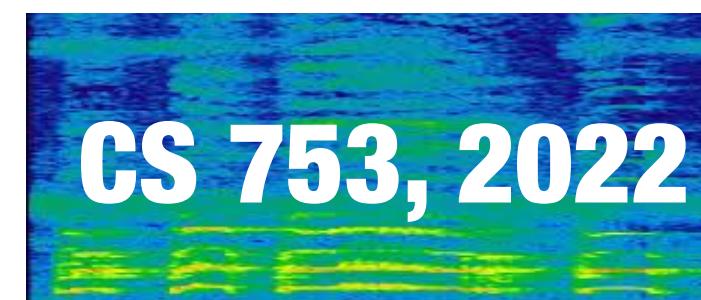


# **Speech Representations**



Instructor: Preethi Jyothi, IITB

# wav2vec

- Algorithm that uses raw audio to learn speech representations (“**self-supervised**” approach)

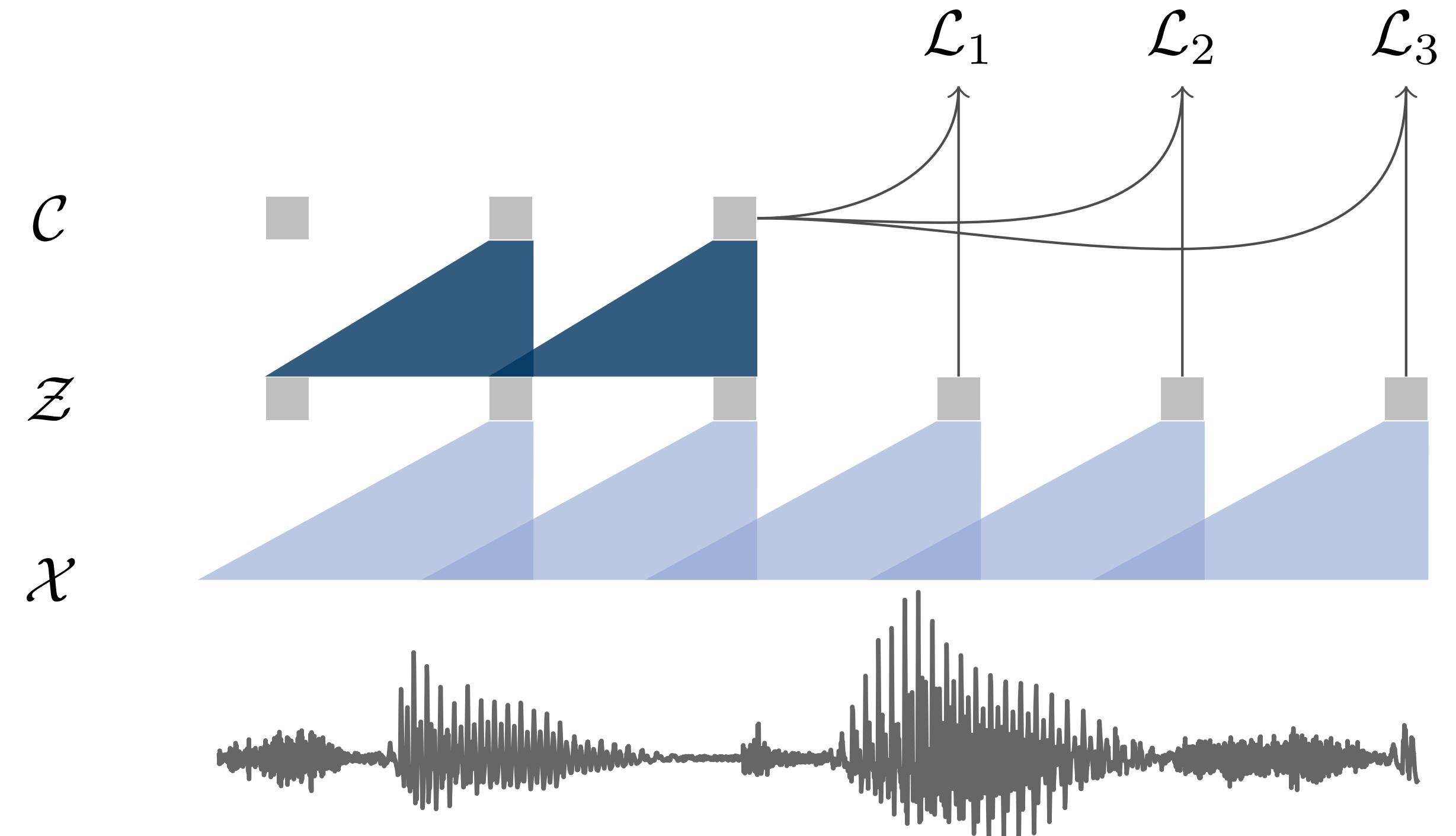
# wav2vec

- Algorithm that uses raw audio to learn speech representations (“**self-supervised**” approach)
  - Encoder network embeds raw audio into a latent representation ( $f : \mathcal{X} \rightarrow \mathcal{Z}$ ) and a context network combines multiple encoded representations into a contextualised embedding ( $g : \mathcal{Z} \rightarrow \mathcal{C}$ )

# wav2vec

- Algorithm that uses raw audio to learn speech representations (“**self-supervised**” approach)

- Encoder network embeds raw audio into a latent representation ( $f : \mathcal{X} \rightarrow \mathcal{Z}$ ) and a context network combines multiple encoded representations into a contextualised embedding ( $g : \mathcal{Z} \rightarrow \mathcal{C}$ )
- Train the model to minimize the following contrastive loss:



# wav2vec

- Algorithm that uses raw audio to learn speech representations (“self-supervised” approach)

- Encoder network embeds raw audio into a latent representation ( $f : \mathcal{X} \rightarrow \mathcal{Z}$ ) and a context network combines multiple encoded representations into a contextualised embedding ( $g : \mathcal{Z} \rightarrow \mathcal{C}$ )

$$\mathcal{L}_k = - \sum_{i=1}^{T-k} \log \frac{\exp(\text{sim}(c_i, z_{i+k}))}{\sum_{\tilde{z}} \exp(\text{sim}(c_i, \tilde{z}))}$$

$$L = \sum_k \mathcal{L}_k$$

$\sigma(c_i^T z_{i+k})$

$N = 1D$

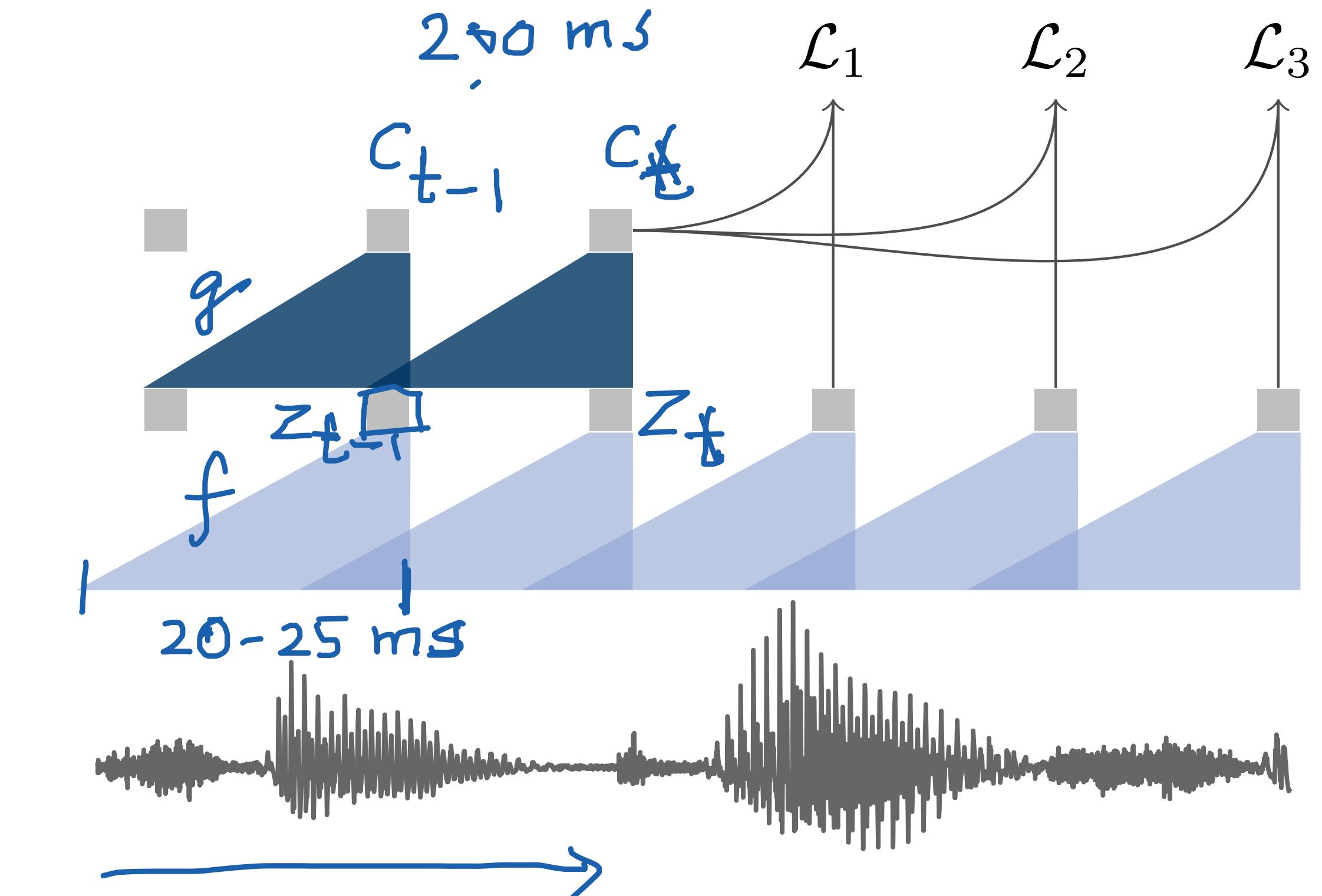
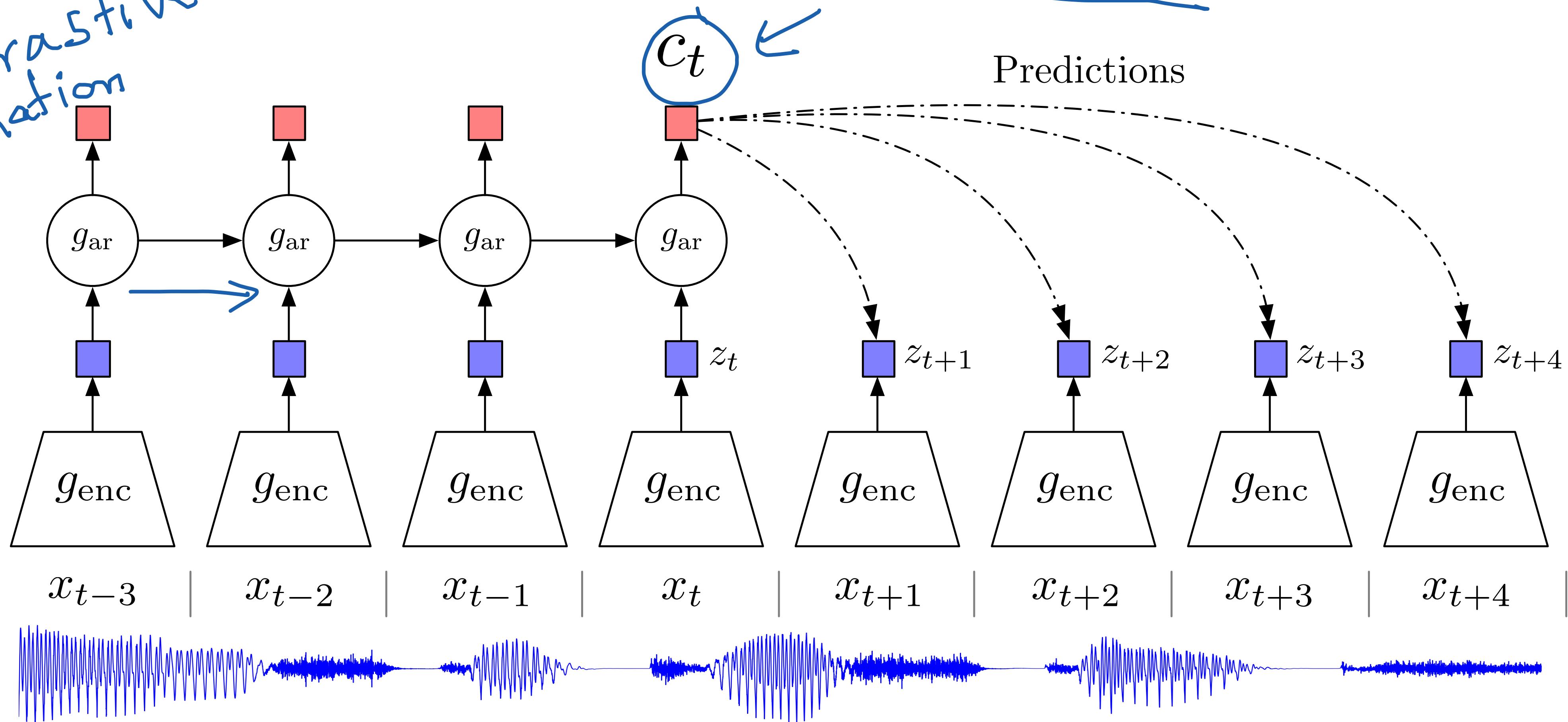


Image from: <https://arxiv.org/pdf/1904.05862.pdf>

# Contrastive Predictive Coding

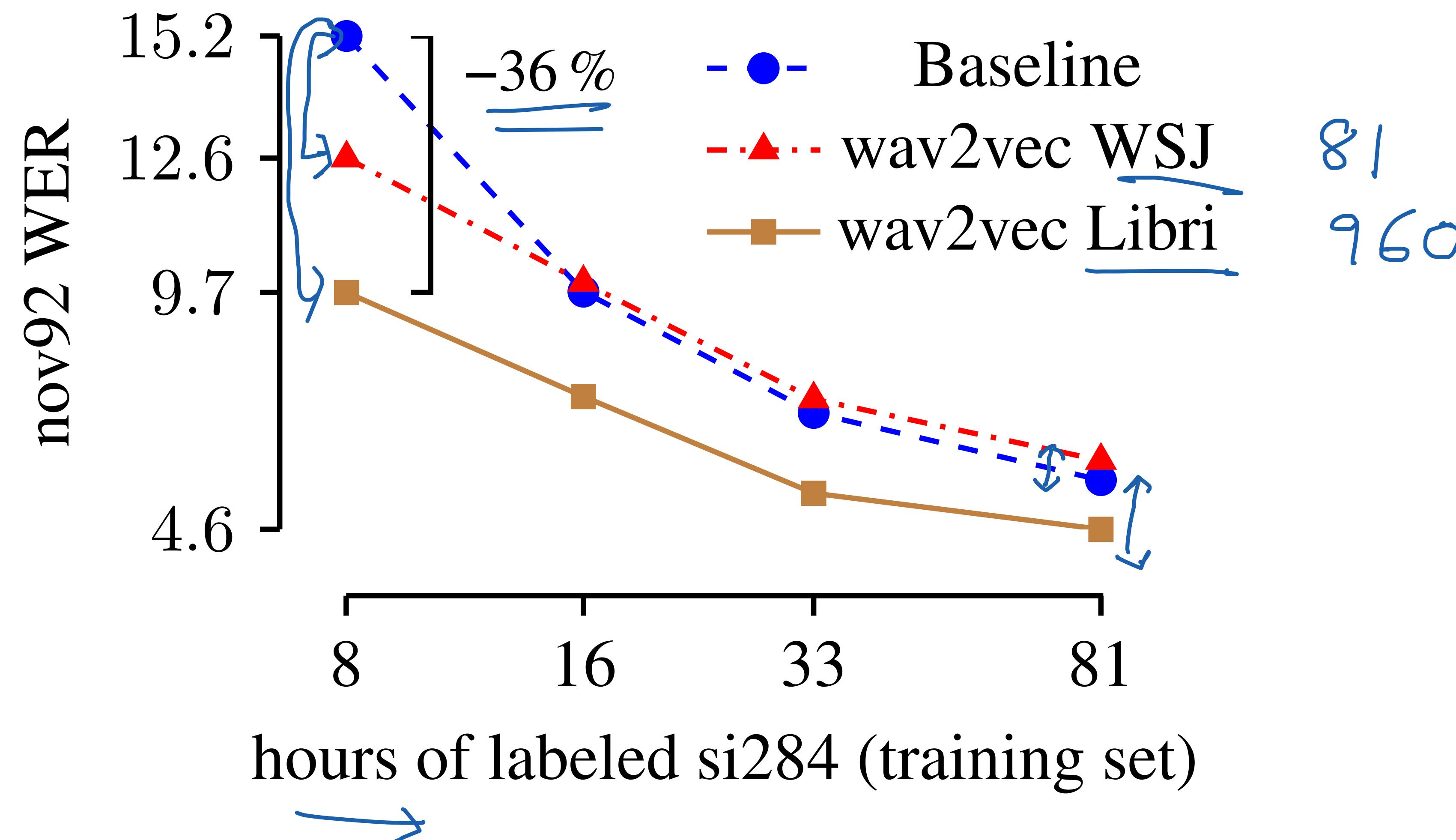
Noise  
Contrastive  
Estimation



$$\mathcal{L}_k = - \sum_{i=1}^{T-k} \log \frac{\exp(\text{sim}(c_i, z_{i+k}))}{\sum_{\tilde{z}} \exp(\text{sim}(c_i, \tilde{z}))}$$

Contrastive Loss

# wav2vec



# wav2vec 2.0: Learning Speech Representations from Raw Audio

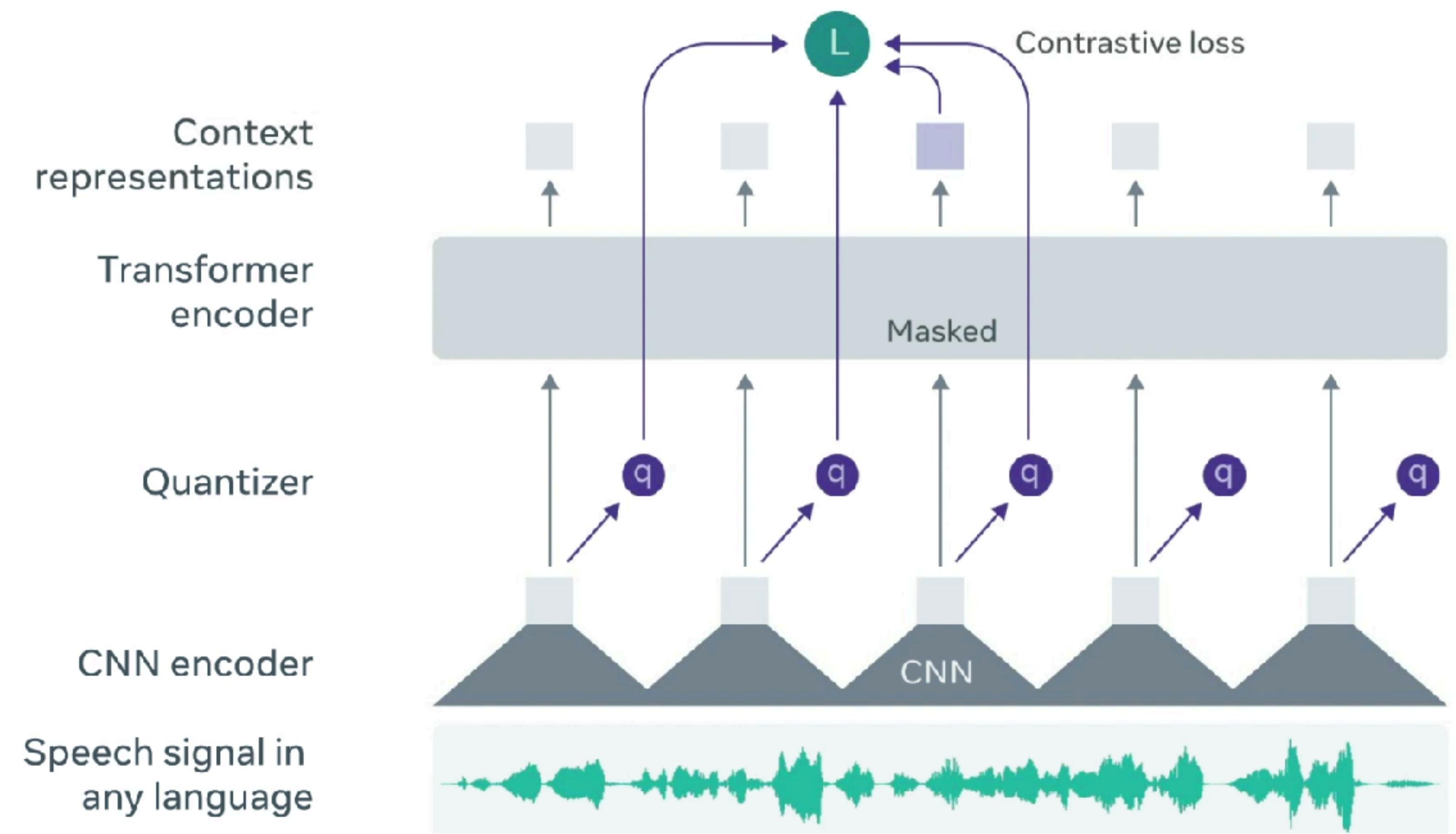


Image from: <https://ai.facebook.com/blog/wav2vec-2-0-learning-the-structure-of-speech-from-raw-audio/>

[1]: <https://arxiv.org/abs/1810.04805>

# wav2vec 2.0: Learning Speech Representations from Raw Audio

- Similar to wav2vec. Outputs from the encoder are further quantized.

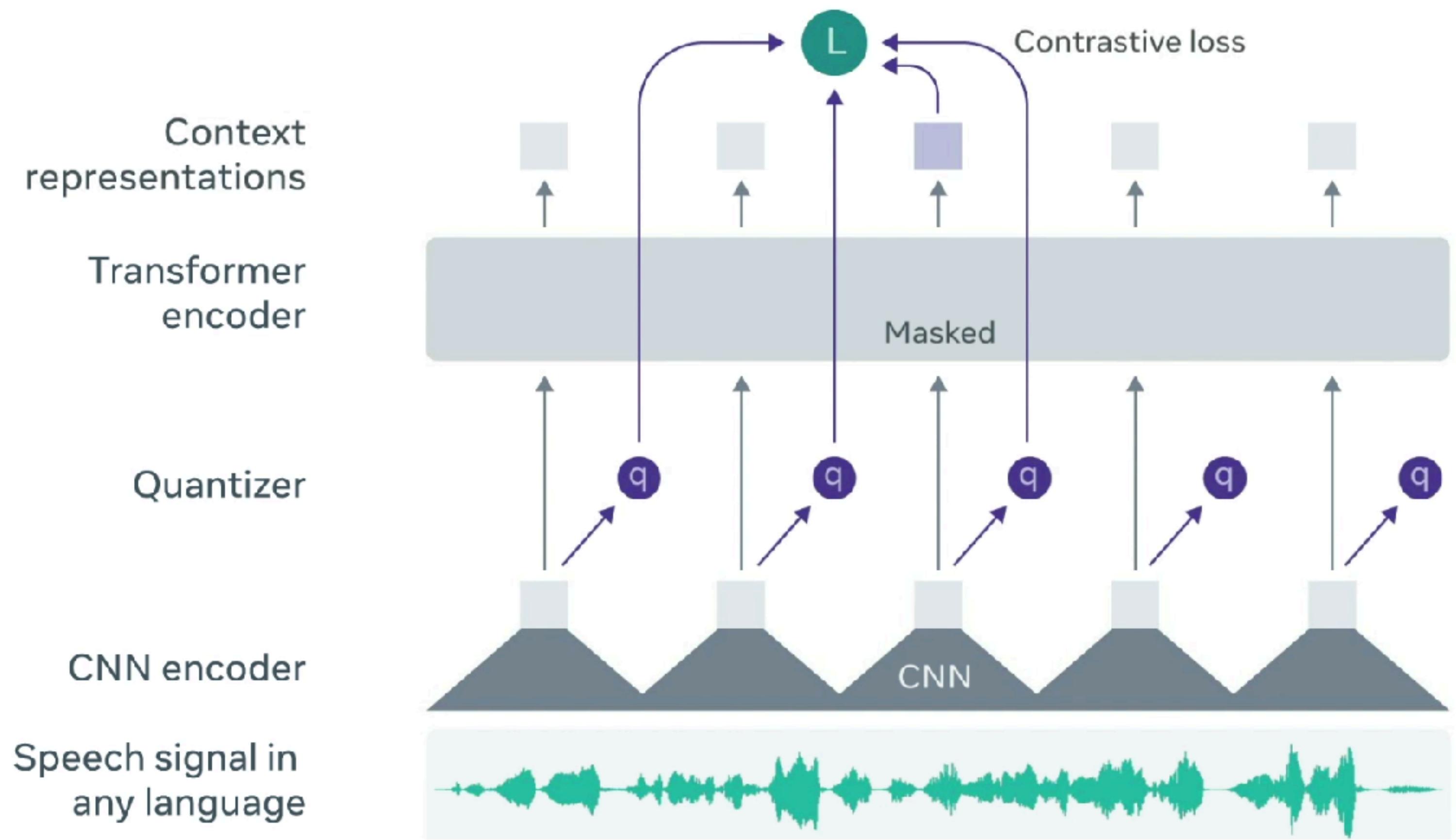


Image from: <https://ai.facebook.com/blog/wav2vec-2-0-learning-the-structure-of-speech-from-raw-audio/>

[1]: <https://arxiv.org/abs/1810.04805>

# wav2vec 2.0: Learning Speech Representations from Raw Audio

- Similar to wav2vec. Outputs from the encoder are further quantized.
- Masks spans of speech representations (as in masked language modelling for BERT [1])

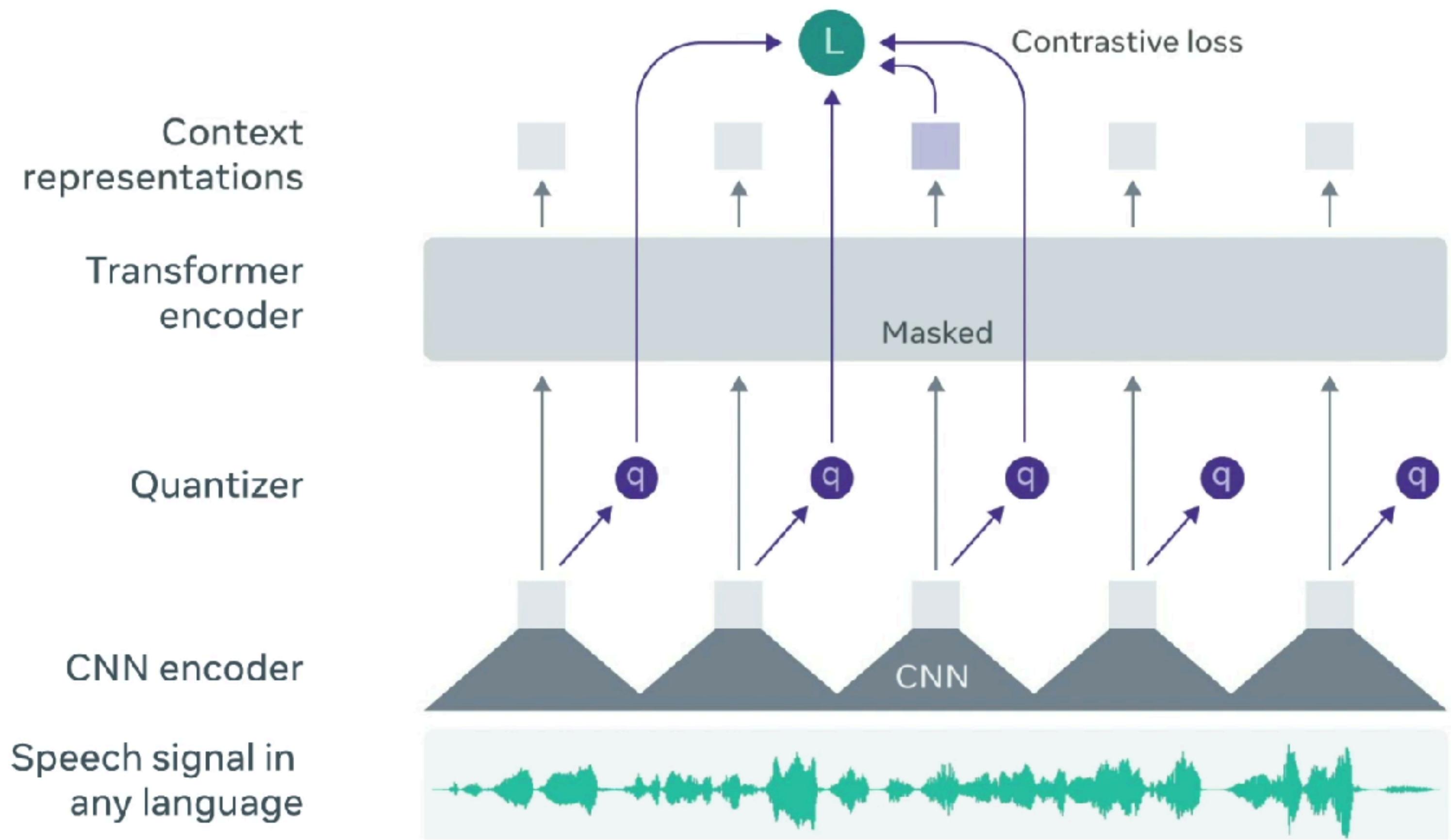


Image from: <https://ai.facebook.com/blog/wav2vec-2-0-learning-the-structure-of-speech-from-raw-audio/>

[1]: <https://arxiv.org/abs/1810.04805>

# wav2vec 2.0: Learning Speech Representations from Raw Audio

vq-wav2vec

- Similar to wav2vec. Outputs from the encoder are further quantized.
- Masks spans of speech representations (as in masked language modelling for BERT [1])
- Training objective is to recover the masked representations among a set of distractors.

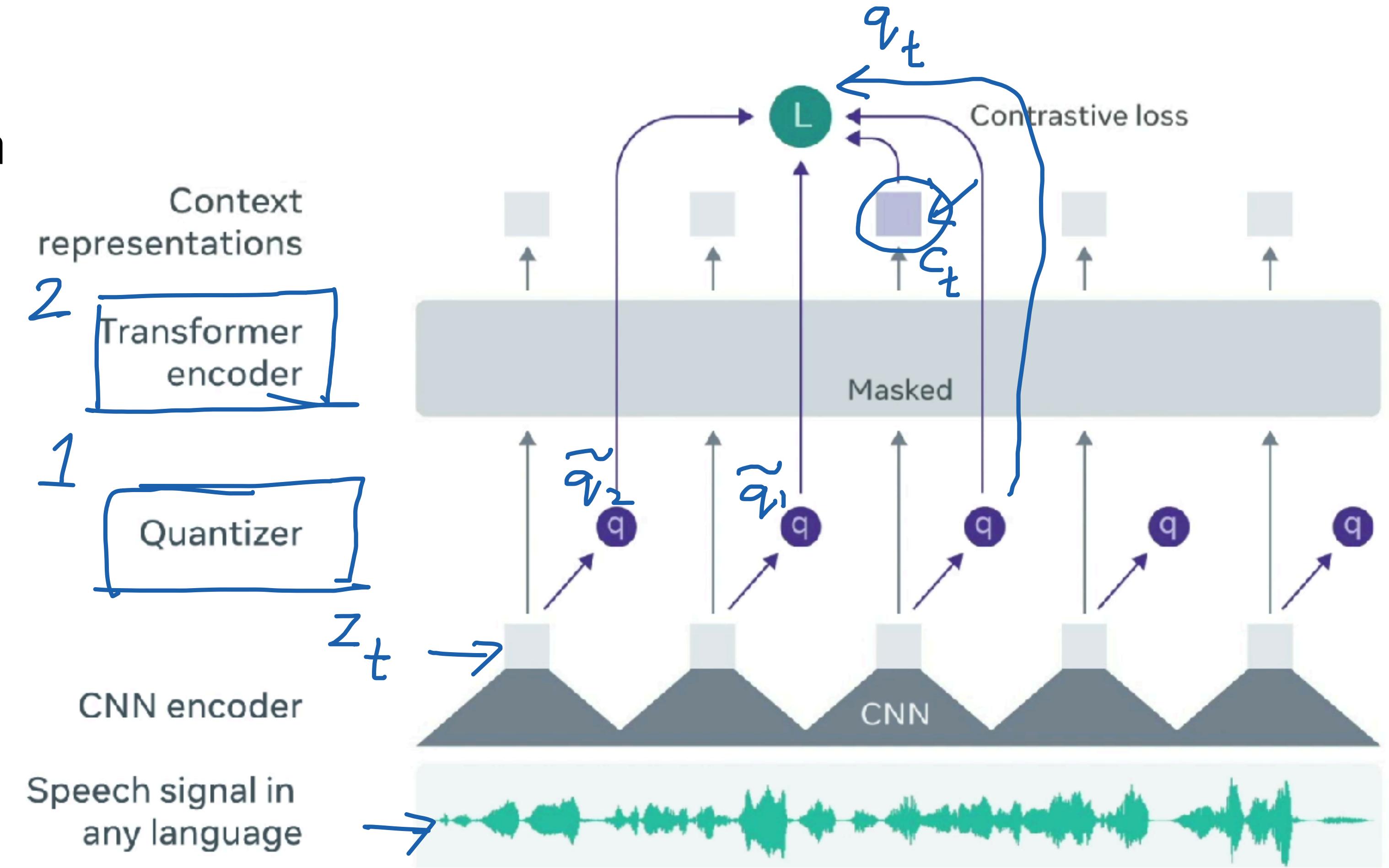
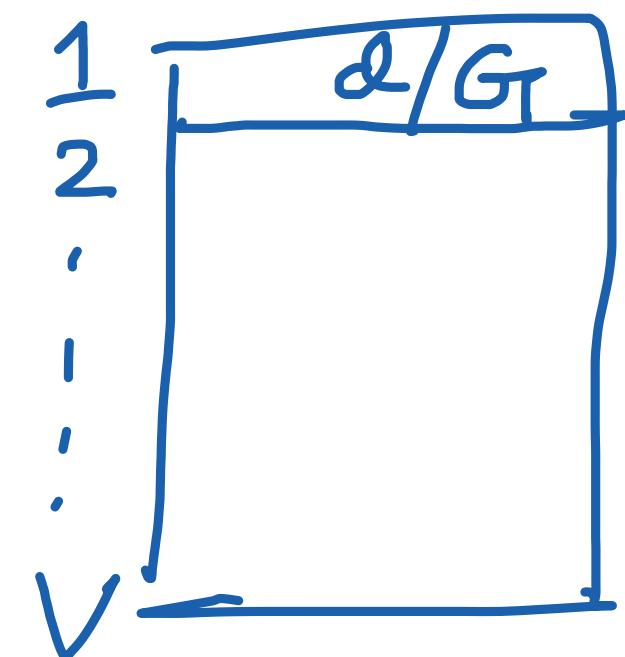


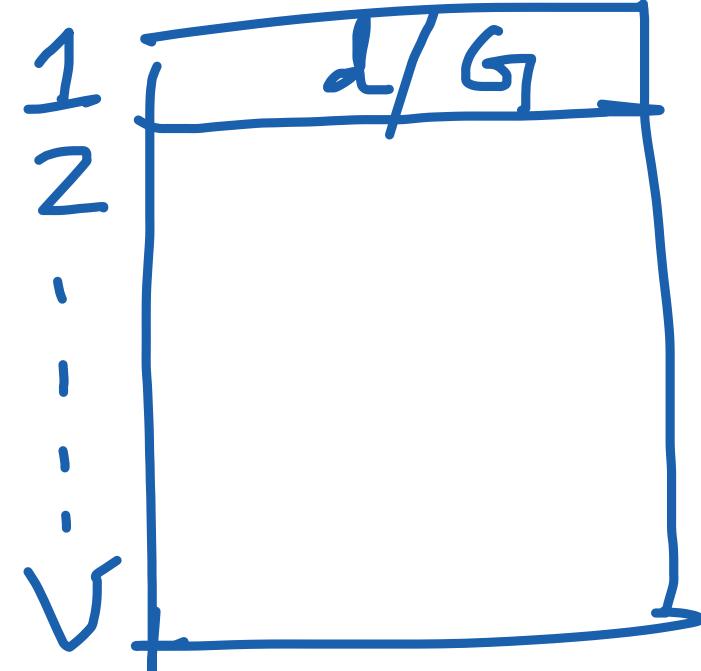
Image from: <https://ai.facebook.com/blog/wav2vec-2-0-learning-the-structure-of-speech-from-raw-audio/>

[1]: <https://arxiv.org/abs/1810.04805>

②  $\underbrace{\text{mapped}}_{\text{using a learnable quantization matrix}} \rightarrow \ell \in \mathbb{R}^{G \times V}$

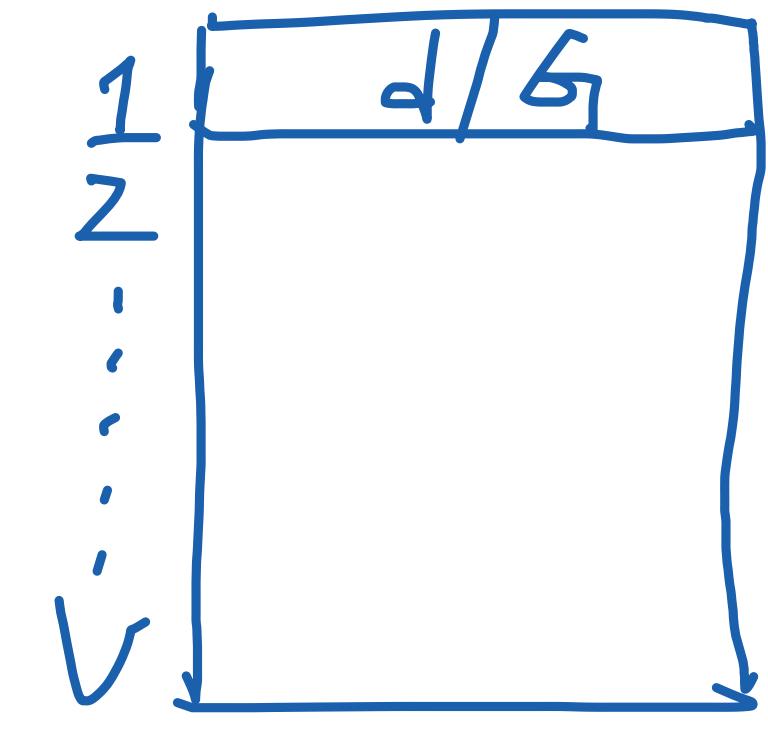


$e_1$



$e_L$

.....



$e_g$

2019

$$P_{g,v} = \frac{\exp(\ell_{g,v} + b_v / \tau)}{\sum_{v=1}^G \exp(\ell_{g,v} + b_v / \tau)}$$

Where  
 $b_v \sim \text{Gumbel}(0, 1)$

$$\tilde{z} = [e_1 | e_2 | \dots | e_g]$$

Gumbel Softmax  
Trick

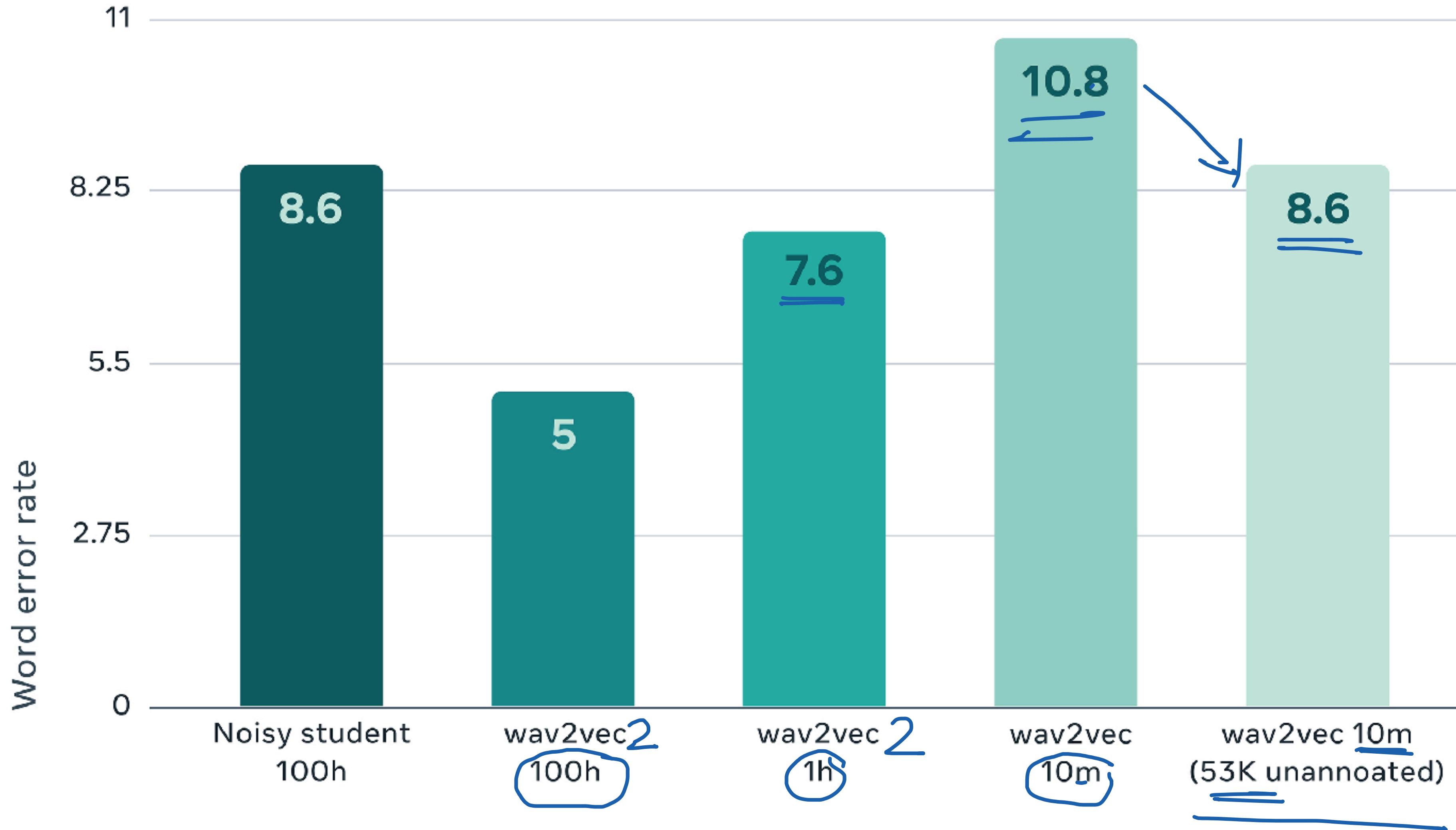
$$\pi_i, \sum_i \pi_i = 1$$

$\Rightarrow$  Sample from categorical distribution  
& make it differentiable

Gumbel trick:  $b = \text{one-hot}(\arg\max_i(g_i + \log \pi_i))$

$$g_i \sim \text{Gumbel}(0, 1) \quad u \sim \text{Unif}(0, 1), g = -\log(-\log u)$$

# wav2vec 2.0: Results on English



# XLSR: Multilingual Self-supervised Speech Model

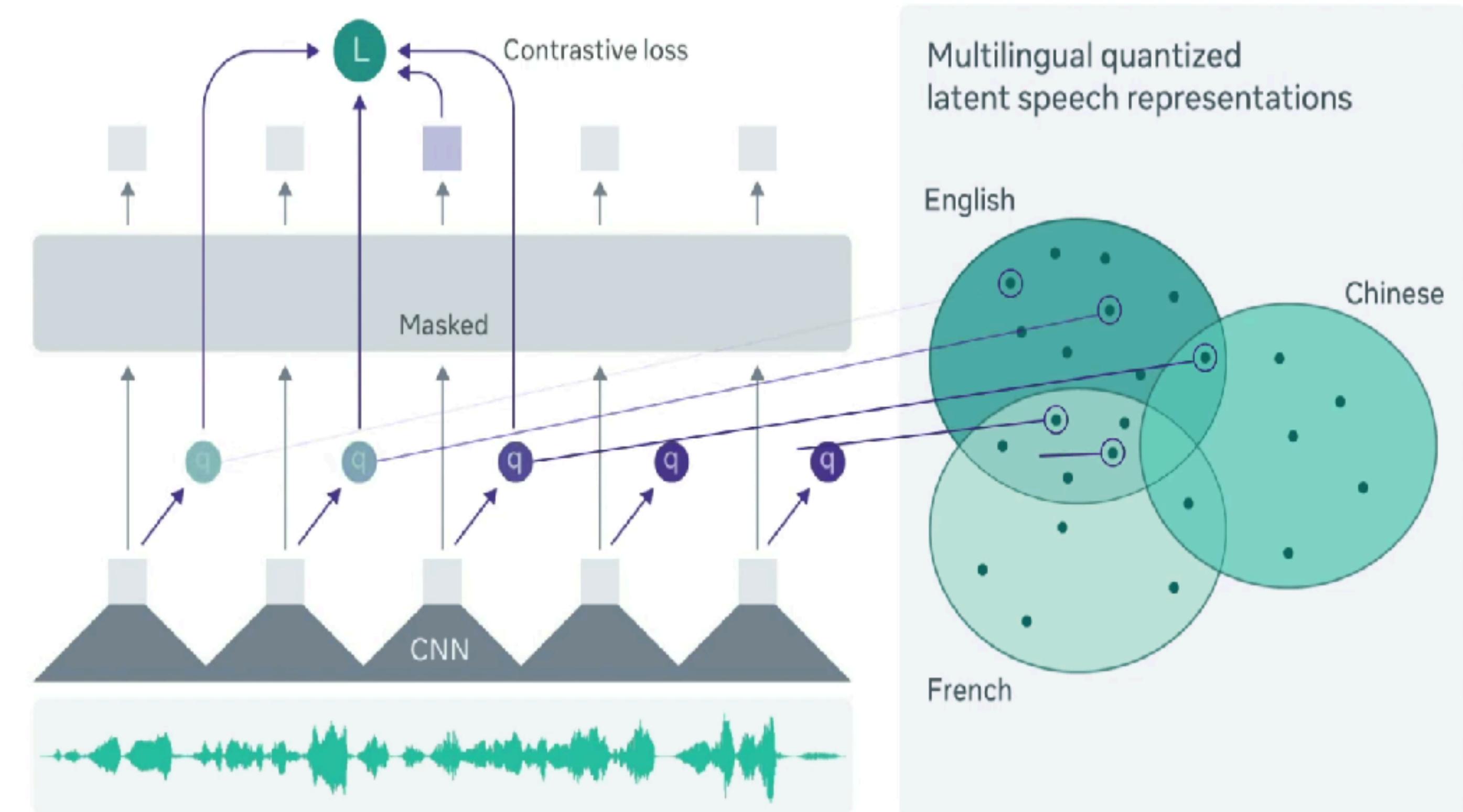
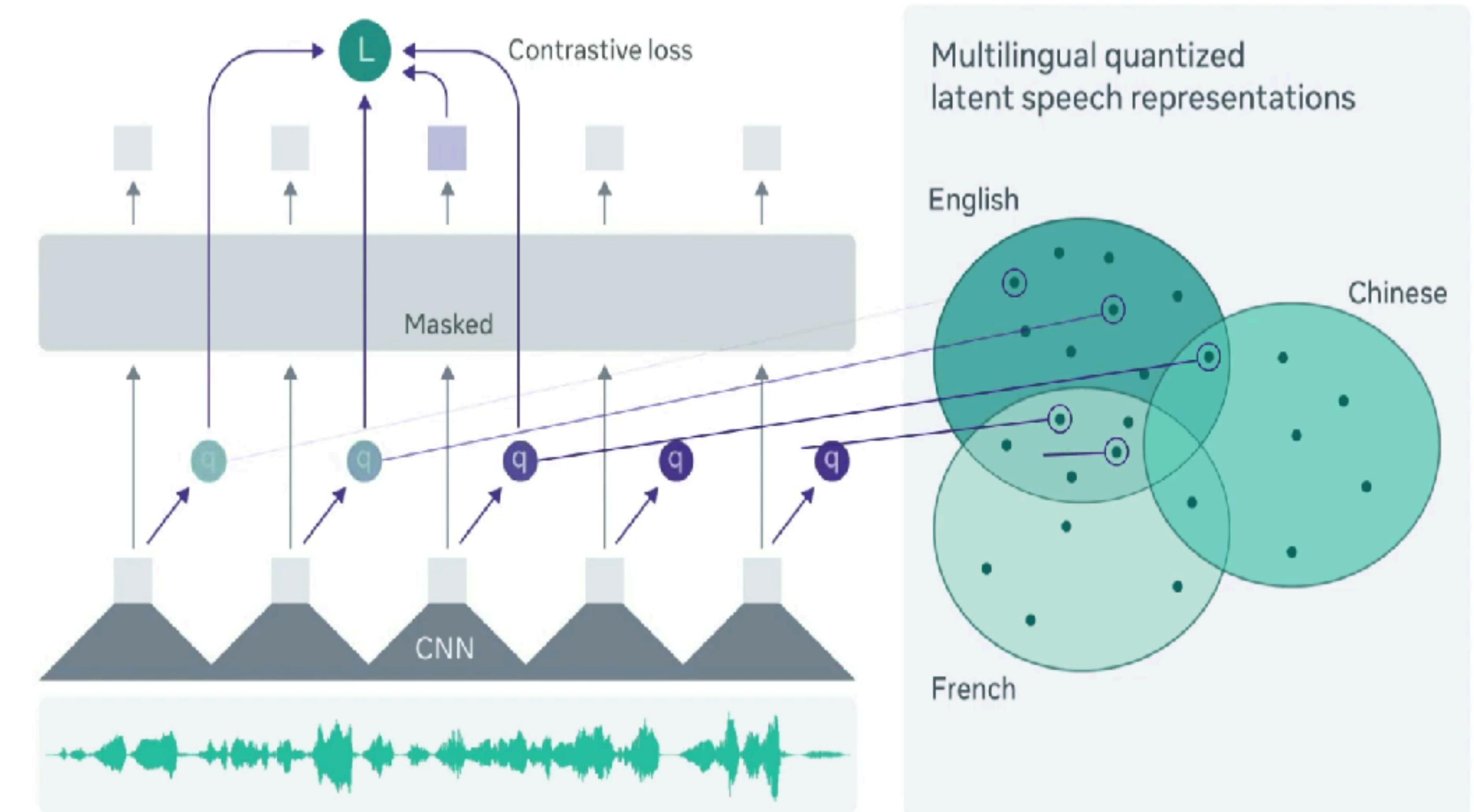


Image from: <https://ai.facebook.com/blog/wav2vec-2-0-learning-the-structure-of-speech-from-raw-audio/>

[1]: <https://arxiv.org/abs/1810.04805>

# XLSR: Multilingual Self-supervised Speech Model

- wav2vec 2.0 model trained on speech in 53 languages



# XLSR: Multilingual Self-supervised Speech Model

- wav2vec 2.0 model trained on speech in 53 languages
- Training objective is to recover the masked representations within a set of distractors

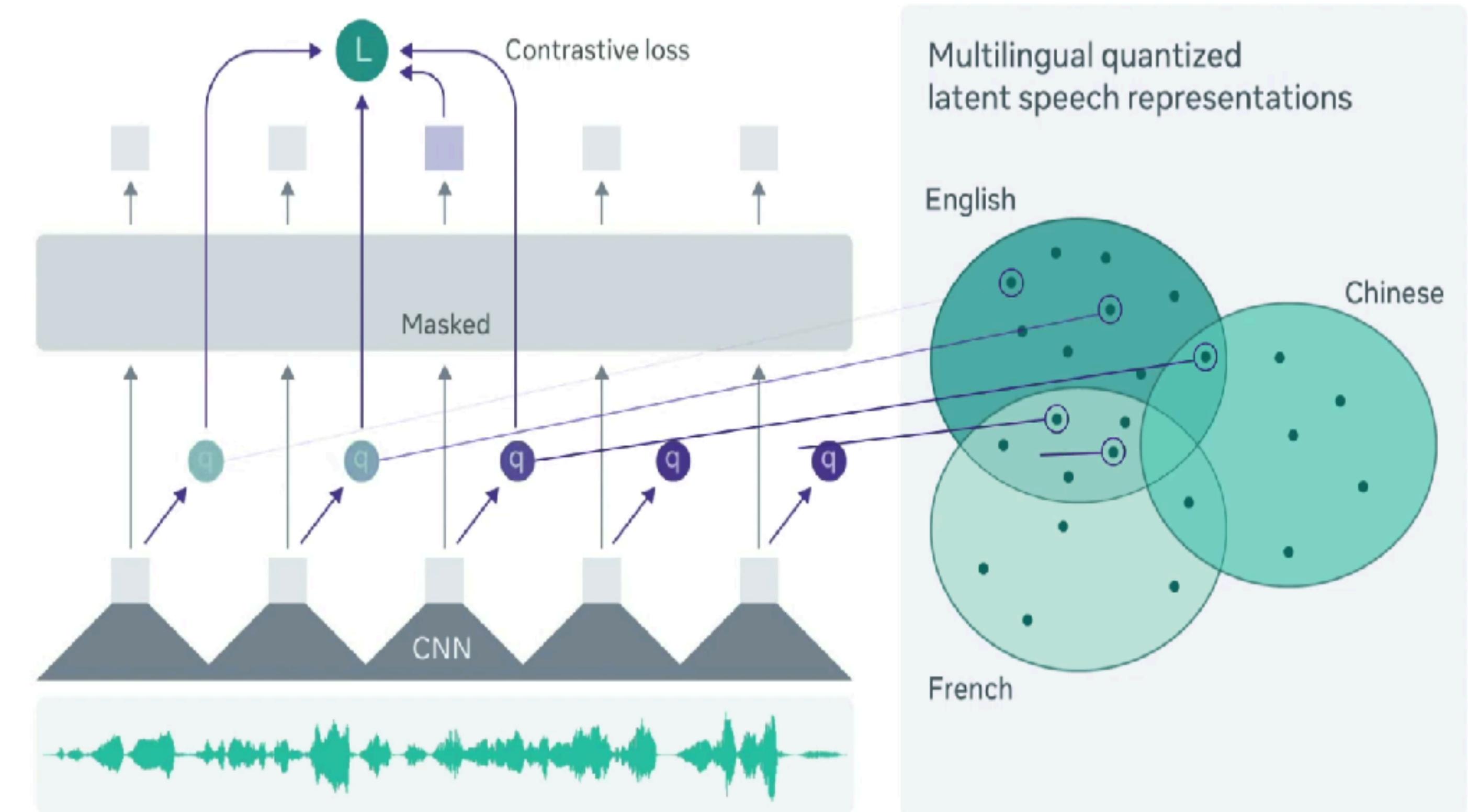


Image from: <https://ai.facebook.com/blog/wav2vec-2-0-learning-the-structure-of-speech-from-raw-audio/>

[1]: <https://arxiv.org/abs/1810.04805>

# XLSR: Multilingual Self-supervised Speech Model

- wav2vec 2.0 model trained on speech in 53 languages
- Training objective is to recover the masked representations within a set of distractors
- Cross-lingual pretraining significantly outperforms monolingual pretraining

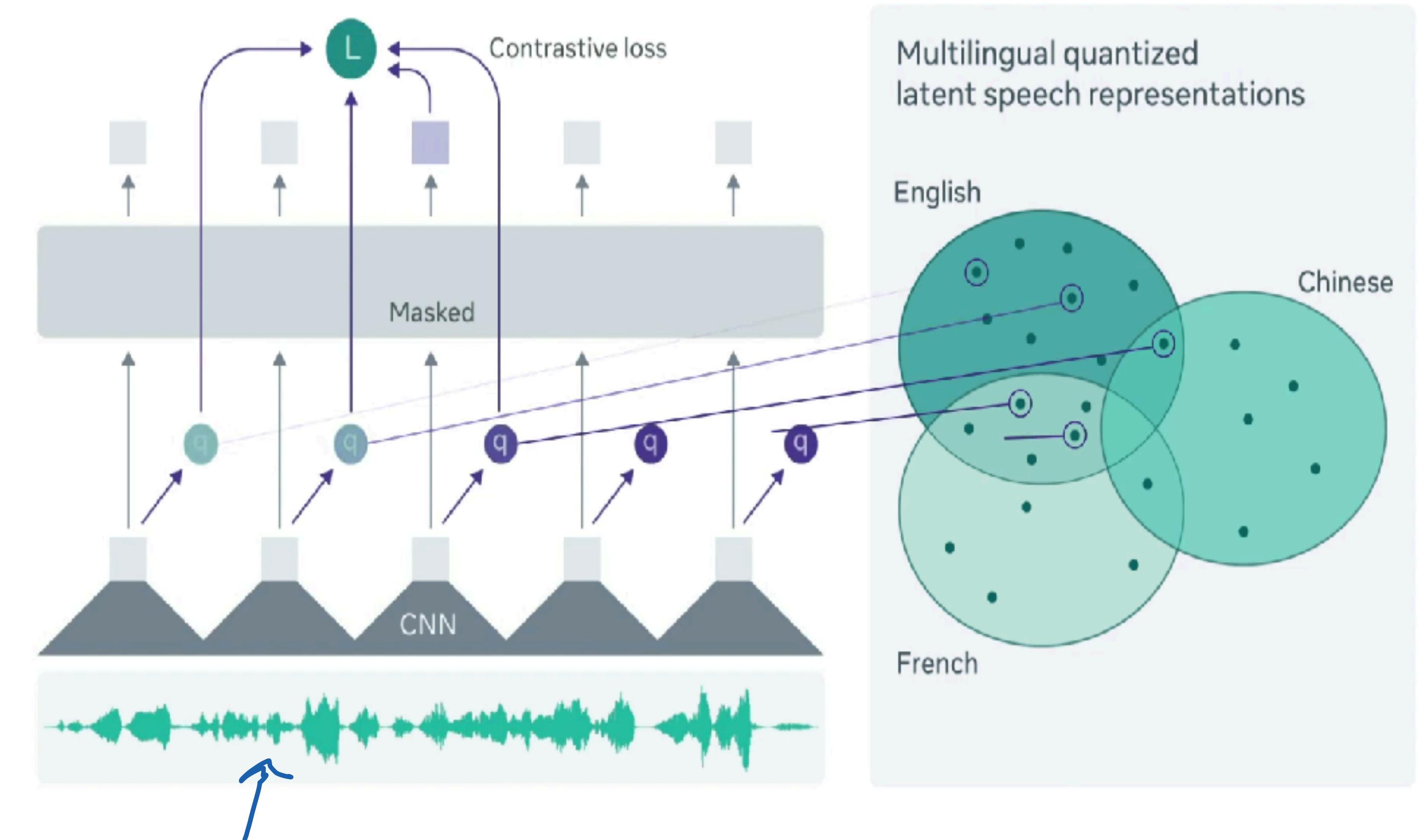


Image from: <https://ai.facebook.com/blog/wav2vec-2-0-learning-the-structure-of-speech-from-raw-audio/>

[1]: <https://arxiv.org/abs/1810.04805>

# XLSR Results

Model	D	#pt	#ft	es	fr	it	ky	nl	ru	sv	tr	tt	zh	Avg
Number of pretraining hours per language				168h	353h	90h	17h	29h	55h	3h	11h	17h	50h	793h
Number of fine-tuning hours per language				1h	10h									
<i>Baselines from previous work</i>														
m-CPC <sup>†</sup> (Rivière et al., 2020)	LS <sub>100h</sub>	10	1	38.7	49.3	42.1	40.7	44.4	45.2	48.8	49.7	44.0	55.5	45.8
m-CPC <sup>†</sup> (Rivière et al., 2020)	LS <sub>360h</sub>	10	1	38.0	47.1	40.5	41.2	42.5	43.7	47.5	47.3	42.0	55.0	44.5
Fer et al. <sup>†</sup> (Fer et al., 2017)	BBL <sub>all</sub>	10	1	36.6	48.3	39.0	38.7	47.9	45.2	52.6	43.4	42.5	54.3	44.9
<i>Our monolingual models</i>														
XLSR-English	CV <sub>en</sub>	1	1	13.7	20.0	19.1	13.2	19.4	18.6	21.1	15.5	11.5	27.1	17.9
XLSR-Monolingual	CV <sub>mo</sub>	1	1	6.8	10.4	10.9	29.6	37.4	11.6	63.6	44.0	21.4	31.4	26.7
<i>Our multilingual models</i>														
XLSR-10 (unbalanced)	CV <sub>all</sub>	10	1	9.7	13.6	15.2	11.1	18.1	13.7	21.4	14.2	9.7	25.8	15.3
XLSR-10	CV <sub>all</sub>	10	1	9.4	14.2	14.1	8.4	16.1	11.0	20.7	11.2	7.6	24.0	13.6
XLSR-10 (separate vocab)	CV <sub>all</sub>	10	10	10.0	13.8	14.0	8.8	16.5	11.6	21.4	12.0	8.7	24.5	14.1
XLSR-10 (shared vocab)	CV <sub>all</sub>	10	10	9.4	13.4	13.8	8.6	16.3	11.2	21.0	11.7	8.3	24.5	13.8
<i>Our multilingual models (Large)</i>														
XLSR-10	CV <sub>all</sub>	10	1	7.9	12.6	11.7	7.0	14.0	9.3	20.6	9.7	7.2	22.8	12.3
XLSR-10 (separate vocab)	CV <sub>all</sub>	10	10	8.1	12.1	11.9	7.1	13.9	9.8	21.0	10.4	7.6	22.3	12.4
XLSR-10 (shared vocab)	CV <sub>all</sub>	10	10	7.7	12.2	11.6	7.0	13.8	9.3	20.8	10.1	7.3	22.3	12.2
<i>Our Large XLSR-53 model pretrained on 56k hours</i>														
XLSR-53	D <sub>53</sub>	53	1	2.9	5.0	5.7	6.1	5.8	8.1	12.2	7.1	5.1	18.3	7.6

# XLSR Results

Model	D	#pt	#ft	es	fr	it	ky	nl	ru	sv	tr	tt	zh	Avg
Number of pretraining hours per language				168h	353h	90h	17h	29h	55h	3h	11h	17h	50h	793h
Number of fine-tuning hours per language				1h	10h									
<i>Baselines from previous work</i>														
m-CPC <sup>†</sup> (Rivière et al., 2020)	LS <sub>100h</sub>	10	1	38.7	49.3	42.1	40.7	44.4	45.2	48.8	49.7	44.0	55.5	45.8
m-CPC <sup>†</sup> (Rivière et al., 2020)	LS <sub>360h</sub>	10	1	38.0	47.1	40.5	41.2	42.5	43.7	47.5	47.3	42.0	55.0	44.5
Fer et al. <sup>†</sup> (Fer et al., 2017)	BBL <sub>all</sub>	10	1	36.6	48.3	39.0	38.7	47.9	45.2	52.6	43.4	42.5	54.3	44.9
<i>Our monolingual models</i>														
XLSR-English	CV <sub>en</sub>	1	1	13.7	20.0	19.1	13.2	19.4	18.6	21.1	15.5	11.5	27.1	17.9
XLSR-Monolingual	CV <sub>mo</sub>	1	1	6.8	10.4	10.9	29.6	37.4	11.6	63.6	44.0	21.4	31.4	26.7
<i>Our multilingual models</i>														
XLSR-10 (unbalanced)	CV <sub>all</sub>	10	1	9.7	13.6	15.2	11.1	18.1	13.7	21.4	14.2	9.7	25.8	15.3
XLSR-10	CV <sub>all</sub>	10	1	9.4	14.2	14.1	8.4	16.1	11.0	20.7	11.2	7.6	24.0	13.6
XLSR-10 (separate vocab)	CV <sub>all</sub>	10	10	10.0	13.8	14.0	8.8	16.5	11.6	21.4	12.0	8.7	24.5	14.1
XLSR-10 (shared vocab)	CV <sub>all</sub>	10	10	9.4	13.4	13.8	8.6	16.3	11.2	21.0	11.7	8.3	24.5	13.8
<i>Our multilingual models (Large)</i>														
XLSR-10	CV <sub>all</sub>	10	1	7.9	12.6	11.7	7.0	14.0	9.3	20.6	9.7	7.2	22.8	12.3
XLSR-10 (separate vocab)	CV <sub>all</sub>	10	10	8.1	12.1	11.9	7.1	13.9	9.8	21.0	10.4	7.6	22.3	12.4
XLSR-10 (shared vocab)	CV <sub>all</sub>	10	10	7.7	12.2	11.6	7.0	13.8	9.3	20.8	10.1	7.3	22.3	12.2
<i>Our Large XLSR-53 model pretrained on 56k hours</i>														
XLSR-53	D <sub>53</sub>	53	1	2.9	5.0	5.7	6.1	5.8	8.1	12.2	7.1	5.1	18.3	7.6

# XLSR Results

Model	D	#pt	#ft	es	fr	it	ky	nl	ru	sv	tr	tt	zh	Avg
Number of pretraining hours per language				168h	353h	90h	17h	29h	55h	3h	11h	17h	50h	793h
Number of fine-tuning hours per language				1h	10h									
<i>Baselines from previous work</i>														
m-CPC <sup>†</sup> (Rivière et al., 2020)	LS <sub>100h</sub>	10	1	38.7	49.3	42.1	40.7	44.4	45.2	48.8	49.7	44.0	55.5	45.8
m-CPC <sup>†</sup> (Rivière et al., 2020)	LS <sub>360h</sub>	10	1	38.0	47.1	40.5	41.2	42.5	43.7	47.5	47.3	42.0	55.0	44.5
Fer et al. <sup>†</sup> (Fer et al., 2017)	BBL <sub>all</sub>	10	1	36.6	48.3	39.0	38.7	47.9	45.2	52.6	43.4	42.5	54.3	44.9
<i>Our monolingual models</i>														
XLSR-English	CV <sub>en</sub>	1	1	13.7	20.0	19.1	13.2	19.4	18.6	21.1	15.5	11.5	27.1	17.9
XLSR-Monolingual	CV <sub>mo</sub>	1	1	6.8	10.4	10.9	29.6	37.4	11.6	63.6	44.0	21.4	31.4	26.7
<i>Our multilingual models</i>														
XLSR-10 (unbalanced)	CV <sub>all</sub>	10	1	9.7	13.6	15.2	11.1	18.1	13.7	21.4	14.2	9.7	25.8	15.3
XLSR-10	CV <sub>all</sub>	10	1	9.4	14.2	14.1	8.4	16.1	11.0	20.7	11.2	7.6	24.0	13.6
XLSR-10 (separate vocab)	CV <sub>all</sub>	10	10	10.0	13.8	14.0	8.8	16.5	11.6	21.4	12.0	8.7	24.5	14.1
XLSR-10 (shared vocab)	CV <sub>all</sub>	10	10	9.4	13.4	13.8	8.6	16.3	11.2	21.0	11.7	8.3	24.5	13.8
<i>Our multilingual models (Large)</i>														
XLSR-10	CV <sub>all</sub>	10	1	7.9	12.6	11.7	7.0	14.0	9.3	20.6	9.7	7.2	22.8	12.3
XLSR-10 (separate vocab)	CV <sub>all</sub>	10	10	8.1	12.1	11.9	7.1	13.9	9.8	21.0	10.4	7.6	22.3	12.4
XLSR-10 (shared vocab)	CV <sub>all</sub>	10	10	7.7	12.2	11.6	7.0	13.8	9.3	20.8	10.1	7.3	22.3	12.2
<i>Our Large XLSR-53 model pretrained on 56k hours</i>														
XLSR-53	D <sub>53</sub>	53	1	2.9	5.0	5.7	6.1	5.8	8.1	12.2	7.1	5.1	18.3	7.6

# XLSR Results

Model	D	#pt	#ft	es	fr	it	ky	nl	ru	sv	tr	tt	zh	Avg
Number of pretraining hours per language				168h	353h	90h	17h	29h	55h	3h	11h	17h	50h	793h
Number of fine-tuning hours per language				1h	10h									
<i>Baselines from previous work</i>														
m-CPC <sup>†</sup> (Rivière et al., 2020)	LS <sub>100h</sub>	10	1	38.7	49.3	42.1	40.7	44.4	45.2	48.8	49.7	44.0	55.5	45.8
m-CPC <sup>†</sup> (Rivière et al., 2020)	LS <sub>360h</sub>	10	1	38.0	47.1	40.5	41.2	42.5	43.7	47.5	47.3	42.0	55.0	44.5
Fer et al. <sup>†</sup> (Fer et al., 2017)	BBL <sub>all</sub>	10	1	36.6	48.3	39.0	38.7	47.9	45.2	52.6	43.4	42.5	54.3	44.9
<i>Our monolingual models</i>														
XLSR-English	CV <sub>en</sub>	1	1	13.7	20.0	19.1	13.2	19.4	18.6	21.1	15.5	11.5	27.1	17.9
XLSR-Monolingual	CV <sub>mo</sub>	1	1	6.8	10.4	10.9	29.6	37.4	11.6	63.6	44.0	21.4	31.4	26.7
<i>Our multilingual models</i>														
XLSR-10 (unbalanced)	CV <sub>all</sub>	10	1	9.7	13.6	15.2	11.1	18.1	13.7	21.4	14.2	9.7	25.8	15.3
XLSR-10	CV <sub>all</sub>	10	1	9.4	14.2	14.1	8.4	16.1	11.0	20.7	11.2	7.6	24.0	13.6
XLSR-10 (separate vocab)	CV <sub>all</sub>	10	10	10.0	13.8	14.0	8.8	16.5	11.6	21.4	12.0	8.7	24.5	14.1
XLSR-10 (shared vocab)	CV <sub>all</sub>	10	10	9.4	13.4	13.8	8.6	16.3	11.2	21.0	11.7	8.3	24.5	13.8
<i>Our multilingual models (Large)</i>														
XLSR-10	CV <sub>all</sub>	10	1	7.9	12.6	11.7	7.0	14.0	9.3	20.6	9.7	7.2	22.8	12.3
XLSR-10 (separate vocab)	CV <sub>all</sub>	10	10	8.1	12.1	11.9	7.1	13.9	9.8	21.0	10.4	7.6	22.3	12.4
XLSR-10 (shared vocab)	CV <sub>all</sub>	10	10	7.7	12.2	11.6	7.0	13.8	9.3	20.8	10.1	7.3	22.3	12.2
<i>Our Large XLSR-53 model pretrained on 56k hours</i>														
XLSR-53	D <sub>53</sub>	53	1	2.9	5.0	5.7	6.1	5.8	8.1	12.2	7.1	5.1	18.3	7.6

# HuggingFace wav2vec Models



Hugging Face

Search models, datasets, users...

Models

Datasets

Spaces

Docs

Solutions

Pricing

Log In

Sign Up

< Back to tag list

Tasks [Clear](#)

Search tags

Natural Language Processing

Fill-Mask Question Answering

Summarization Table Question Answering

Text Classification Text Generation

Text2Text Generation Token Classification

Translation Zero-Shot Classification

Sentence Similarity Conversational

Feature Extraction

Audio

Text-to-Speech

Automatic Speech Recognition

Audio-to-Audio Audio Classification

Voice Activity Detection

Computer Vision

Image Classification Object Detection

Image Segmentation Text-to-Image

Models 1,898

Search Models

↑↓ Sort: Most Downloads

facebook/wav2vec2-base-960h

Automatic Speech Recognition • Updated 7 days ago • ↓ 522k • ❤ 30

facebook/hubert-large-ls960-ft

Automatic Speech Recognition • Updated 7 days ago • ↓ 67k • ❤ 11

facebook/wav2vec2-large-960h-lv60-self

Automatic Speech Recognition • Updated 7 days ago • ↓ 50.9k • ❤ 11

facebook/wav2vec2-large-960h

Automatic Speech Recognition • Updated 7 days ago • ↓ 32k • ❤ 3

vumichien/wav2vec2-large-xlsr-japanese-hiragana

Automatic Speech Recognition • Updated Jun 18, 2021 • ↓ 31k • ❤ 2

hf-internal-testing/processor\_with\_lm

Automatic Speech Recognition • Updated Jan 18 • ↓ 21.4k

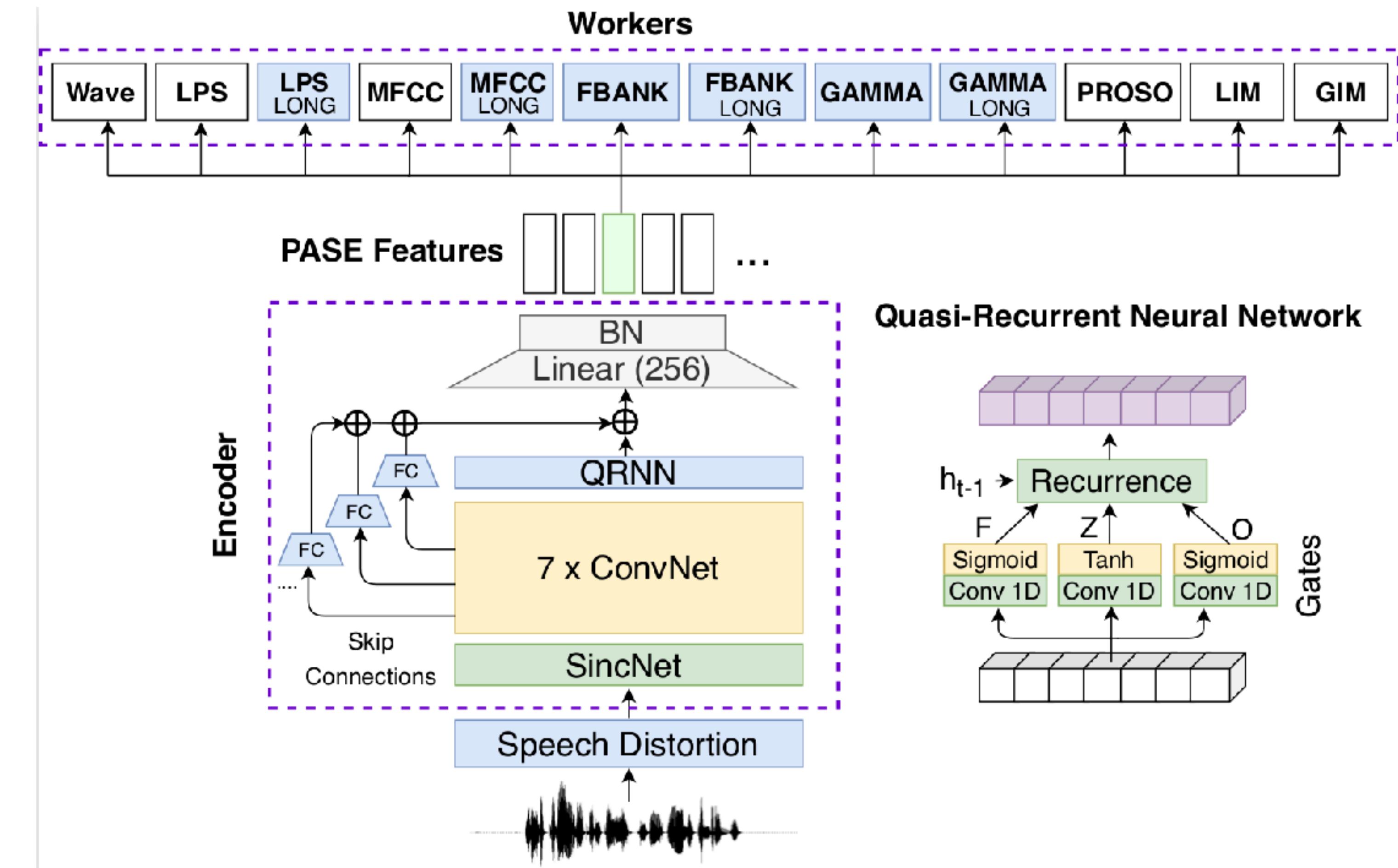
jonatasgrosman/wav2vec2-large-xlsr-53-english

Automatic Speech Recognition • Updated 20 days ago • ↓ 18.7k • ❤ 8

theainerd/wav2vec2-large-xlsr-53-odia

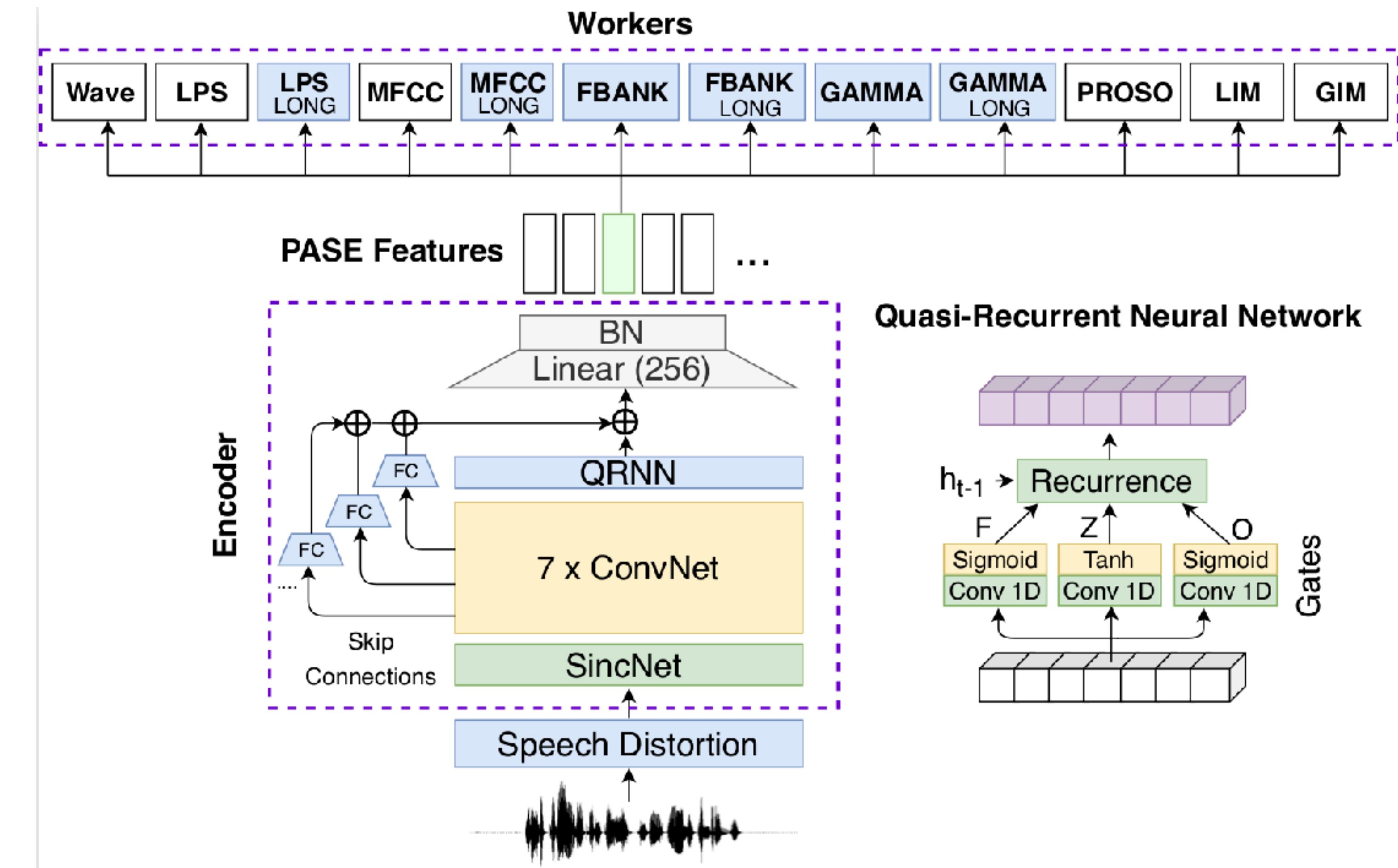
Automatic Speech Recognition • Updated Mar 24, 2021 • ↓ 13.1k

# Multi-task: PASE/PASE+



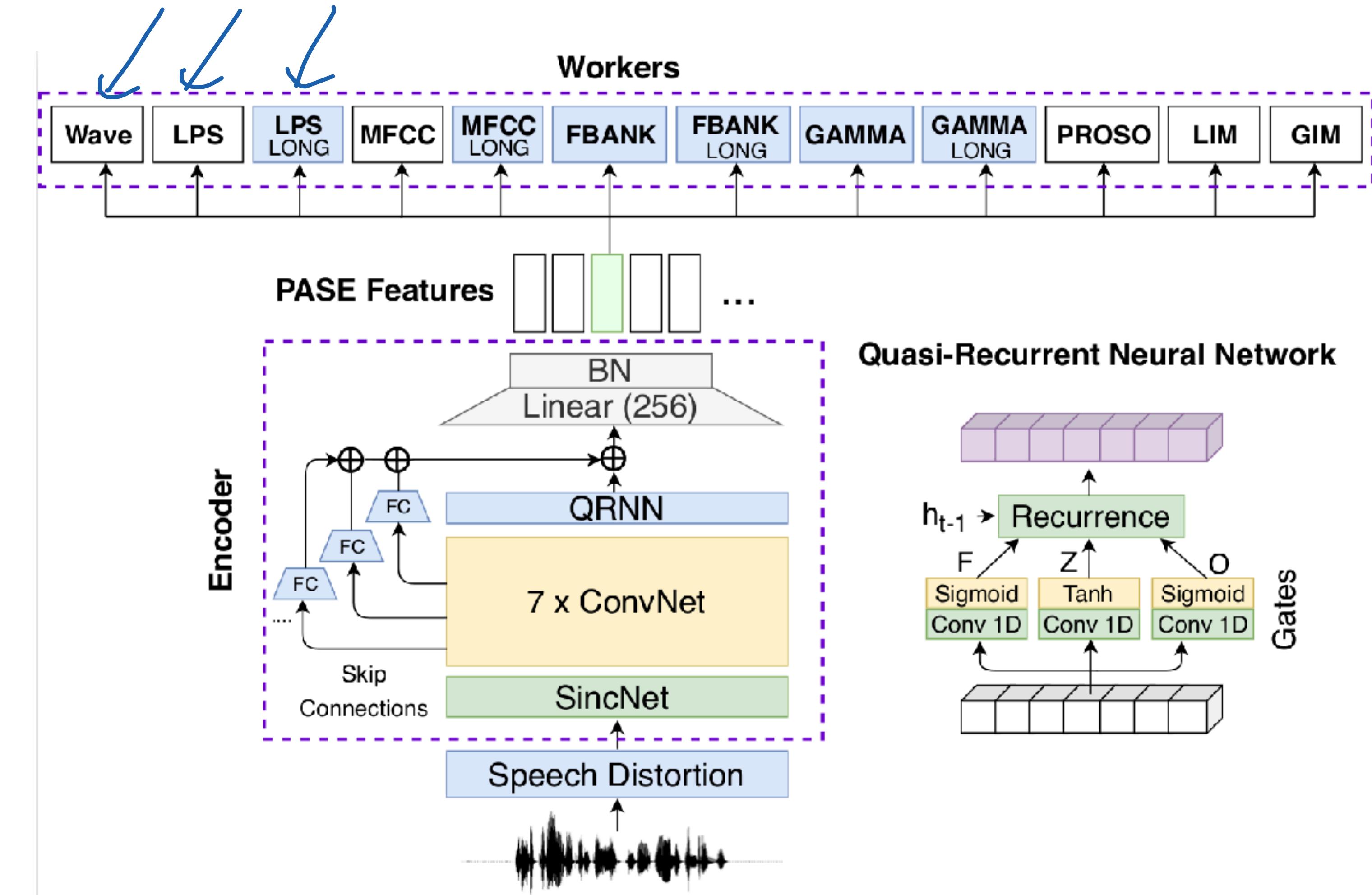
# Multi-task: PASE/PASE+

- PASE consists of a speech encoder, followed by an army of workers with self-supervised objectives trained in a multi-task setup.

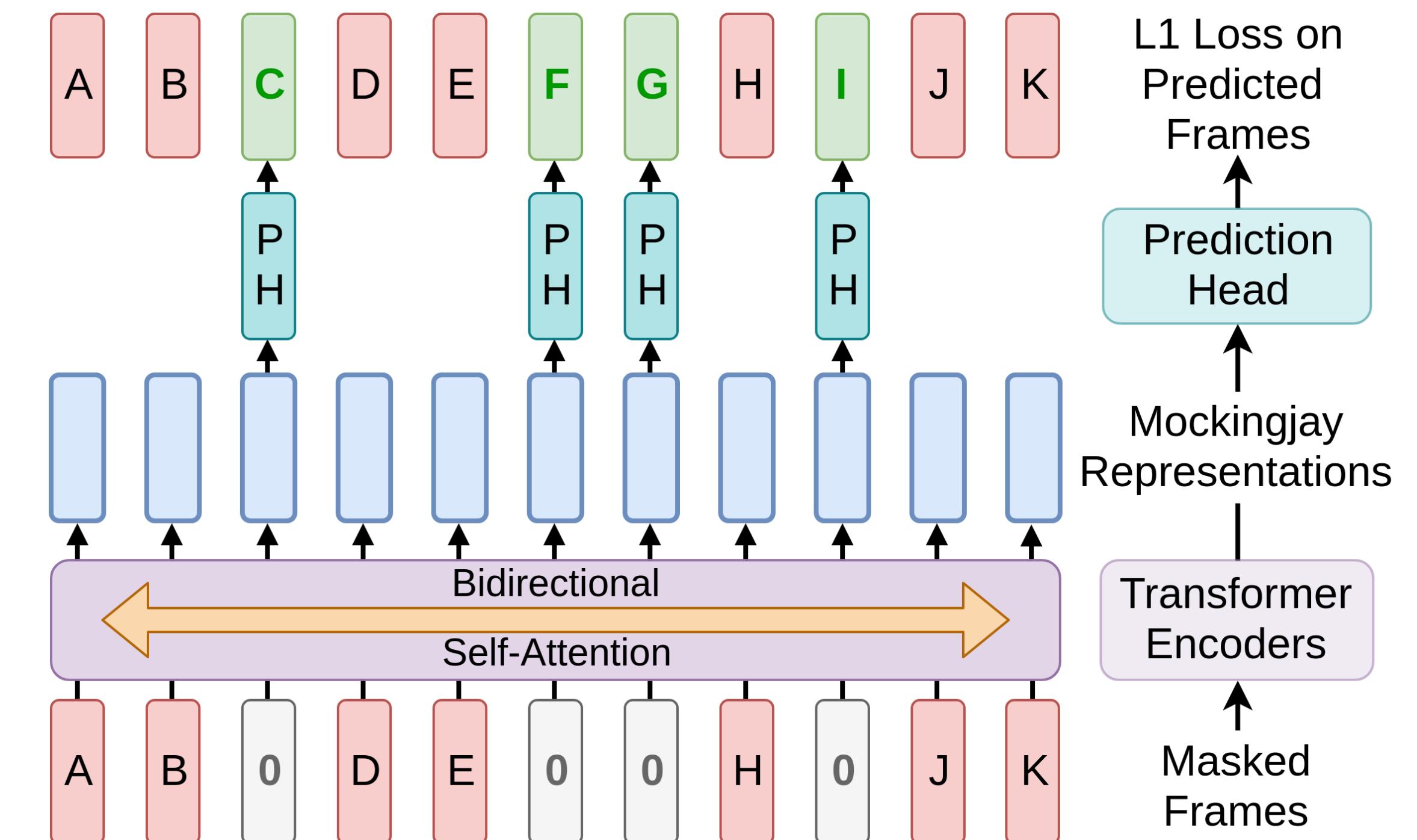


# Multi-task: PASE/PASE+

- PASE consists of a speech encoder, followed by an army of workers with self-supervised objectives trained in a multi-task setup.
- Self-supervised tasks are either regression tasks (MFCC, LPS, etc.) or classification tasks (LIM, GIM, etc.).

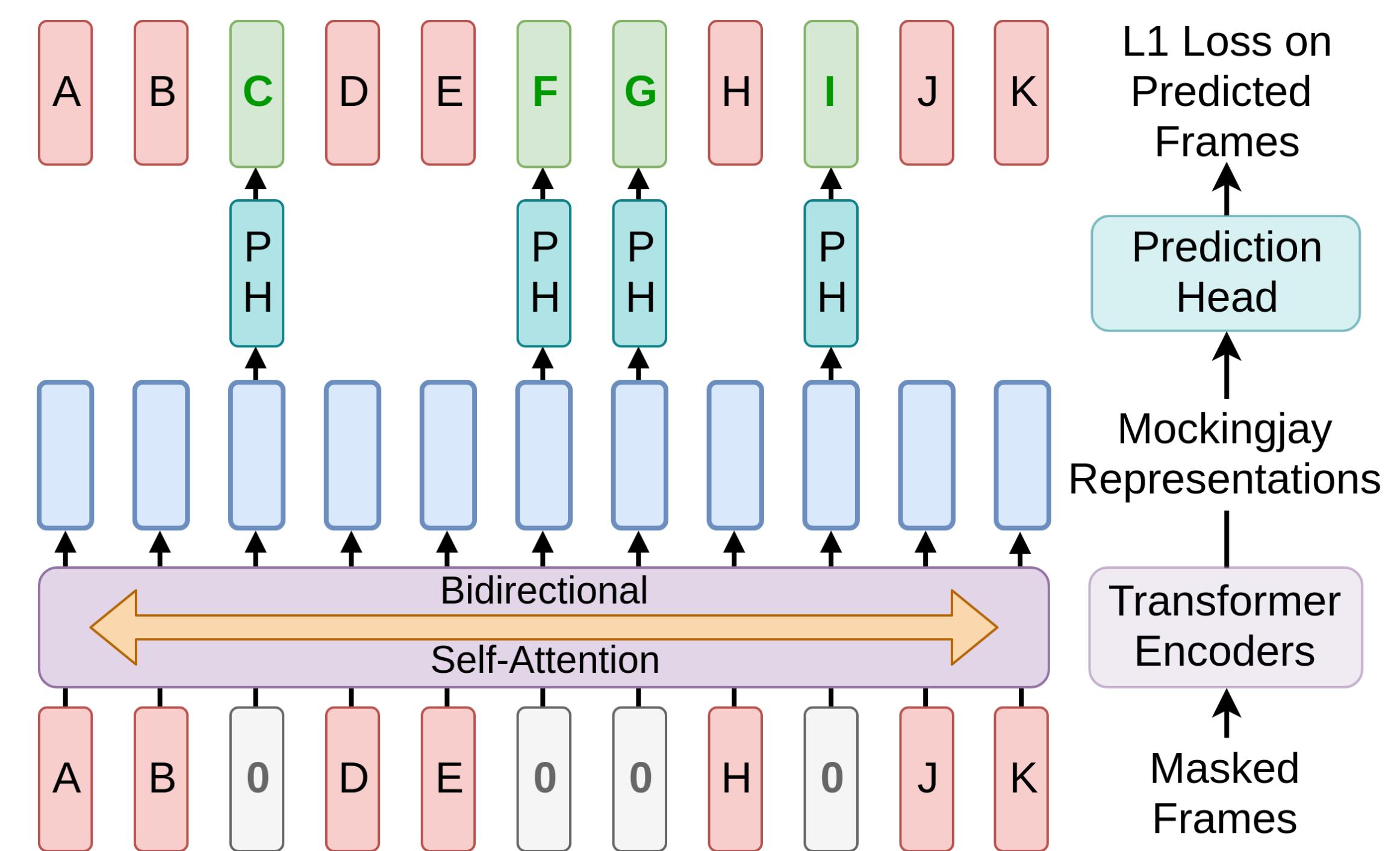


# Masked Reconstruction: Mockingjay



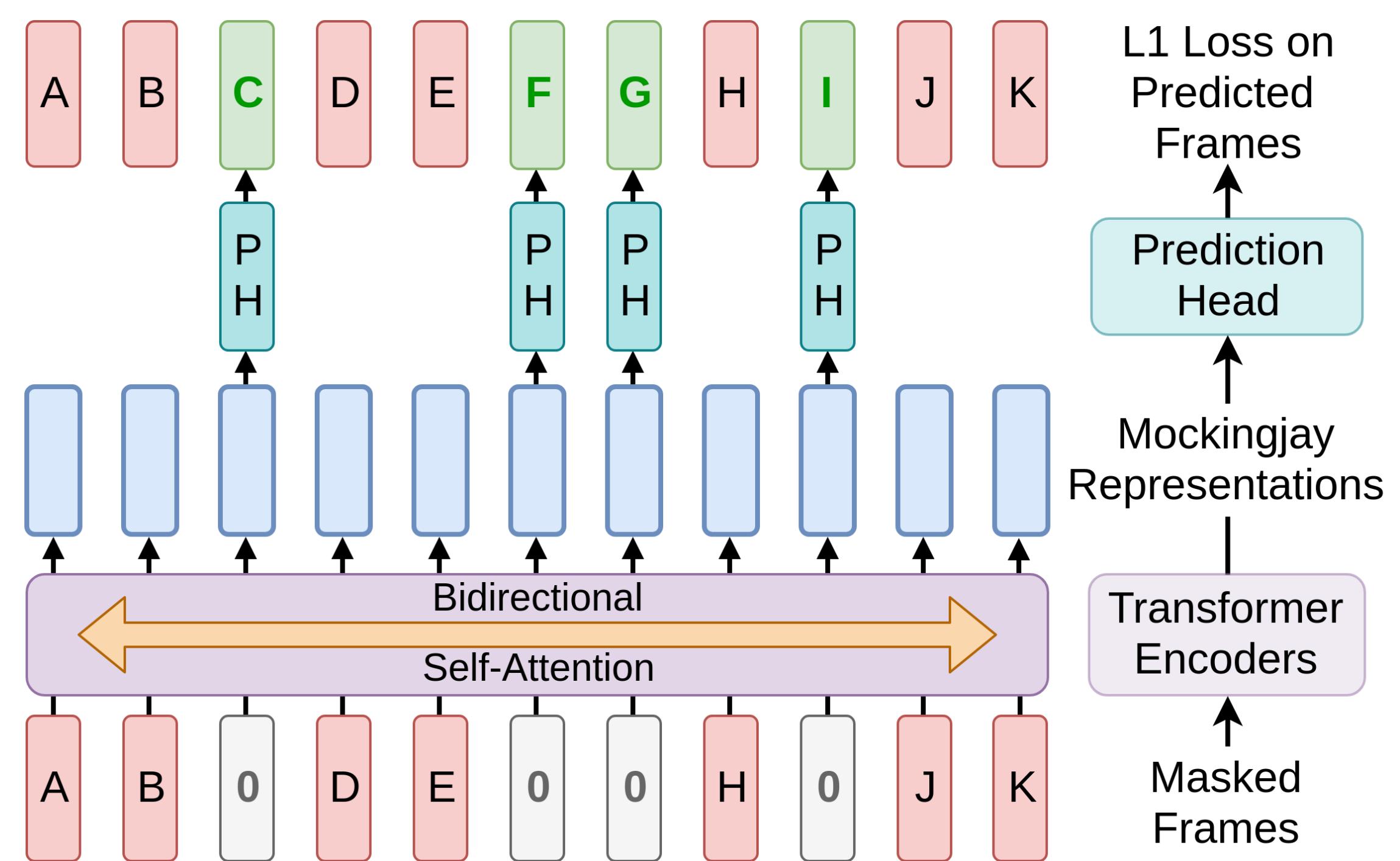
# Masked Reconstruction: Mockingjay

- Predict the current frame by jointly conditioning on both past and future contexts



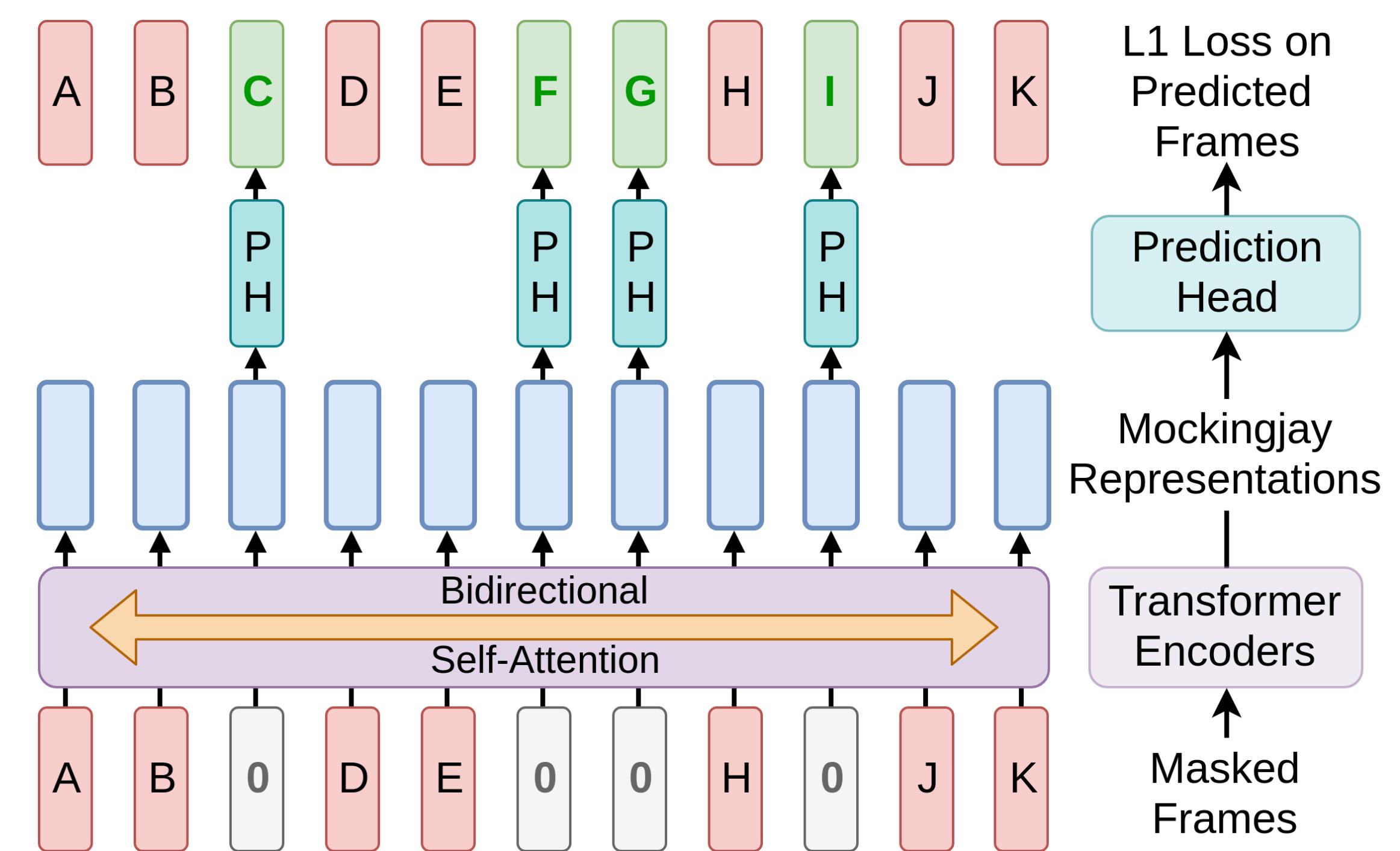
# Masked Reconstruction: Mockingjay

- Predict the current frame by jointly conditioning on both past and future contexts
- Minimize reconstruction error between predictions and ground-truth on the masked frames

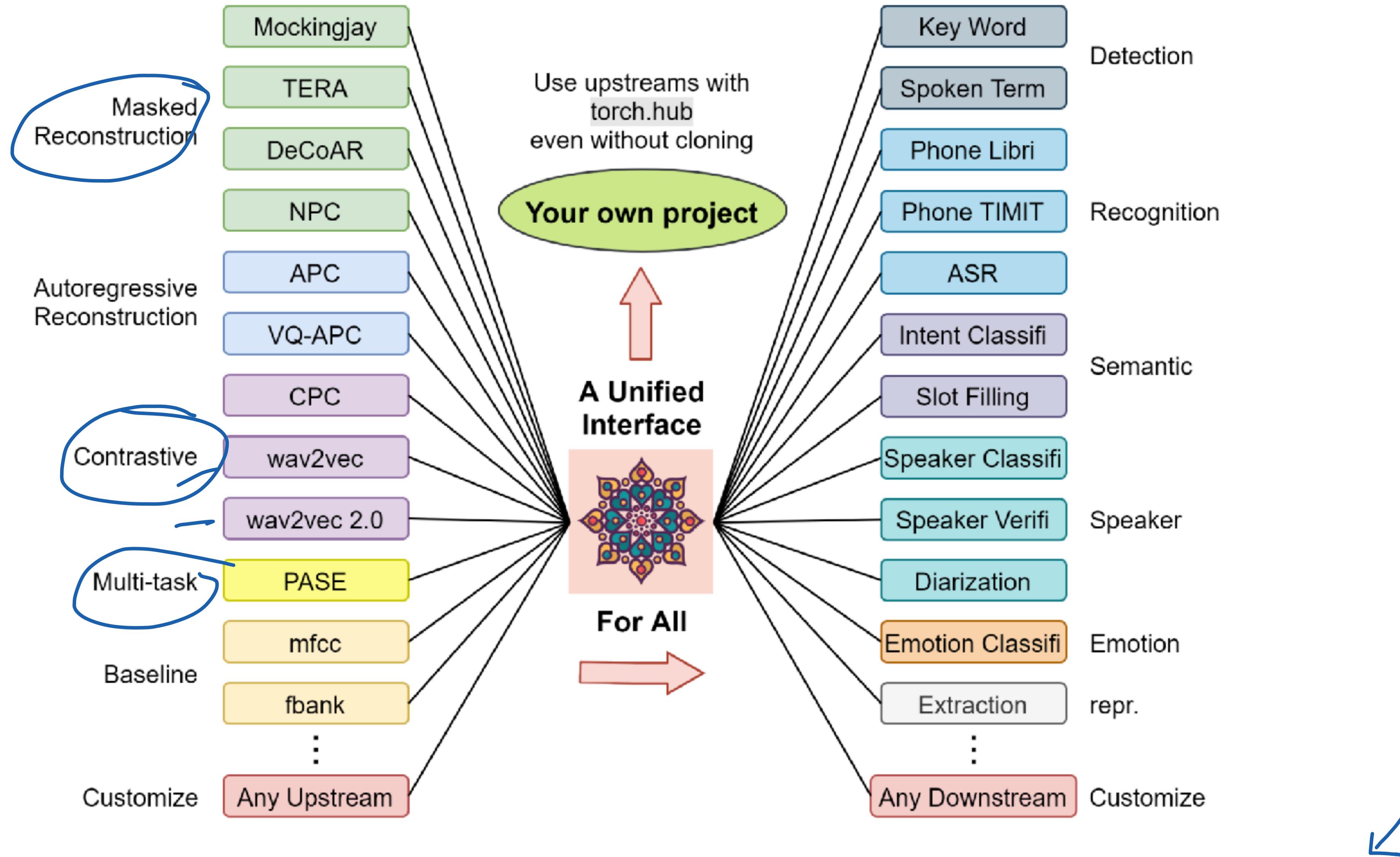


# Masked Reconstruction: Mockingjay

- Predict the current frame by jointly conditioning on both past and future contexts
- Minimize reconstruction error between predictions and ground-truth on the masked frames
- Final representations are the Transformer encoder's hidden states

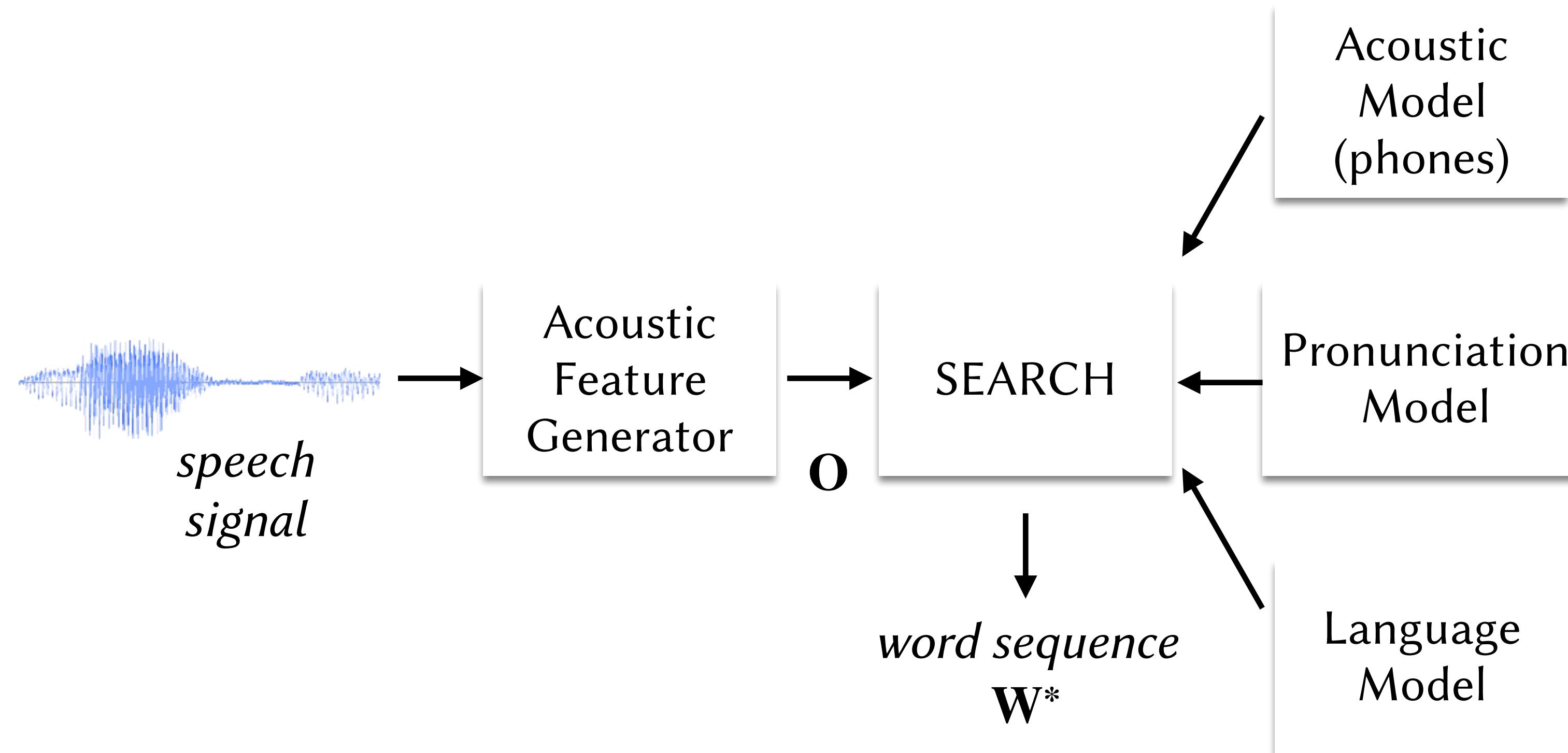


# S3PRL Toolkit

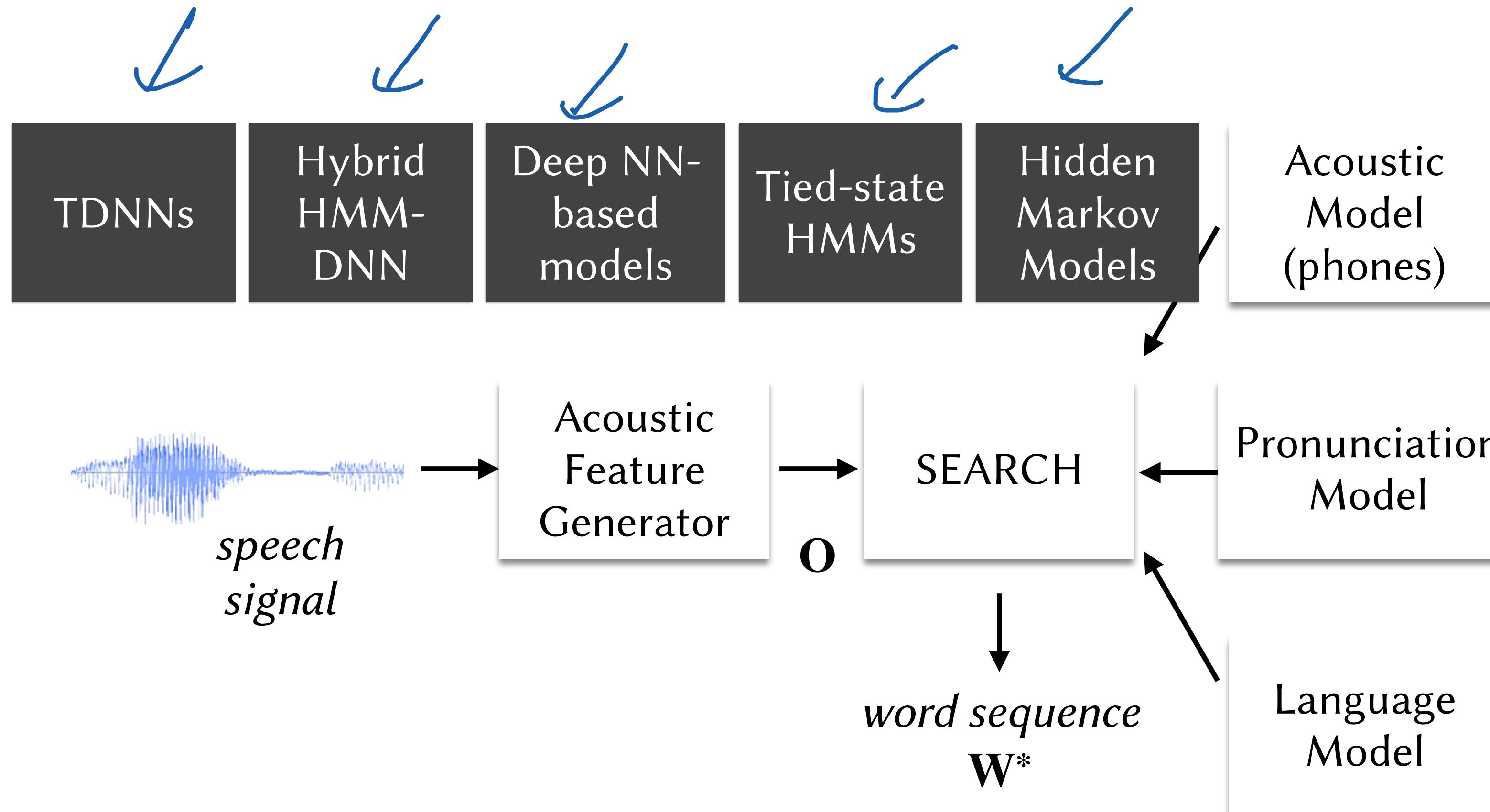


# **CS753 Concluding Remarks**

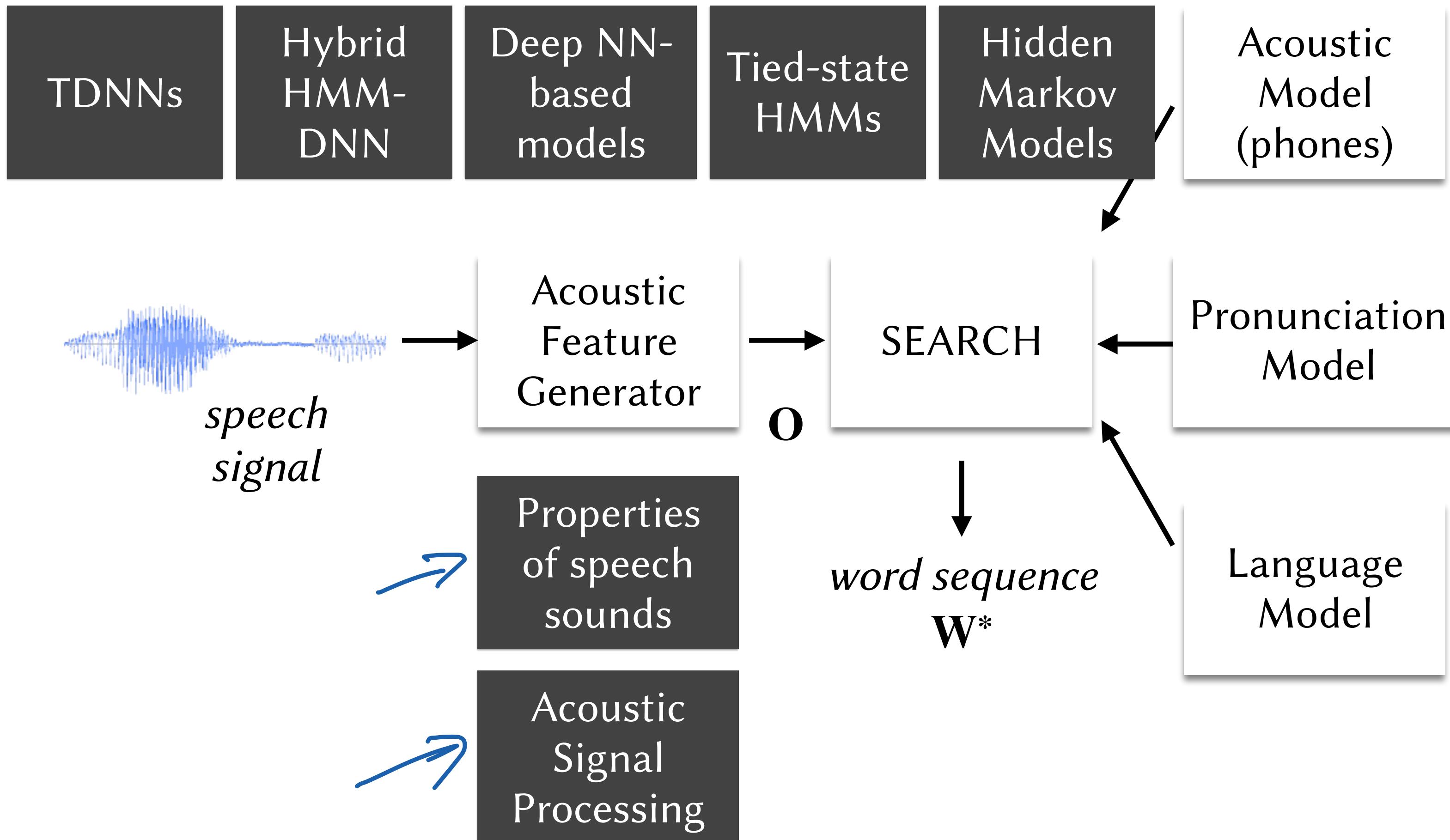
# Topics covered



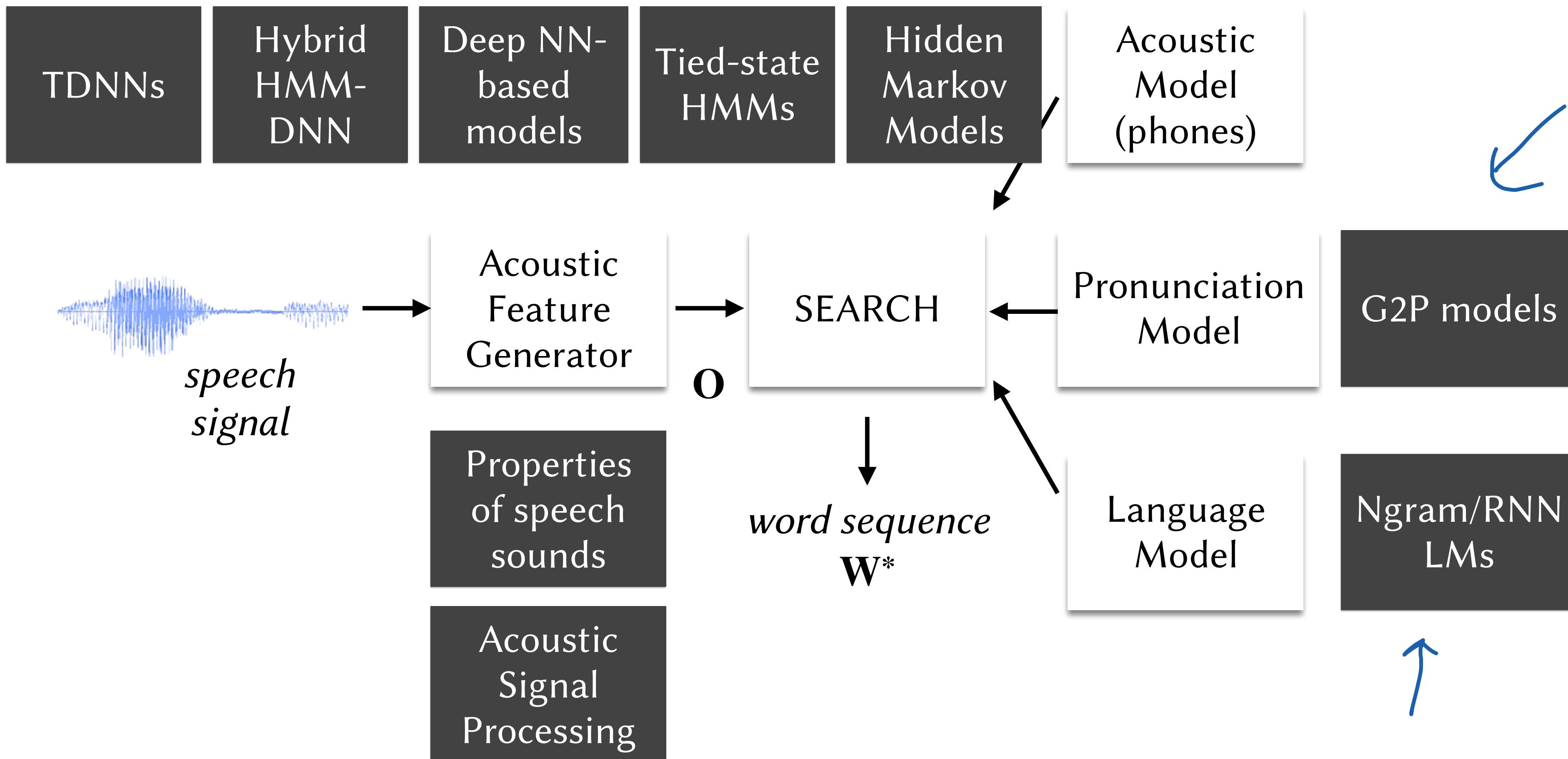
# Topics covered



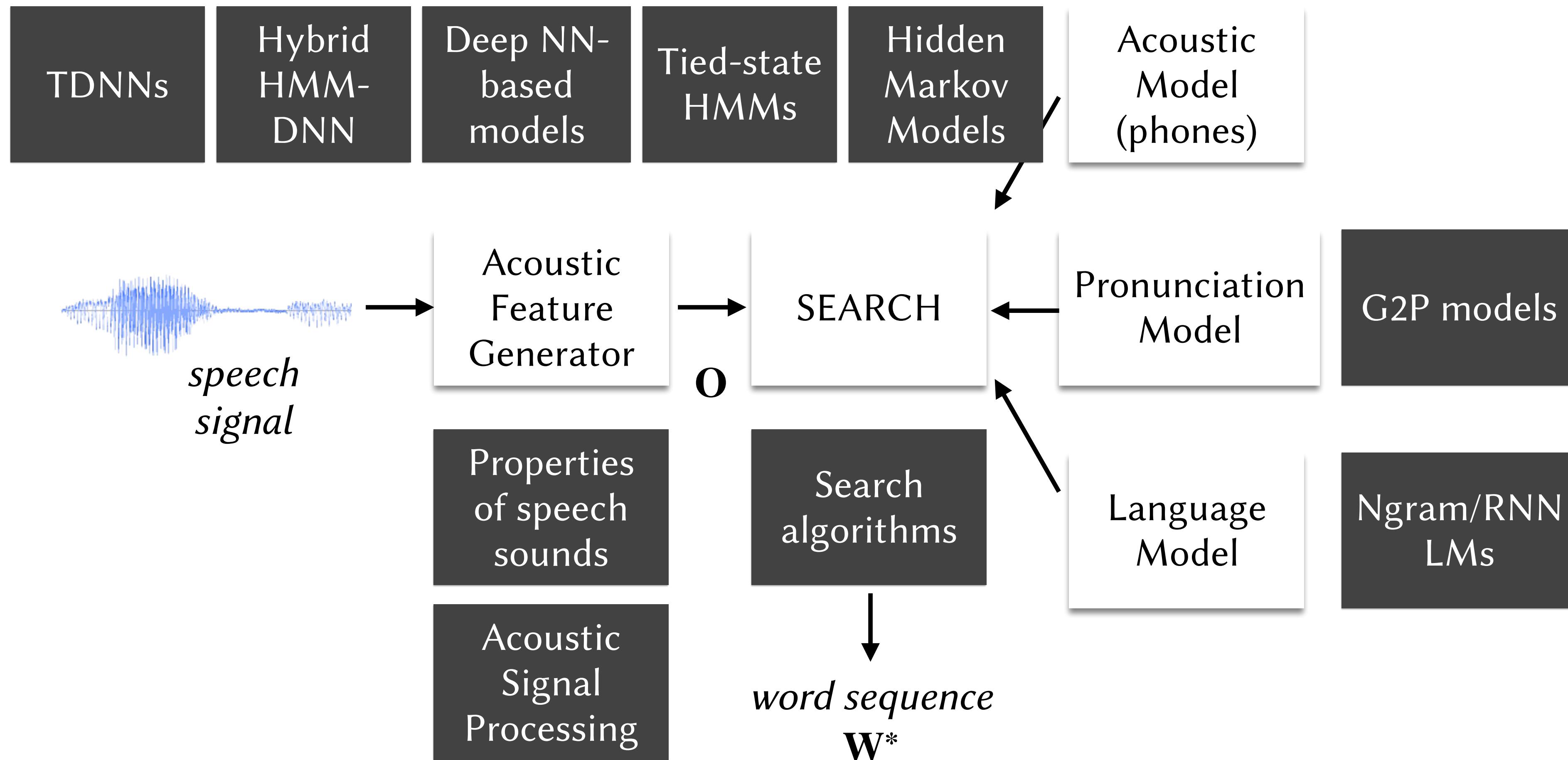
# Topics covered



# Topics covered



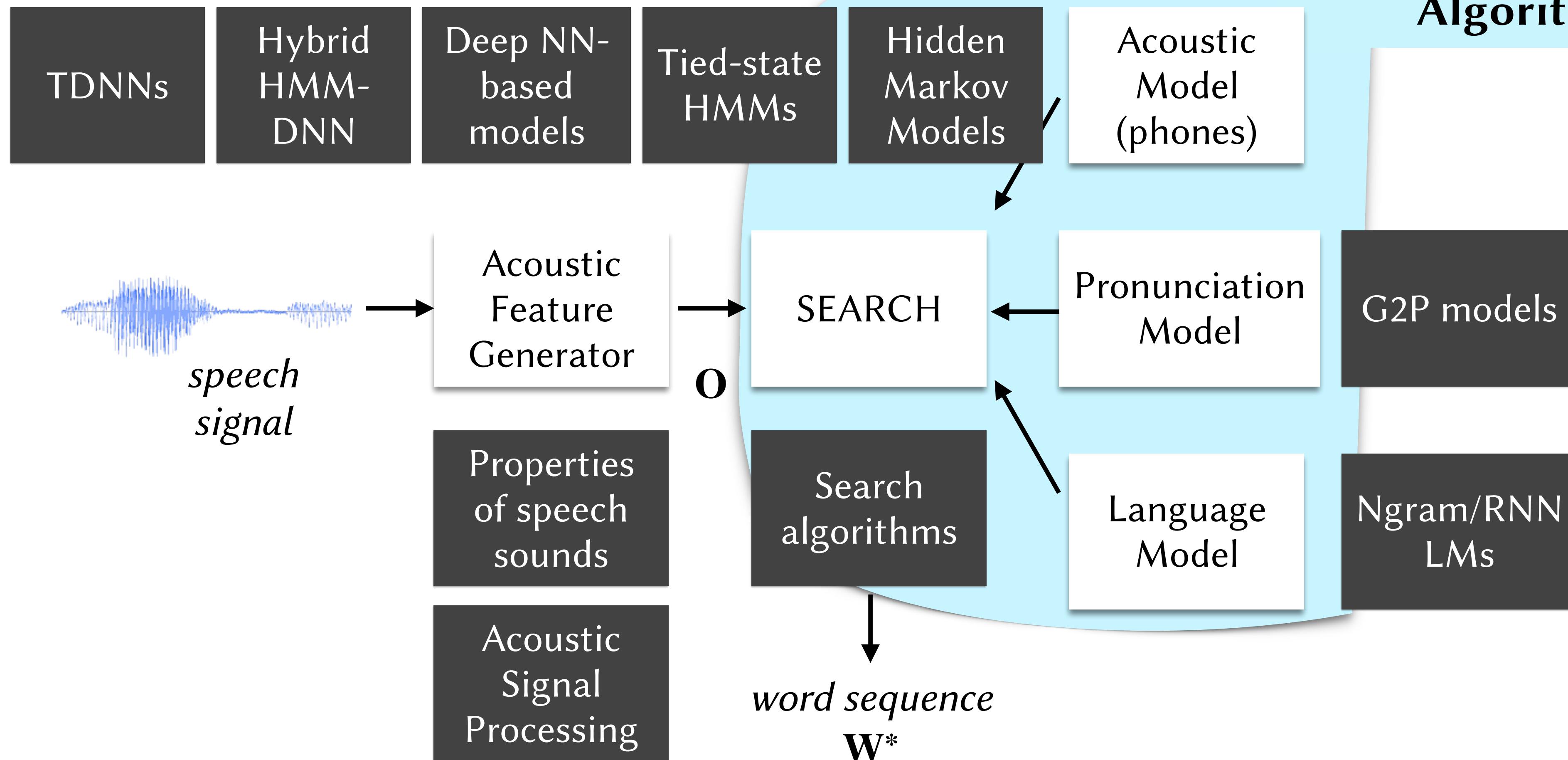
# Topics covered



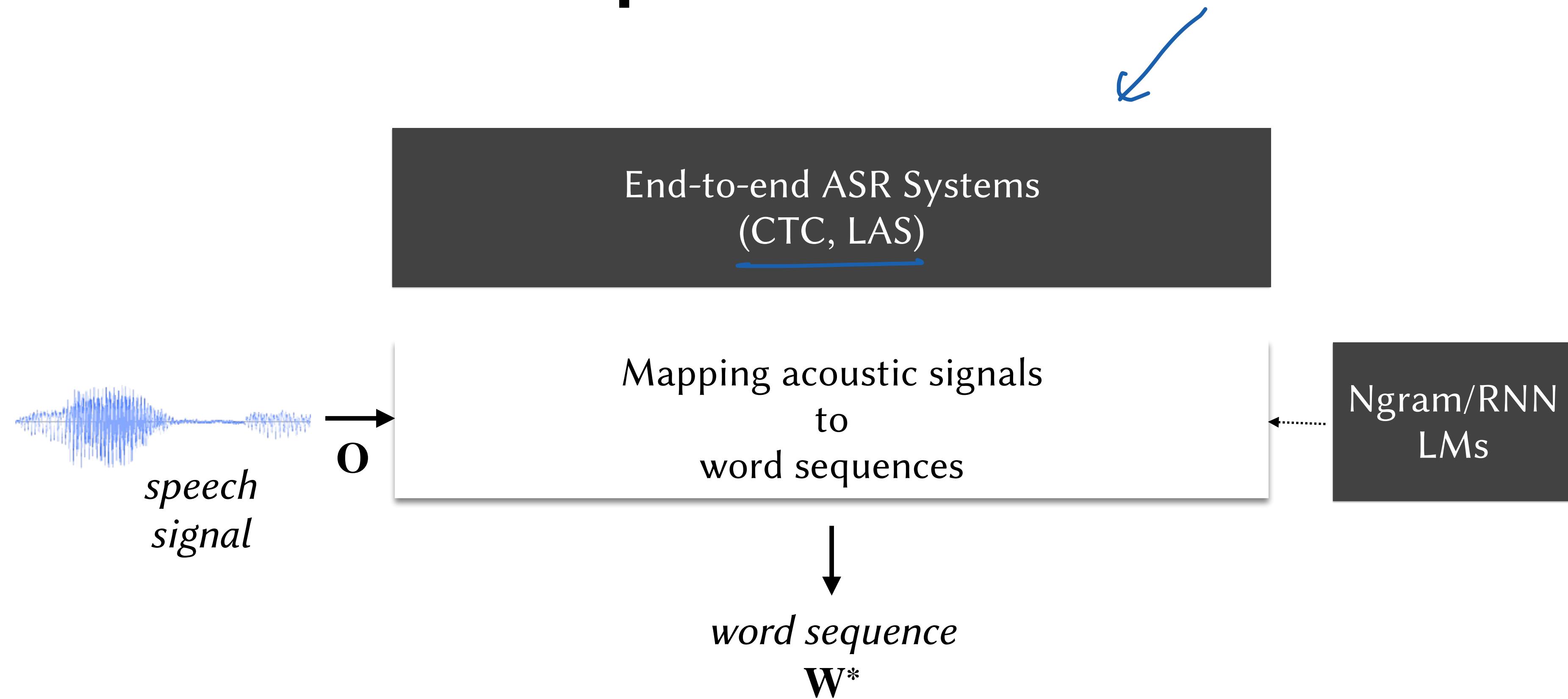
# Topics covered



Formalism: Finite State  
Transducers (WFST  
Algorithms)



# Topics covered



# Seminar Topics

EMFORMER

Blocksync Beam

Voice Conv

Simul Speech

Stream Trans ASR 1,2

Auto VC

AV Dereverb

Relthink ASR

MaskCTC

Audio Adversary

Quaternion

VoiceSep

Rhythm Transfer

wav2vec U

Right2Talk

NeuralSynth

BERT ASR

AV Robustness

SpeakQL

SoundSep Video

Diarization

Visual Voice

Large Margin SV

Efficient Vocalize

Multiling ASR

Style Tokens

ImpercepAdv ASR

Look2Speech

CluedAVVP

GE2E

TriBERT

Speech2SQL

LatticeInputs

Spoken Translation

Move2Hear

RNNT Decoding

AVID

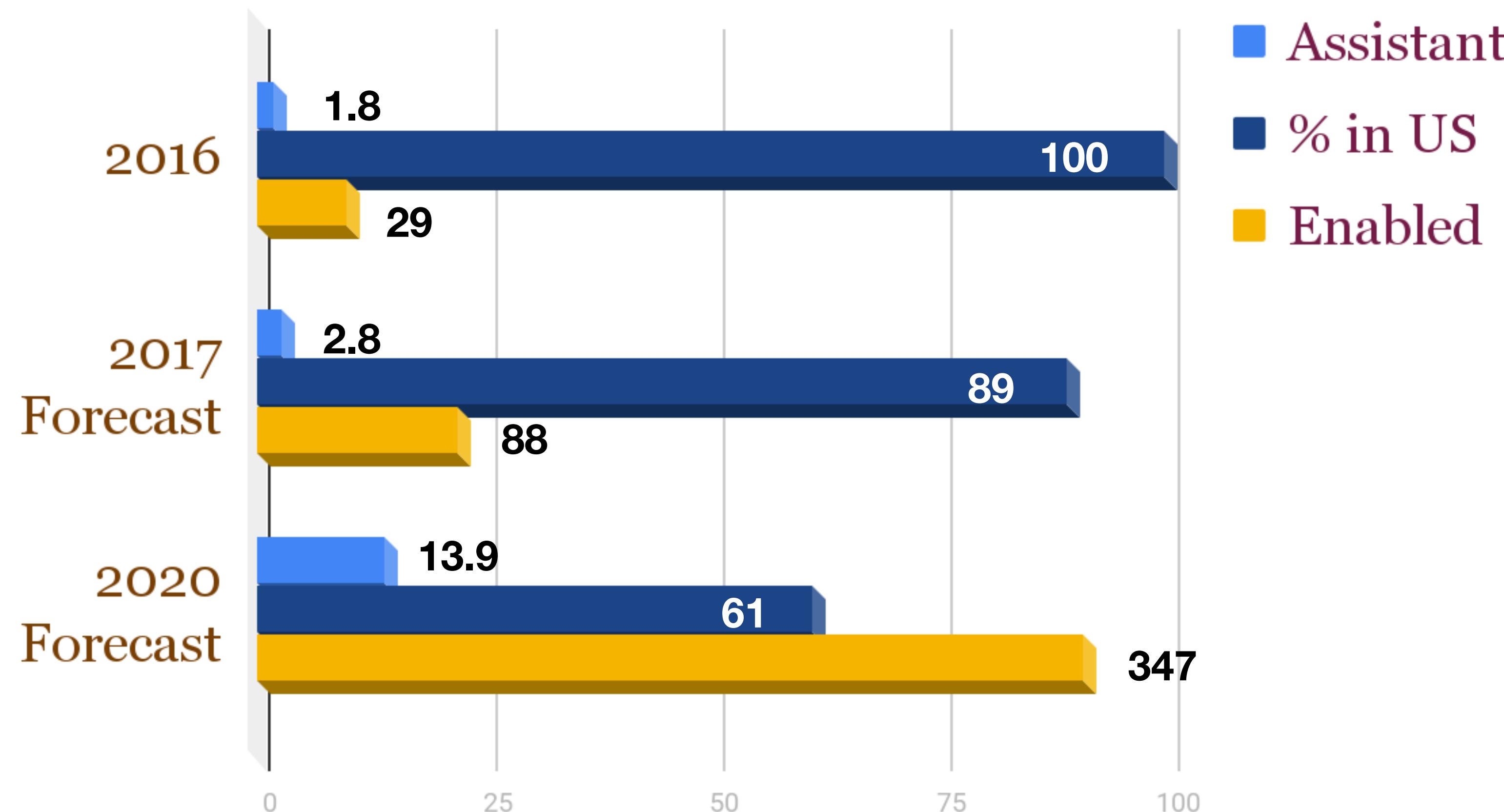
Word-level E2E ASR

Deepfakes

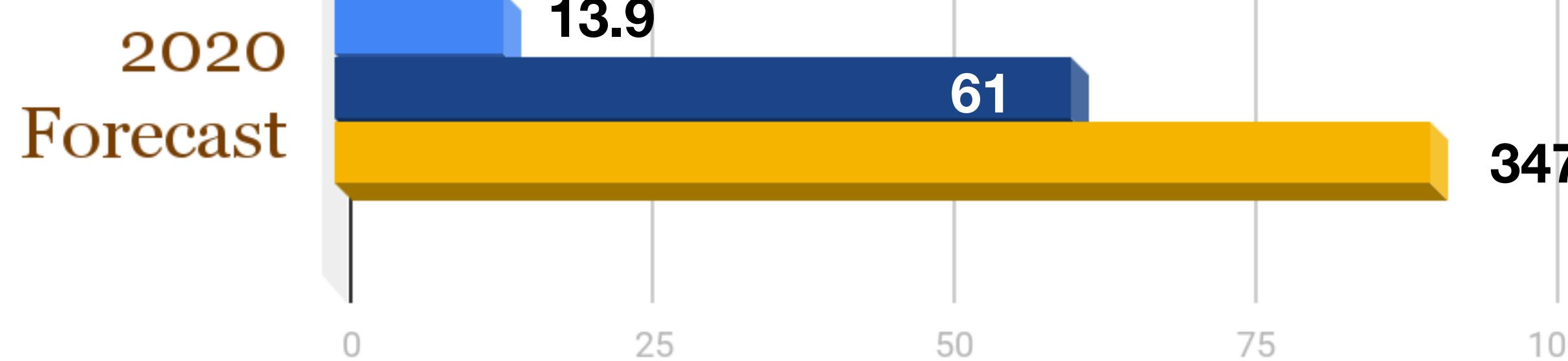
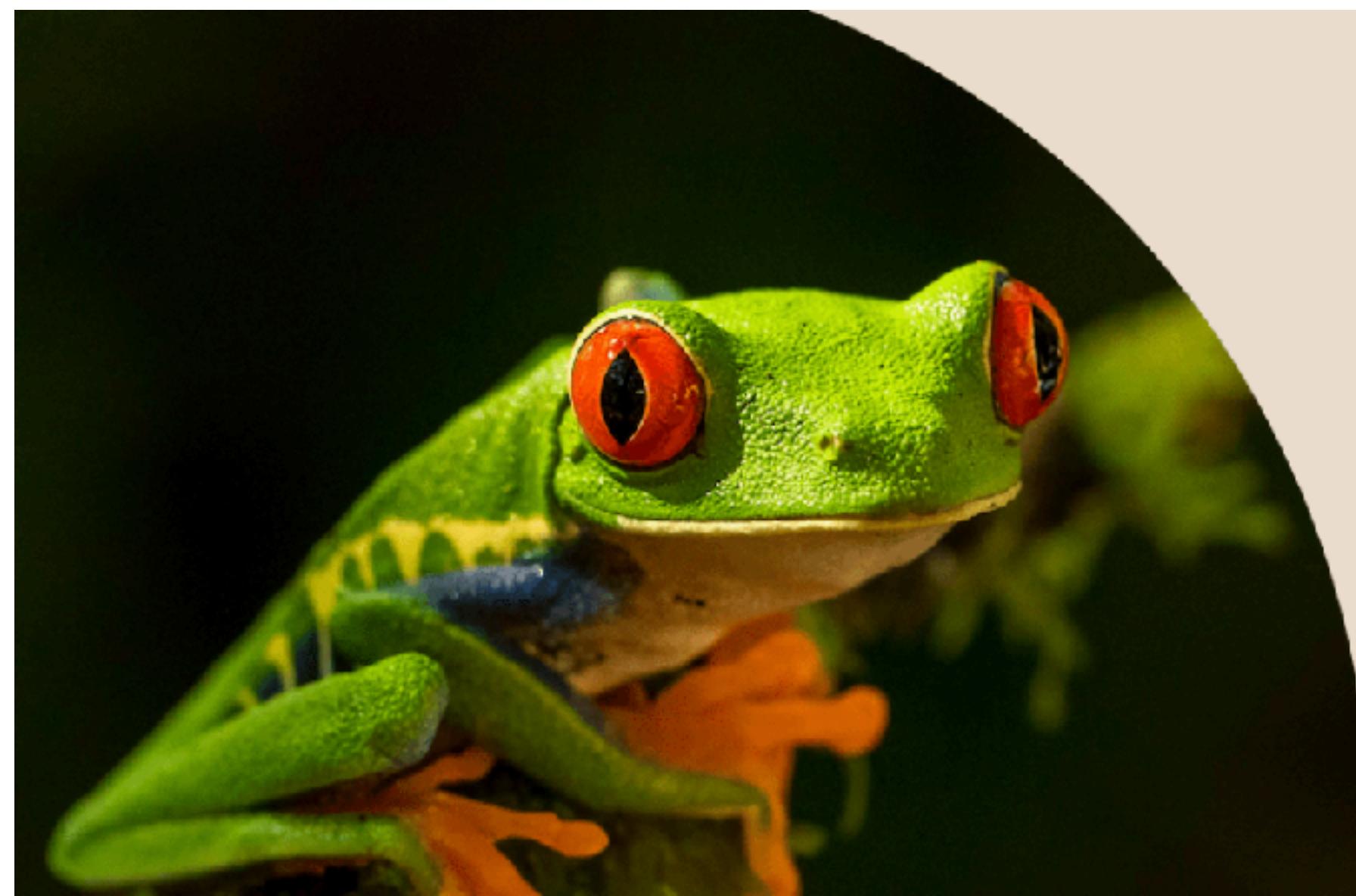
Metric GAN

# Exciting time to do speech research

## Market for Voice



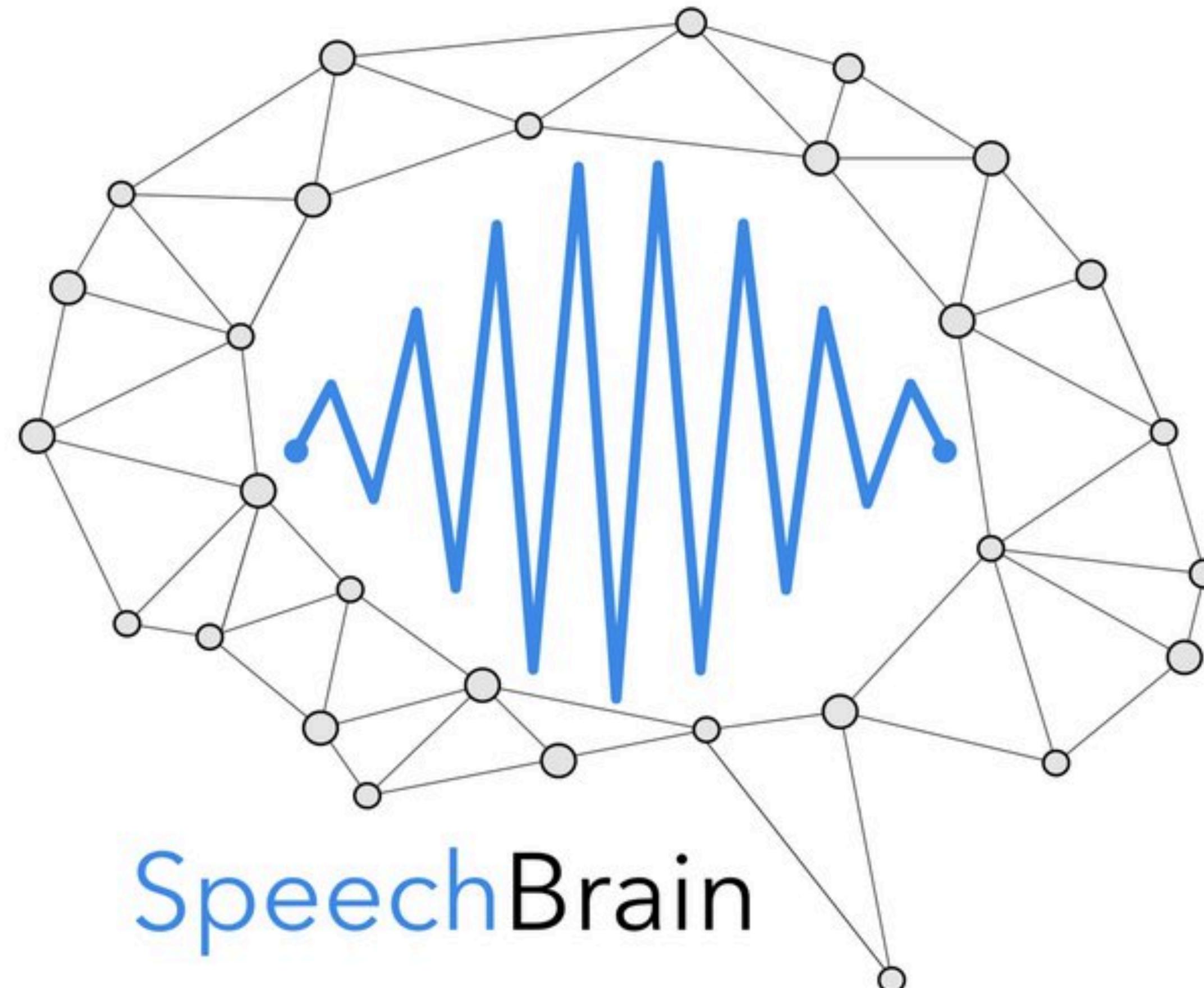
# Exciting time to do speech research



## Coqui, Freeing Speech

Coqui, a startup providing open speech tech for everyone 🐸  
Sign up with your email address to receive the Coqui newsletter.

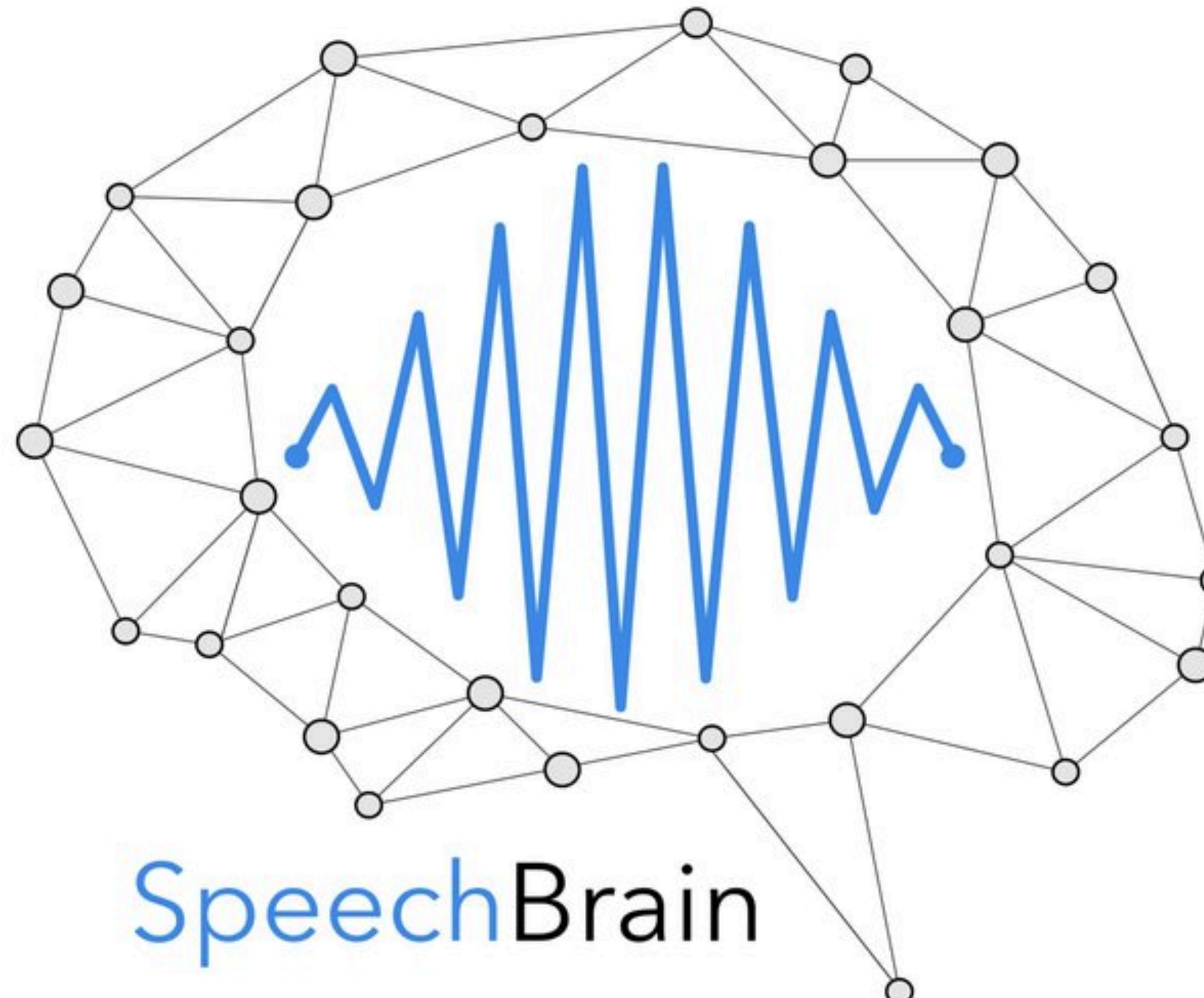
# Exciting time to do speech research



## Coqui, Freeing Speech

Coqui, a startup providing open speech tech for everyone 🐸  
Sign up with your email address to receive the Coqui newsletter.

# Exciting time to do speech research



## Coqui, Freeing Speech

Coqui, a startup providing open speech tech for everyone 🐸  
Sign up with your email address to receive the Coqui newsletter.

Over  
**80%+**

Of respondents  
are actively using  
ASR to transcribe  
speech data.

However...

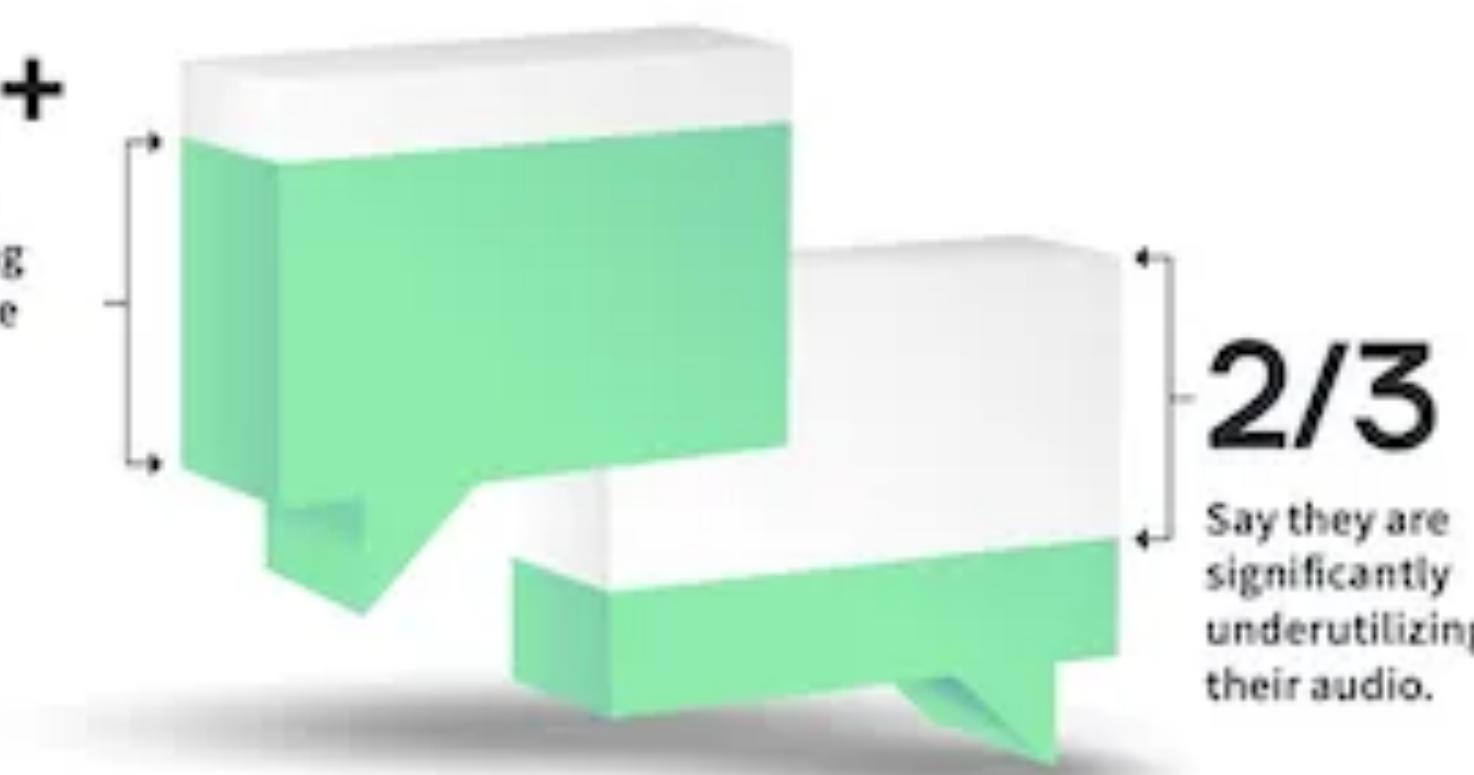


Image from: <https://coqui.ai/>  
<https://speechbrain.github.io/>

# What's next?

## Need to do more...

- Robust to variations in age, accent and ability
- Handling noisy real-life settings with many speakers (e.g., meetings, parties)
- Handling pronunciation variability
- Handling new languages/ dialects

# Accented Speech Recognition

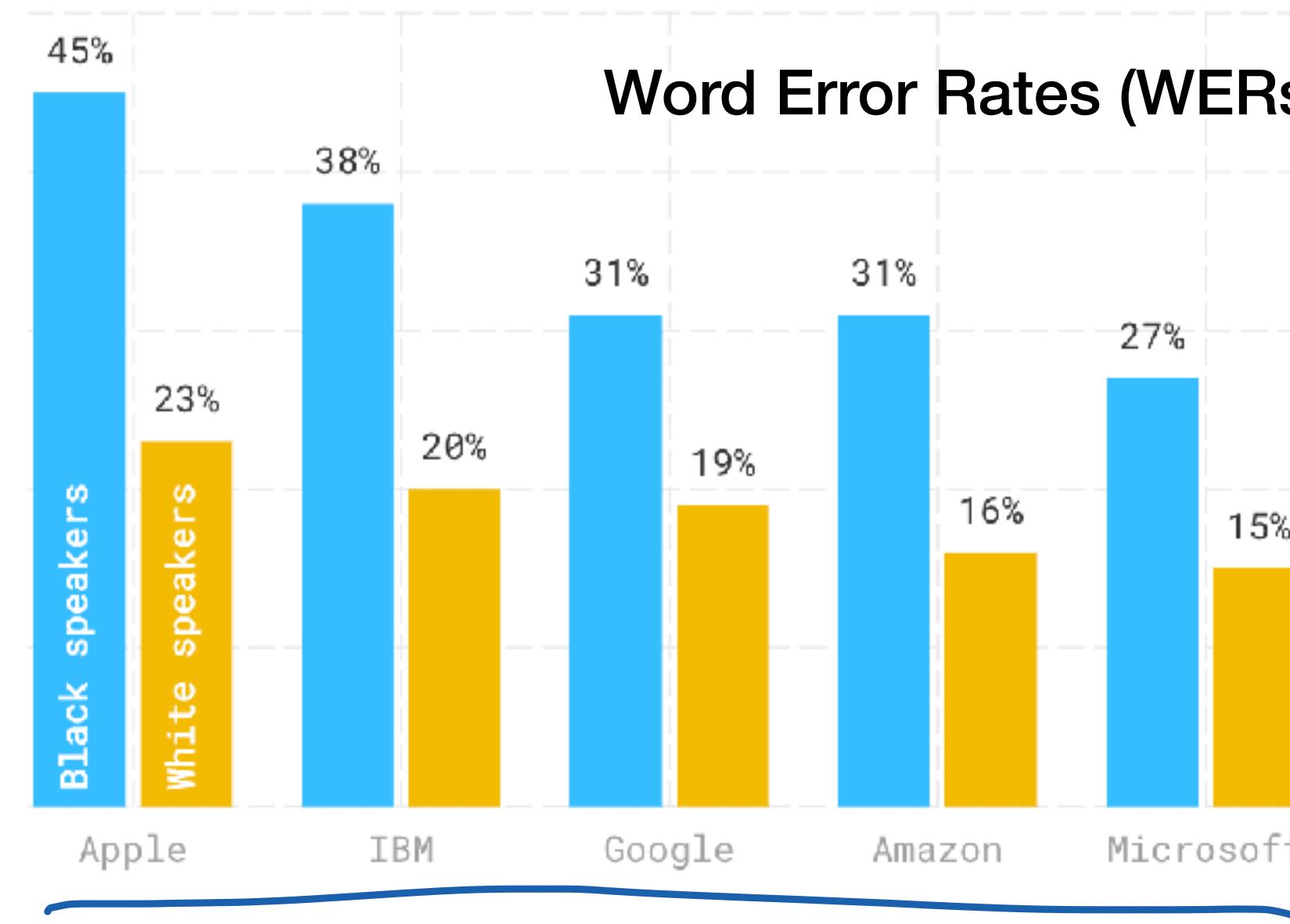
---

- “Voice is the next big platform, unless you have an accent” Wired'17

# Accented Speech Recognition

- “Voice is the next big platform, unless you have an accent” Wired’17

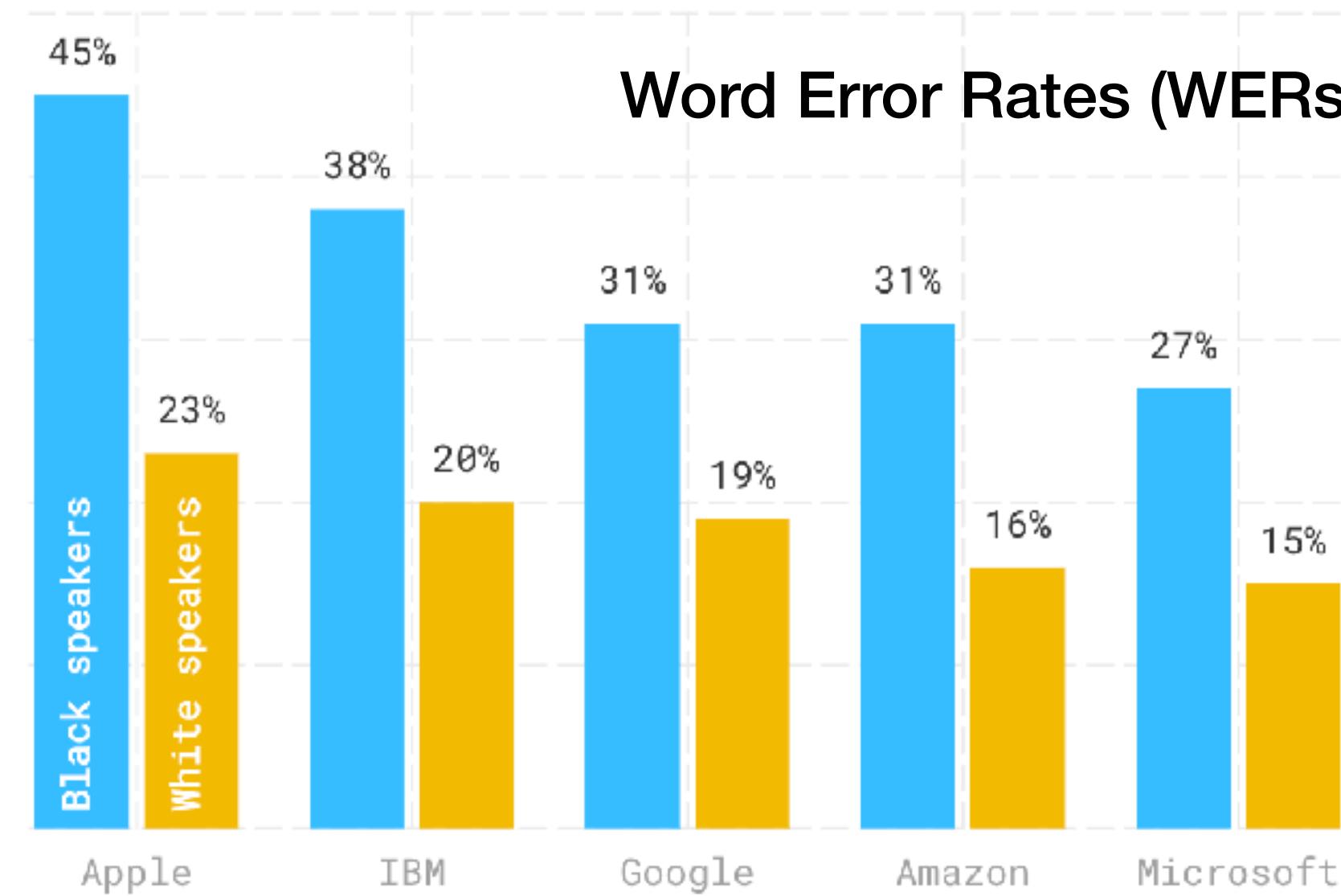
Non-native accents still pose a significant challenge to state-of-the-art ASR systems



# Accented Speech Recognition

- “Voice is the next big platform, unless you have an accent” Wired’17

Non-native accents still pose a significant challenge to state-of-the-art ASR systems



- WERs range from 11% to 53% across different Indian accented English speech test sets, using a state-of-the-art ASR system AKSJ’21

Wired’17 <https://www.wired.com/2017/03/voice-is-the-next-big-platform-unless-you-have-an-accent/>

Image from <https://fairspeech.stanford.edu/>, 2020

# What's next?

## Need to do more...

- Robust to variations in age, accent and ability
- Handling noisy real-life settings with many speakers (e.g., meetings, parties)
- Handling pronunciation variability
- Handling new languages/ dialects

## ... with less

- Fast (real-time) decoding using limited computational power/ memory
- Faster training algorithms
- Reduce duplicated effort across domains/languages
- Reduce dependence on language-specific resources
- Train with less labeled data

# Participation Points

Quiz	Points	# of responses
1	5	193
2	7	154
3	6	156
4	5	143

# Participation Points

- Four Moodle-quizzes

Quiz	Points	# of responses
1	5	193
2	7	154
3	6	156
4	5	143

# Participation Points

- Four Moodle-quizzes
- $\geq 50\%$  points across 4 quizzes gets 5 participation points

Quiz	Points	# of responses
1	5	193
2	7	154
3	6	156
4	5	143

# Participation Points

- Four Moodle-quizzes
- $\geq 50\%$  points across 4 quizzes gets 5 participation points
- [40-50) – 4  
[30-40) – 3  
[20-30) – 2  
[10-20) – 1  
 $< 10$  – 0

Quiz	Points	# of responses
1	5	193
2	7	154
3	6	156
4	5	143

# Participation Points

- Four Moodle-quizzes
  - $\geq 50\%$  points across 4 quizzes gets 5 participation points
  - [40-50) – 4  
[30-40) – 3  
[20-30) – 2  
[10-20) – 1  
 $< 10$  – 0
  - Subtract  $x$  from your participation points if you do not satisfy the seminar attendance requirement
    - $x = 1$  if you attended 3 or 4
    - $x = 2$  if you attended 1 or 2
    - $x = 3$  if you attended none
- | Quiz | Points | # of responses |
|------|--------|----------------|
| 1    | 5      | 193            |
| 2    | 7      | 154            |
| 3    | 6      | 156            |
| 4    | 5      | 143            |

# Final Exam (35 points)

1. WFST algorithms
2. WFSTs used in ASR
3. HMM algorithms/EM/Tied-state Triphone models
4. DNN-based acoustic models/TDNNs
5. N-gram LMs/Smoothing/RNN LMs
6. End-to-end ASR (CTC, LAS)
7. MFCC feature extraction
8. Search & Decoding
9. Speech Representations

# Final Exam (35 points)

1. WFST algorithms
2. WFSTs used in ASR
3. HMM algorithms/EM/Tied-state Triphone models
4. DNN-based acoustic models/TDNNs
5. N-gram LMs/Smoothing/RNN LMs
6. End-to-end ASR (CTC, LAS)
7. MFCC feature extraction
8. Search & Decoding
9. Speech Representations

*Questions can be asked on any of the 9 topics listed above. This will be an open-notes exam.*

# Final Exam (35 points)

*Date and Time: April 28th, 17:30 pm - 20:30 pm*



1. WFST algorithms
2. WFSTs used in ASR
3. HMM algorithms/EM/Tied-state Triphone models
4. DNN-based acoustic models/TDNNs
5. N-gram LMs/Smoothing/RNN LMs
6. End-to-end ASR (CTC, LAS)
7. MFCC feature extraction
8. Search & Decoding
9. Speech Representations

*Questions can be asked on any of the 9 topics listed above. This will be an open-notes exam.*

# **Seminar: Presentation + Review + Scientific Article**

# **Seminar: Presentation + Review + Scientific Article**

- Presentation scores (out of 8) will be released this week

# Seminar: Presentation + Review + Scientific Article

- Presentation scores (out of 8) will be released this week
- Best presentation: **PhoeniTIX (Efficient Vocalize): T Pavan Kalyan, Pranamya Prashant Kulkarni, Arif Ashfaque Ahmad**

# Seminar: Presentation + Review + Scientific Article

- Presentation scores (out of 8) will be released this week
- Best presentation: **PhoeniTIX (Efficient Vocalize): T Pavan Kalyan, Pranamya Prashant Kulkarni, Arif Ashfaque Ahmad**
- Honorable Mentions (in no particular order):  
AudioAdversary (Unmuters) /Visual Voice (DeepSpeech) /Mask CTC (angry\_nerds)/  
AVID (Tri-Vachan) /SpeakQL (235161) /Large margin SV (Speech Hackers)

# Seminar: Presentation + Review + Scientific Article

- Presentation scores (out of 8) will be released this week
- Best presentation: **PhoeniTIX (Efficient Vocalize): T Pavan Kalyan, Pranamya Prashant Kulkarni, Arif Ashfaque Ahmad**
- Honorable Mentions (in no particular order):  
AudioAdversary (Unmuters) /Visual Voice (DeepSpeech) /Mask CTC (angry\_nerds)/  
AVID (Tri-Vachan) /SpeakQL (235161) /Large margin SV (Speech Hackers)
- Review will be due on Moodle on or before midnight of April 15th.

# Seminar: Presentation + Review + Scientific Article

- Presentation scores (out of 8) will be released this week
- Best presentation: **PhoeniTIX (Efficient Vocalize): T Pavan Kalyan, Pranamya Prashant Kulkarni, Arif Ashfaque Ahmad**
- Honorable Mentions (in no particular order):  
AudioAdversary (Unmuters) /Visual Voice (DeepSpeech) /Mask CTC (angry\_nerds)/  
AVID (Tri-Vachan) /SpeakQL (235161) /Large margin SV (Speech Hackers)
- Review will be **due on Moodle on or before midnight of April 15th.**
- Scientific article will be due first week of May (as per poll). Submission portal will be open this week for those who want to submit by April 17th.

# **Seminar: Poster Presentation + Hacker**

# Seminar: Poster Presentation + Hacker

- Poster presentation will be held online on May 6, 7, 8. Specifics will be conveyed via Moodle.
- Hacker presentations will also be held on the same 3 dates (May 6, 7, 8). For the Hacker role:
  - ✓ Share a link to your GitHub repository.
  - ✓ Each team describes what they have implemented and how it relates to the paper they were assigned.
  - ✓ Short viva-voce.
  - ✓ Note there is no report or write-up associated with the hacker role.