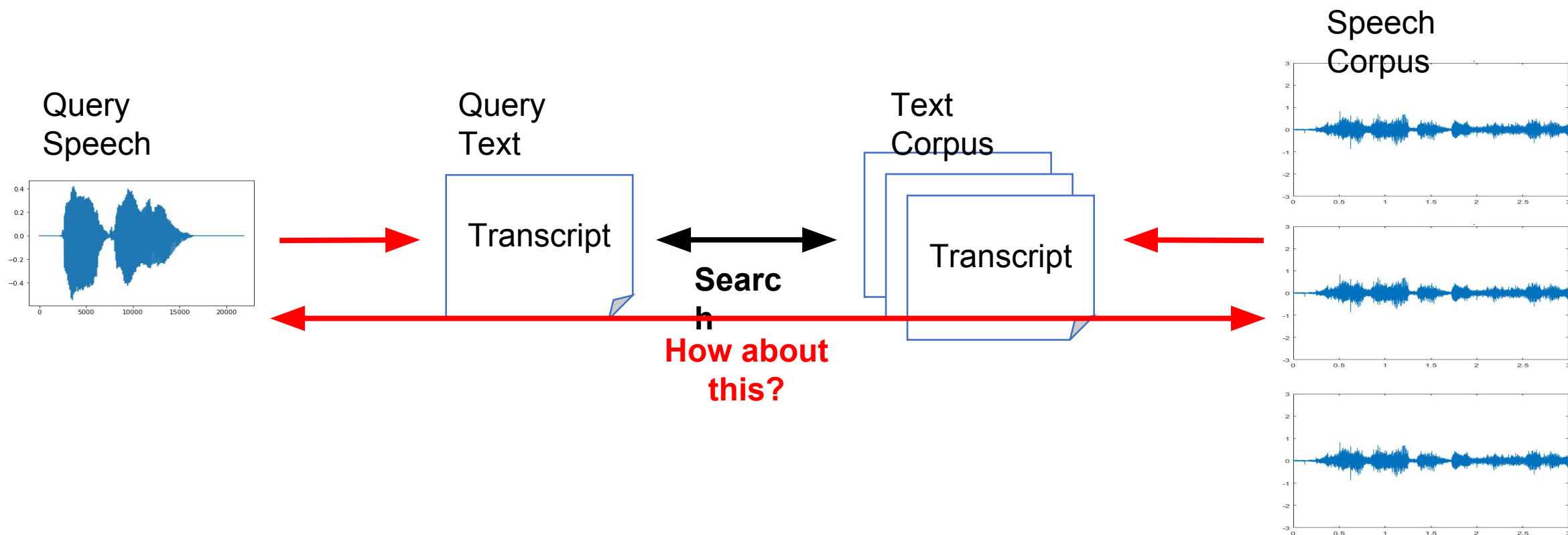
Text
Corpus



Speech
Corpus



Task Scope



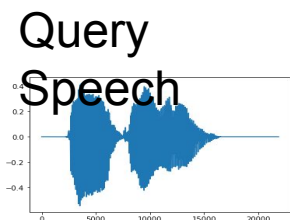
Dataset

- Librispeech (dev-clean version)
 - 2700 files (337 M)
 - 2 – 33 seconds for each file
 - Transcript available for each file
 - Audio books for fiction

Ground Truth

- Setting
 - Classical information retrieval problem
 - Text-to-text search
- Search engine library: [Whoosh](#)
 - Ranking by Okapi BM25 (Tf-Idf): top-k documents as ground truth
 - Default parsing: no stemming etc.
- Measurement
 - Precision
 - Recall

Search with ASR



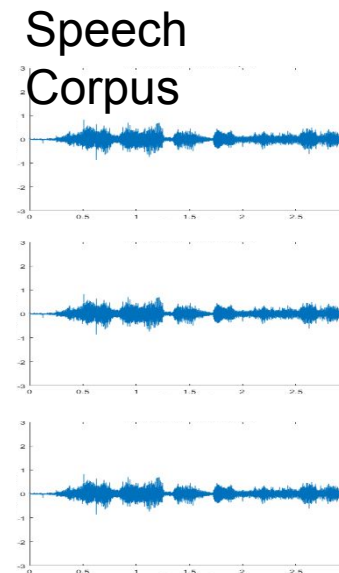
ASR

Query Text

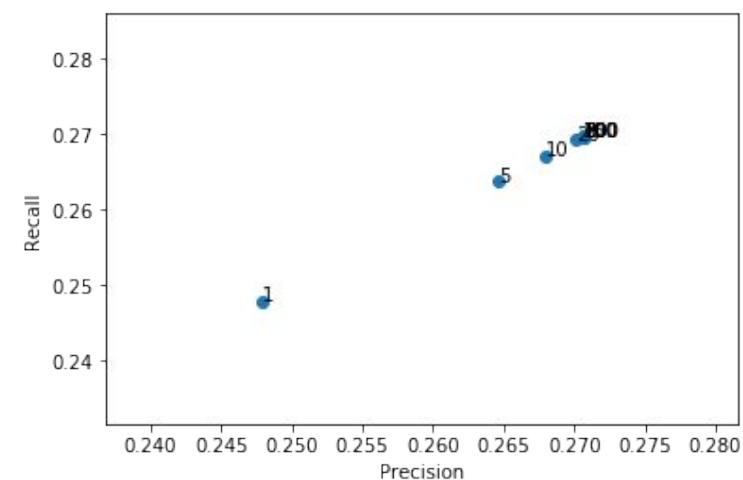
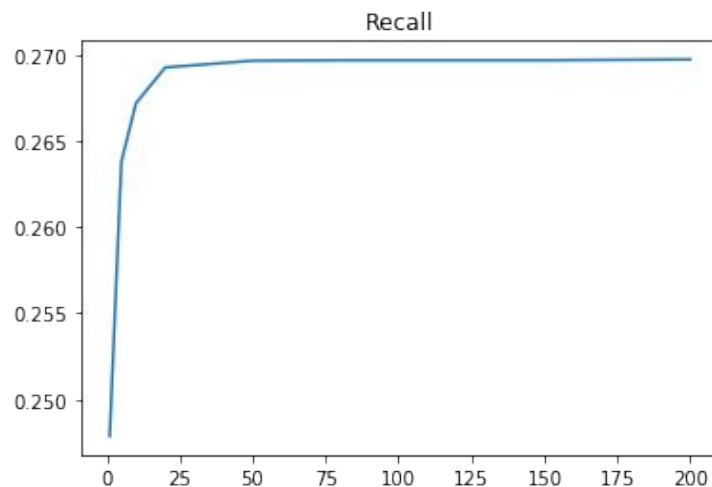
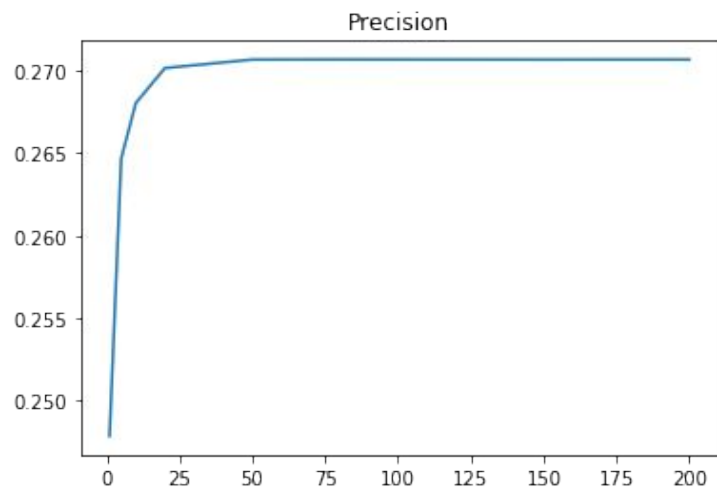
Search

Text Corpus

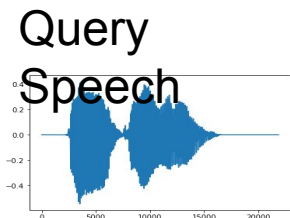
ASR



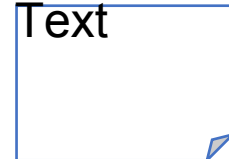
Performance: precision, recall



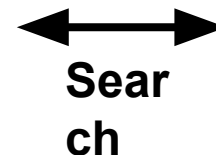
Search with ASR



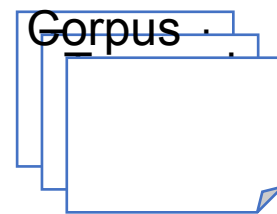
Query Text



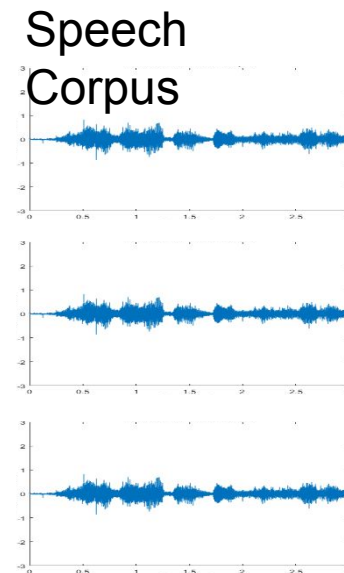
A blue-outlined rectangular box representing the query text output.



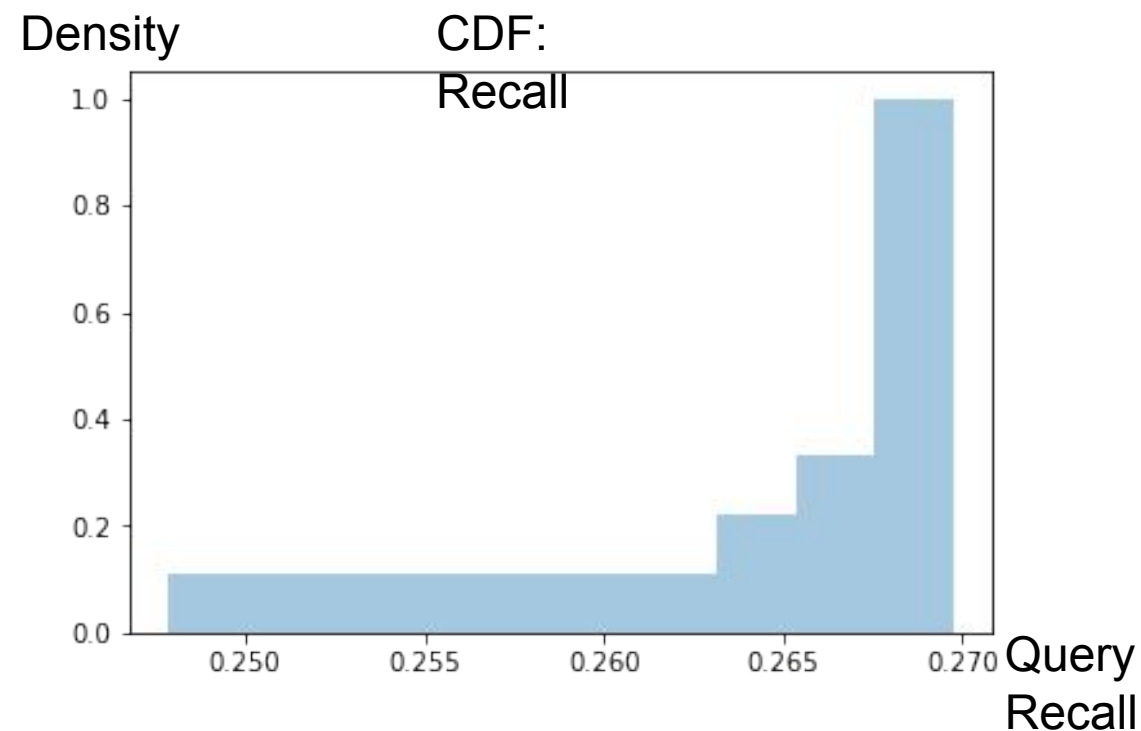
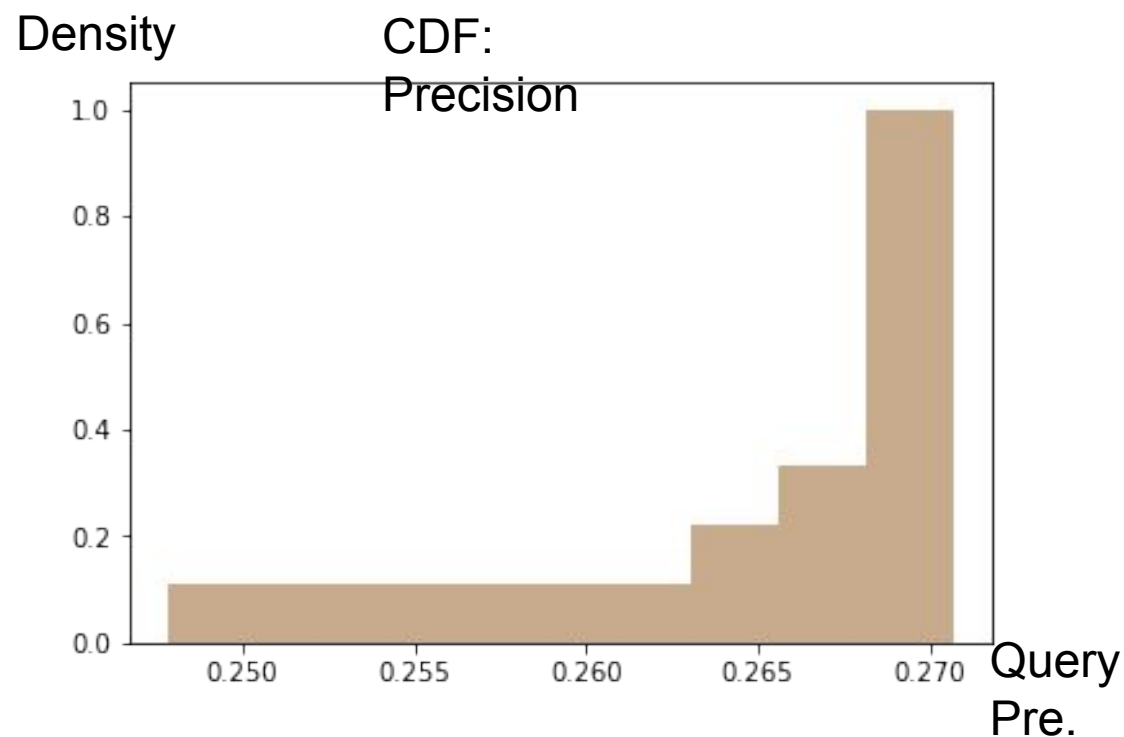
Text Corpus



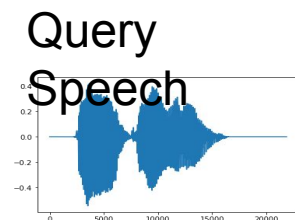
Three overlapping blue-outlined rectangular boxes representing the text corpus.



Performance distribution

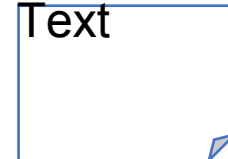


Search with ASR



AS
R

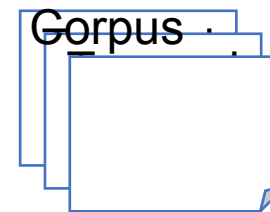
Query
Text



A rectangular box with a blue border and a folded bottom-right corner, representing a text document.

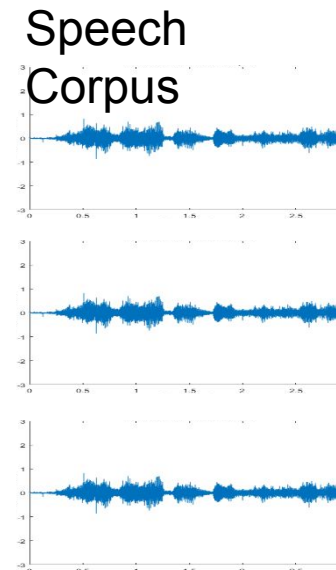
Search

Text
Corpus



Three overlapping rectangular boxes with blue borders and folded bottom-right corners, representing a collection of text documents.

AS
R



- Problem

- ASR Error: from both sided!!

- Short token conversion

- OOV

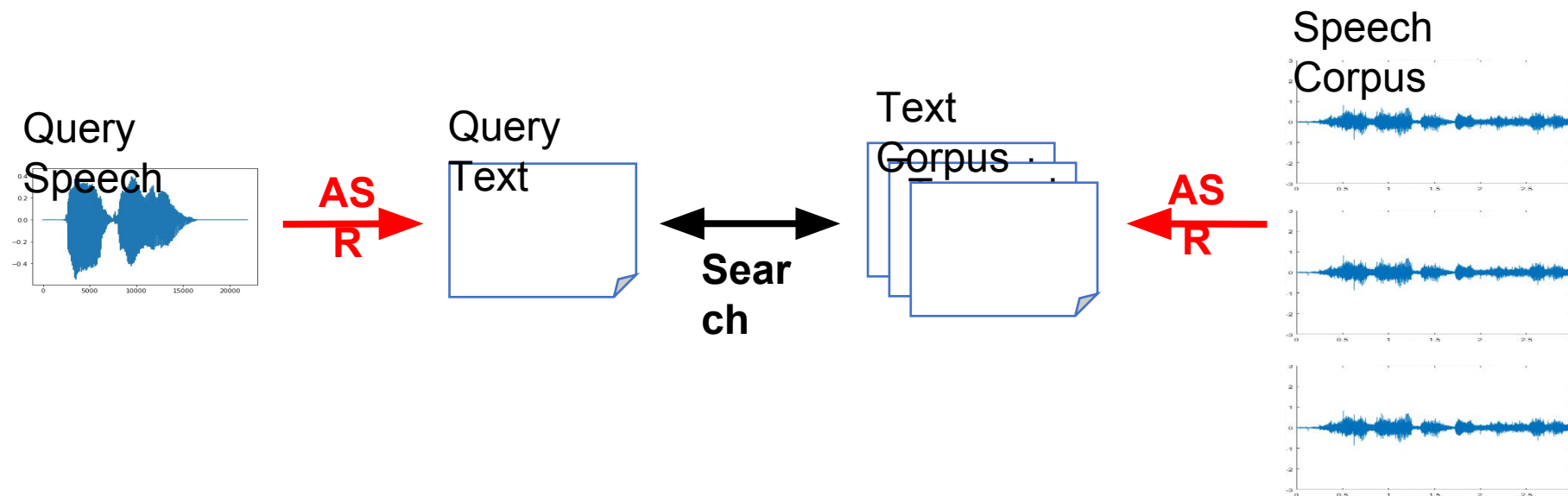
Antonia -> Anthony

ph d -> PhD

rangitata (place) -> rangi

Hansel -> cancel

Search with ASR

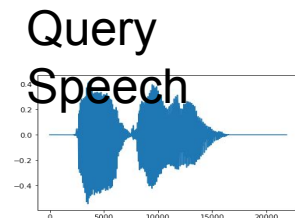


Can we reduce the errors propagated by ASR?

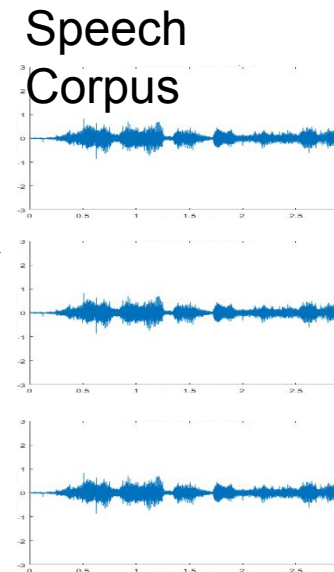
Task Setting

- Find embeddings of audio files
- Nearest neighbor search for each embedding

Search with MFCC

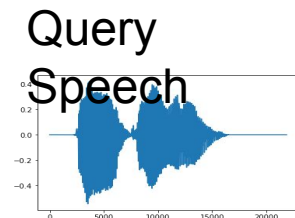


MFCC
representation
Search

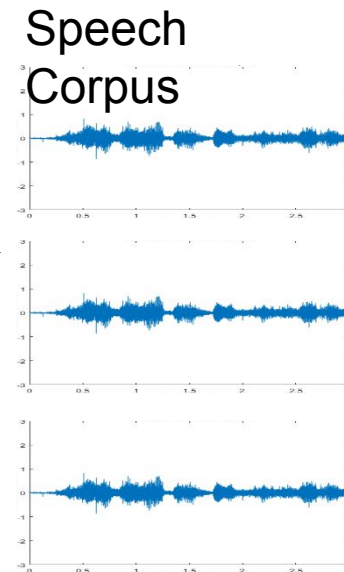


- MFCC:
 - The most common feature used in ASR system
 - Mimics people perception system and extract lower frequency features

Search with MFCC



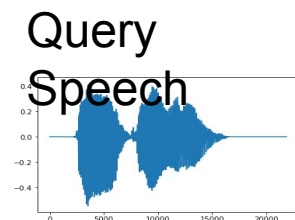
MFCC
representation
Search



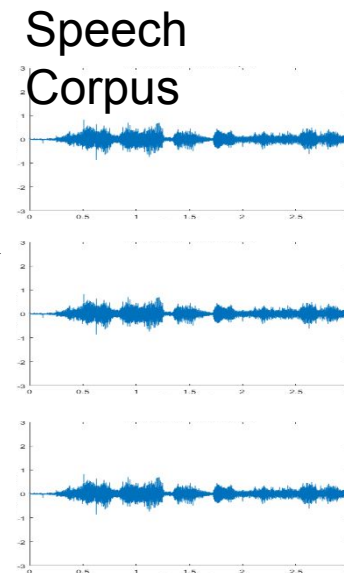
- How do we search?
 - Each query has ~100 MFCC frames
 - Retrieve top-k NN for each query and then use majority vote and distance to rank

Search with MFCC

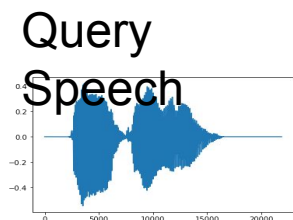
- Recall
- Precision



MFCC
representation
Search



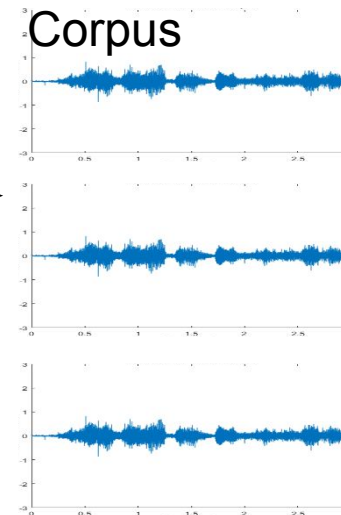
Search with MFCC



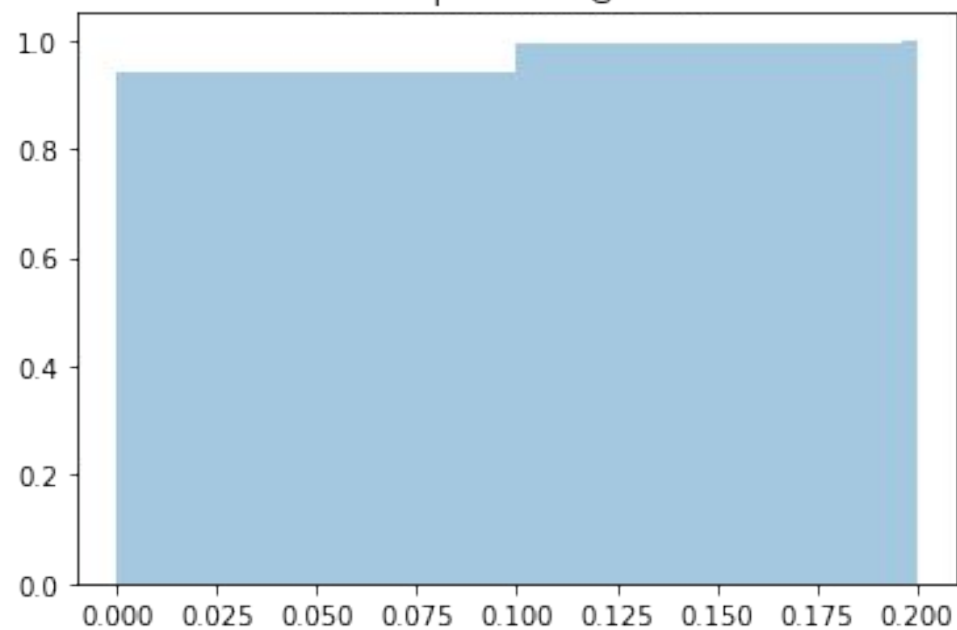
MFCC
representation

Search

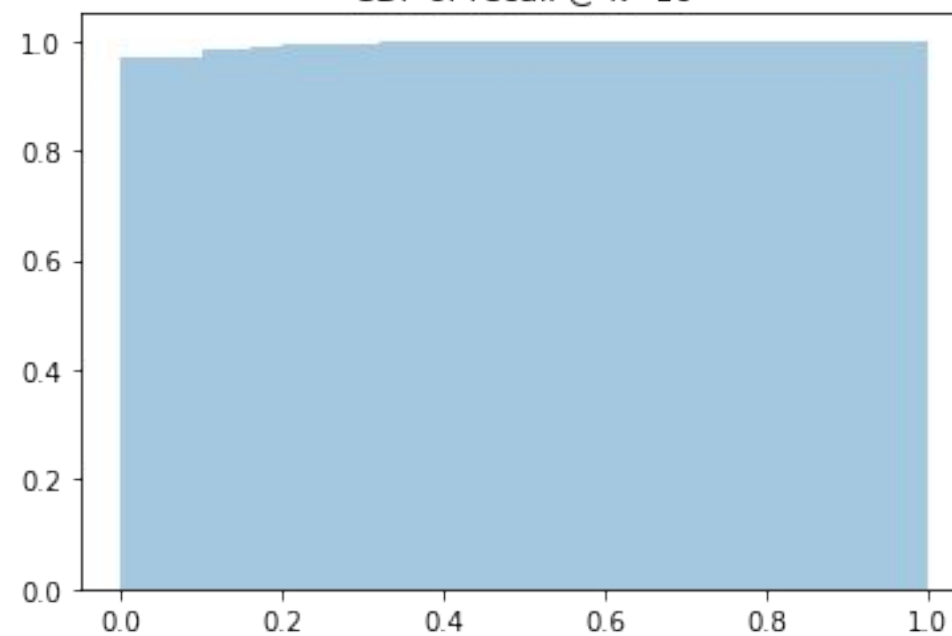
Speech
Corpus



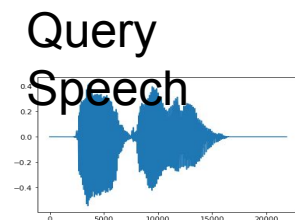
CDF of precision @ k=10



CDF of recall @ k=10

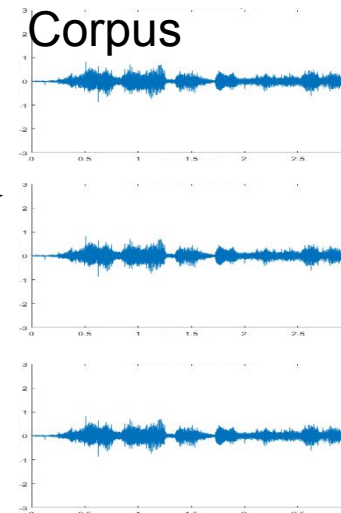


Search with MFCC



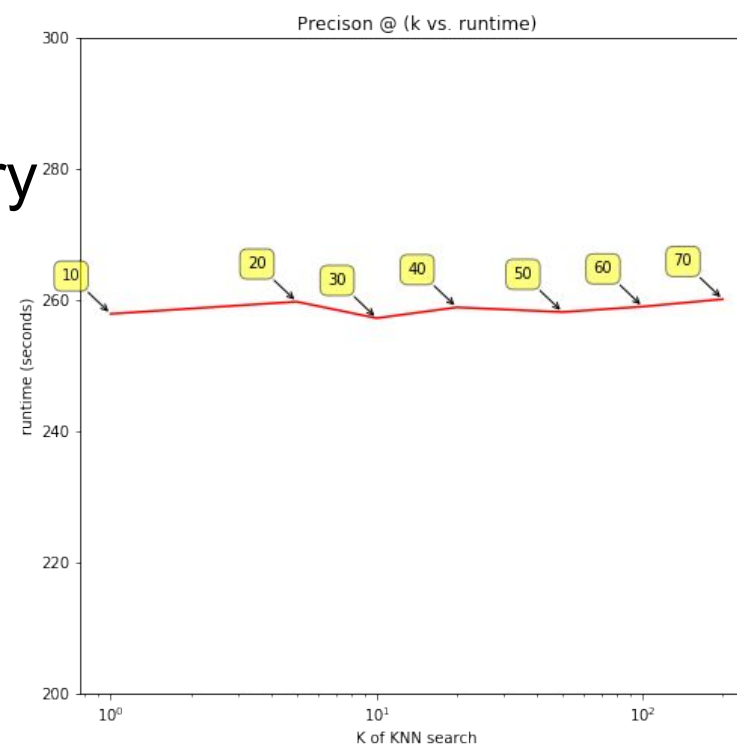
MFCC
representation
Search

Speech
Corpus

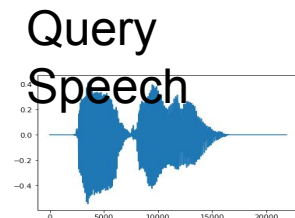


- Runtime analysis

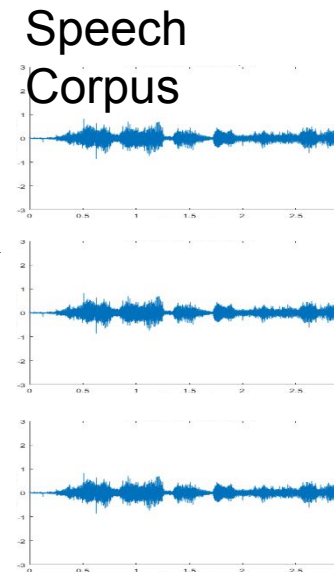
- Spark cluster with 4 workers
- Total 32 cores and 240 G memory
- K varies



Search with MFCC

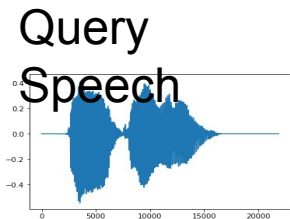


MFCC
representation
Search



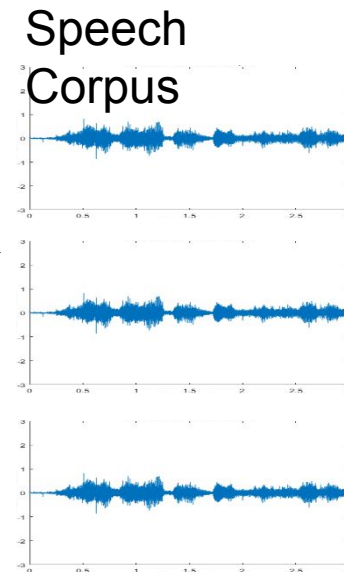
- Observation:
 - Not effective even with sliding window
 - Some queries have recall of 1.00 !!

Search with Deep Autoencoder

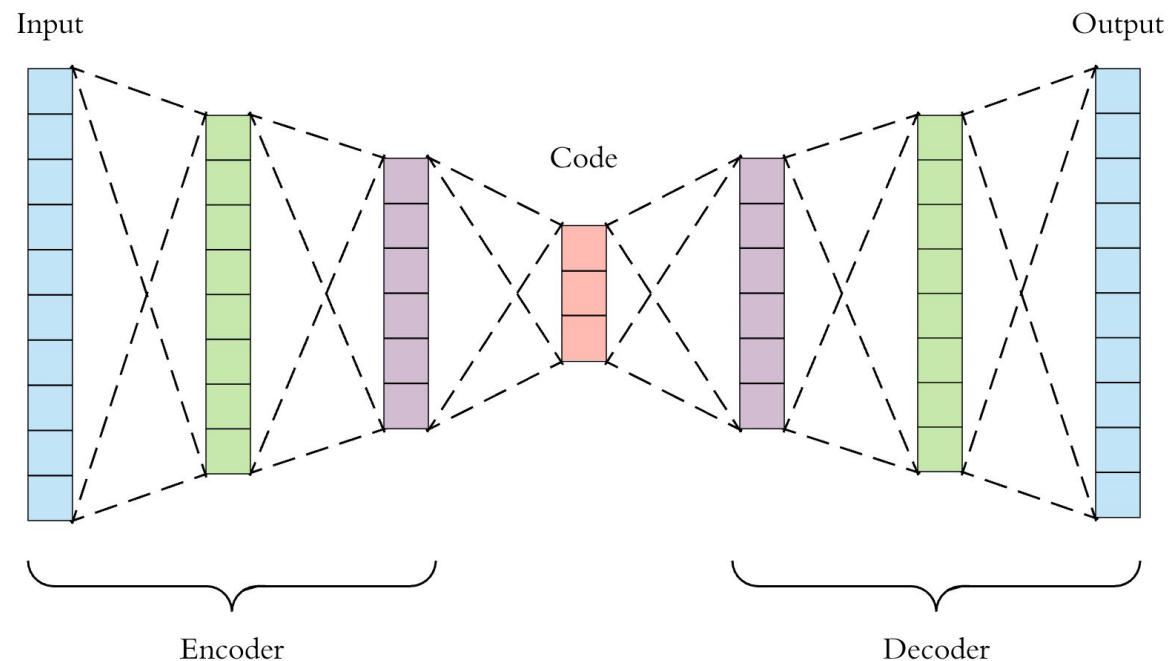


Autoencoder representation

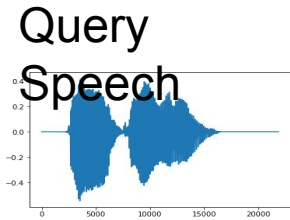
Search



- Autoencoder: ANN that generates output similar to the input
 - 3 Layer: [500, 180, 120]
 - Pretrain
 - Input: 100 MFCC frames = 1200
 - Sliding windows with stride = 5



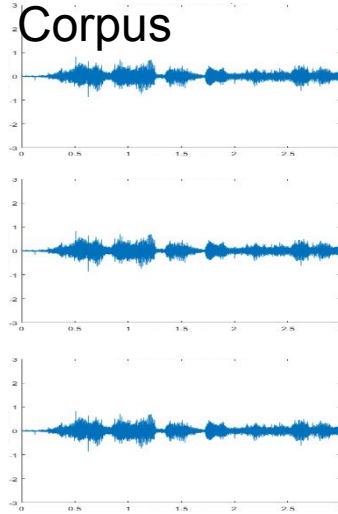
Search with Deep Autoencoder



Autoencoder representation

Search

Speech Corpus



- Recall
- Precision
- Problem

Takeaway

- Signal processing is hard

Future work

- Spectrogram:
 - Time-versus-freq. features
- Last layer encoding of ASR

