

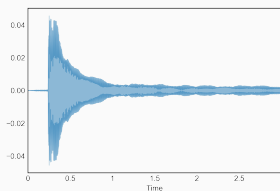
End to End ASR and Self Supervised Learning

Anirudh Gupta



Automatic Speech Recognition

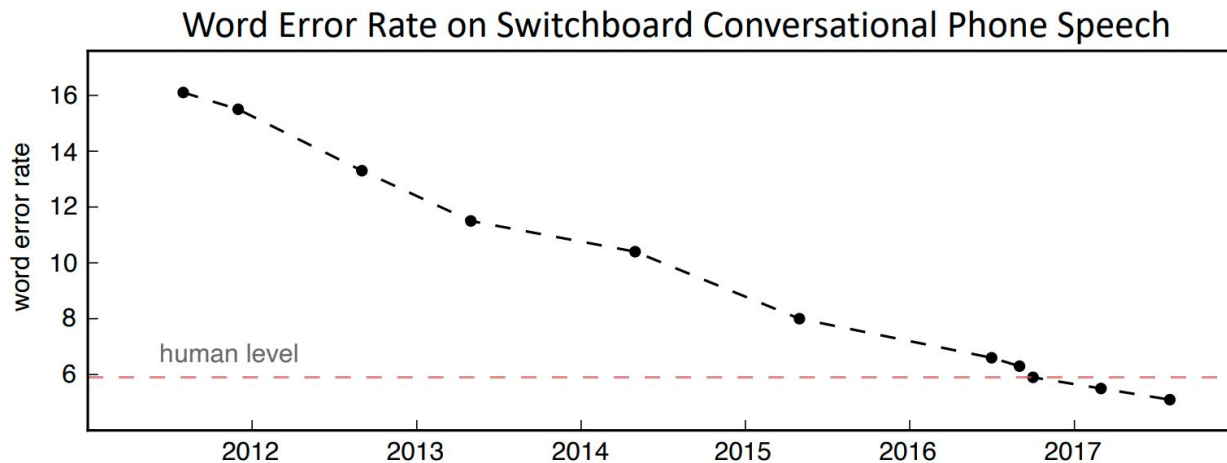
Goal: Input Speech -----> Output Transcription



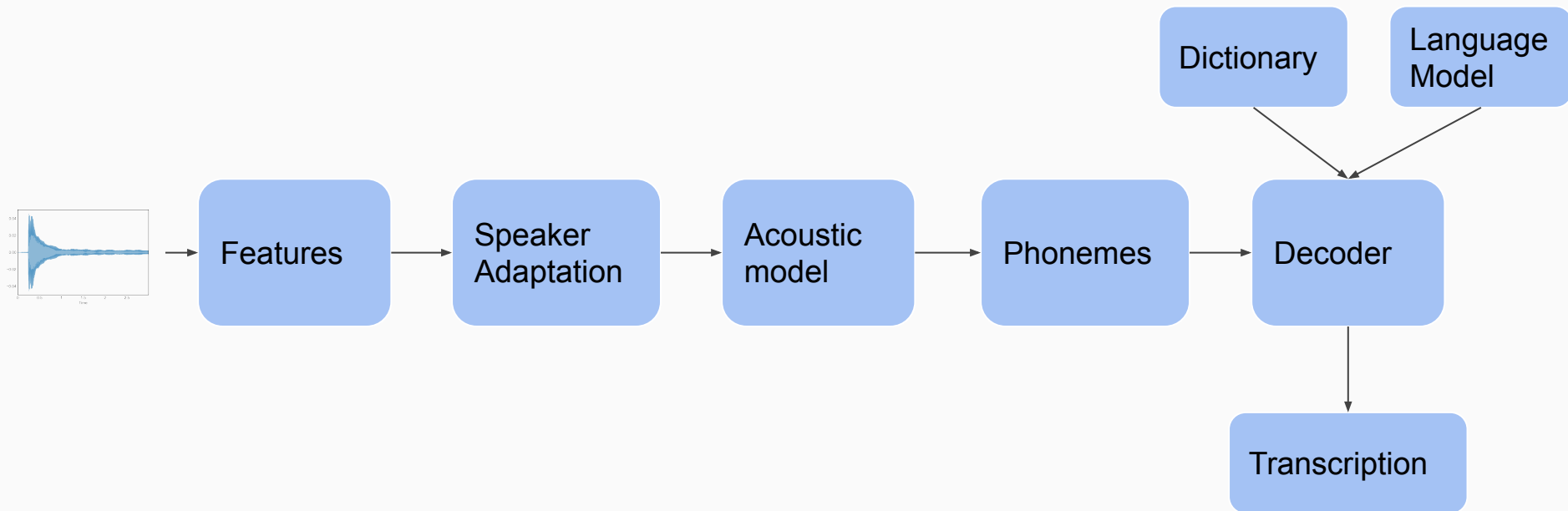
Automatic Speech
Recognizer

The quick brown fox jumps
over the lazy dog

Modern Speech Recognition



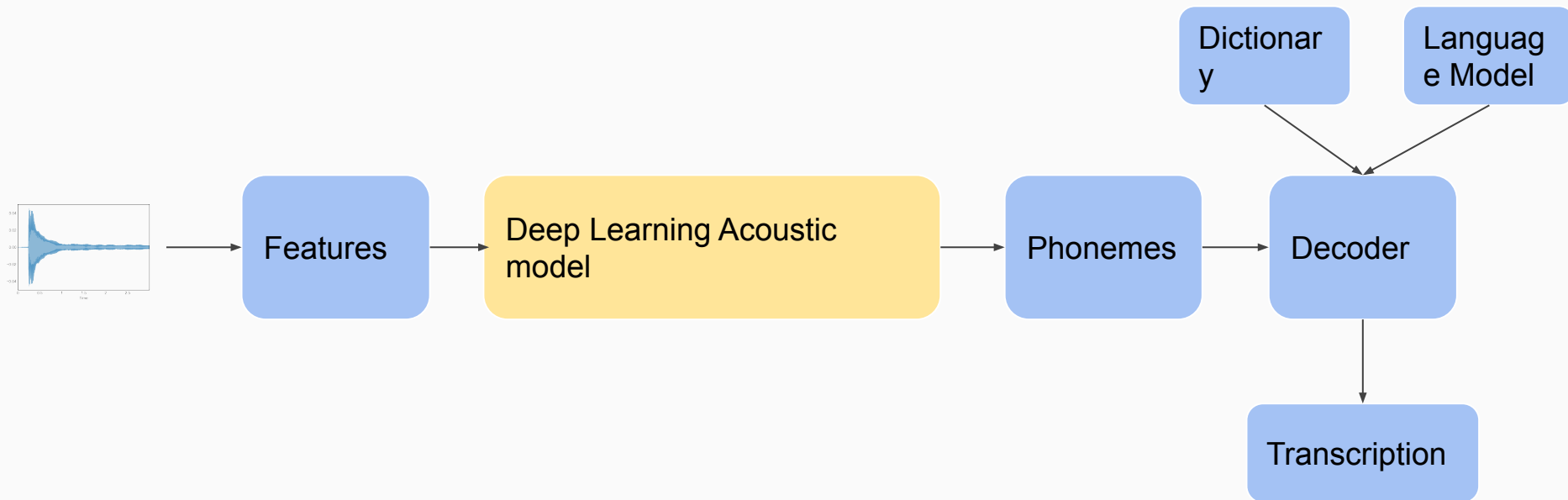
Traditional ASR Systems



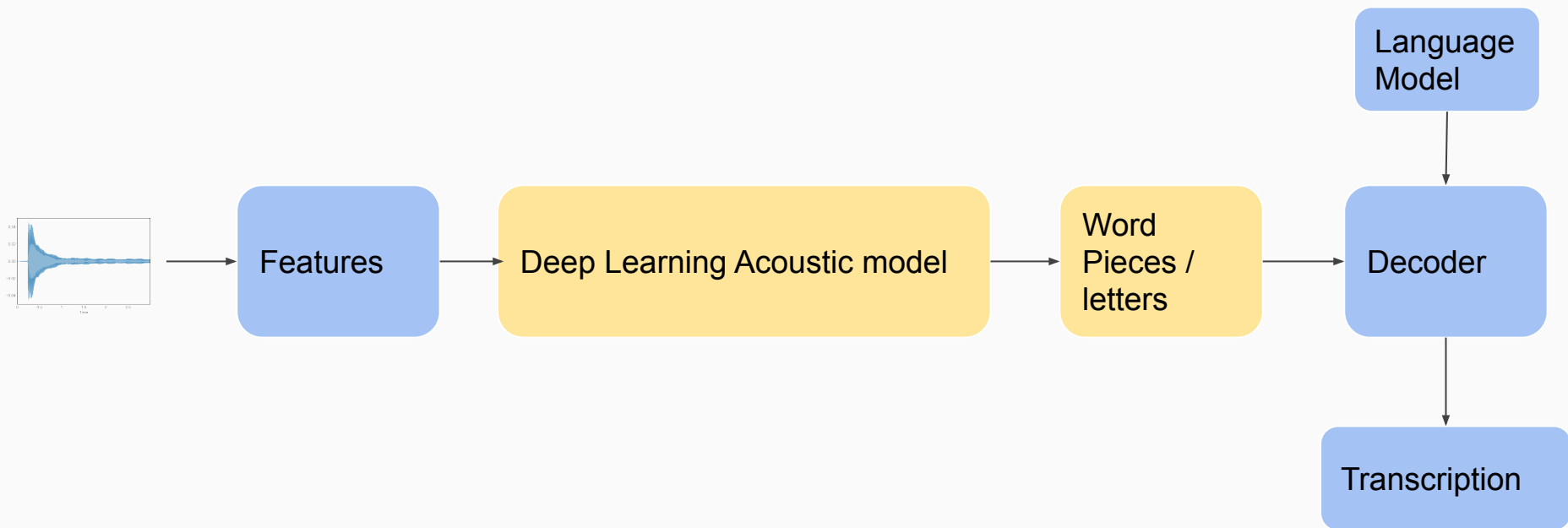
Traditional ASR Systems

- Traditional pipeline is highly tweak-able, but also hard to get working well.
- Historically, each part of the system has its own set of challenges.
Eg - choosing feature representation
- Cascading errors due to different individual components.

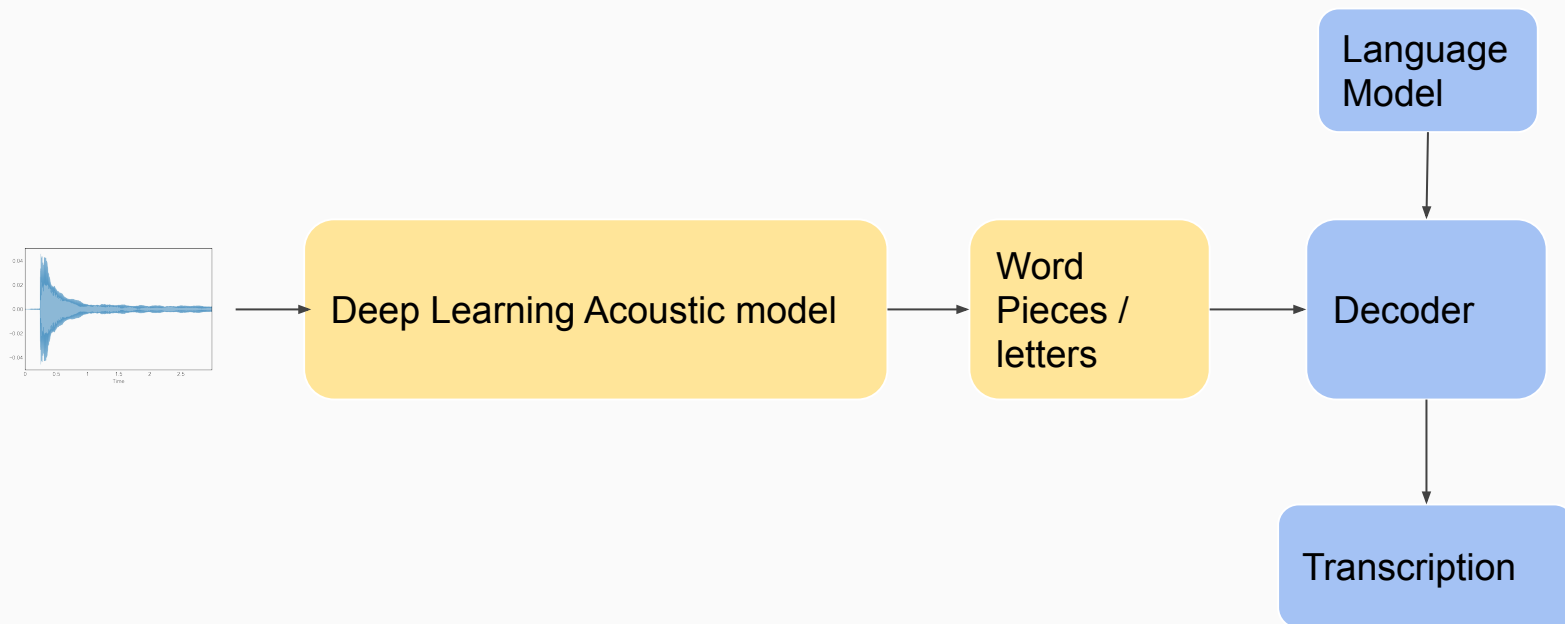
DL based acoustic models



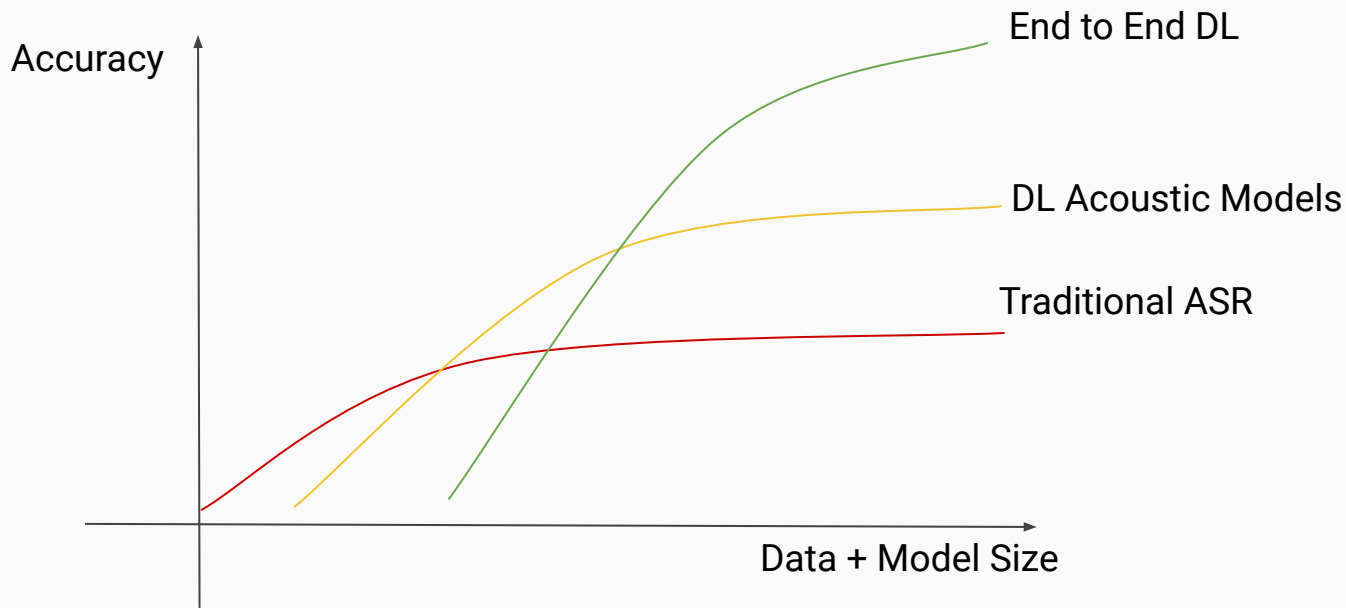
End to end networks



End to end networks

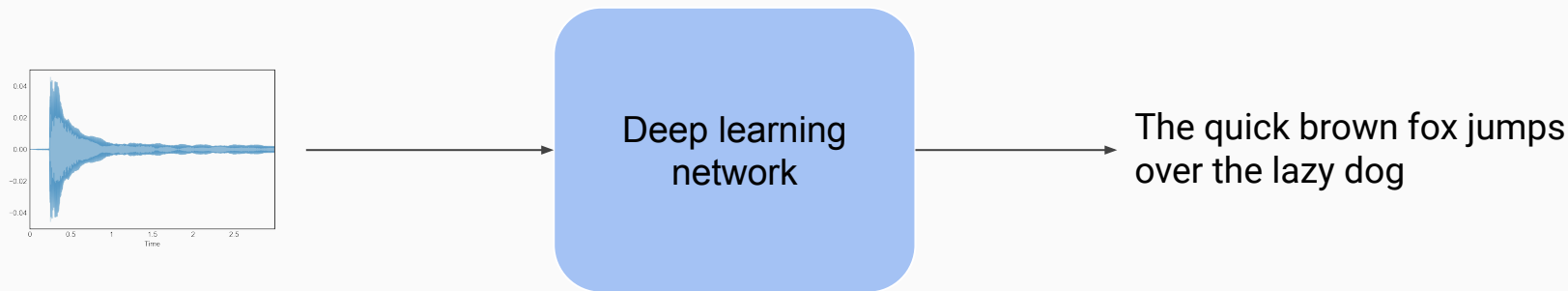


Moving towards end to end networks



Supervised E2E Networks

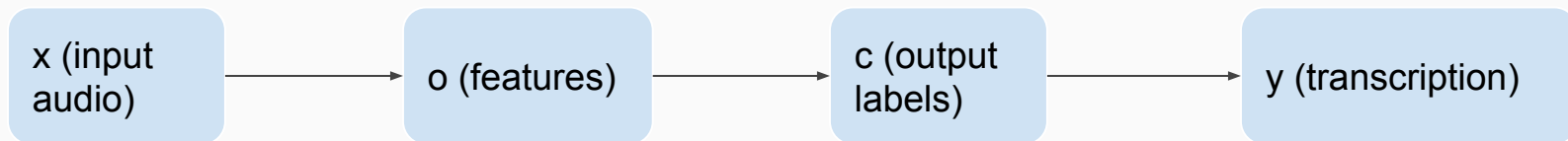
Towards end to end networks



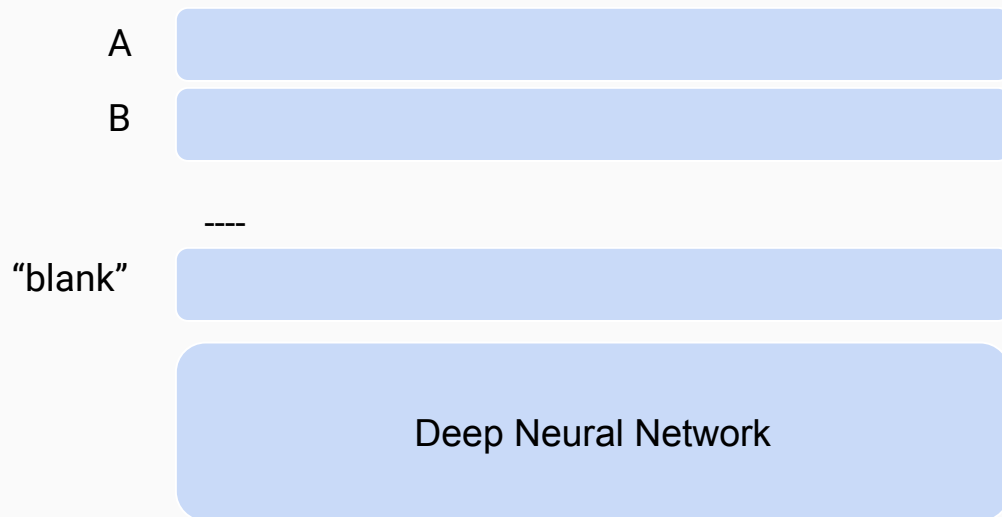
- Main issue: $\text{length}(x) \neq \text{length}(y)$
 - Difficult to understand how y maps to frames of audio. (eg. same word over different time scales)
- Initial fix: Train a network that would predict the sound at every frame. (alignment)

Connectionist Temporal Classification (CTC)

- A network outputs c , a distribution over symbols. ($\text{length}(c) = \text{length}(x)$)
 $c \in \{A, B, C, \dots, Z, \text{blank}\}$
- Define a mapping $\alpha : c \rightarrow y$
- Maximize likelihood of y^* under this model.



Connectionist Temporal Classification (CTC)



Connectionist Temporal Classification (CTC)

Assuming independence a distribution can be defined over the whole character sequences.

$$P(c|x) = \prod_{i=1}^N P(c_i|x)$$

Eg. $P(c = \text{HHH_E_LL_LO_}_ _ _ | x) = P(c_1 = \text{H}|x) P(c_2 = \text{H}|x) \dots P(c_{15} = \text{blank}|x)$

Connectionist Temporal Classification (CTC)

Define a mapping $\alpha : c \rightarrow y$

Given a specific character sequence c , squeeze out duplicates and blanks to yield transcription.

$y = \alpha(c) = \alpha(\text{HHH_E_LL_LO_}) = \text{"HELLO"}$

$$P(y|x) = \sum_{c:\alpha(c)=y} P(c|x)$$

Connectionist Temporal Classification (CTC)

$$\begin{aligned}\theta^* &= \operatorname{argmax}_{\theta} \sum_i \log P(y^{*(i)} | x^{(i)}) \\ &= \operatorname{argmax}_{\theta} \sum_i \log \sum_{c: \alpha(c)=y^{*(i)}} P(c | x^{(i)})\end{aligned}$$

[Graves et. al 2006] provides an efficient dynamic programming algorithm to compute the summation and its gradient.

DeepSpeech

- Trained on almost ~10,000 hours of audio data.
- Does not need hand designed components to handle noise, reverberation or speaker variation.
- Does not need a phoneme dictionary.
- Achieved 16% WER on Switchboard Hub5'00

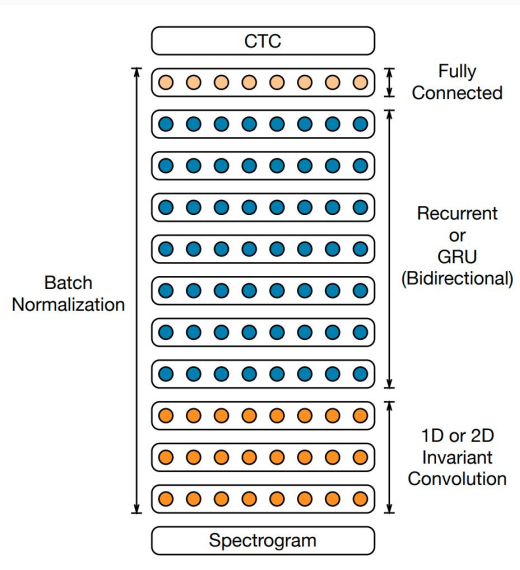


Fig: Network Architecture

Types of Speech Data

- Find data that matches our goals.

Styles of Speech: Read, conversational, Spontaneous

Issues: Disfluency, noise, mic quality, reverb, echo

Applications: Dictation, meeting transcriptions, call centres etc.

Speech Data

- Transcribing speech data isn't cheap.
- Typical speech benchmarks offer 100s to few 1000s hours of data.
 - LibriSpeech
 - LDC corpora
 - LDCIL

Why Self Supervised Learning?

- AI has made tremendous progress in developing systems that can learn from massive amounts of carefully labelled data.
- However, this is a bottleneck for building more generalist models that can perform multiple tasks without massive amounts of data.
- It is practically impossible to label everything in the world!

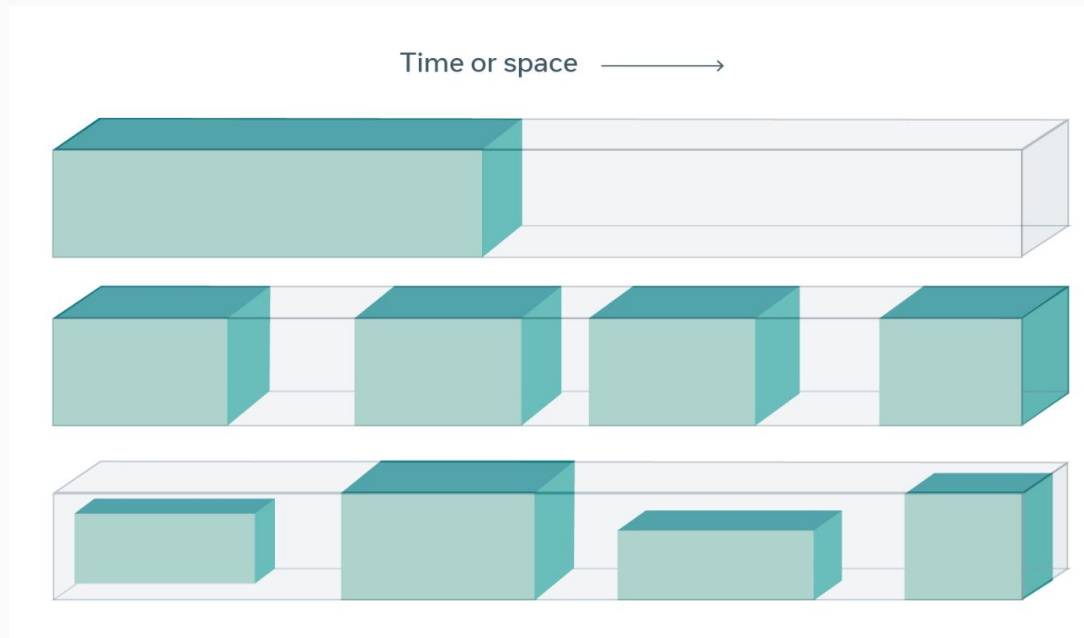
“Alternate” Supervision

- Getting labelled data is extremely difficult and expensive. (Switchboard)
- Labelling does not scale very well.

Semi automatic process: Obtain data using forced alignment and other techniques.

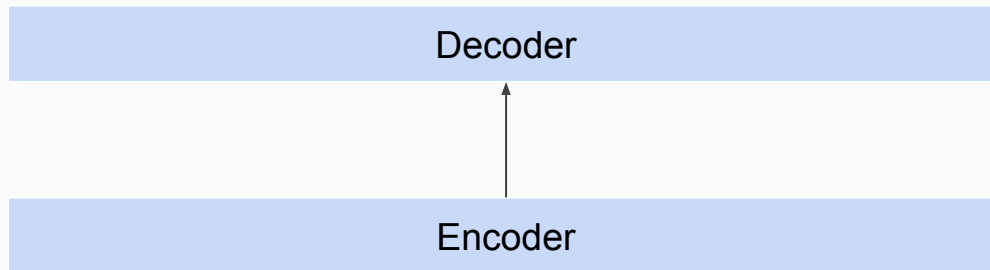
Using the data itself: **“self”- supervised**

Self Supervision as predictive learning



Self Supervised Learning in NLP

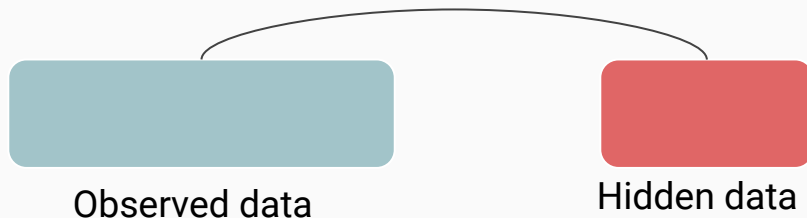
OUTPUT : The quick brown fox jumps over the lazy dog



INPUT : The quick [.....] fox jumps [.....] the lazy [.....]

Contrastive Learning

- Technique to learn general features of a dataset without labels by teaching the model which data points are similar or different.
- Given a set of input labels $\{x_i\}$, each corresponding to a label $y_i \in \{1, \dots, L\}$ among L classes.
- We would like to learn a function $f_{\theta}(\cdot) : X \rightarrow \mathbb{R}^d$ that encodes x_i into an embedding vector space such that examples from same classes have similar embeddings and sampled from different classes have very different ones.



Wav2vec 2.0

- Learn audio representations in a pretraining task solely from audio data as a contrastive task.
- Fine tune on **some** amount of audio data for speech recognition task using CTC loss function.

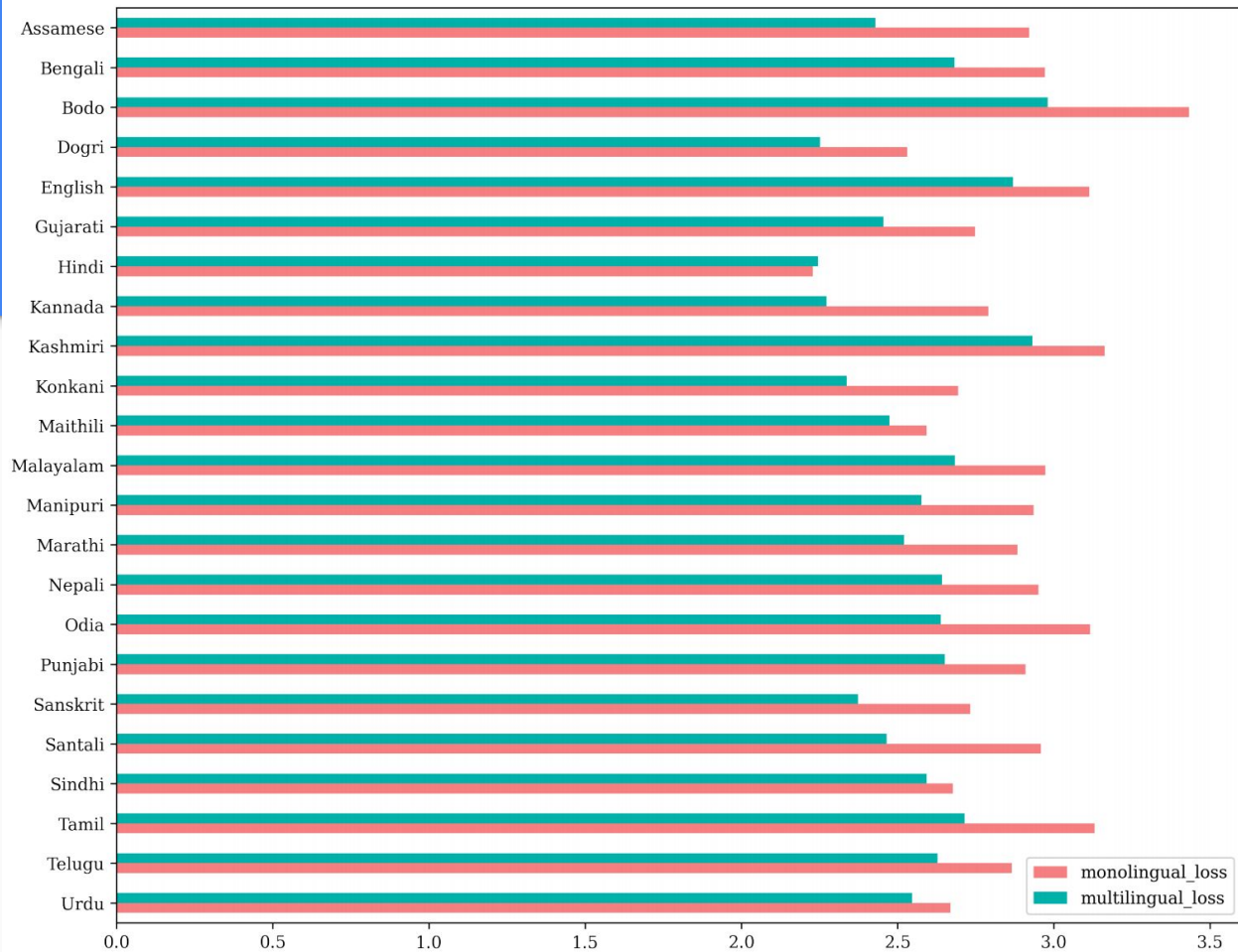
Wav2vec 2.0

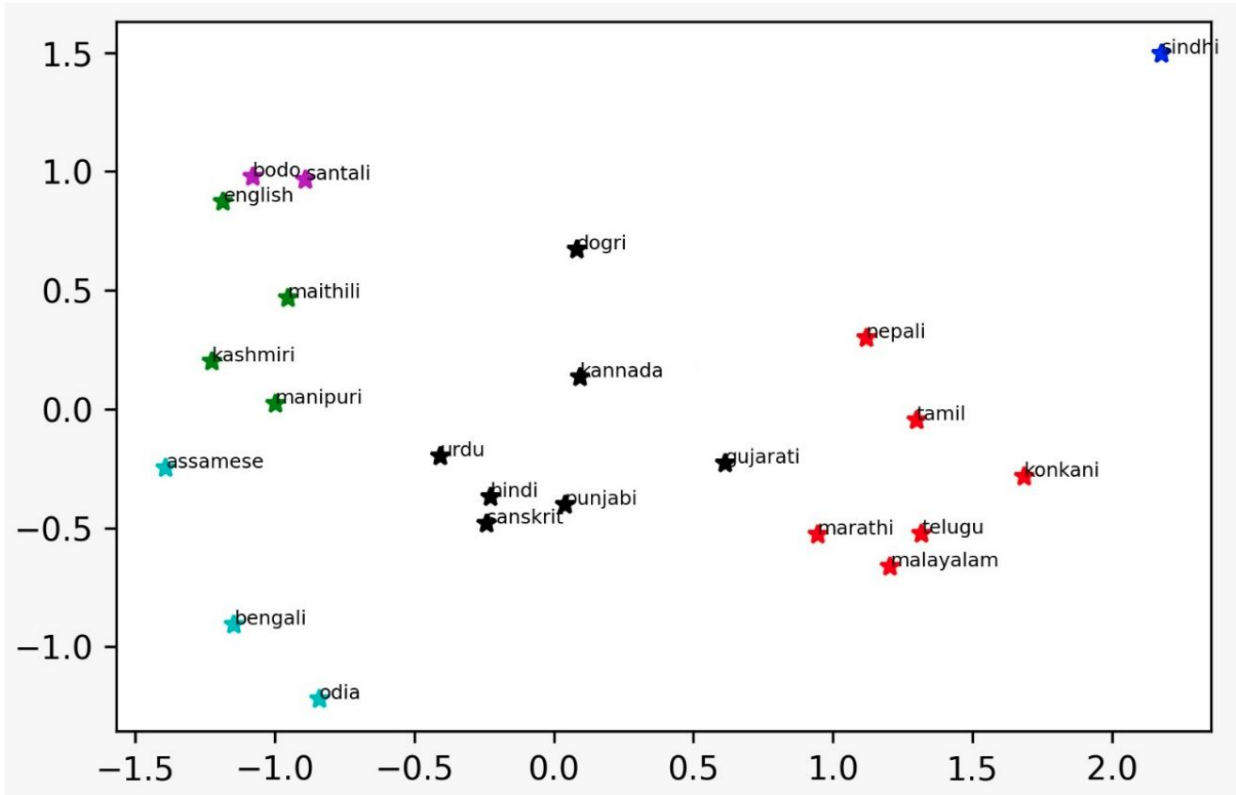
- Compensate for amount of fine-tuning data with the amount of pre-training data.

Model	Unlabeled data	LM	dev		test	
			clean	other	clean	other
10 min labeled						
BASE	LS-960	None	46.1	51.5	46.9	50.9
		4-gram	8.9	15.7	9.1	15.6
		Transf.	6.6	13.2	6.9	12.9
LARGE	LS-960	None	43.0	46.3	43.5	45.3
		4-gram	8.6	12.9	8.9	13.1
		Transf.	6.6	10.6	6.8	10.8
LARGE	LV-60k	None	38.3	41.0	40.2	38.7
		4-gram	6.3	9.8	6.6	10.3
		Transf.	4.6	7.9	4.8	8.2
1h labeled						
BASE	LS-960	None	24.1	29.6	24.5	29.7
		4-gram	5.0	10.8	5.5	11.3
		Transf.	3.8	9.0	4.0	9.3
LARGE	LS-960	None	21.6	25.3	22.1	25.3
		4-gram	4.8	8.5	5.1	9.4
		Transf.	3.8	7.1	3.9	7.6
LARGE	LV-60k	None	17.3	20.6	17.2	20.3
		4-gram	3.6	6.5	3.8	7.1
		Transf.	2.9	5.4	2.9	5.8
10h labeled						
BASE	LS-960	None	10.9	17.4	11.1	17.6
		4-gram	3.8	9.1	4.3	9.5
		Transf.	2.9	7.4	3.2	7.8
LARGE	LS-960	None	8.1	12.0	8.0	12.1
		4-gram	3.4	6.9	3.8	7.3
		Transf.	2.9	5.7	3.2	6.1
LARGE	LV-60k	None	6.3	9.8	6.3	10.0
		4-gram	2.6	5.5	3.0	5.8
		Transf.	2.4	4.8	2.6	4.9
100h labeled						
BASE	LS-960	None	6.1	13.5	6.1	13.3
		4-gram	2.7	7.9	3.4	8.0
		Transf.	2.2	6.3	2.6	6.3
LARGE	LS-960	None	4.6	9.3	4.7	9.0
		4-gram	2.3	5.7	2.8	6.0
		Transf.	2.1	4.8	2.3	5.0
LARGE	LV-60k	None	3.3	6.5	3.1	6.3
		4-gram	1.8	4.5	2.3	4.6
		Transf.	1.9	4.0	2.0	4.0

Cross lingual representations

- Pre-training with multiple languages facilitates learning cross lingual speech representations.
- Data with high resource languages can be used to aid learning of low resource languages.
- Cross lingual pre-training significantly outperforms monolingual pre-training in speech recognition tasks.





Still not solved!

- Conversational speech with multiple speakers.
- Performance in noisy environments.
- Performance on esoteric accents and model bias.
- Goal: works on some speakers sometimes ----> works for all people all the time.
- Looking into semantic errors as well while evaluating WER.
- Human-error rate on benchmarks is often quite high since the evaluation is done in a context free environment.



The screenshot shows a Twitter thread with four tweets. The first tweet is by Thomas G. Dietterich (@tdietterich) from 22 Jul 21, discussing automated transcription of @icmlconf talks and mentioning the 'Mushroom learning community'. It has 8 retweets, 1 quote tweet, and 117 likes. The second tweet is a reply by Thomas G. Dietterich 20h ago, noting that 'decoder' was transcribed as 'Dakota' and 'vocal tract' as 'vocal Trump', with 14 likes. The third tweet is a reply by Piotr Żelasko (@PiotrZelasko) 55m ago, stating 'People tell me ASR is a solved problem' with a sad face emoji, and 1 like. The fourth tweet is a reply by Matt Post (@mjpost) 1h ago, mentioning Microsoft's success at NAACL in 2016, with 1 like. All tweets have icons for replies, retweets, likes, and shares.

Thomas G. Dietterich @tdietterich · 22 Jul 21 · Twitter Web App

It's a lot of fun watching the automated transcription of these @icmlconf talks. My favorite one: "Mushroom learning community"

9:16 · 22 Jul 21 · Twitter Web App

8 Retweets 1 Quote Tweet 117 Likes

Thomas G. Dietterich @tdietterich · 20h

"decoder" is consistently rendered as "Dakota" and the "vocal tract" came out as "vocal Trump"

14 Likes

Piotr Żelasko @PiotrZelasko · 55m

Replying to @tdietterich @stanfordnlp and @icmlconf

People tell me ASR is a solved problem 😞

1 Like

Thomas G. Dietterich @tdietterich · 53m

On the other hand, AI-based entertainment is stronger than ever!

1 Like

Show replies

Matt Post @mjpost · 1h

Replying to @tdietterich @stanfordnlp and @icmlconf

The only time I've seen this done right was when Microsoft provided it, for NAACL a few years ago, 2016 I think. They got all the PDFs ahead of time and adapted their LM.

Thank You