



# Part 1: ASR - Automatic Speech Recognition

Building Conversational AI Applications

# About Me

Instructor Name – Instructor Contact (e.g. LinkedIn / email)



- About instructor 1
- About instructor 2
- About instructor 3



## Course Timeline

8:50	—	9:00	Welcome
9:00	—	10:00	Lecture 1
10:00	—	11:00	Lab 1
11:00	—	11:15	Break
11:15	—	12:15	Lecture 2
12:15	—	13:15	Lab 2
13:15	—	14:00	break
14:00	—	15:00	Lecture 3
15:00	—	16:00	Lab 3
16:00	—	17:00	Assessment support



# Course Agenda

## Part 1: ASR - Automatic Speech Recognition

Lecture: Challenges of conversational AI applications, building blocks, ASR (automatic speech recognition) focus

Lab: Investigate ASR pipeline with NeMo and Riva

## Part 2: TTS – Text To Speech and Conversational AI Customization

Lecture: How can we take advantage of conversational AI platforms to build custom applications?

Lab: Investigate the TTS pipeline with NeMo and Riva; then build a full custom conversational AI pipeline.

## Part 3: Production Deployment

Lecture: Discussion of production deployment considerations including an overview of NVIDIA Triton Inference Server and TensorRT.

Lab: Hands-on deployment of an example application using Helm and Kubernetes; includes profiling the application



# ASR (Part 1)

## Lecture

- What is Conversational AI? Complexity of Deploying Conversational AI Applications
- Deep Learning Fundamentals
- Speech and Text Processing
- Automatic Speech Recognition
- ASR Tools: Riva, NGC, and NeMo
- Lab Overview

## Lab

- Investigate ASR Pipeline with NeMo and Riva



Shannon's



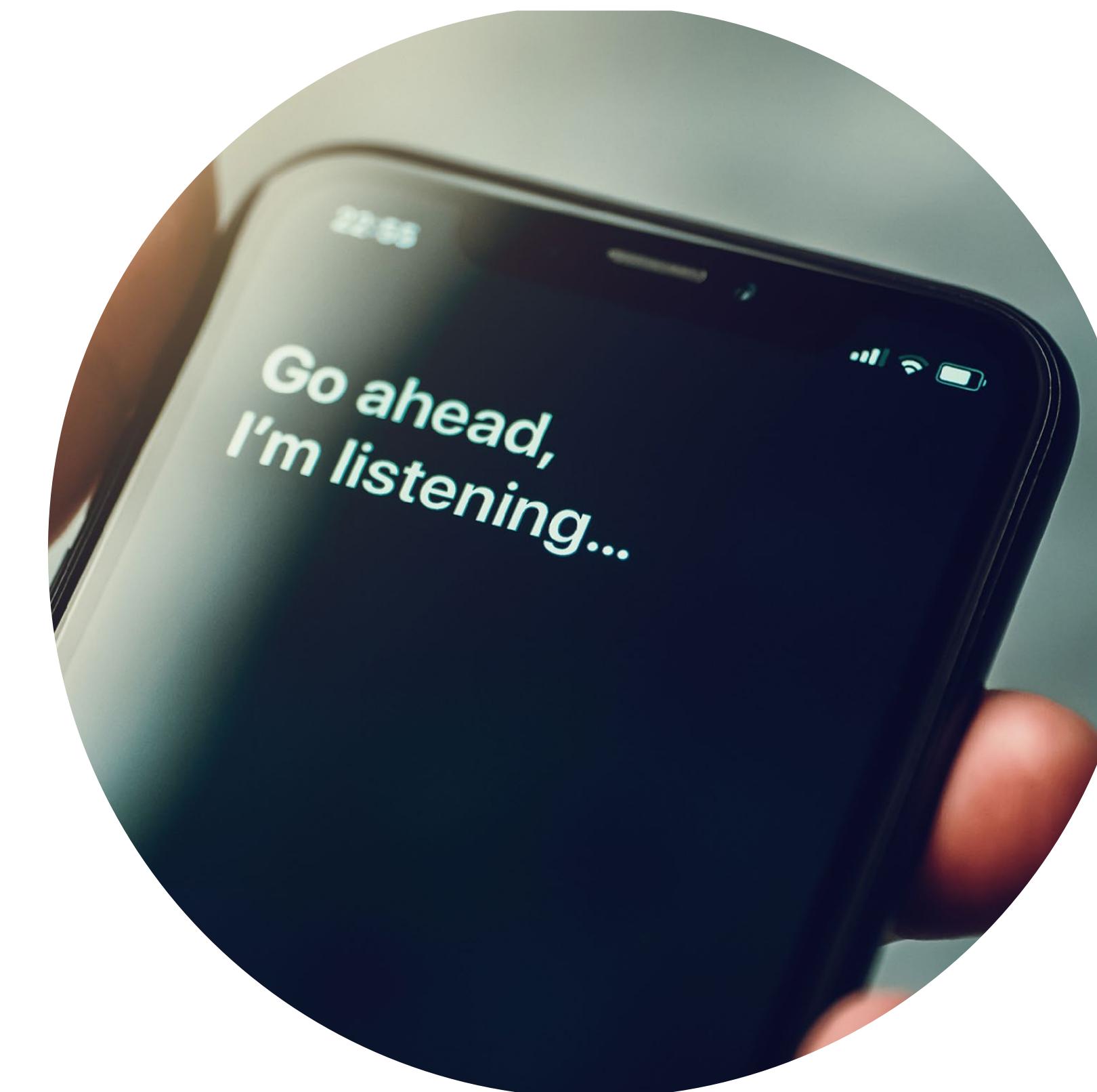
# Automatic Speech Recognition Is Everywhere!

Hundreds of billions of minutes of speech generated daily



Call Center

500M Calls Daily



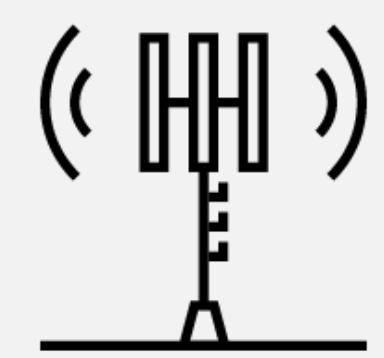
Consumer Apps

1.8B Daily Minutes



Online Meetings

200M Daily



Telecom



Finance



Healthcare



Retail



Building Conversational AI Applications  
is Not Trivial

# Challenges

New models , domain-specific accuracy, complex pipelines, real-time responses

## HIGH ACCURACY

Lack of access to state-of-the-art speech algorithms and models

## REAL-TIME RESULTS

Require low-latency for natural interactions

## SCALABLE DEPLOYMENT

- Limited support for deployment platforms
- Limited support for scaling to multiple client requests

## CUSTOMIZATION

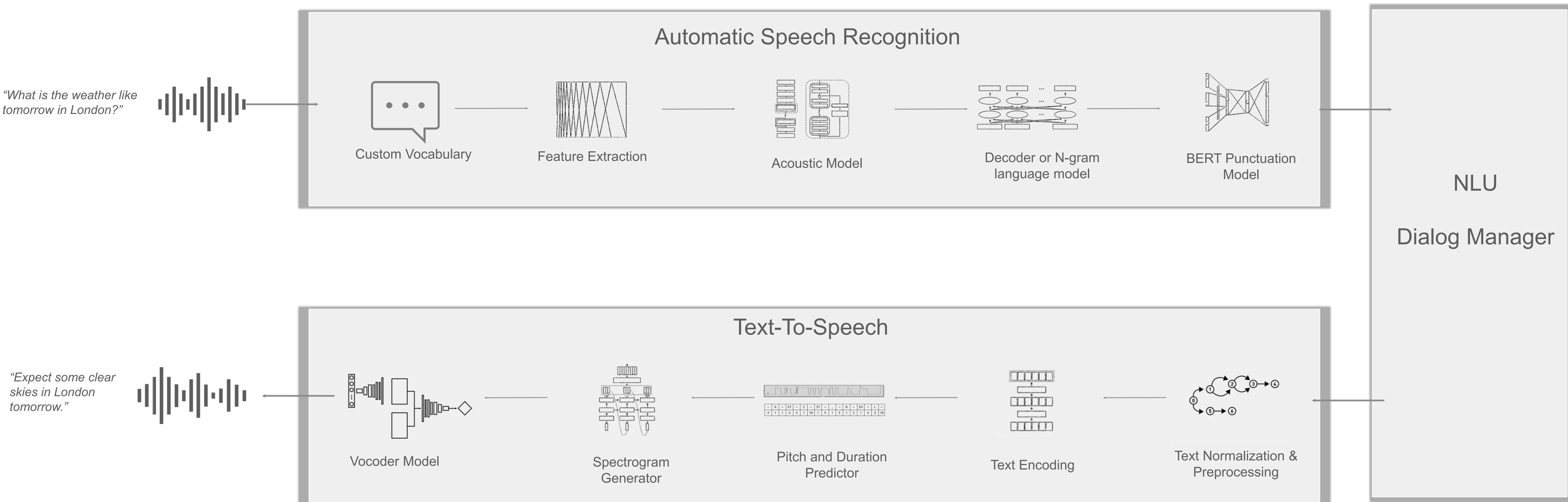
Limited model, domain, and pipeline customization

## DATA SOVEREIGNTY

Limited control over the proprietary data used for training and inference

# Challenges

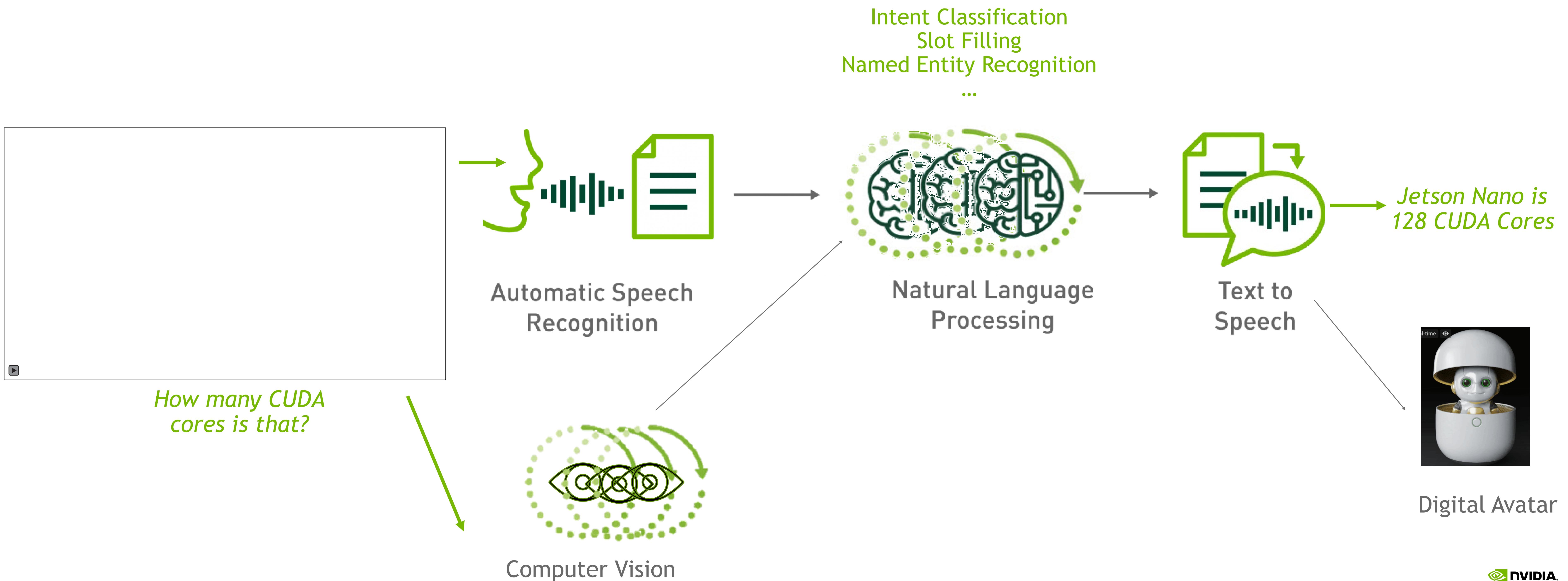
New models , domain-specific accuracy, complex pipelines, real-time responses



Real-time applications need to deliver latency <300 ms

# Challenges

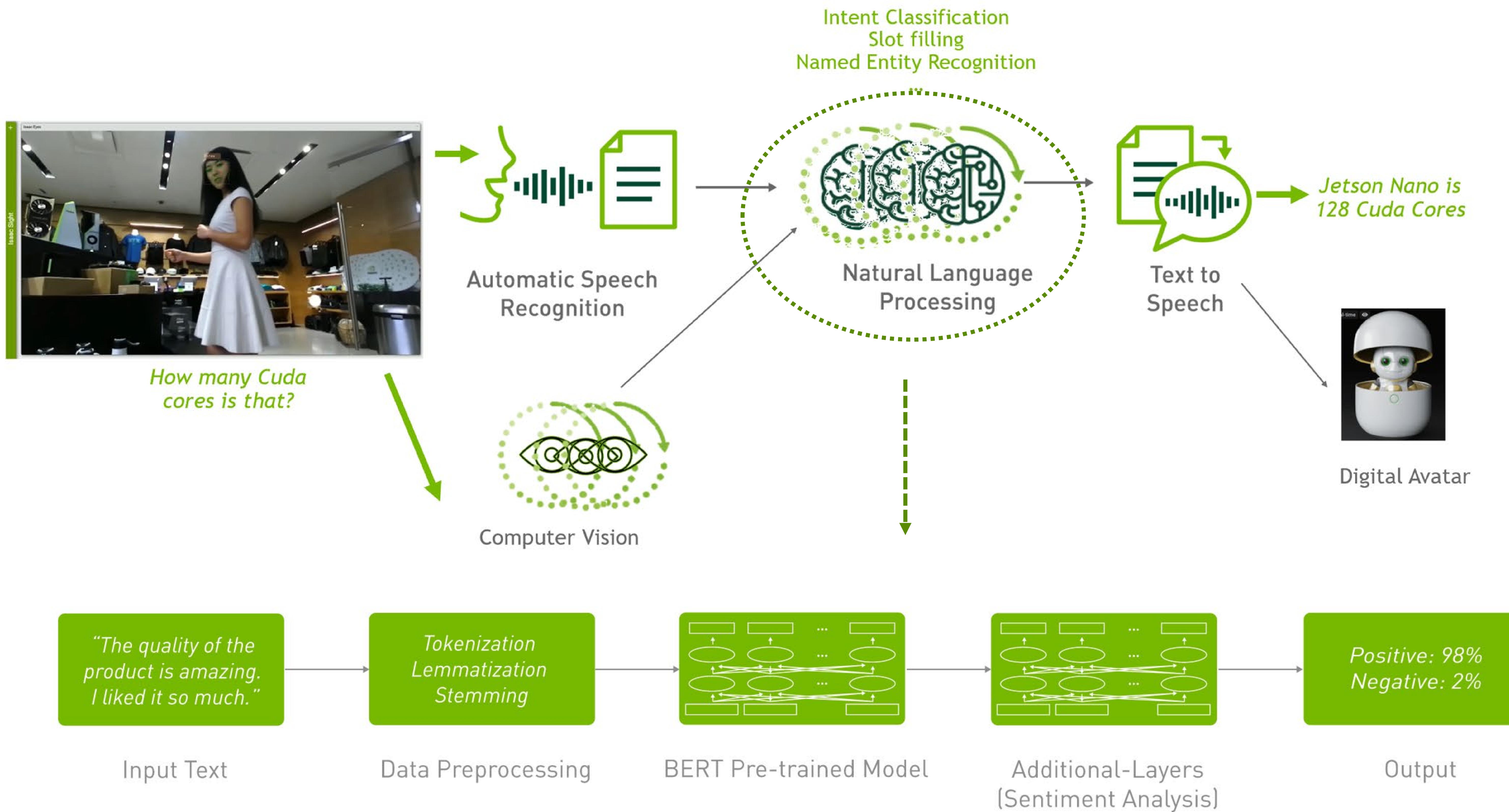
New models , domain-specific accuracy, complex pipelines, real-time responses



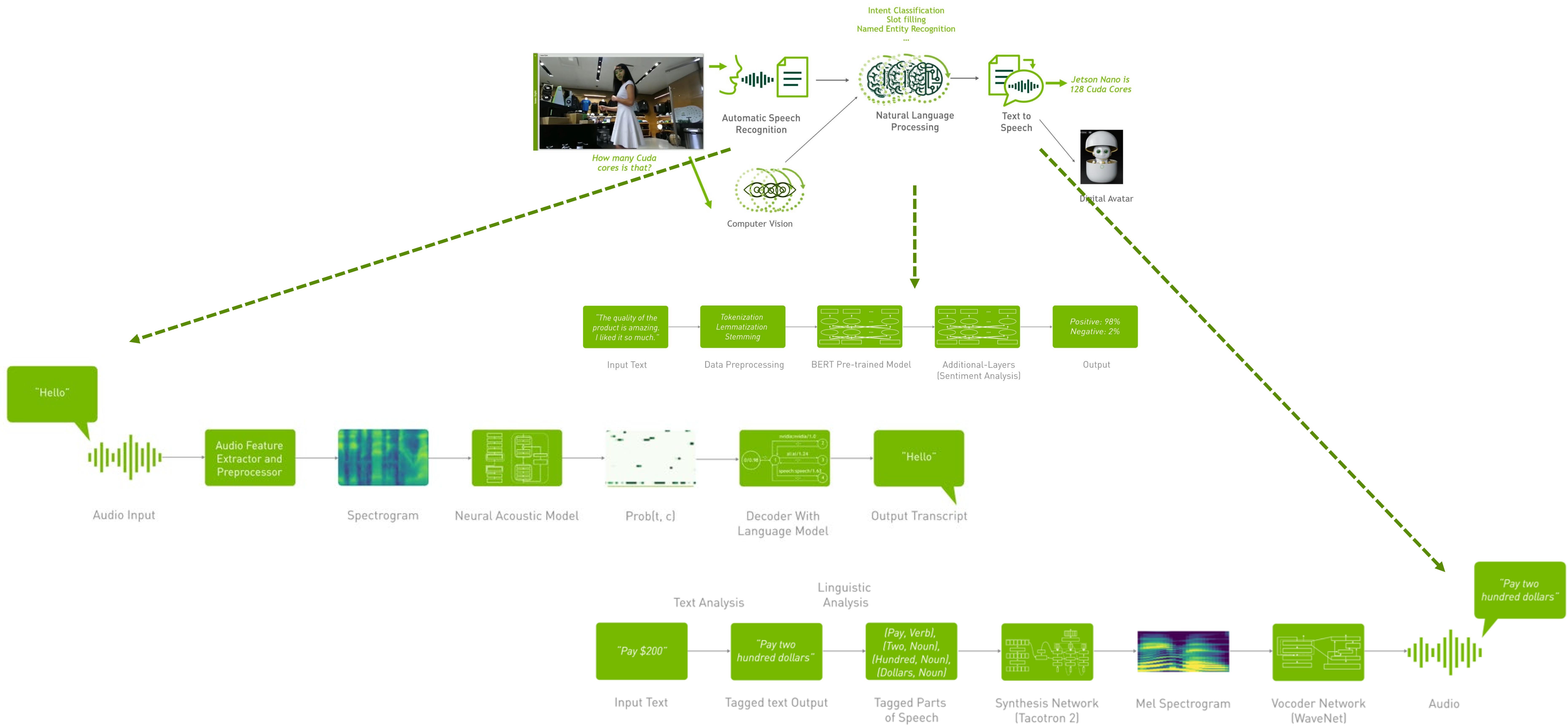
!

This is still an oversimplification of the  
problem

# What About Additional Logic?

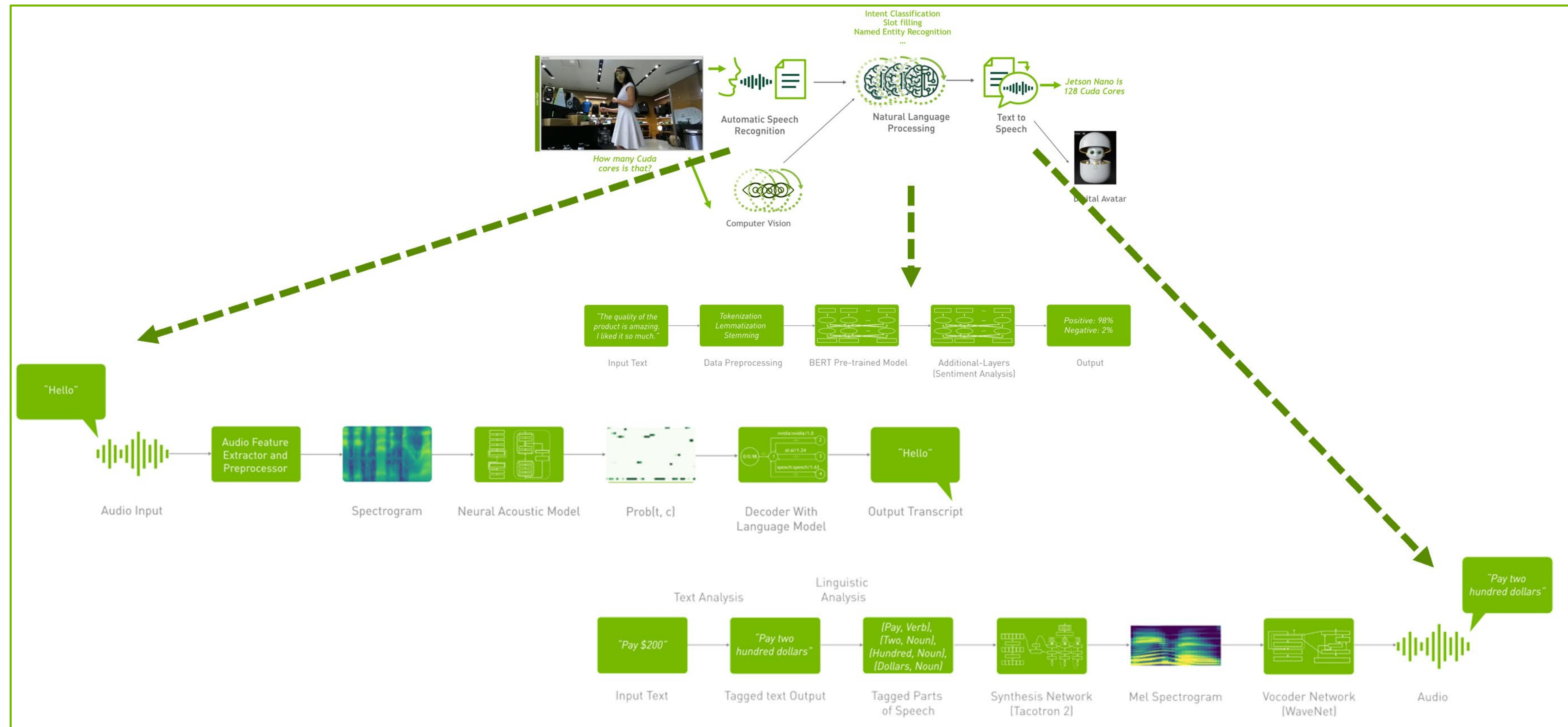


# Not Something Unique to NLP



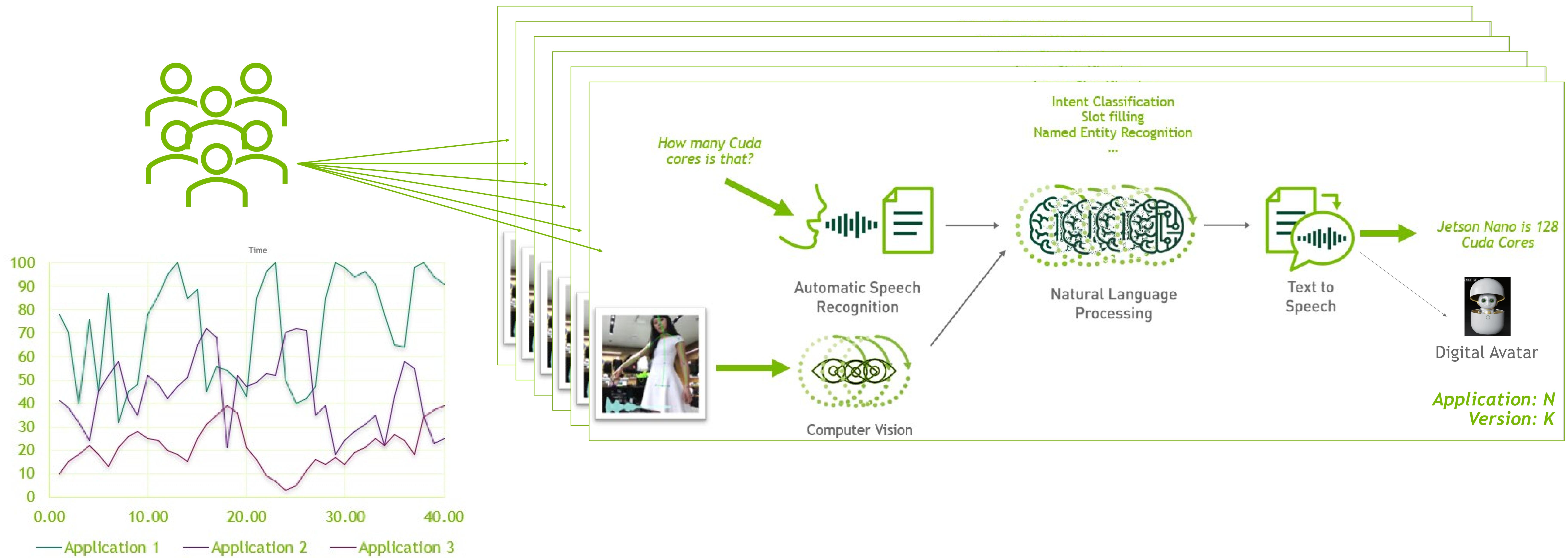
# This Was Just a Single Application

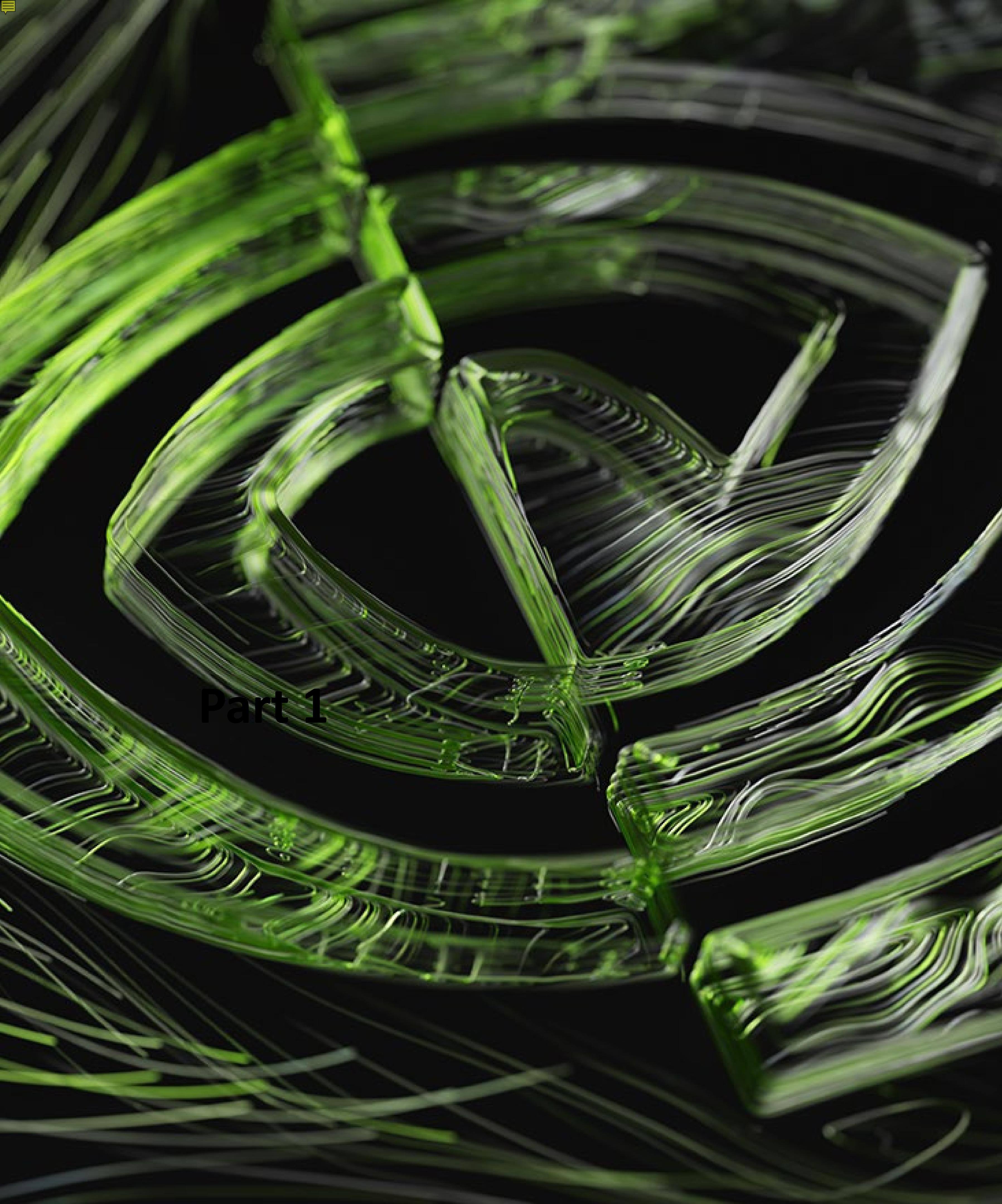
## Application 1



# Challenges

Latency, throughput, utilization, capital investment, ongoing costs





# ASR (Part 1)

## Lecture

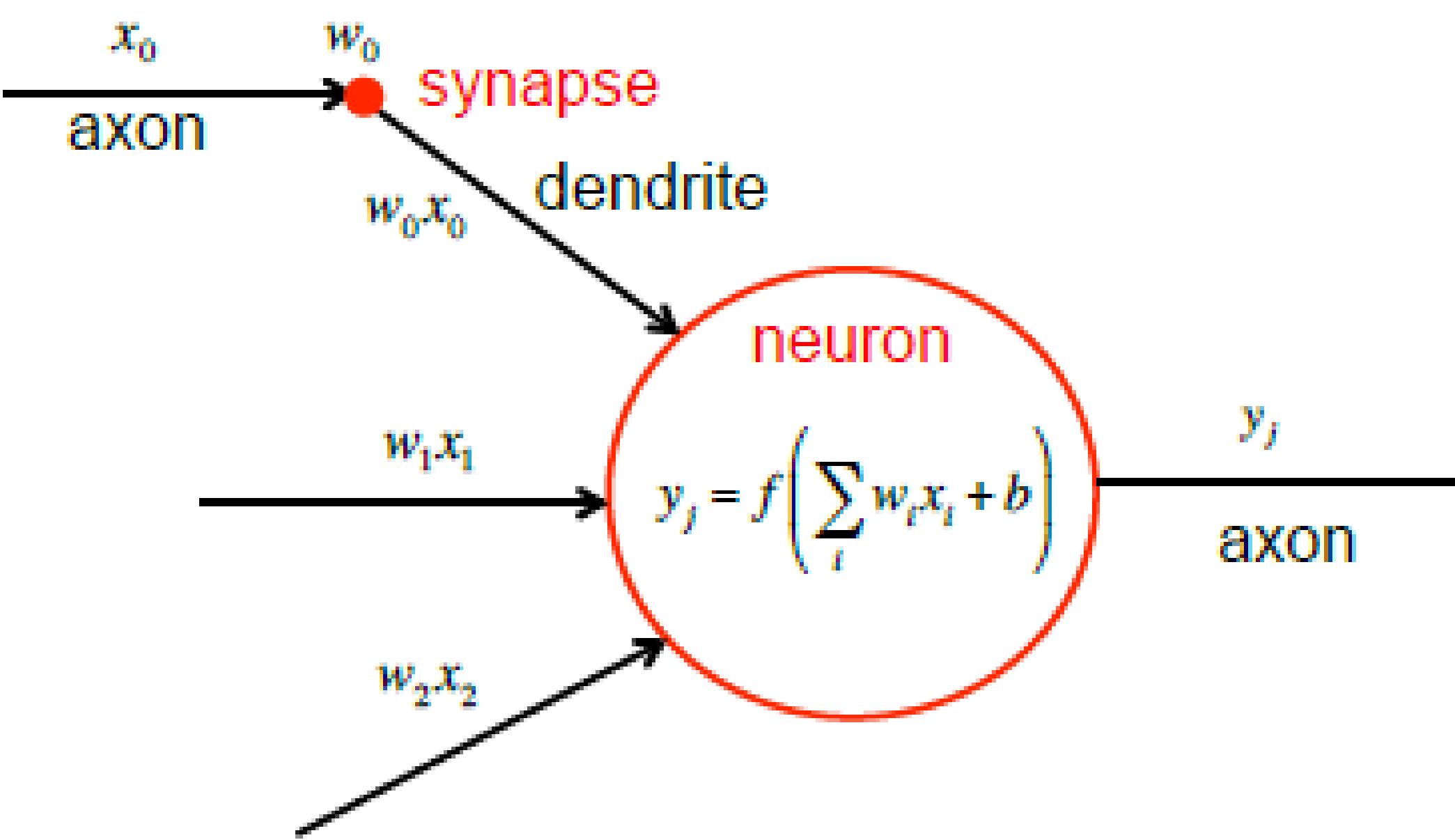
- What is Conversational AI? Complexity of Deploying Conversational AI Applications
- Deep Learning Fundamentals
- Speech and Text Processing
- Automatic Speech Recognition
- ASR Tools: Riva, NGC, and NeMo
- Lab Overview

## Lab

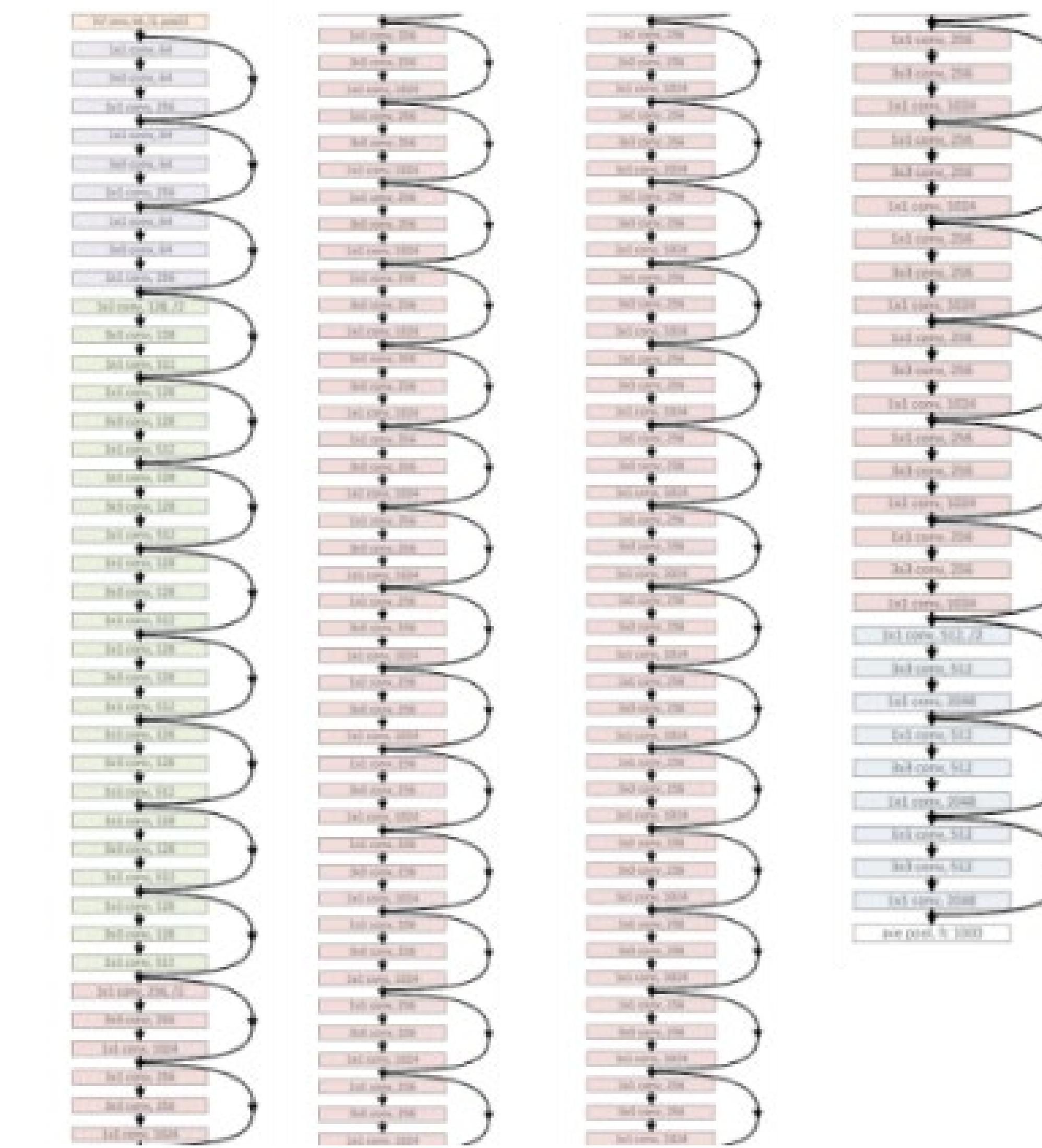
- Investigate ASR Pipeline with NeMo and Riva

# Neural Networks are Not New

They are surprisingly simple as an algorithm



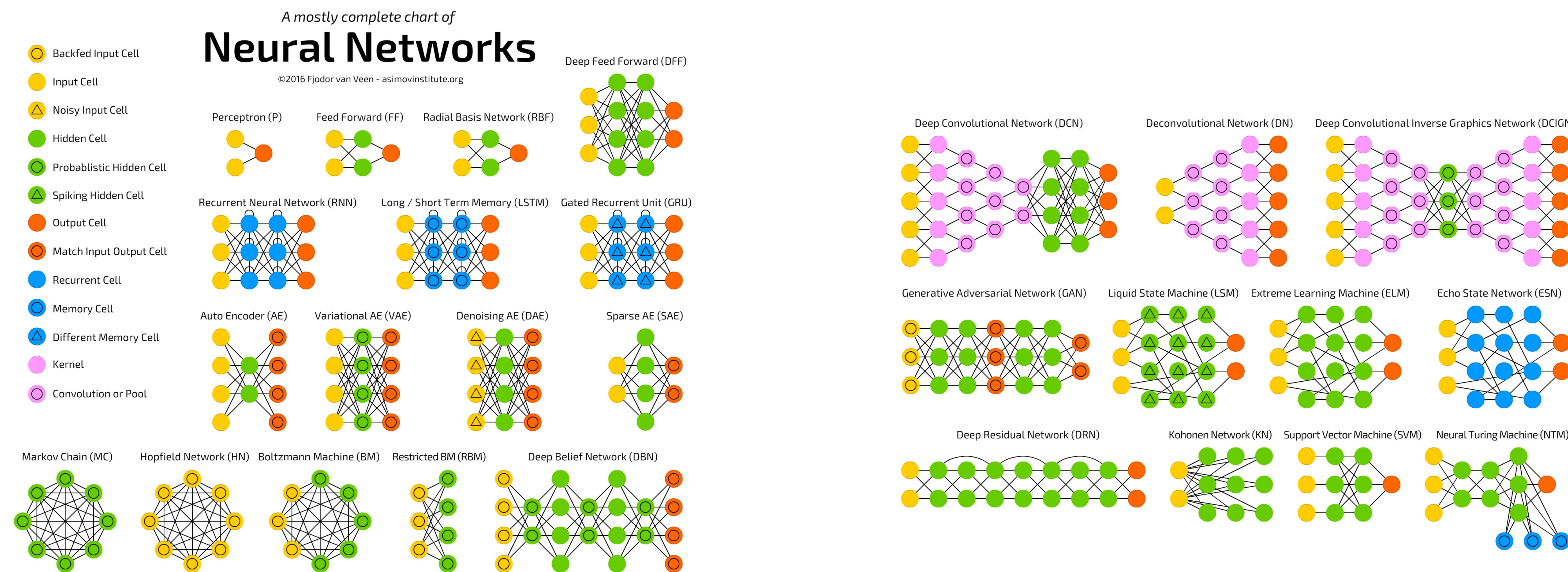
Neural Network - Neuron



Deep Neural Network

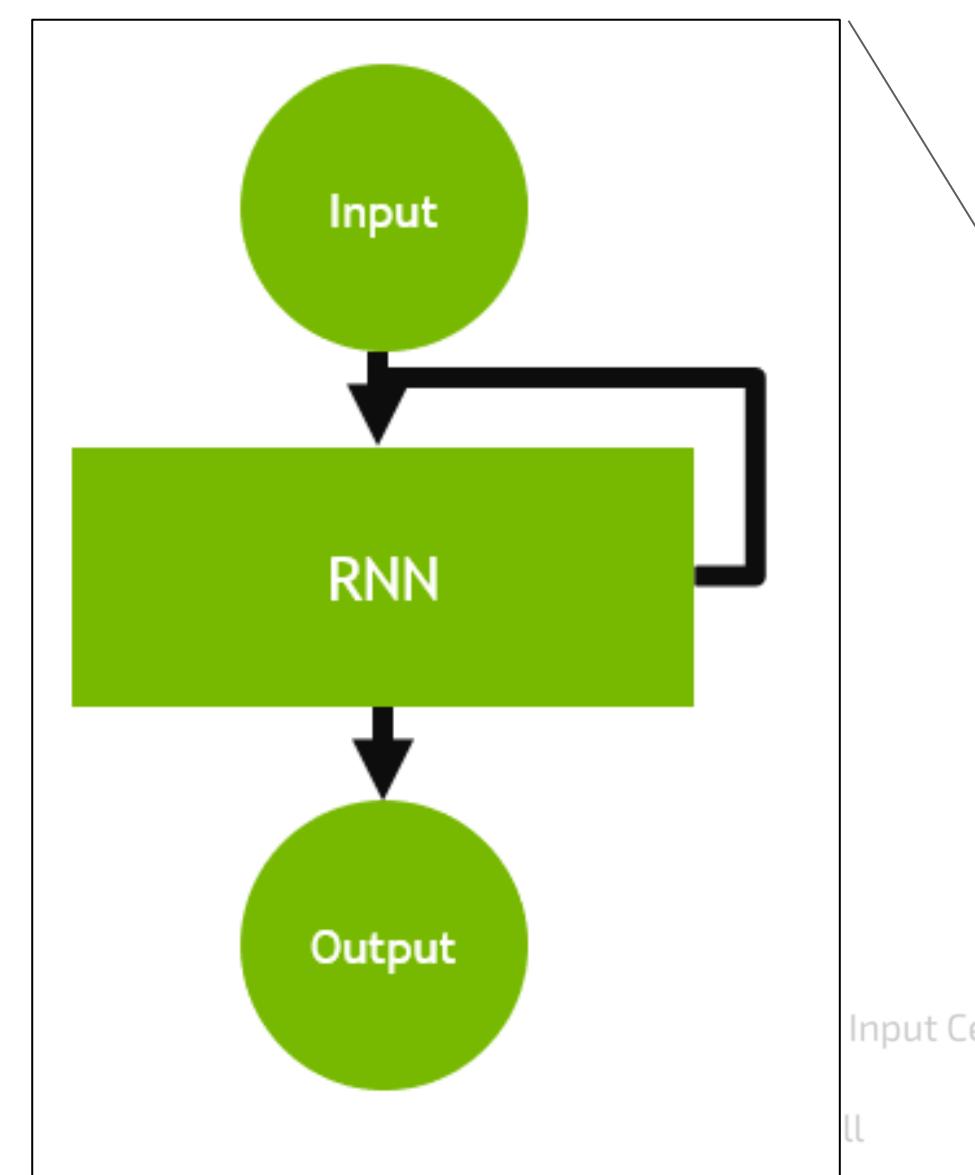
# Neural Networks

Convolutional and recurrent neural networks



# Neural Networks

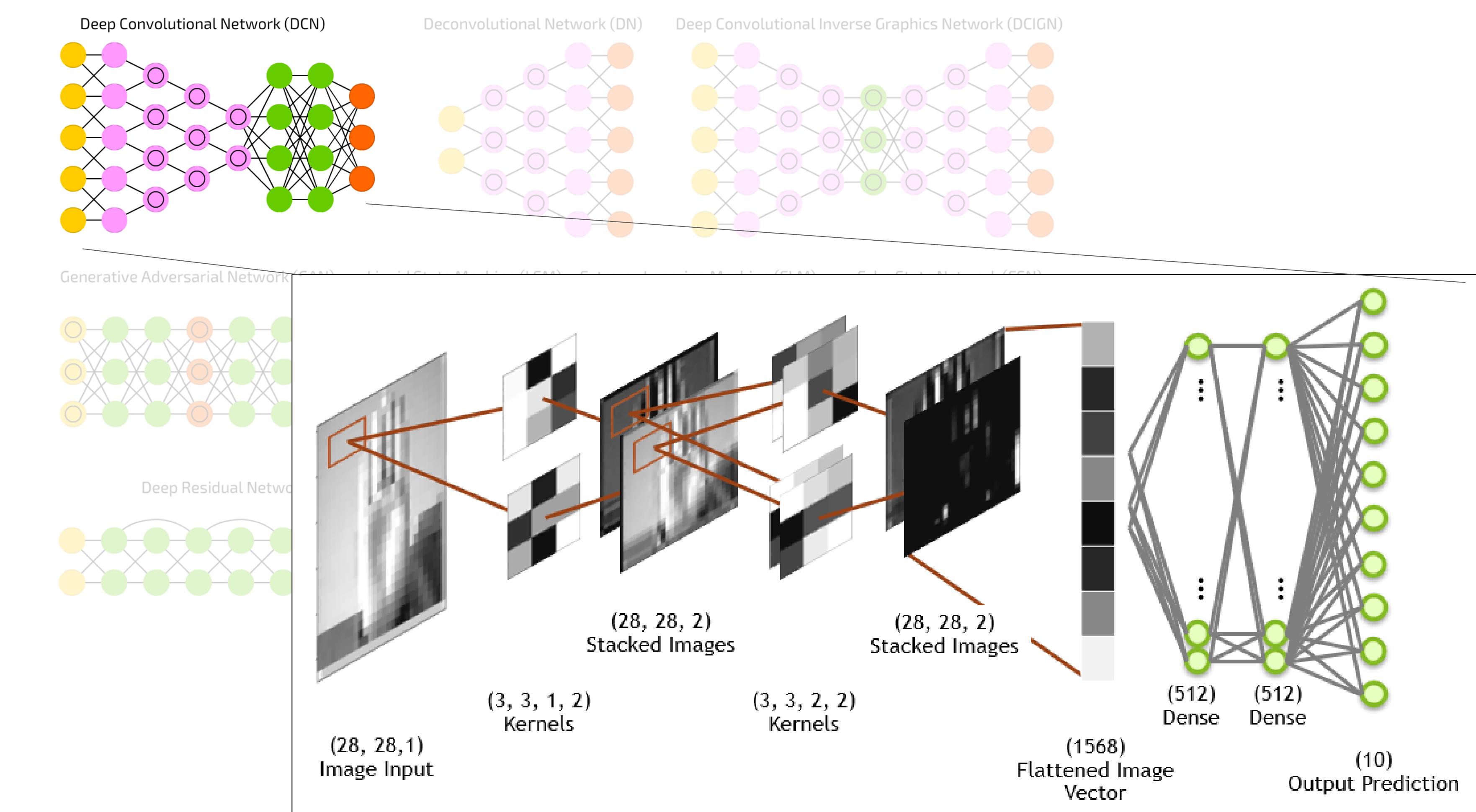
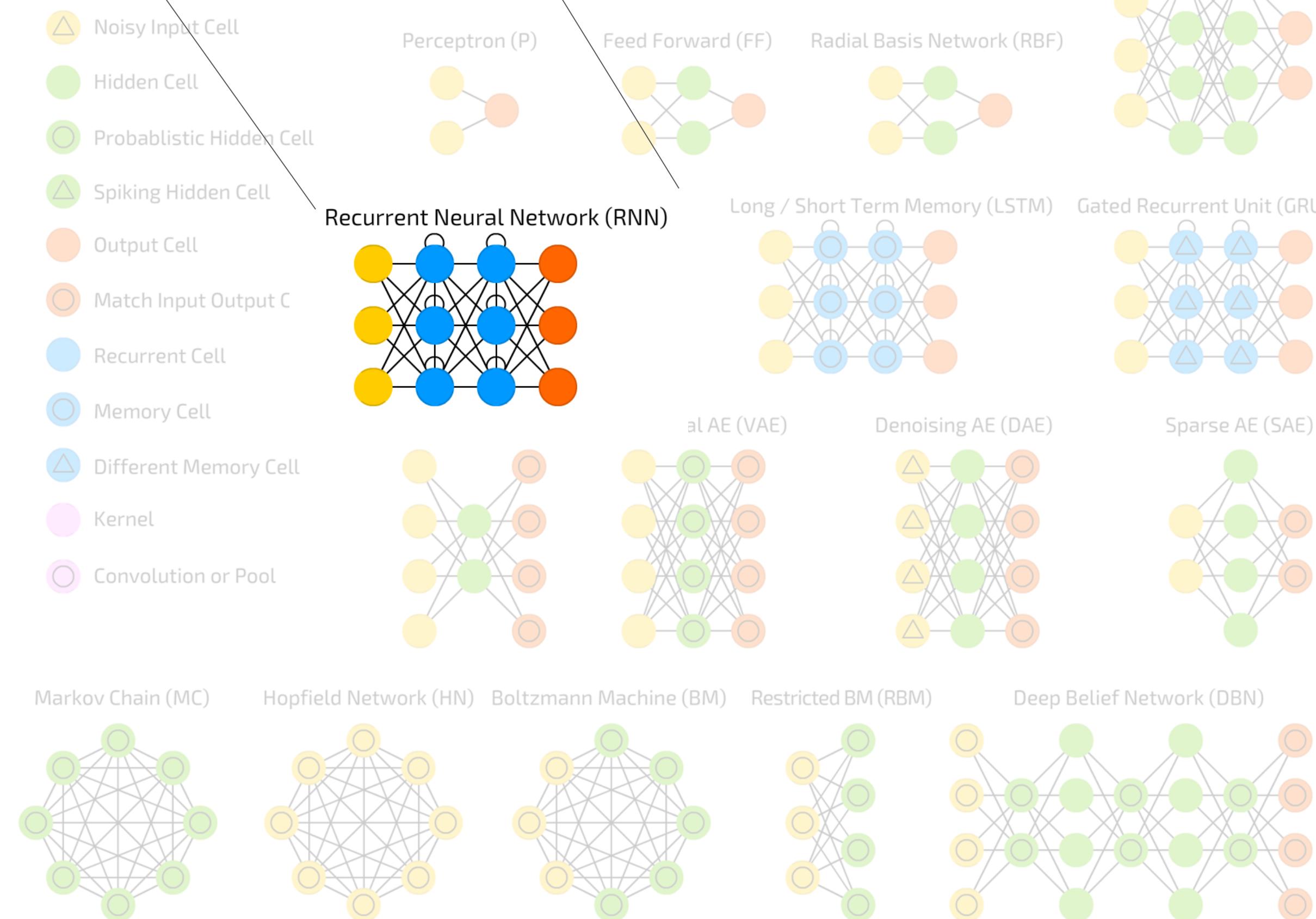
## Convolutional and recurrent neural networks



### A mostly complete chart of Neural Networks

Legend:

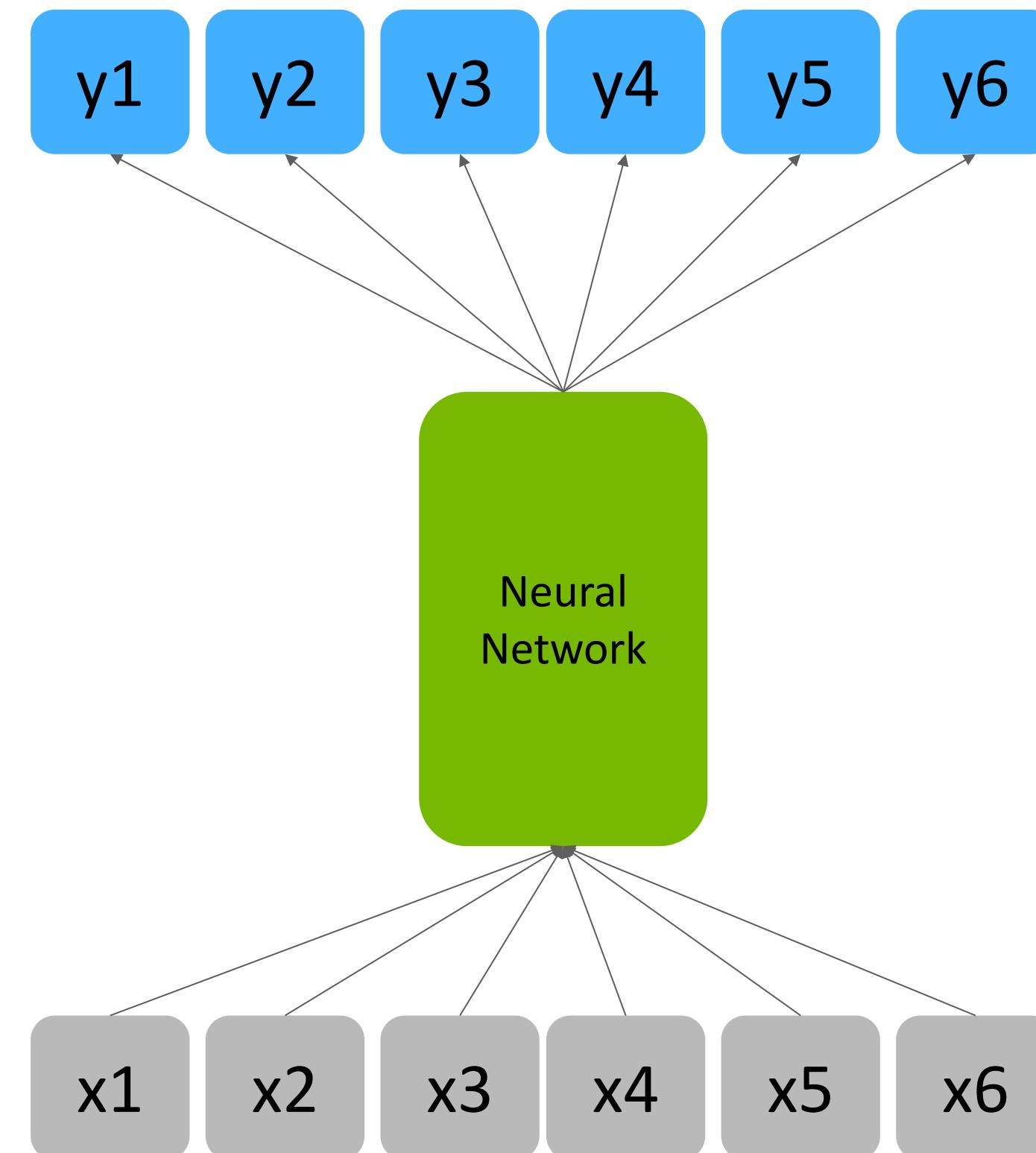
- Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Different Memory Cell
- Kernel
- Convolution or Pool



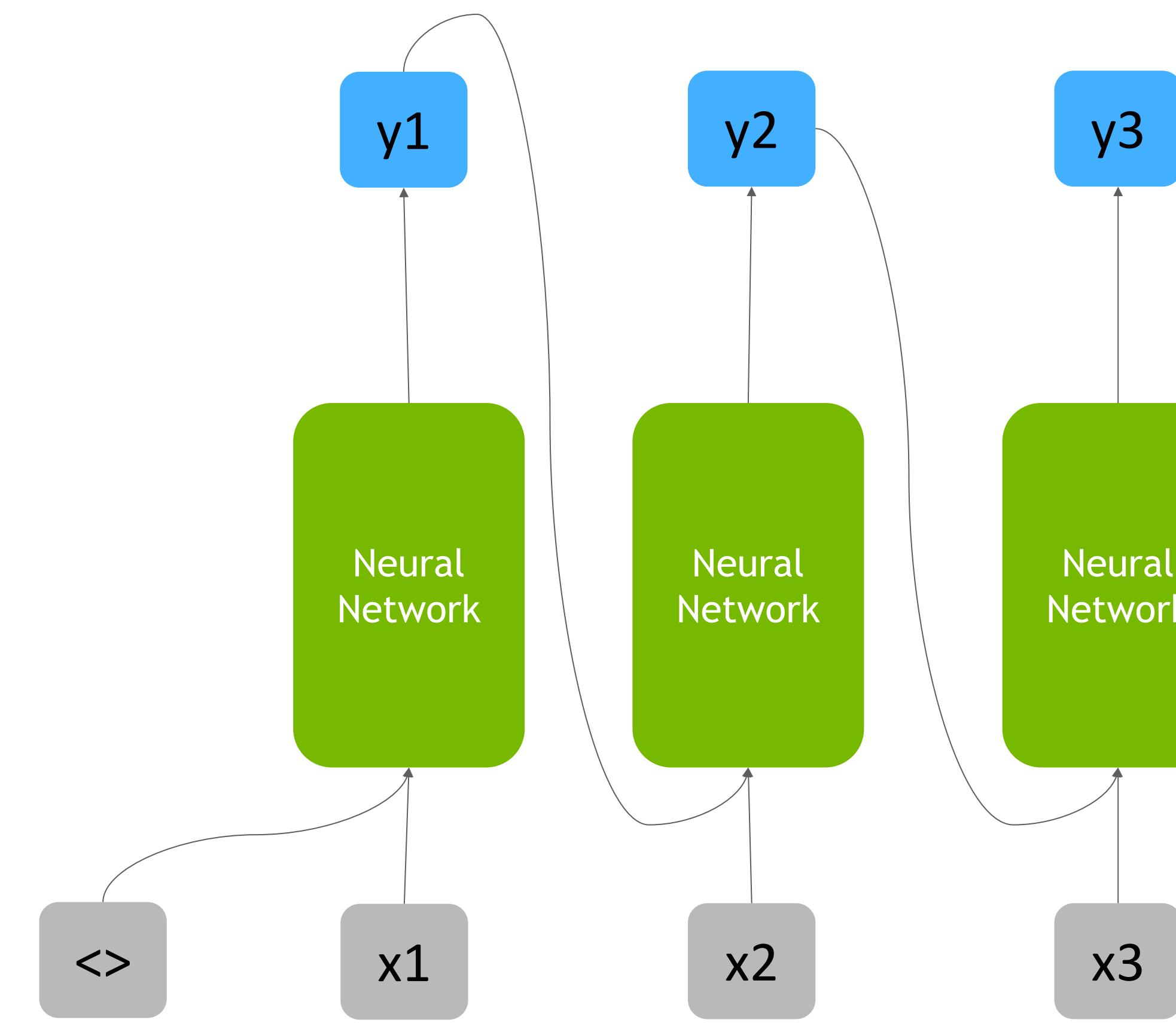
# Neural Networks

## Autoregressive and non-autoregressive

Non Autoregressive Models



Autoregressive Models

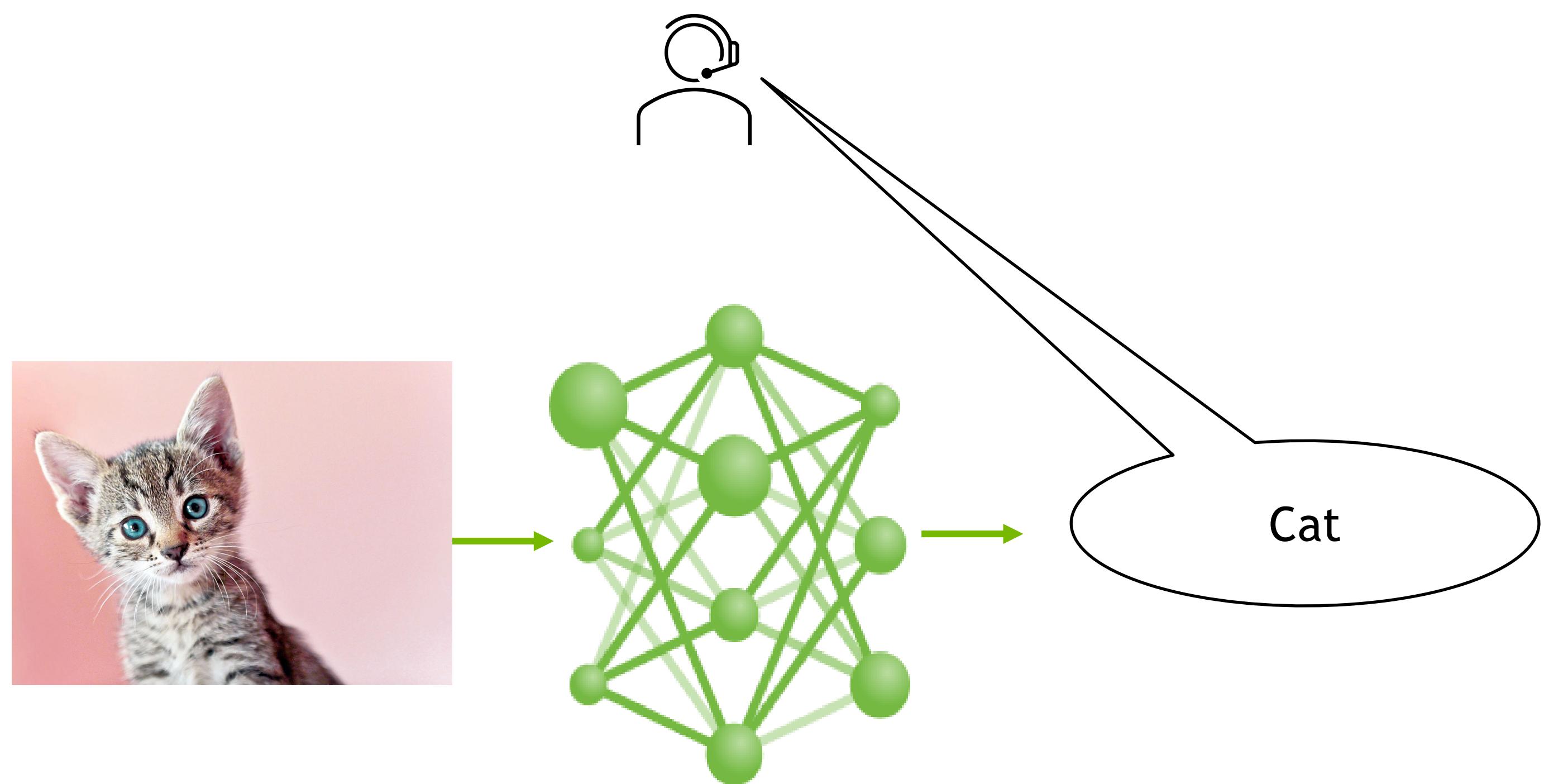


- Predictions of a Non-autoregressive model at the current time step does not depend on the predictions of the previous time-steps
- Fast at inference time
- Called also parallel

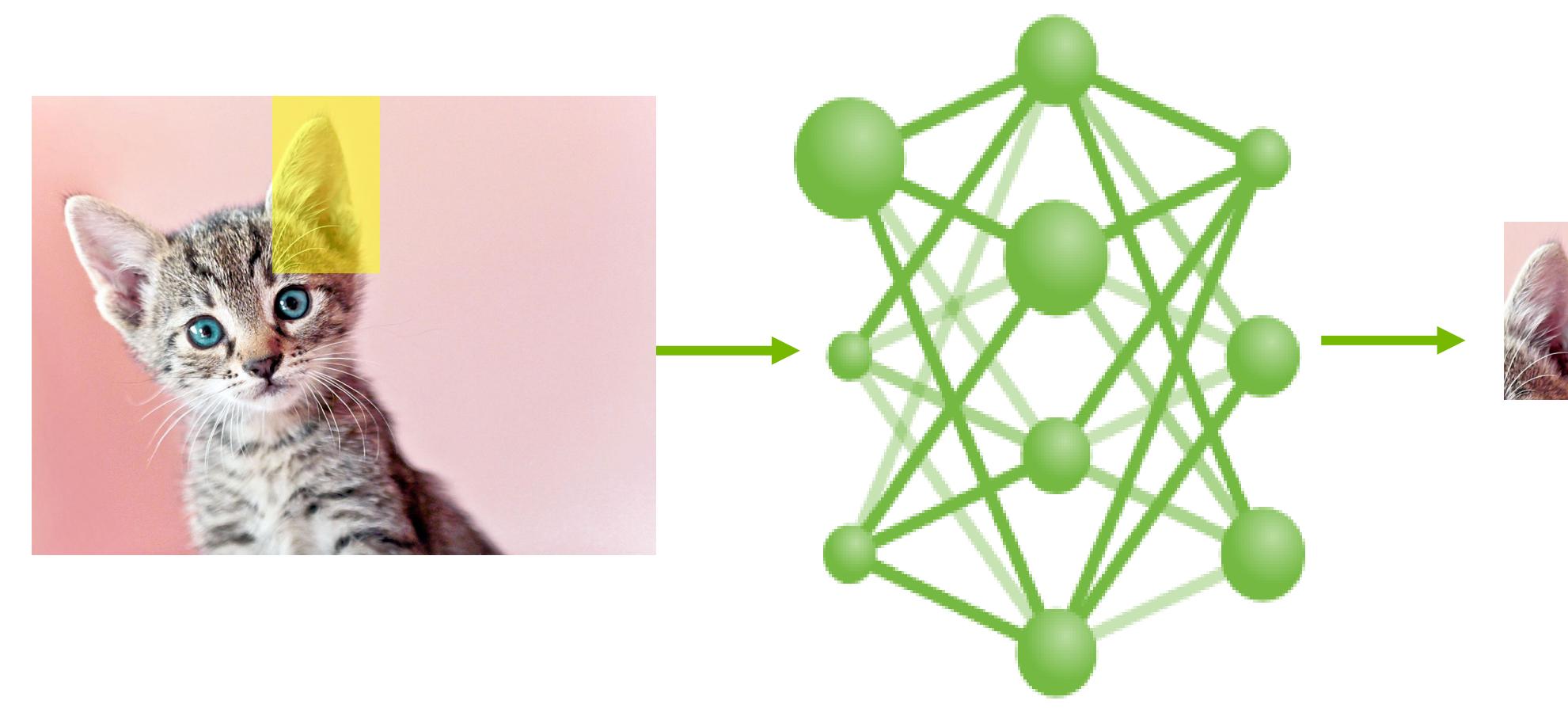
- Predictions of an autoregressive model at the current time step depends on the predictions of the previous time-steps
- Usually higher accuracy compared to non-autoregressive models due to temporal dependencies
- Slow at inference time

# Neural Networks

Supervised, Self-Supervised



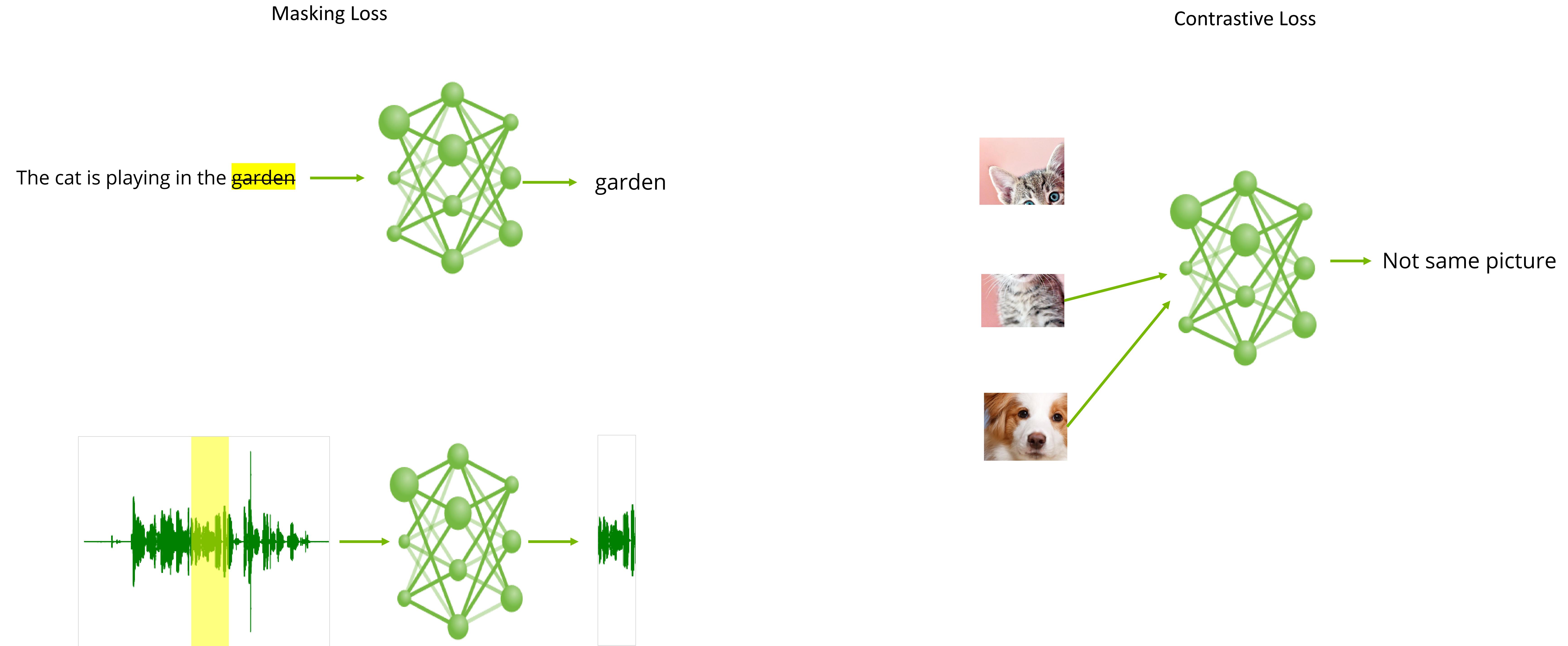
Supervised Training



Self-Supervised Training

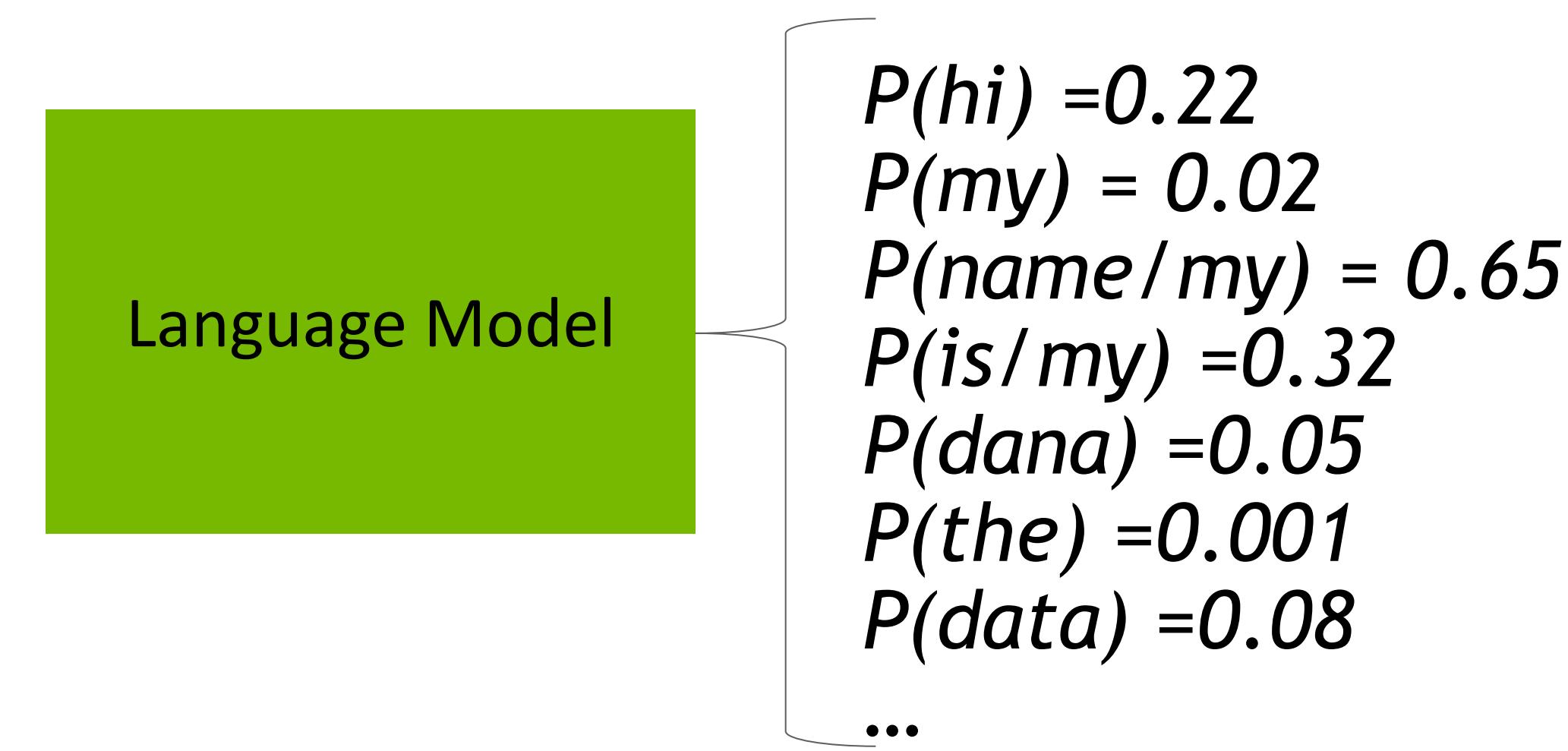
# Neural Networks

## Self-Supervised



# Neural Language Models

## Language Model



$P(hi \ my \ name \ is \ dana) ?$

Probability of a sentence is decomposed into conditional probabilities of each word given its previous context:

$$P(hi \ my \ name \ is \ dana) = p(hi) * p(my/hi) * p(name/hi my) * p(is/hi my name) * p(dana/hi my name is)$$



# ASR (Part 1)

## Lecture

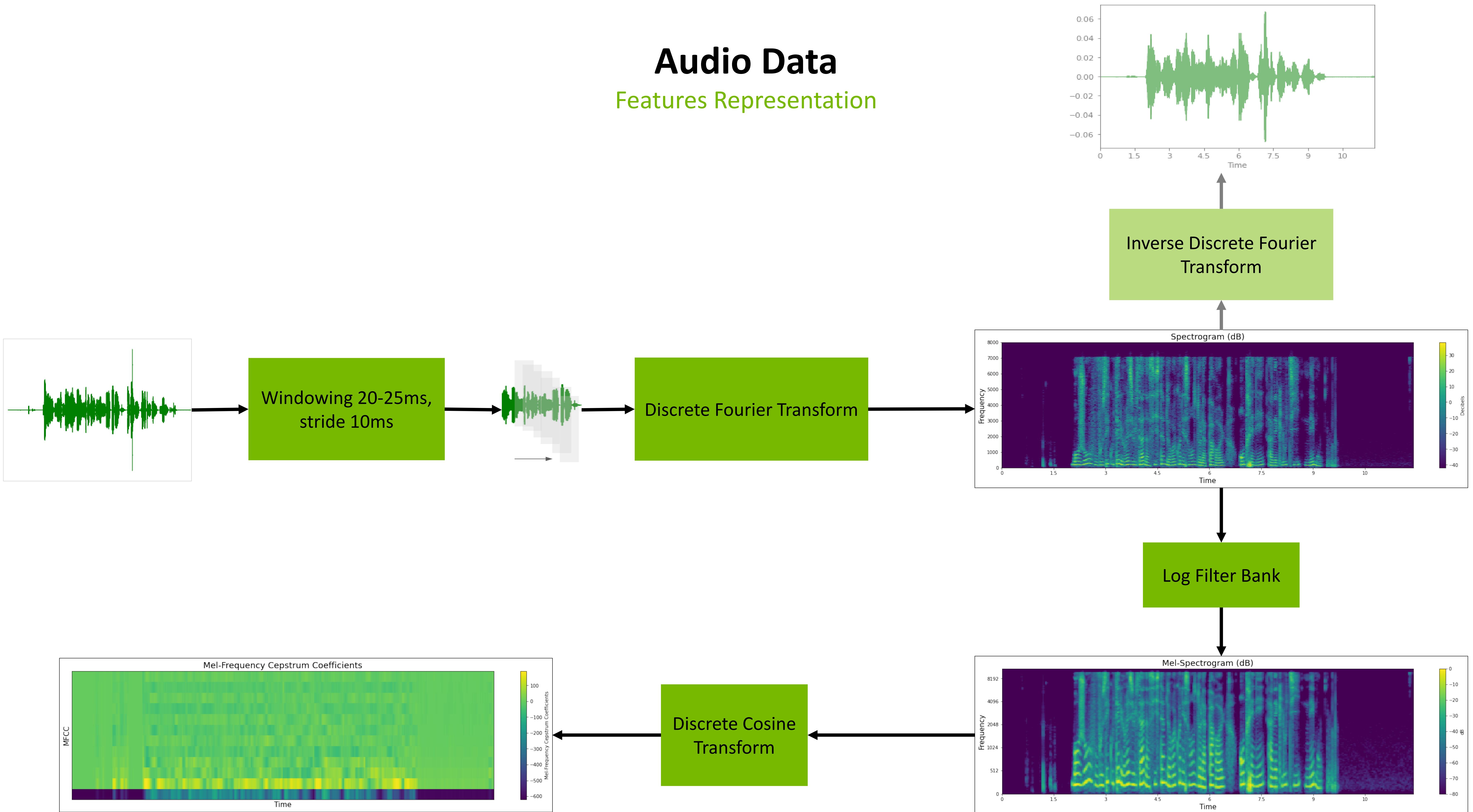
- What is Conversational AI? Complexity of Deploying Conversational AI Applications
- Deep Learning Fundamentals
  - Speech and Text Processing
- Automatic Speech Recognition
- ASR Tools: Riva, NGC, and NeMo
- Lab Overview

## Lab

- Investigate ASR Pipeline with NeMo and Riva

# Audio Data

## Features Representation



# Text Data

char/graphemes/words → phonemes/tokens

Text	Definition	Examples
Char	Alphabet Characters	a, b, c, d, ...n,...
Graphemes	Small unit of a <u>written</u> language	n, nn, kn, gn, pn
Words	Words	Hello, world, ..

Text Representation	Definition	Examples	Standard	Comments
Phonemes	Small unit of spoken <u>sound</u>	/n/	Grapheme to phonemes <u>International Phonetic Alphabet (IPA):</u> 'H', 'E', 'L', 'L', 'O', ' ', ' ', 'w', 'ə', 'l', 'd' <u>ARPABET:</u> HH, AH0, L, OW1, ,W, ER1, L, D	Non Phonetic Languages: <ul style="list-style-type: none"> <li>Multiple graphemes for a phoneme (eg. "f" "ph")</li> <li>Different pronunciations for same spelling (e.g. "I read book" in present and past)</li> </ul>
Tokens	Subwords	_use_ful	Tokenizers: Byte-Pair Encoding, WordPiece or SentencPiece <u>bert-base-uncased:</u> 'i', 'work', 'for', 'n', '#vid', '#ia' 1045, 2147, 2005, 1050, 17258, 2401 <u>xlnet-base-cased:</u> ' ', 'i', '_work', '_for', ' ', 'n', 'vid', 'ia' 17, 150, 154, 28, 17, 180, 7791, 780	<ul style="list-style-type: none"> <li>ML Algorithms deal with indexes instead of tokens</li> <li>Tokenizers uses a fixed size of vocabulary</li> <li>Out of Vocabulary OOV:               <ul style="list-style-type: none"> <li>Chars are part of the vocabulary</li> <li>Special token UNK</li> </ul> </li> </ul>



# ASR (Part 1)

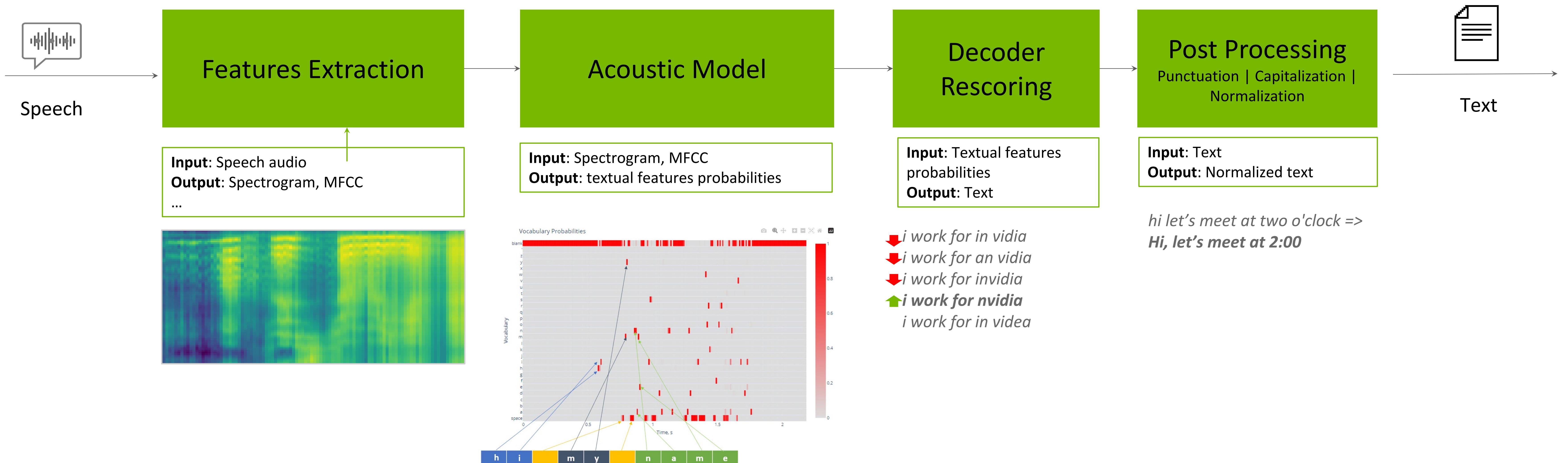
## Lecture

- What is Conversational AI? Complexity of Deploying Conversational AI Applications
- Deep Learning Fundamentals
- Speech and Text Processing
- **Automatic Speech Recognition**
- ASR Tools: Riva, NGC, and NeMo
- Lab Overview

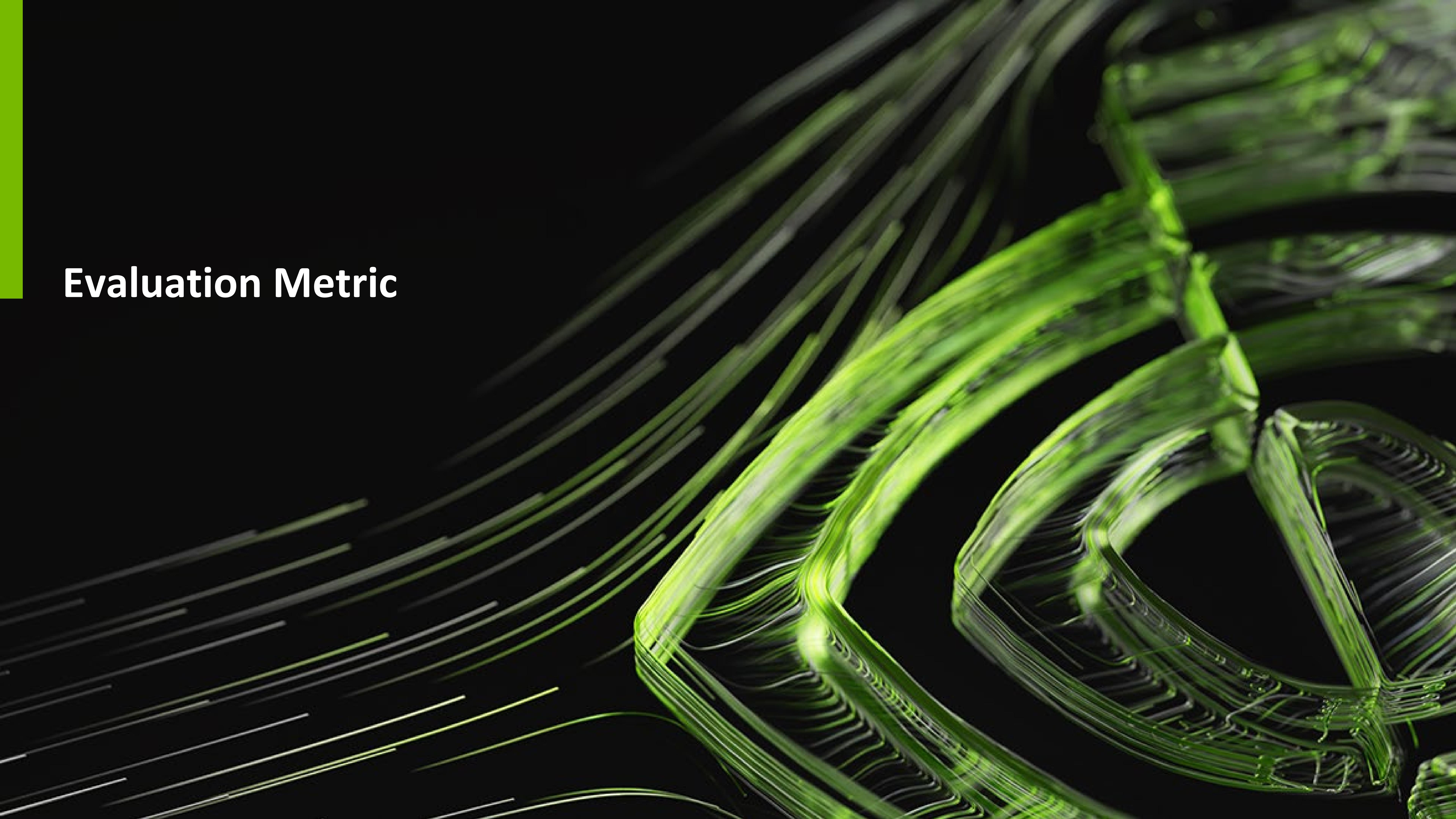
## Lab

- Investigate ASR Pipeline with NeMo and Riva

# ASR Pipeline



# Evaluation Metric



# Measuring Success of ASR

WER | CER

$$\text{Word Error Rate (WER)} = \frac{S + D + I}{N}$$

- $S$ : Number of words substituted
- $I$ : Number of words inserted
- $D$  Number of words deleted
- $N$ : Number of words in the reference

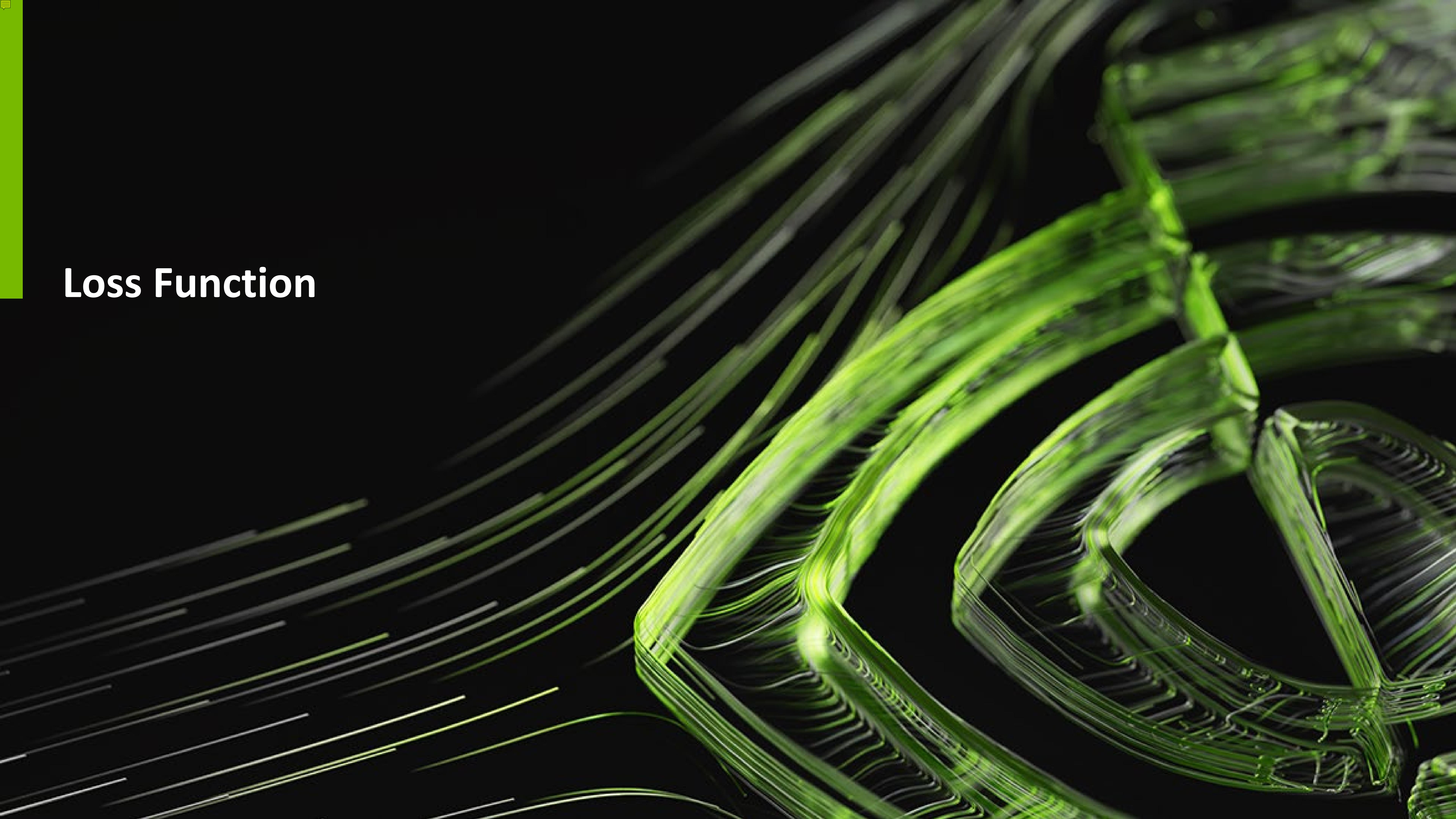
Reference	the <b>cat</b> sat on <b>the</b> mat
Transcription	the <b>cats</b> at on mat yeah
WER=0.66	D=1, S=2, I=1, N=6

$$\text{Character Error Rate (CER)} = \frac{S + D + I}{N}$$

- $S$ : Number of chars substituted
- $I$ : Number of chars inserted
- $D$  Number of chars deleted
- $N$ : Number of chars in the reference

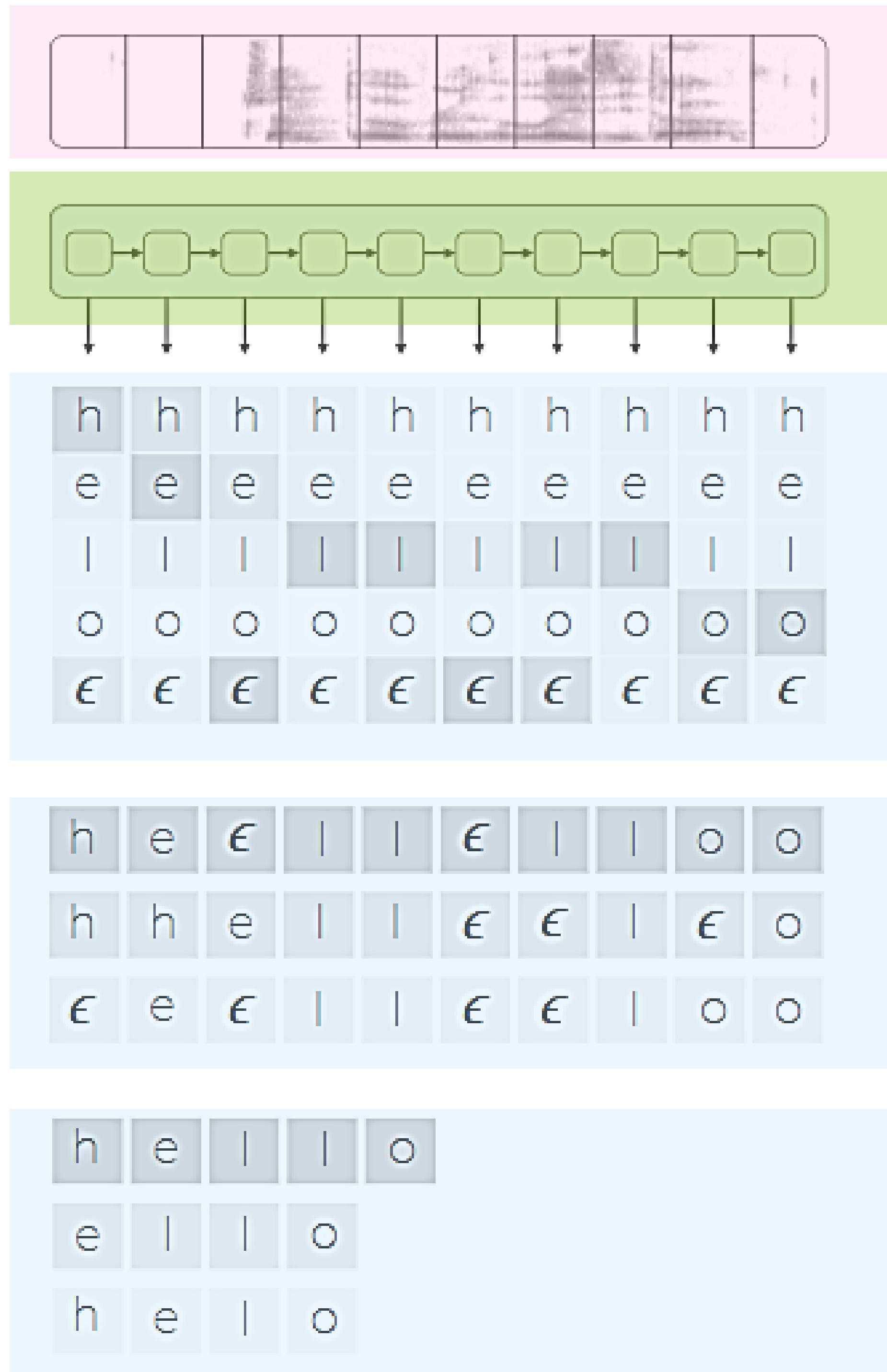
Reference	't', 'h', 'e', ' ', 'c', 'a', 't', ' ', 's', 'a', 't', ' ', 'o', 'n', ' ', 't', 'h', 'e', ' ', 'm', 'a', 't'
Transcription	't', 'h', 'e', ' ', 'c', 'a', 't', ' ', 's', 'a', 't', ' ', 'o', 'n', ' ', 'n', ' ', 'a', ' ', 't'
CER=0.2272	D=4, S=1, I=0, N=22

# Loss Function



# Acoustic Models

## Connectionist Temporal Classification (CTC)



We start with an input sequence,  
like a spectrogram of audio.

The input is fed into an RNN,  
for example.

The network gives  $p_t(a | X)$ ,  
a distribution over the outputs  
 $\{h, e, l, o, \epsilon\}$  for each input step.

With the per time-step output  
distribution, we compute the  
probability of different sequences

By marginalizing over alignments,  
we get a distribution over outputs.

To be precise, the CTC objective for a single  $(X, Y)$  pair is:

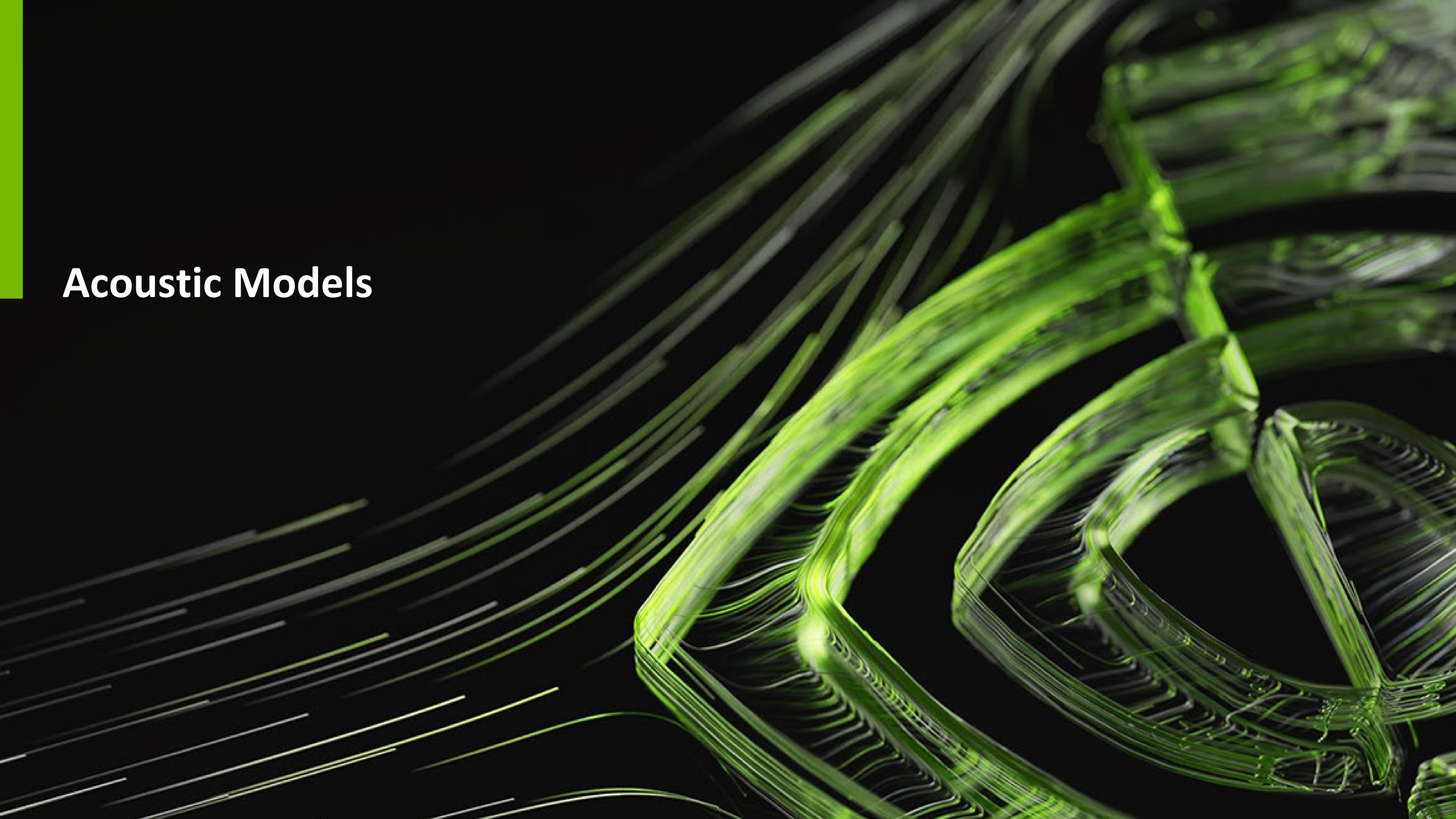
$$p(Y | X) = \sum_{A \in \mathcal{A}_{X,Y}} \prod_{t=1}^T p_t(a_t | X)$$

The CTC conditional  
probability

marginalizes over the  
set of valid alignments

computing the probability for a  
single alignment step-by-step.

# Acoustic Models



# Acoustic Models

RNN-based approaches: DeepSpeech and DeepSpeech2

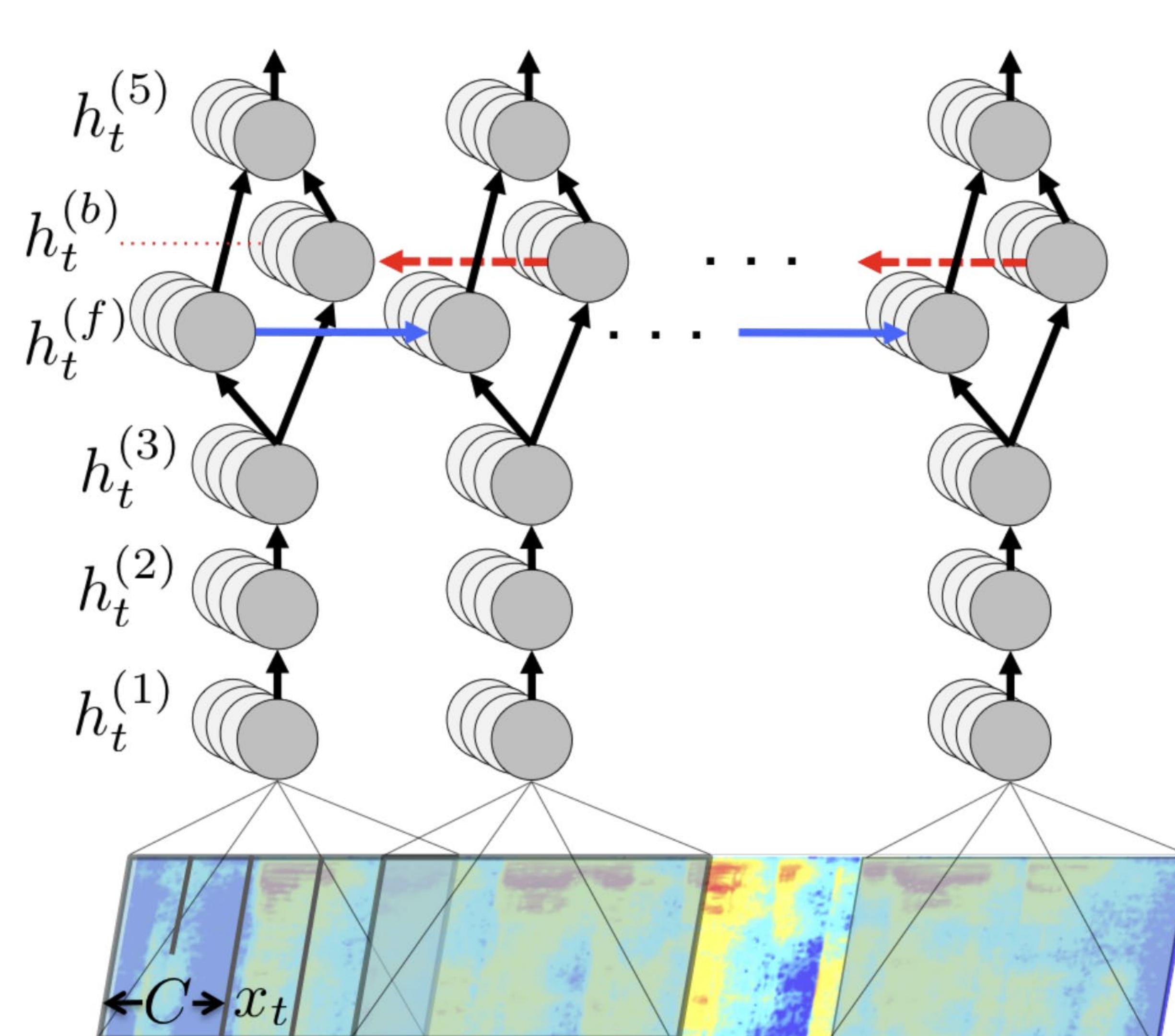


Figure 1: Structure of our RNN model and notation.

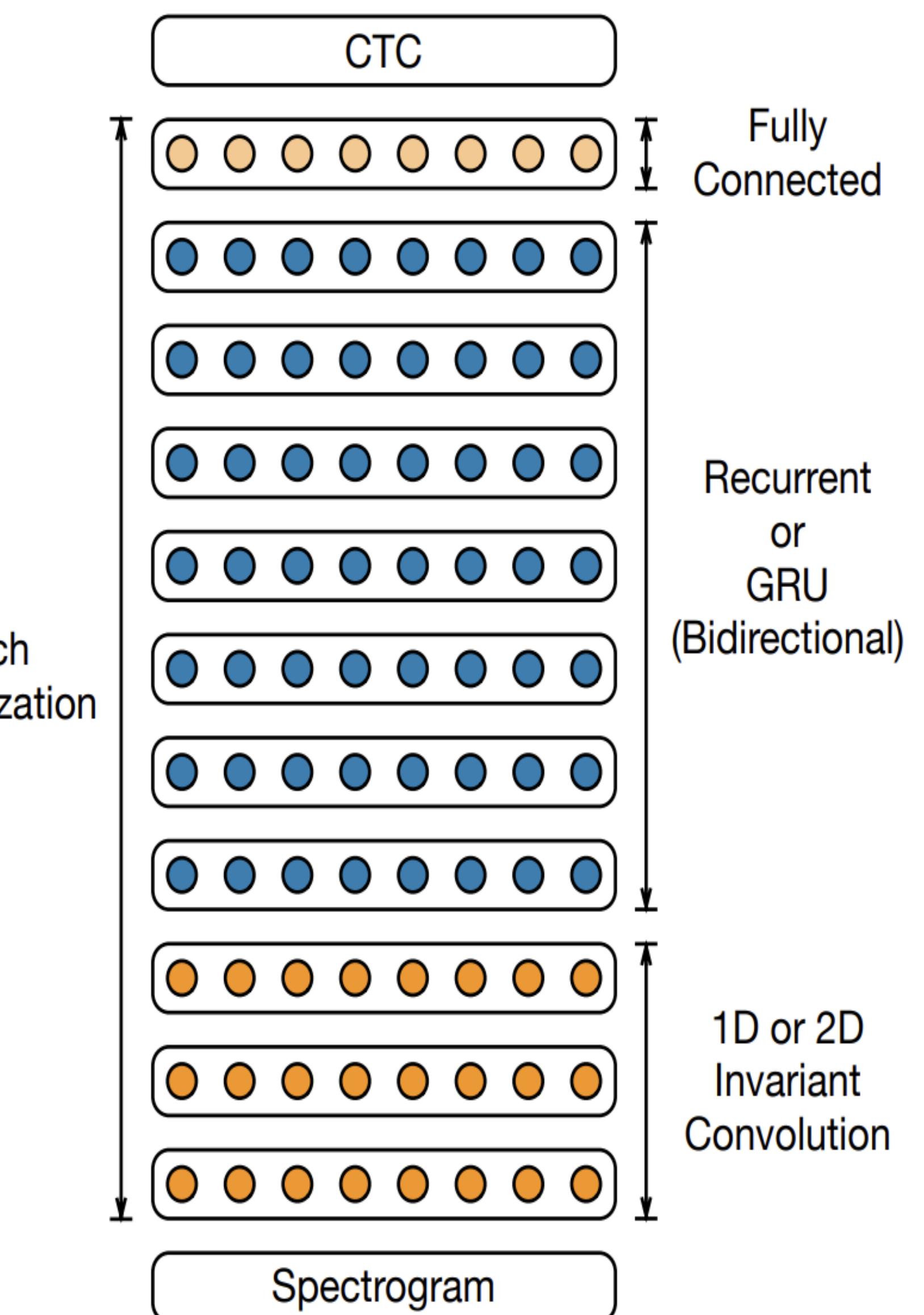


Figure 1: Architecture of the DS2 system used to train on both English and Mandarin speech. We explore variants of this architecture by varying the number of convolutional layers from 1 to 3 and the number of recurrent or GRU layers from 1 to 7.

# Acoustic Models

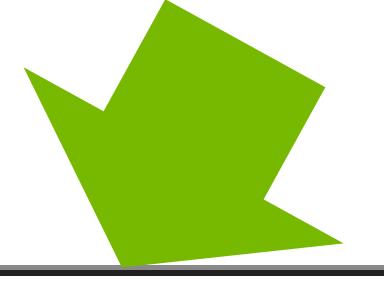
## RNN-based approaches: DeepSpeech and DeepSpeech2

RNN output	Decoded Transcription
what is the weather like in bostin right now	what is the weather like in boston right now
prime miniter nerernr modi	prime minister narendra modi
arther n tickets for the game	are there any tickets for the game

Table 1: Examples of transcriptions directly from the RNN (left) with errors that are fixed by addition of a language model (right).

System	Clean (94)	Noisy (82)	Combined (176)
Apple Dictation	14.24	43.76	26.73
Bing Speech	11.73	36.12	22.05
Google API	6.64	30.47	16.72
wit.ai	7.94	35.06	19.41
<b>Deep Speech</b>	<b>6.56</b>	<b>19.06</b>	<b>11.85</b>

Table 4: Results (%WER) for 5 systems evaluated on the original audio. Scores are reported *only* for utterances with predictions given by all systems. The number in parentheses next to each dataset, e.g. Clean (94), is the number of utterances scored.



GPU	OpenMPI all-reduce	Our all-reduce	Performance Gain
4	55359.1	2587.4	21.4
8	48881.6	2470.9	19.8
16	21562.6	1393.7	15.5
32	8191.8	1339.6	6.1
64	1395.2	611.0	2.3
128	1602.1	422.6	3.8

Table 7: Comparison of two different all-reduce implementations. All times are in seconds. Performance gain is the ratio of OpenMPI all-reduce time to our all-reduce time.

Language	Architecture	CPU CTC Time	GPU CTC Time	Speedup
English	5-layer, 3 RNN	5888.12	203.56	28.9
Mandarin	5-layer, 3 RNN	1688.01	135.05	12.5

Table 8: Comparison of time spent in seconds in computing the CTC loss function and gradient in one epoch for two different implementations. Speedup is the ratio of CPU CTC time to GPU CTC time.

# Acoustic Models

## Simplifying with 1D convolutional layers: Wave2Letter

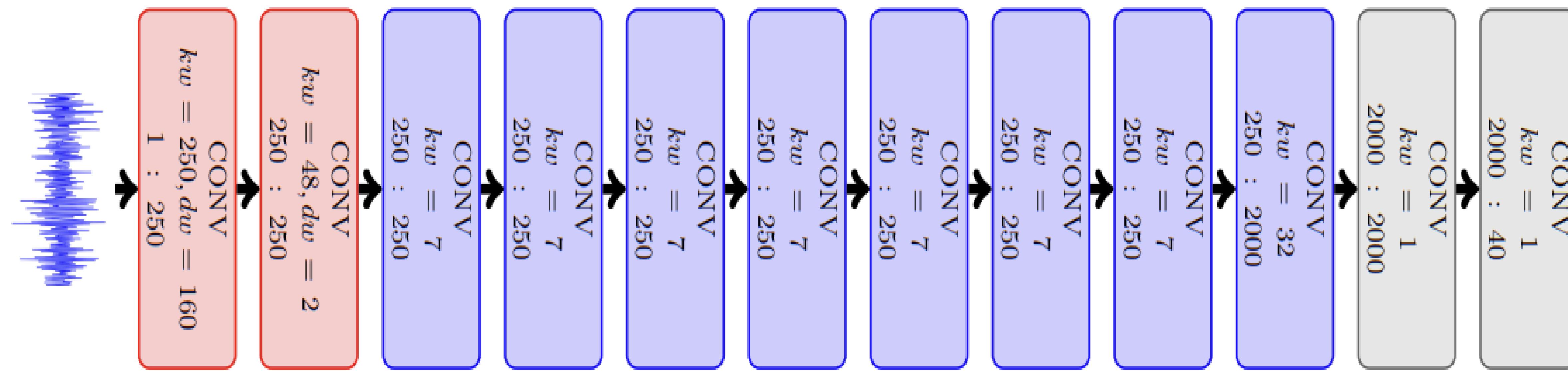
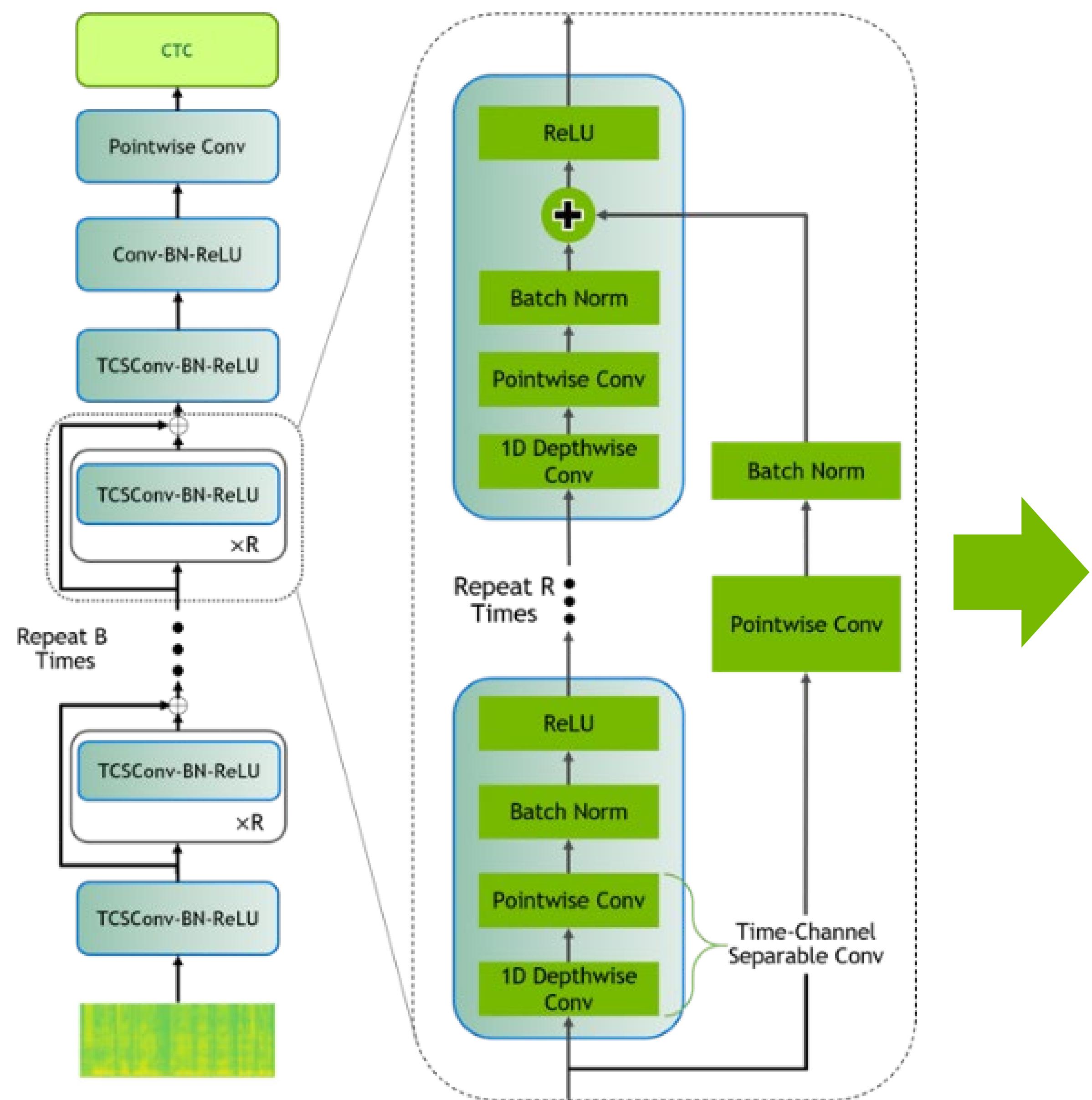


Table 2: LER/WER of the best sets of hyper-parameters for each feature types.

	MFCC		PS		Raw	
	LER	WER	LER	WER	LER	WER
dev-clean	6.9		9.3		10.3	
test-clean	6.9	7.2	9.1	9.4	10.6	10.1

# Acoustic Models

## QuartzNet



## Citrinet

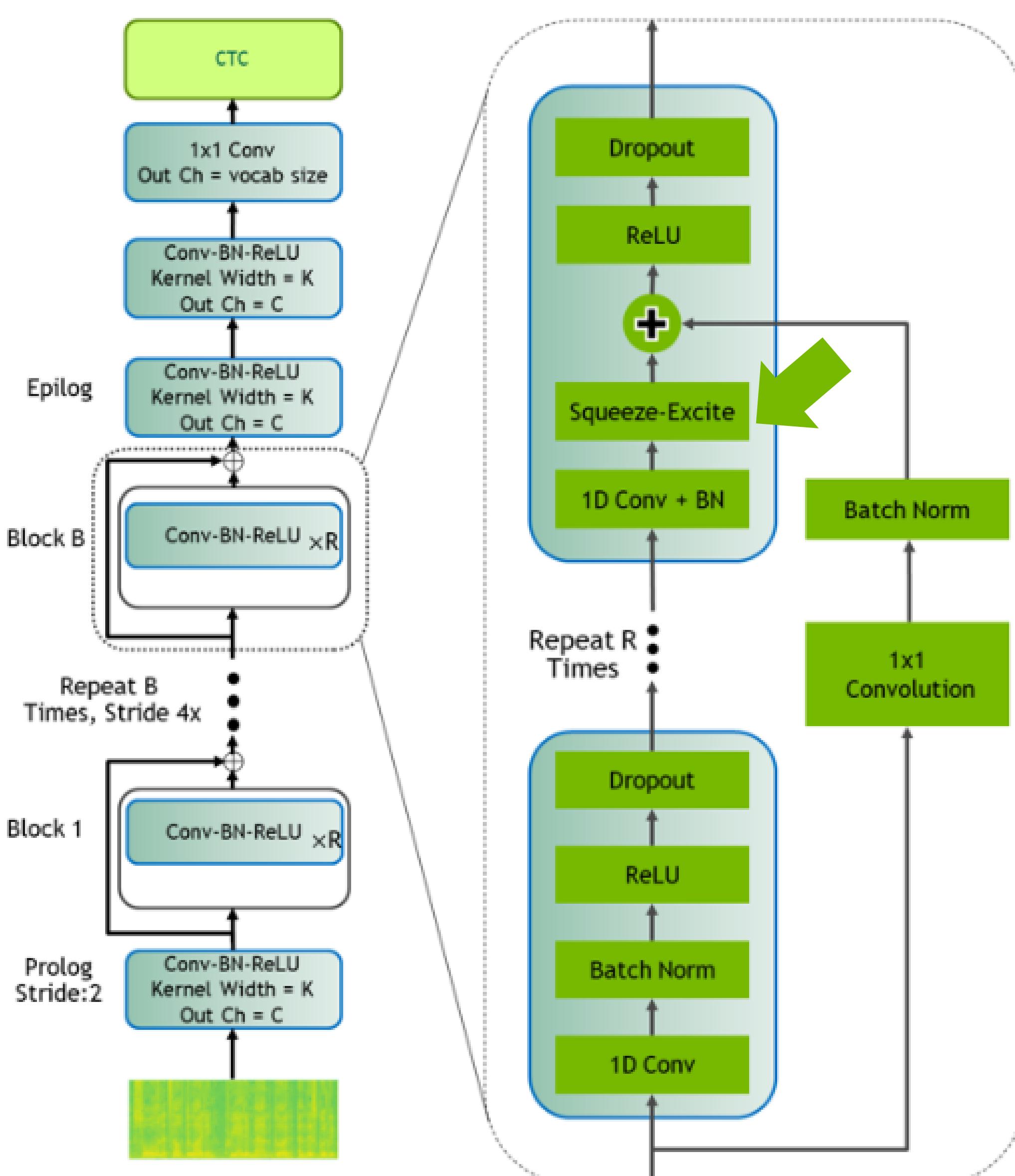


Table 3: *LibriSpeech: Citrinet vs Transducers, WER(%)*

Model	LM	Test clean	Test other	Params, M
ContextNet-L [14]	-	2.10	4.60	112.7
	RNN	1.90	4.10	
Conformer-L[15]	-	2.10	4.30	118
	RNN	1.90	3.90	
Citrinet-256	-	3.78	9.6	
	6-gram Transf	3.65	8.06	9.8
Citrinet-384	-	3.20	7.90	
	6-gram Transf	2.94	6.71	21.0
Citrinet-512	-	3.11	7.82	
	6-gram Transf	2.40	6.08	36.5
Citrinet-768	-	2.57	6.35	
	6-gram Transf	2.15	5.11	81
Citrinet-1024	-	2.52	6.22	
	6-gram Transf	2.10	5.06	142
		2.00	4.69	

# Acoustic Models

## Conformer and Conformer-CTC

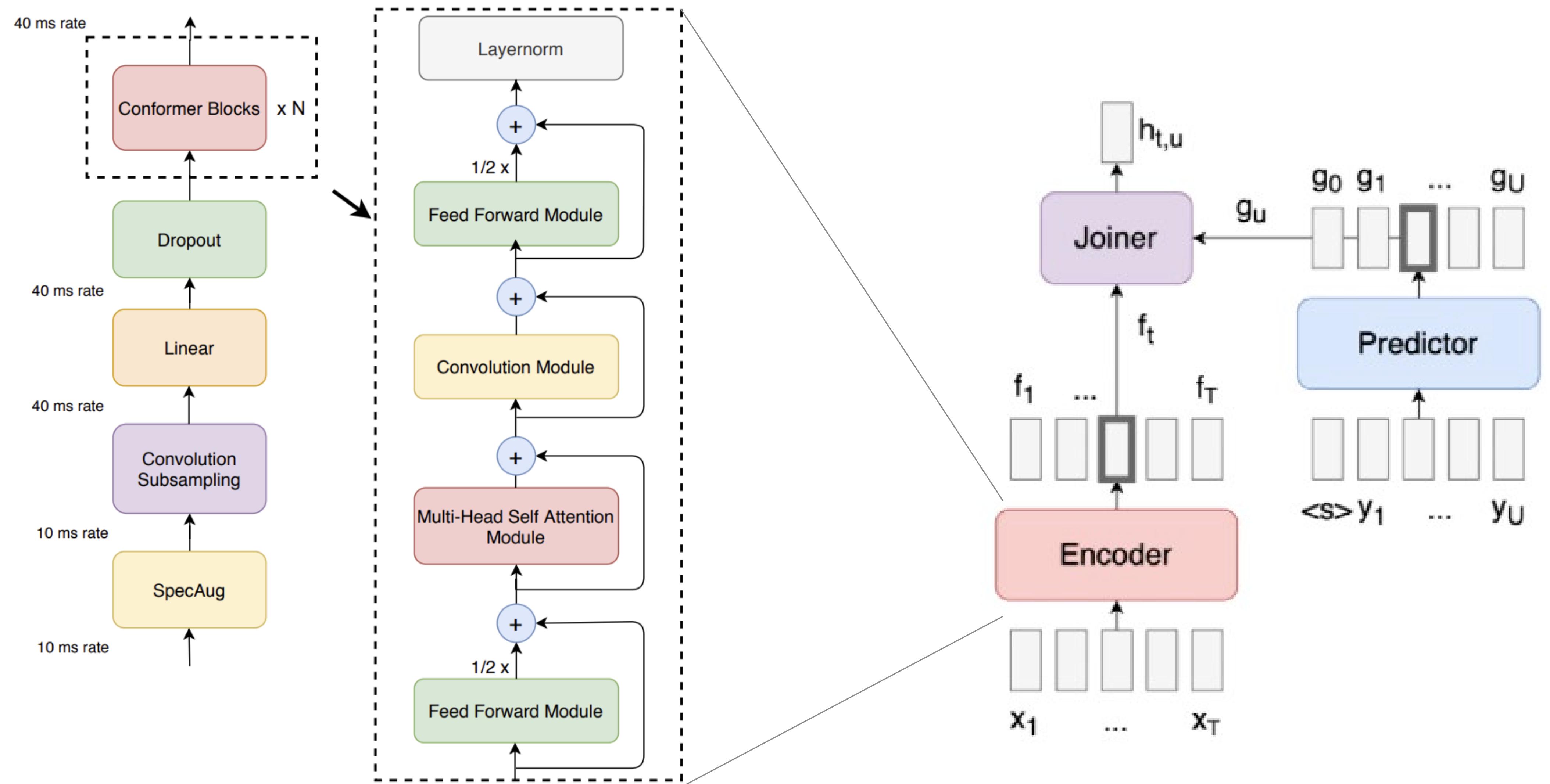


Figure 1: **Conformer encoder model architecture.** Conformer comprises of two macaron-like feed-forward layers with half-step residual connections sandwiching the multi-headed self-attention and convolution modules. This is followed by a post layernorm.

Method	#Params (M)	WER Without LM		WER With LM	
		testclean	testother	testclean	testother
<b>Hybrid</b>					
Transformer [33]	-	-	-	2.26	4.85
<b>CTC</b>					
QuartzNet [9]	19	3.90	11.28	2.69	7.25
<b>LAS</b>					
Transformer [34]	270	2.89	6.98	2.33	5.17
Transformer [19]	-	2.2	5.6	2.6	5.7
LSTM	360	2.6	6.0	2.2	5.2
<b>Transducer</b>					
Transformer [7]	139	2.4	5.6	2.0	4.6
ContextNet(S) [10]	10.8	2.9	7.0	2.3	5.5
ContextNet(M) [10]	31.4	2.4	5.4	<b>2.0</b>	4.5
ContextNet(L) [10]	112.7	<b>2.1</b>	4.6	<b>1.9</b>	4.1
<b>Conformer (Ours)</b>					
Conformer(S)	10.3	<b>2.7</b>	<b>6.3</b>	<b>2.1</b>	<b>5.0</b>
Conformer(M)	30.7	<b>2.3</b>	<b>5.0</b>	<b>2.0</b>	<b>4.3</b>
Conformer(L)	118.8	<b>2.1</b>	<b>4.3</b>	<b>1.9</b>	<b>3.9</b>

LSTM-layer decoder in all our models.

# Acoustic Models

## Conformer-CTC

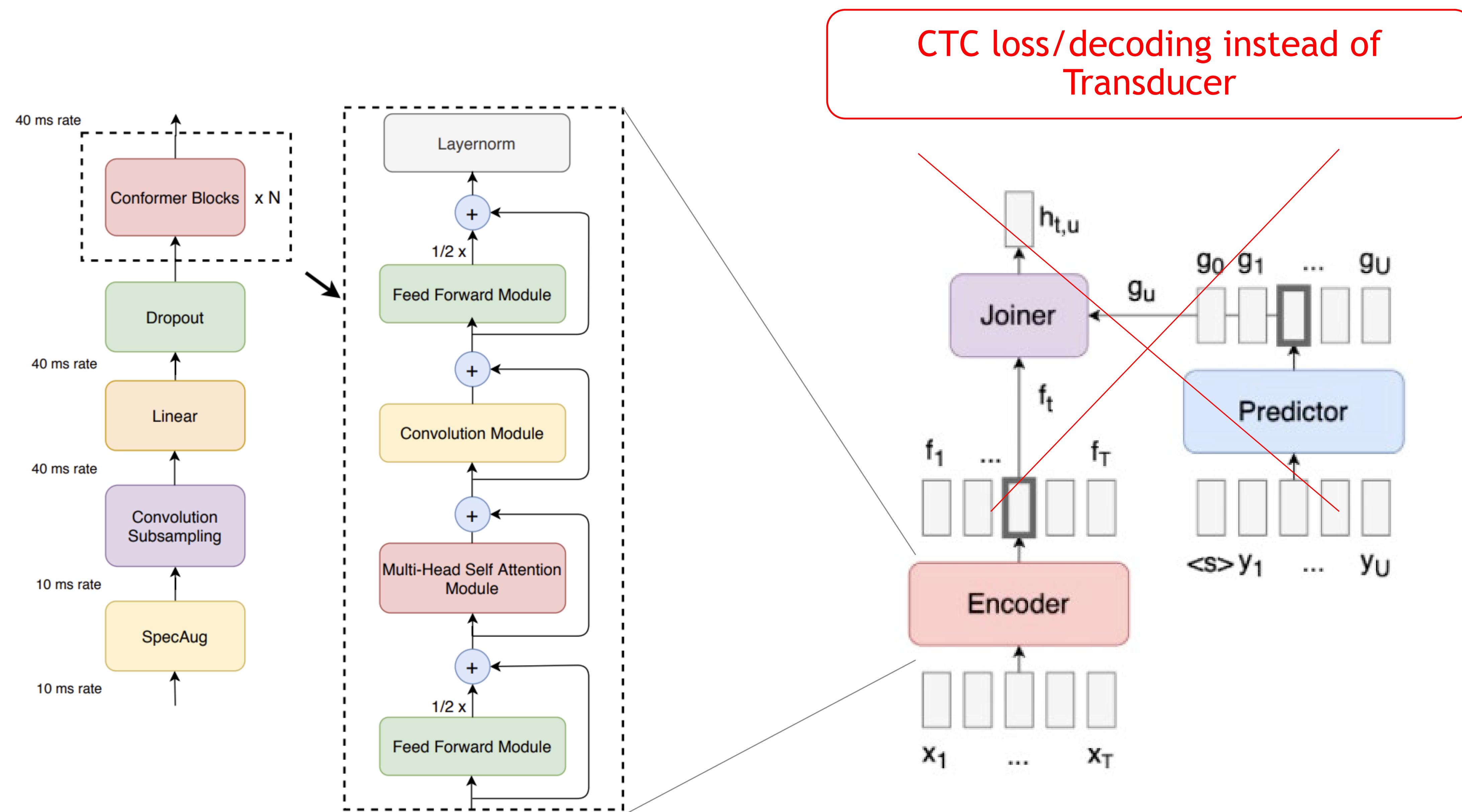
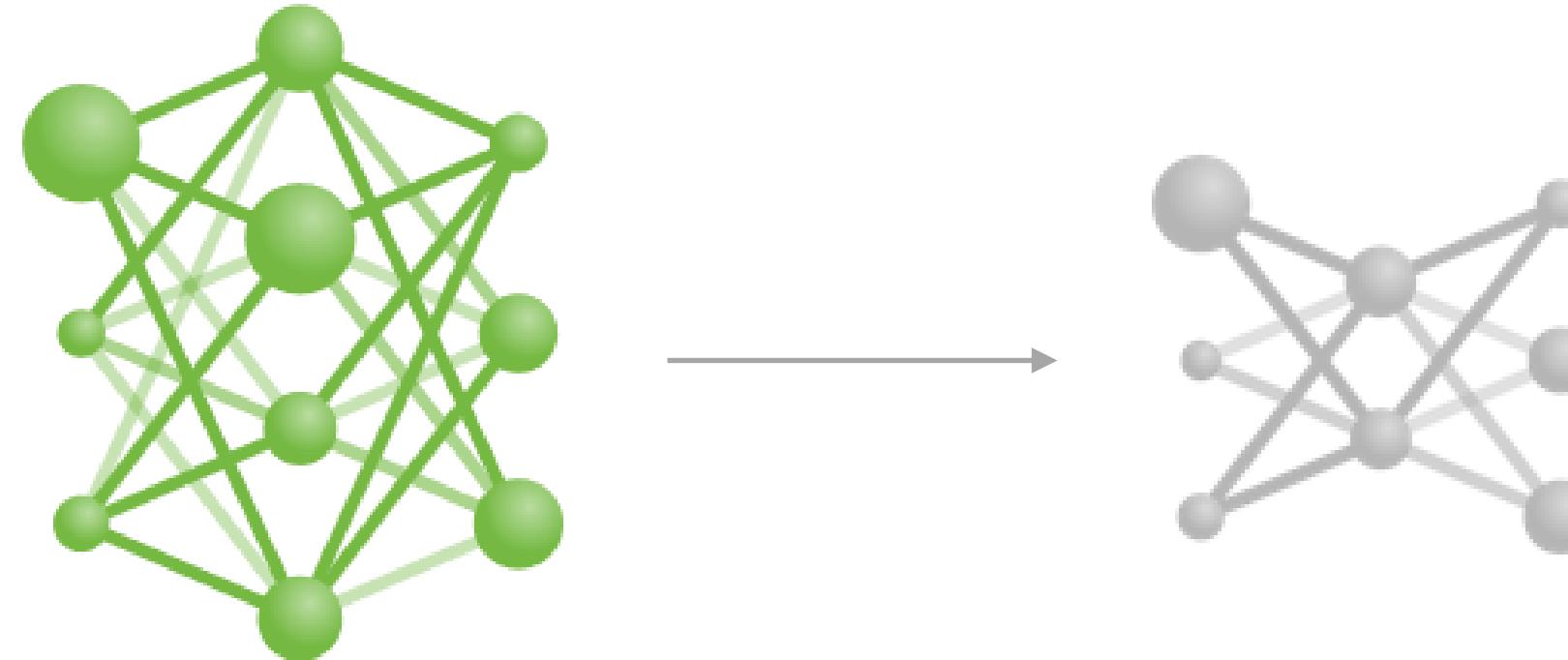


Figure 1: **Conformer encoder model architecture.** Conformer comprises of two macaron-like feed-forward layers with half-step residual connections sandwiching the multi-headed self-attention and convolution modules. This is followed by a post layernorm.

# Acoustic Models

## Transfer learning

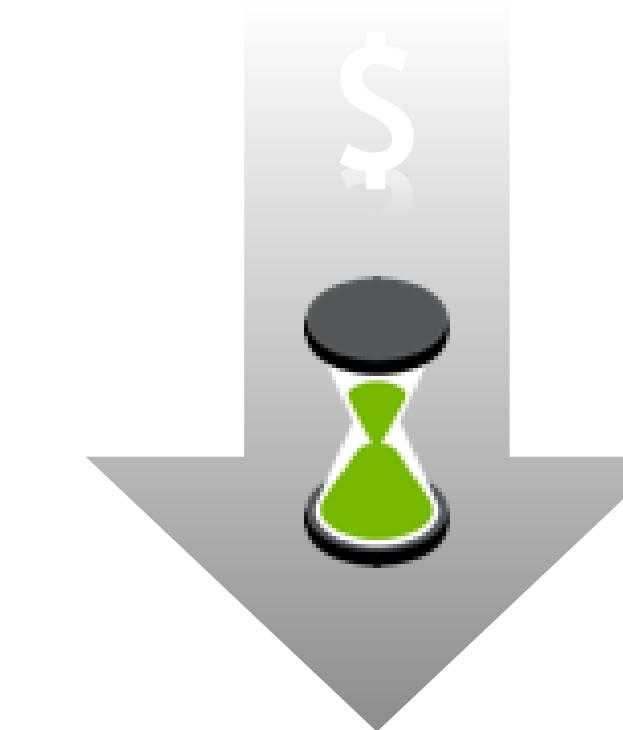
*“Transfer learning is a process of transferring learned features from one model to another”*



### Key Benefits



Less Data Required to Train Accurately

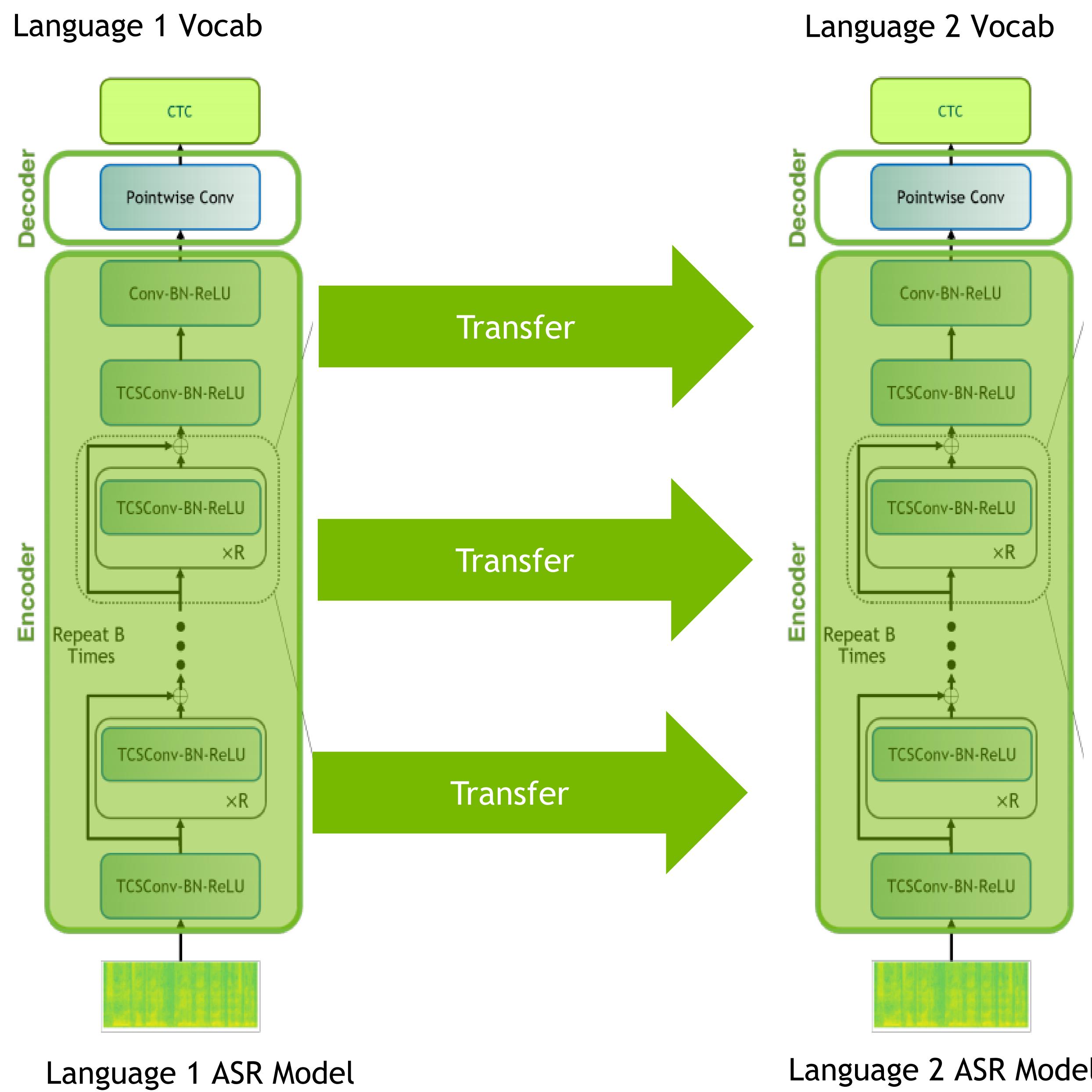


Reduce Training Time and Cost

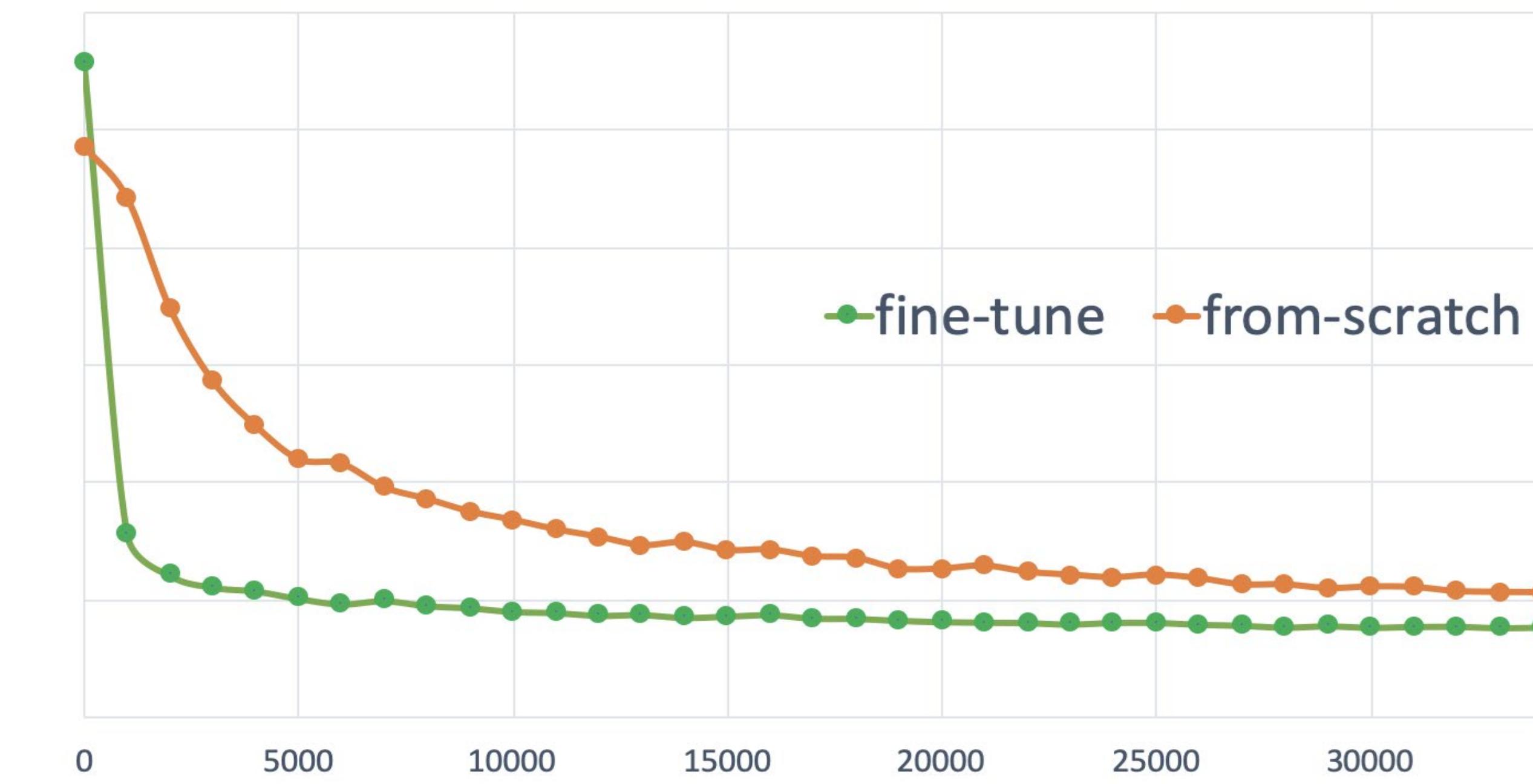
<https://blogs.nvidia.com/blog/2019/02/07/what-is-transfer-learning/>

# Acoustic Models

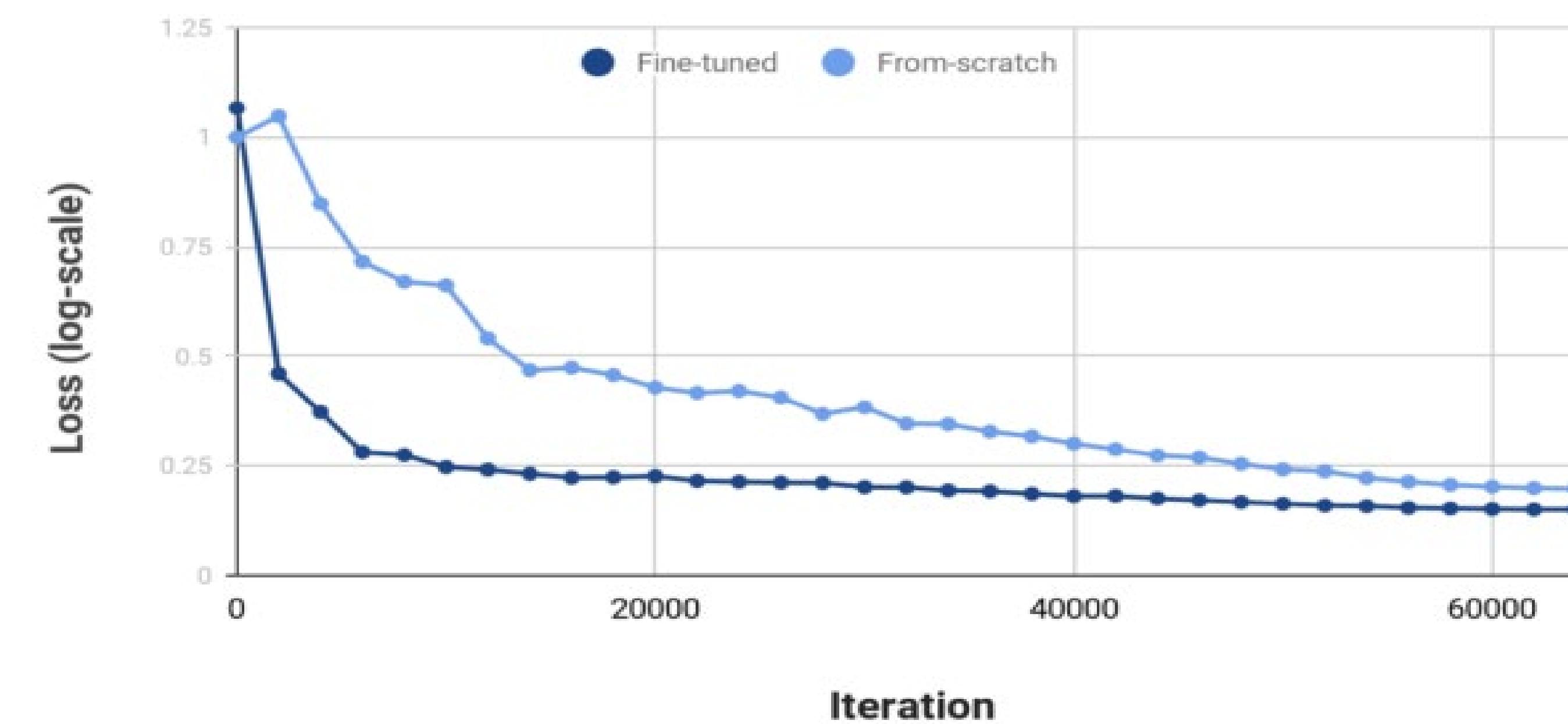
## Transfer learning



Evaluation WER vs Iteration

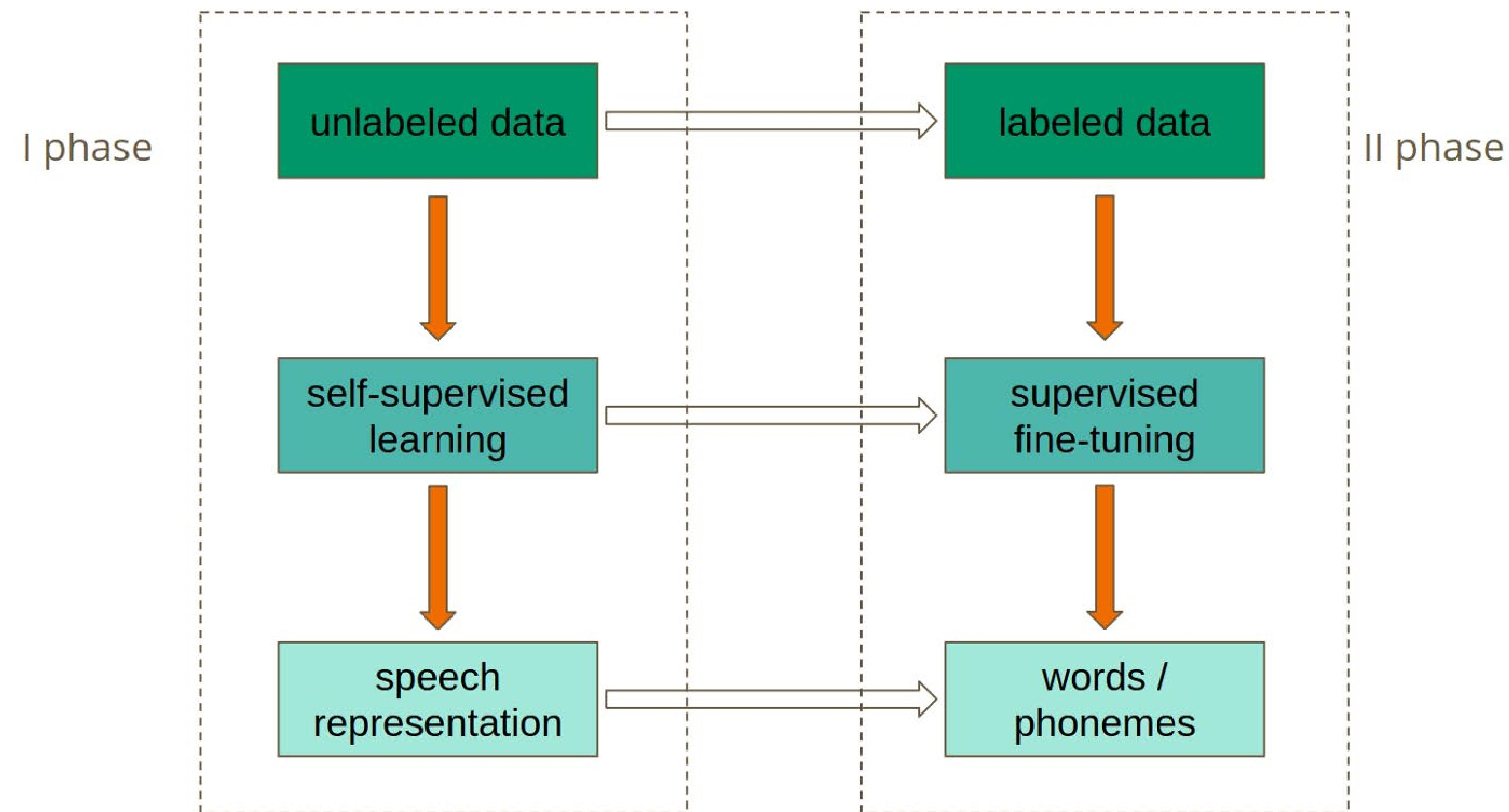


Evaluation WER on Spanish dataset



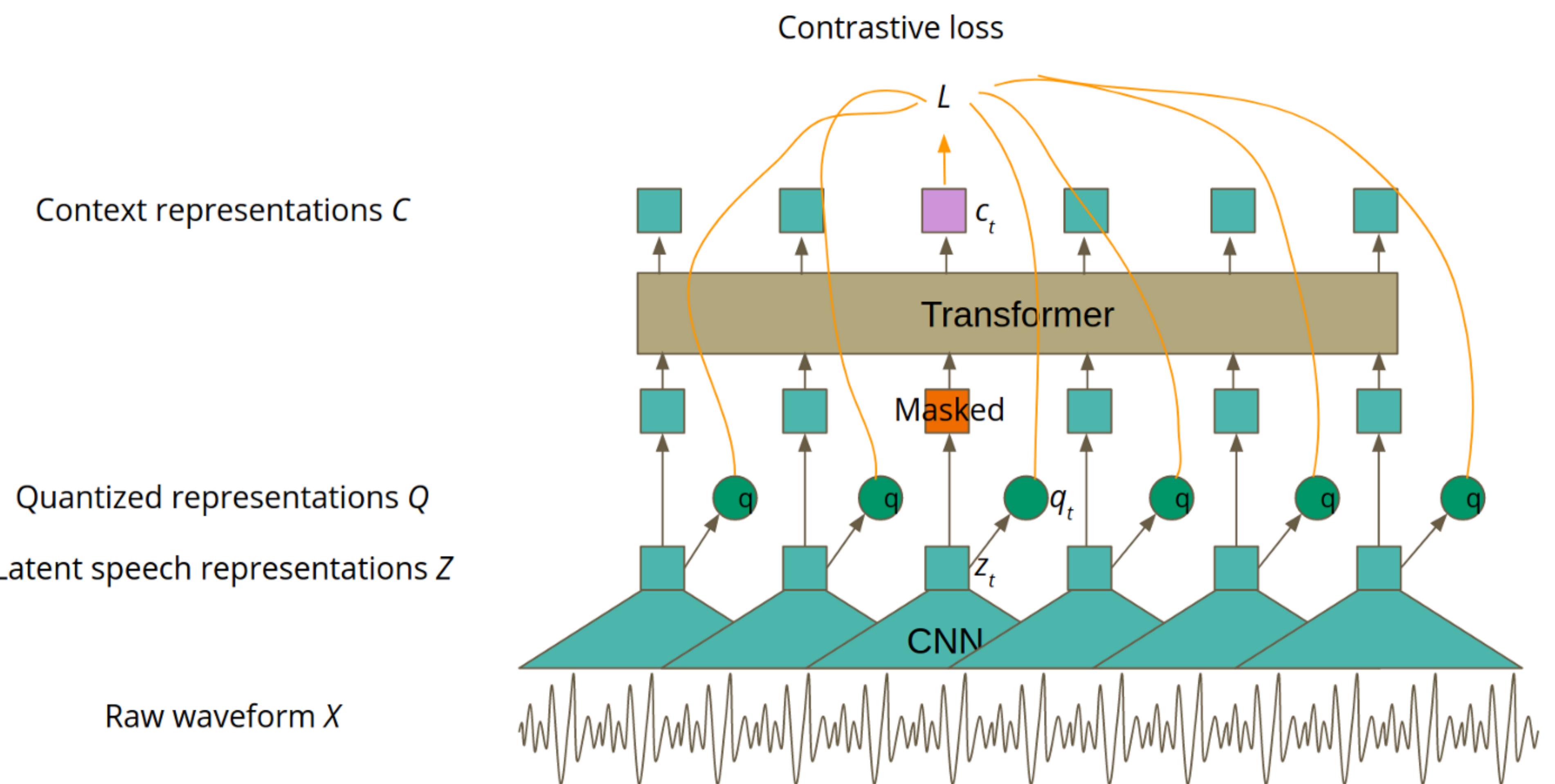
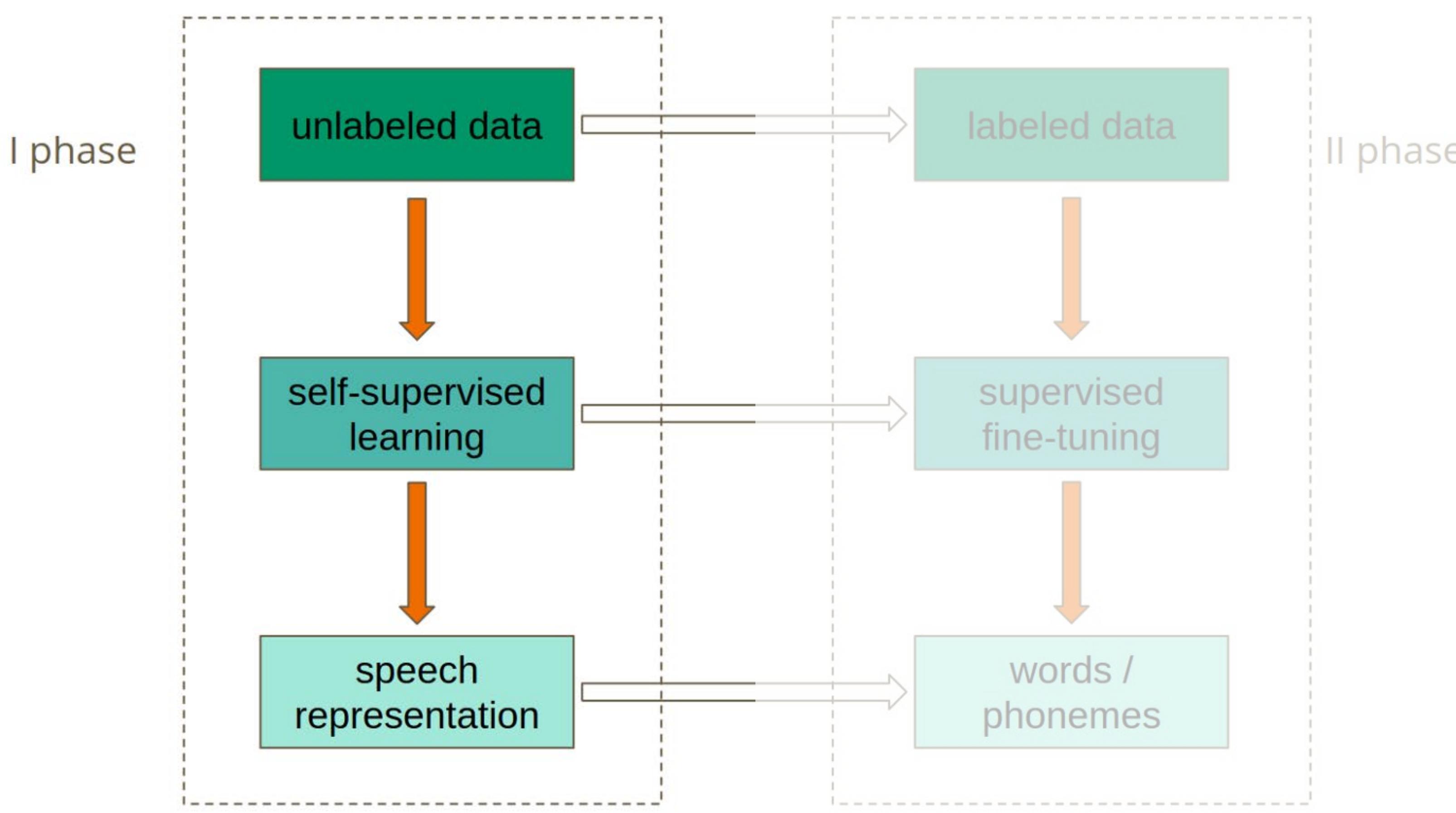
# Acoustic Models

Wave2vec 2.0



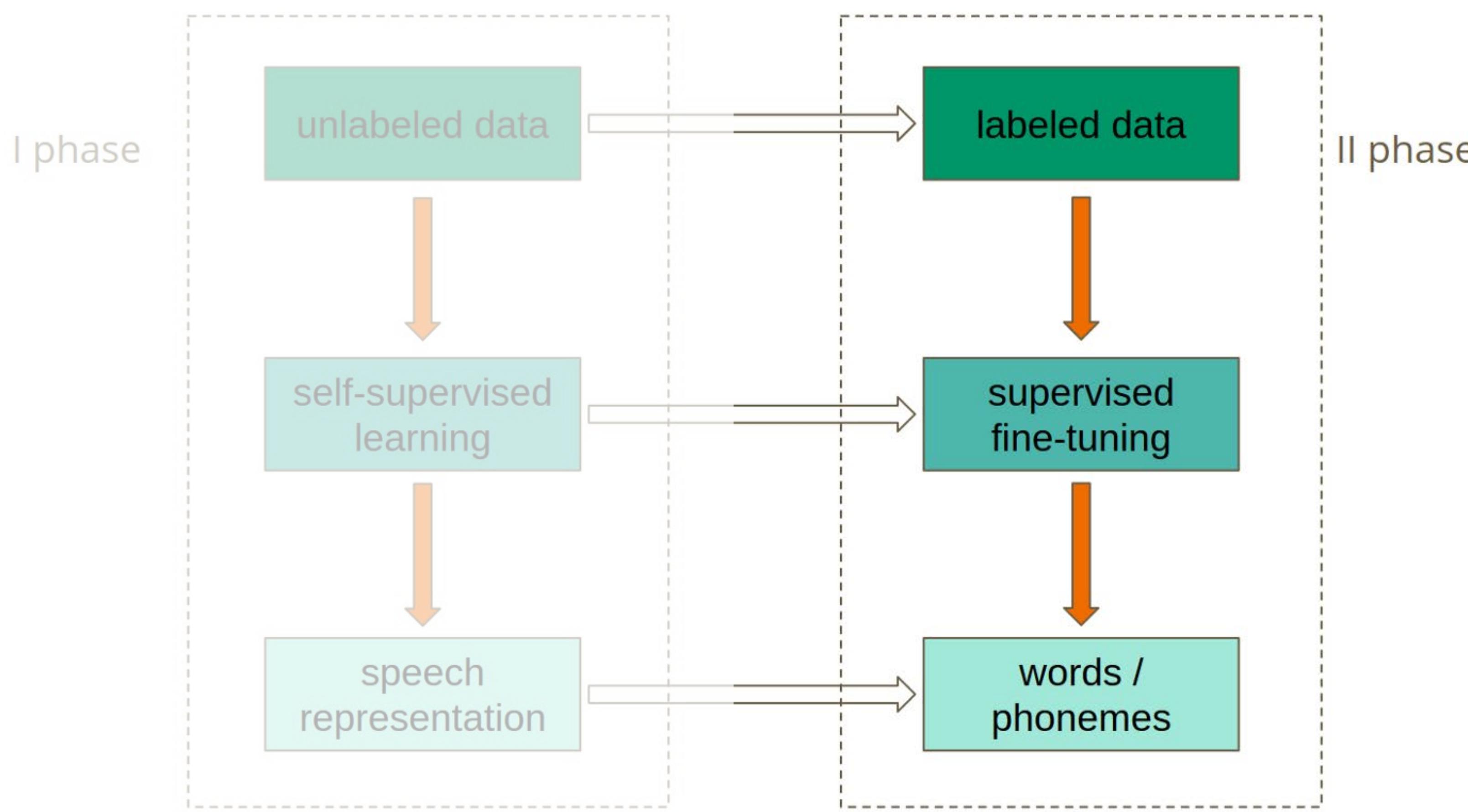
# Acoustic Models

Wave2vec 2.0

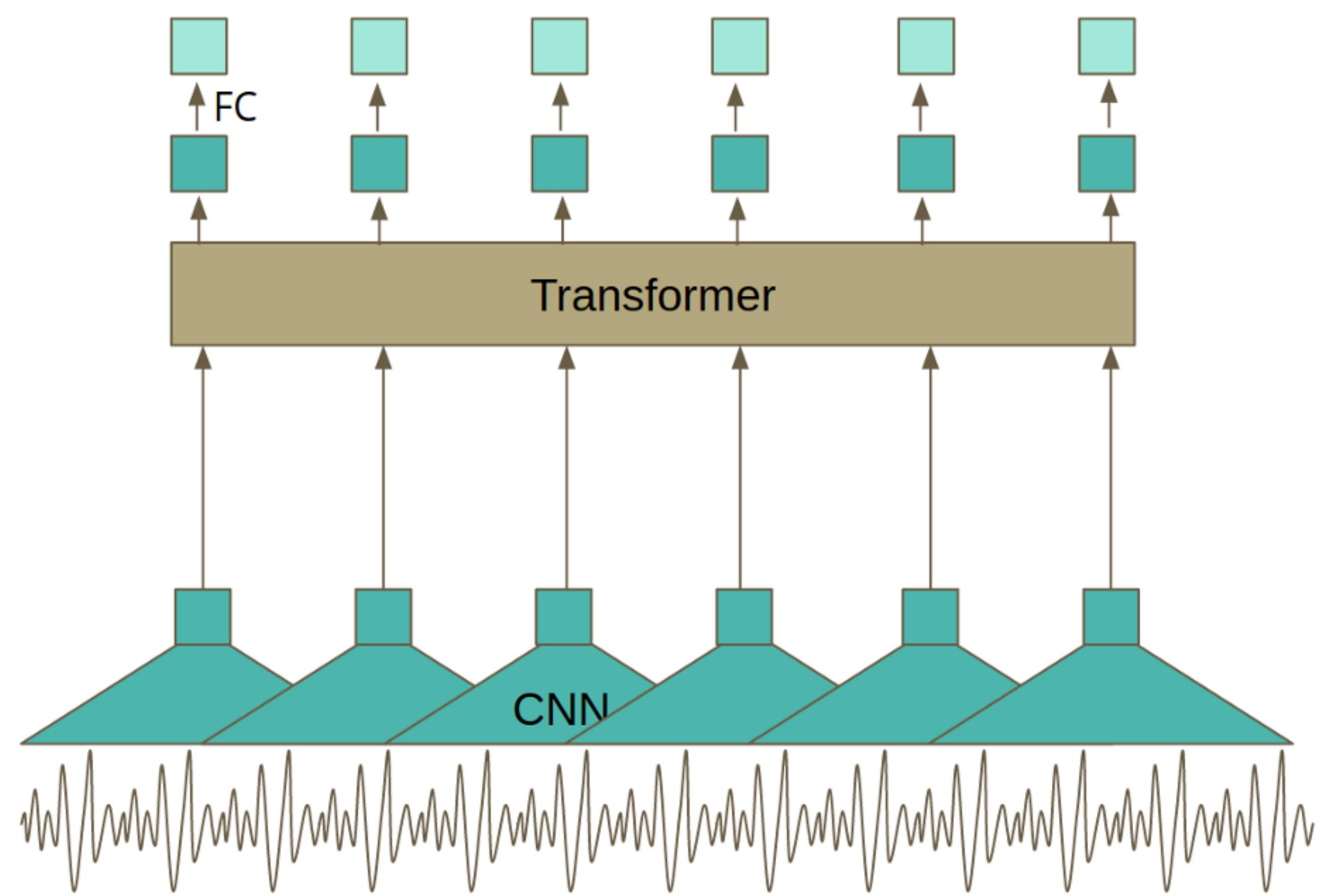


# Acoustic Models

Wave2vec 2.0



Output  $Y$   
Context representations  $C$   
Latent speech representations  $Z$   
Raw waveform  $X$



# Acoustic Models

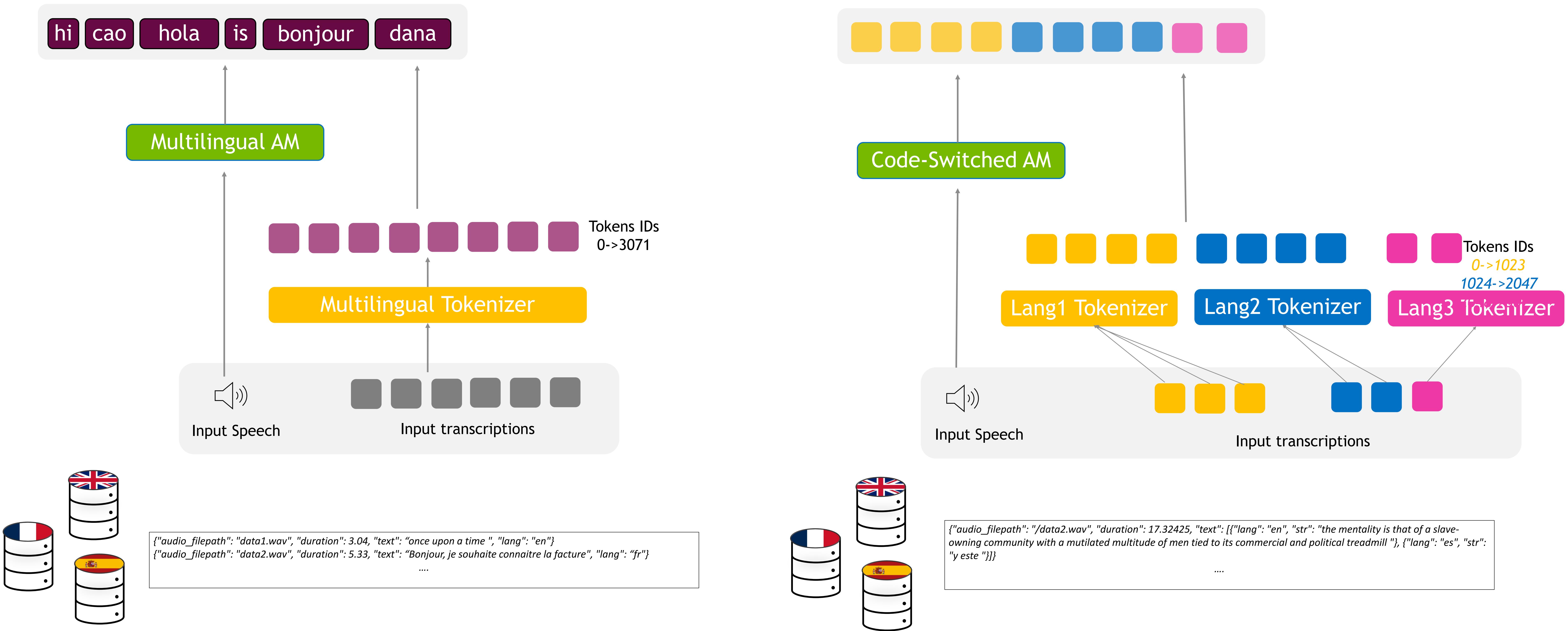
Wave2vec 2.0

*“This demonstrates the feasibility of speech recognition with limited amounts of labeled data.”*

*“When lowering the amount of labeled data to one hour, wav2vec 2.0 outperforms the previous state of the art on the 100 hour subset while using 100 times less labeled data. ”*

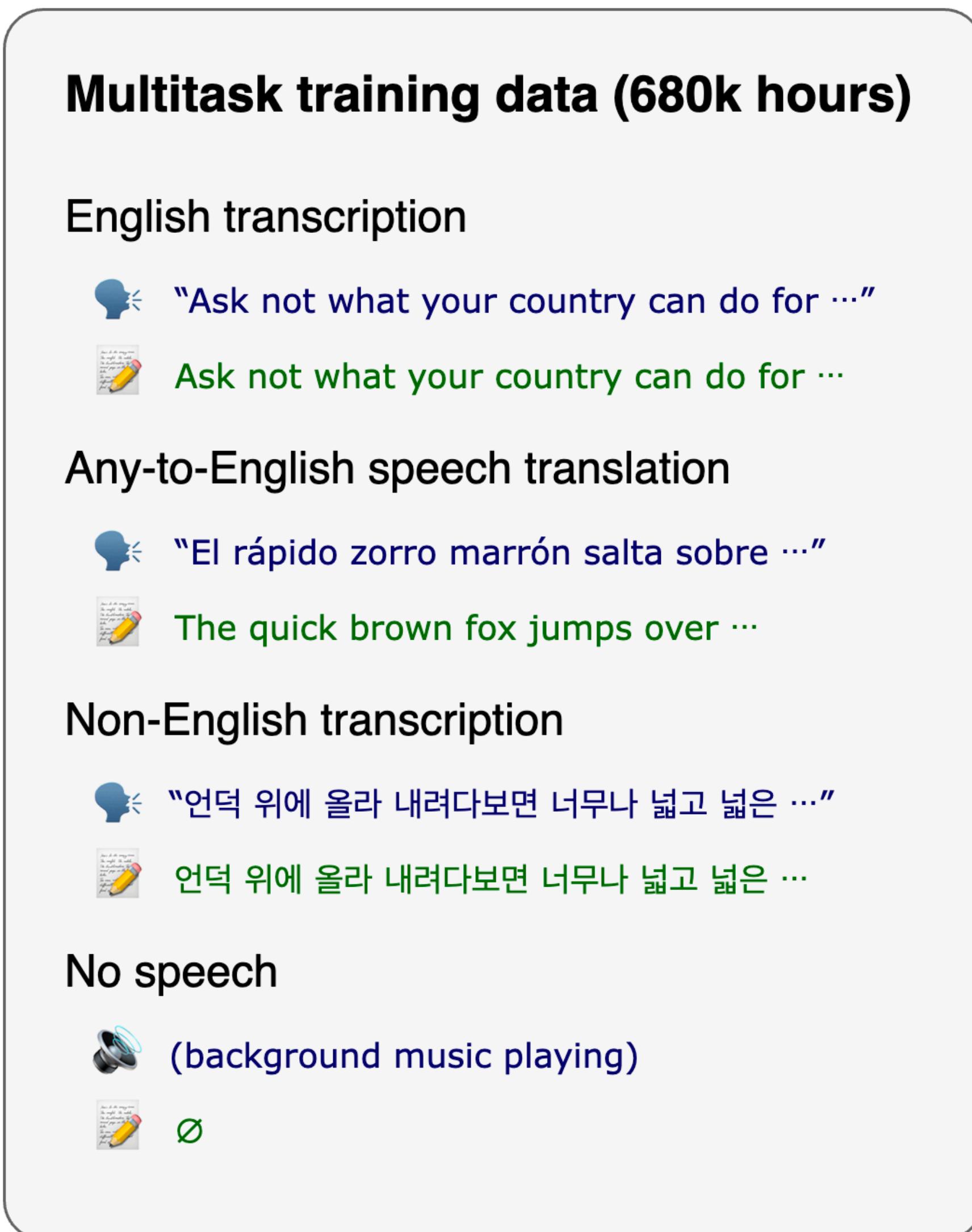
*“ Using just ten minutes of labeled data and pre-training on 53k hours of unlabeled data still achieves 4.8/8.2 WER.*

# Multilingual vs Code-Switched Acoustic Models

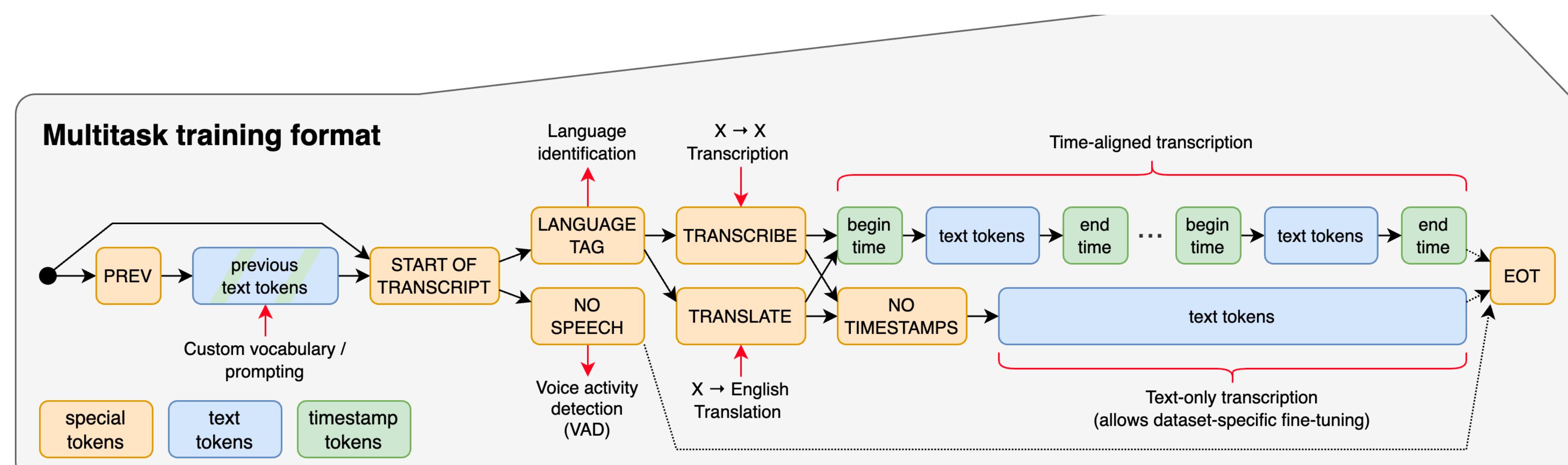
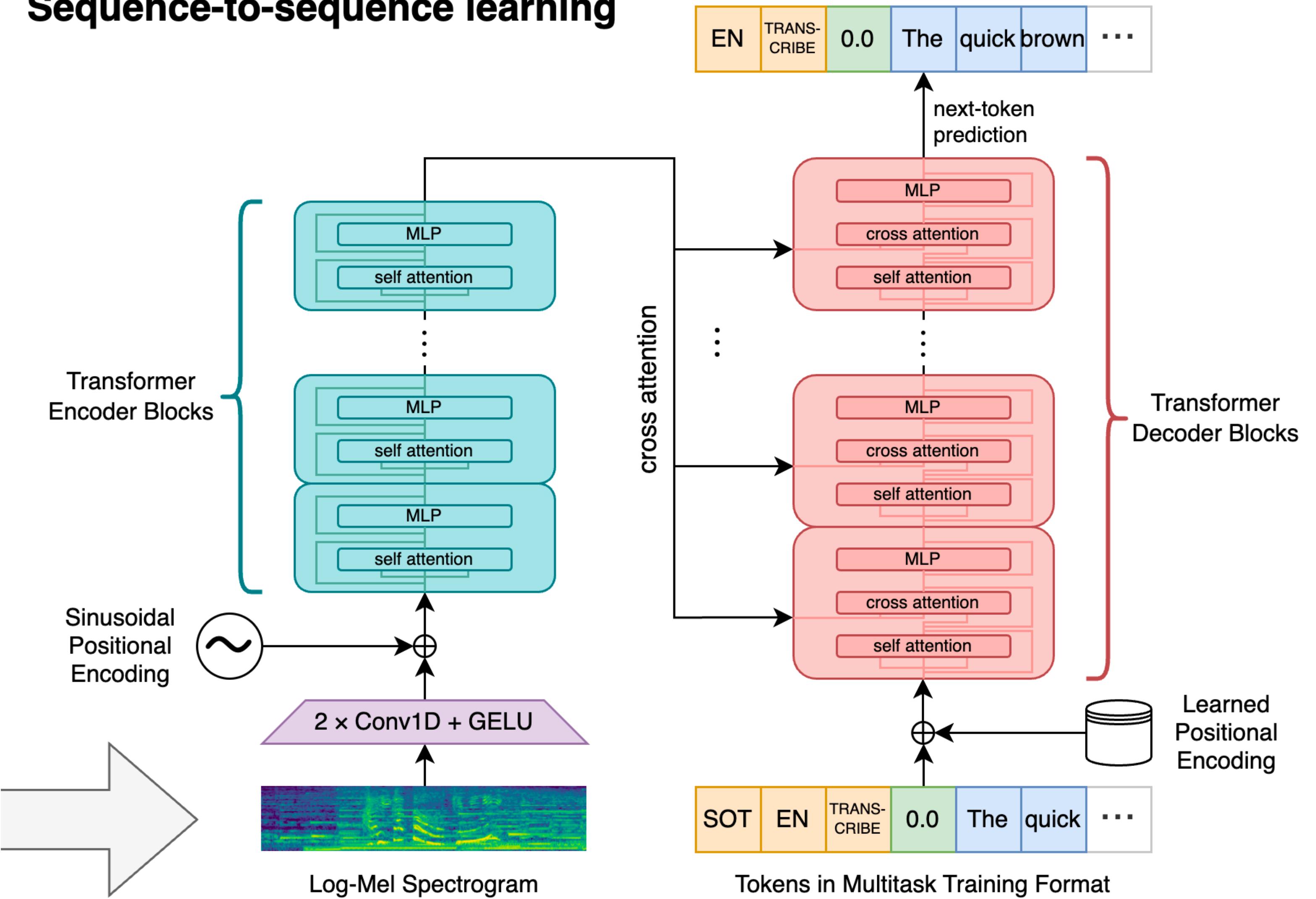


# Multilingual Acoustic Models

Whisper



## Sequence-to-sequence learning

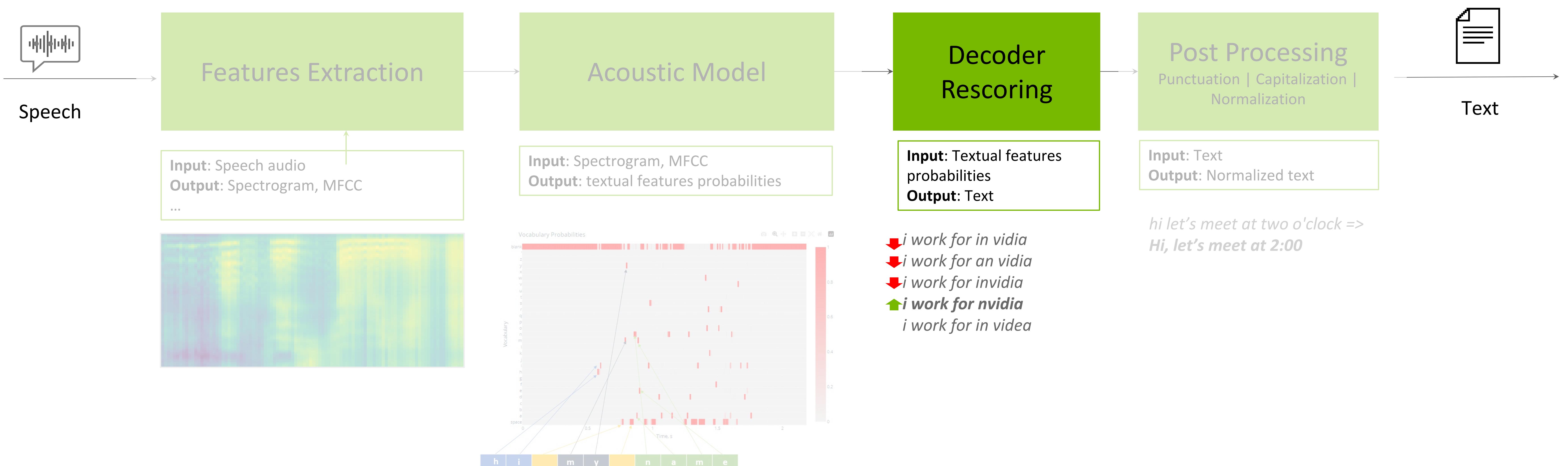


# Acoustic Models

## Recap

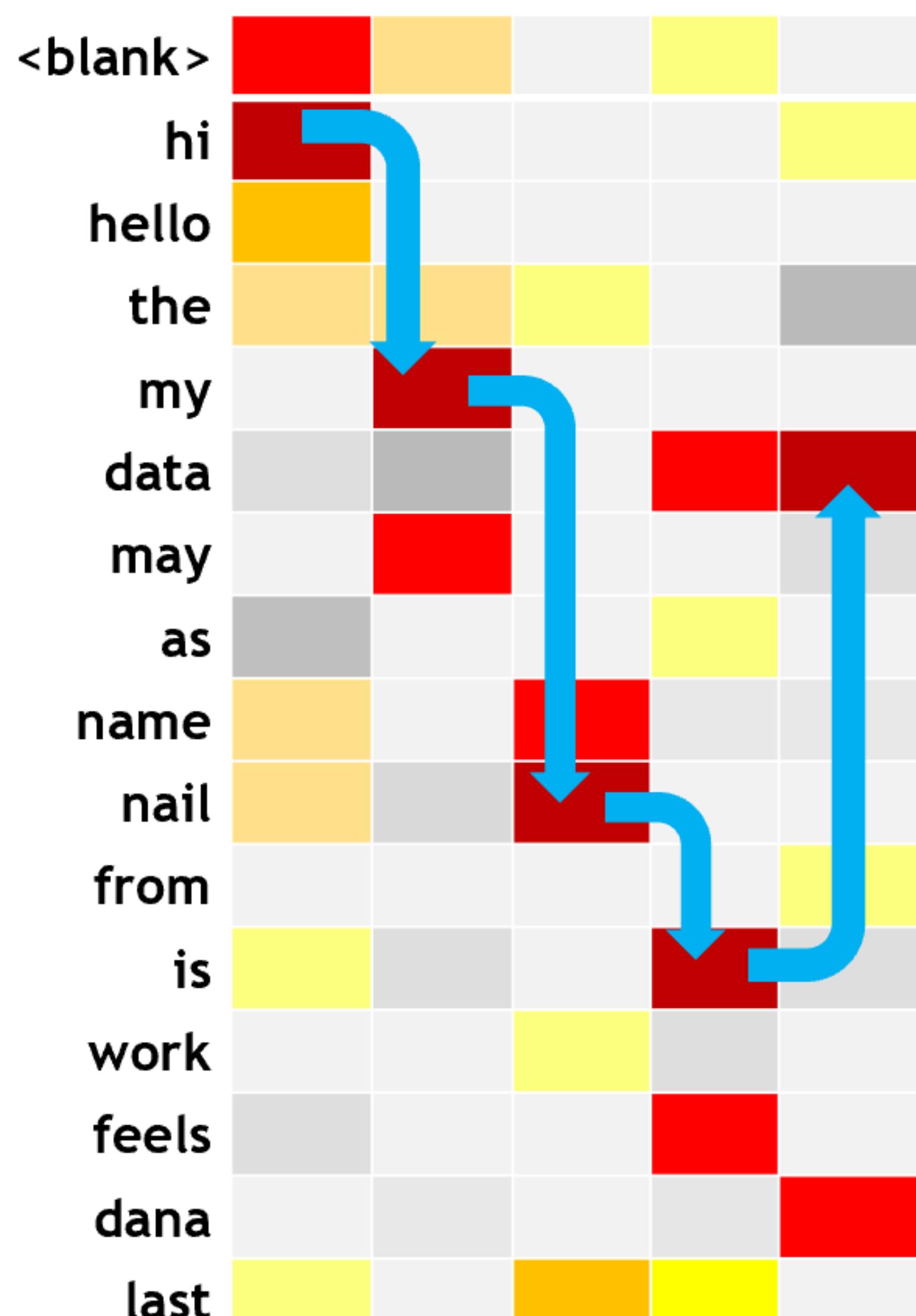
Models	Tasks	Approach	Training	Comments
DeepSpeech	ASR	RNN-based + CTC	Supervised	Autoregressive
DeepSpeech2	ASR	RNN-based + CTC	Supervised	Autoregressive
Wave2vec 2.0		Transformer + Convolutional	Self-Supervised + Supervised finetuning	Non Autoregressive
Citrinet	ASR	Convolutional + CTC	Supervised	Non Autoregressive
Conformer-CTC	ASR	Transformer + Convolutional + CTC	Supervised	Non Autoregressive
Conformer-RNNTransducer	ASR	Transformer + Convolutional + LSTM	Supervised	Autoregressive
Wave2Letter	ASR	Convolutional + CTC	Supervised	Non Autoregressive
Whisper	Multilingual Multitask	Transformer + Convolutional	Supervised	Autoregressive

# ASR Pipeline

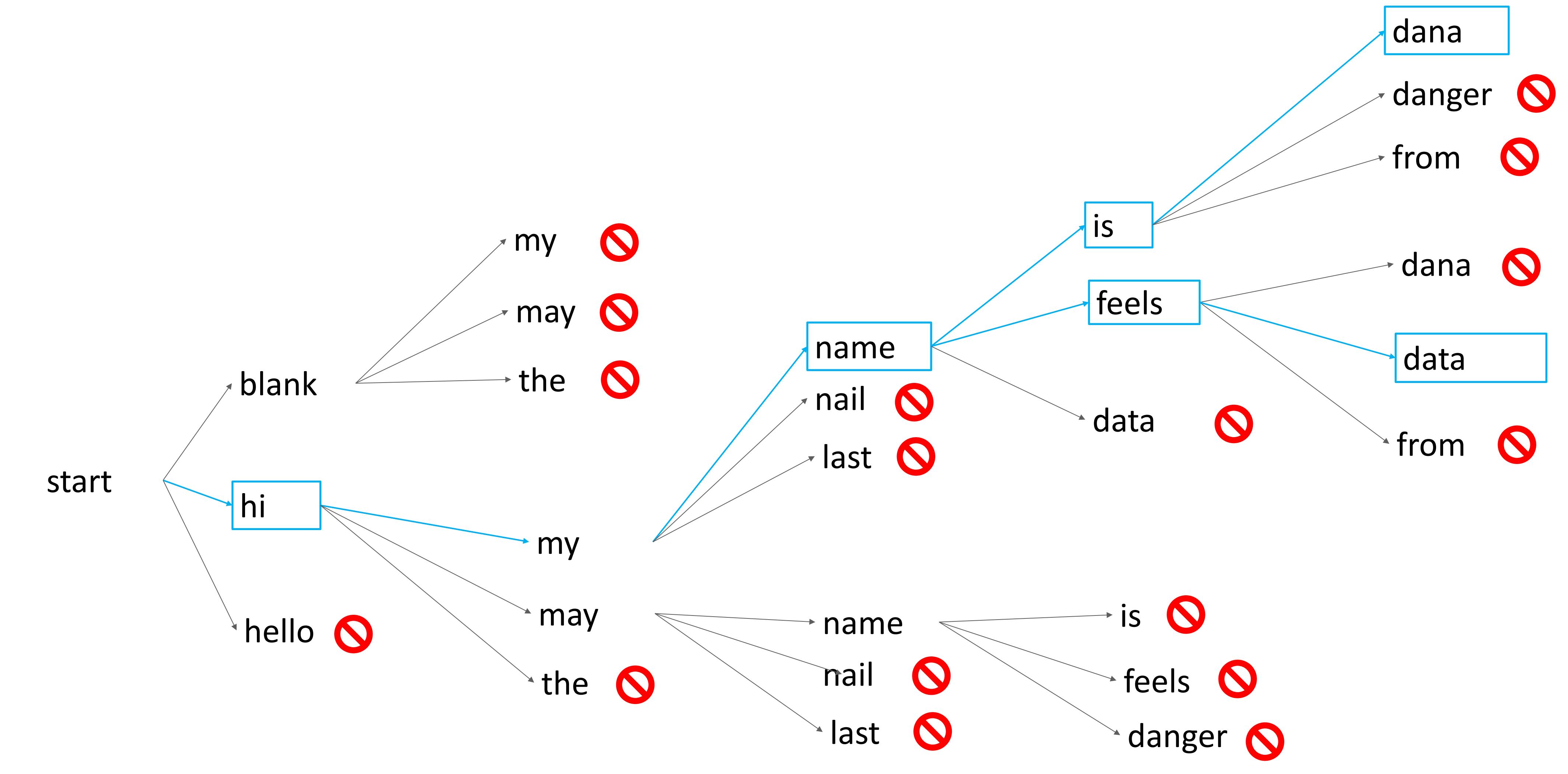


# Decoder Rescoring

Greedy decoder | beam search decoder



*Hi my nail is data*



*hi my name feels data | score 0.78*  
*hi my name is dana | score = 0.72*

$$P(\text{hi my name is dana}) = p(\text{hi}) * p(\text{my}/\text{hi}) * p(\text{name}/\text{hi my}) * p(\text{is}/\text{hi my name}) * p(\text{dana}/\text{hi my name is})$$

# Decoder Rescoring

Beam search decoder with a language model

- Language Model: The probability of a word depends its previous words.

$$P(\text{hi my name is dana}) = p(\text{hi}) * p(\text{my}/\text{hi}) * p(\text{name}/\text{hi my}) * p(\text{is}/\text{hi my name}) * p(\text{dana}/\text{hi my name is})$$

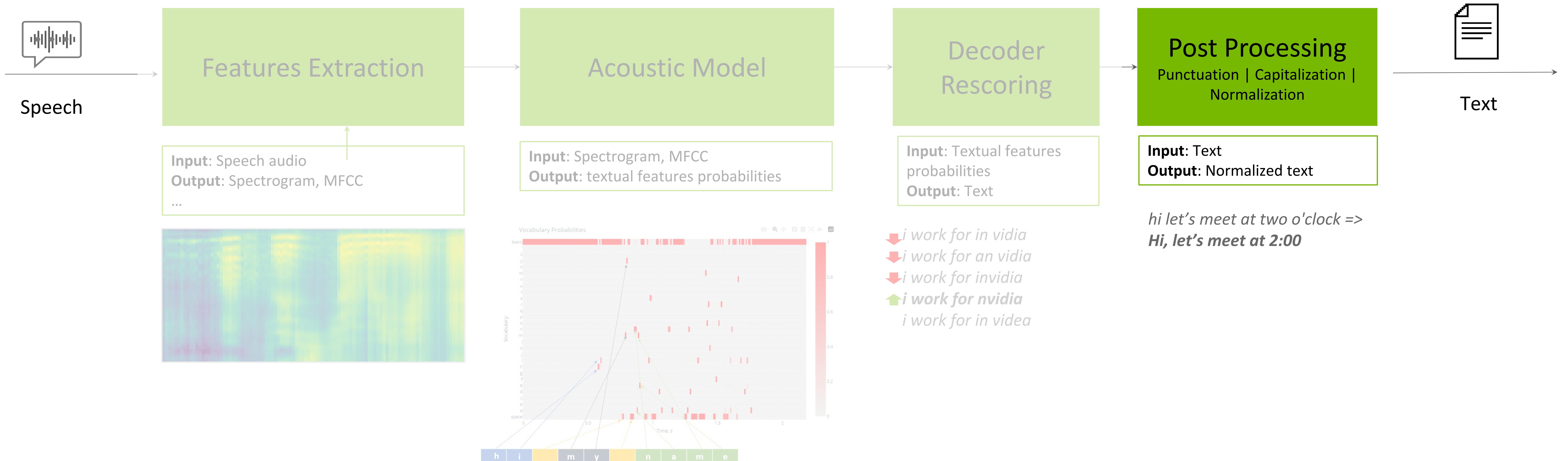
- N-gram LM: The probability of a word depends on a fixed  $N$  previous words.

Example 2-gram:  $P(\text{hi my name is dana}) = p(\text{hi}) * p(\text{my}/\text{hi}) * p(\text{name}/\text{hi my}) * p(\text{is}/\text{hi my name}) * p(\text{dana}/\text{hi my name is})$

$$\text{final\_score} = \text{acoustic\_score} + \alpha \text{lm\_score} + \beta \text{seq\_length}$$

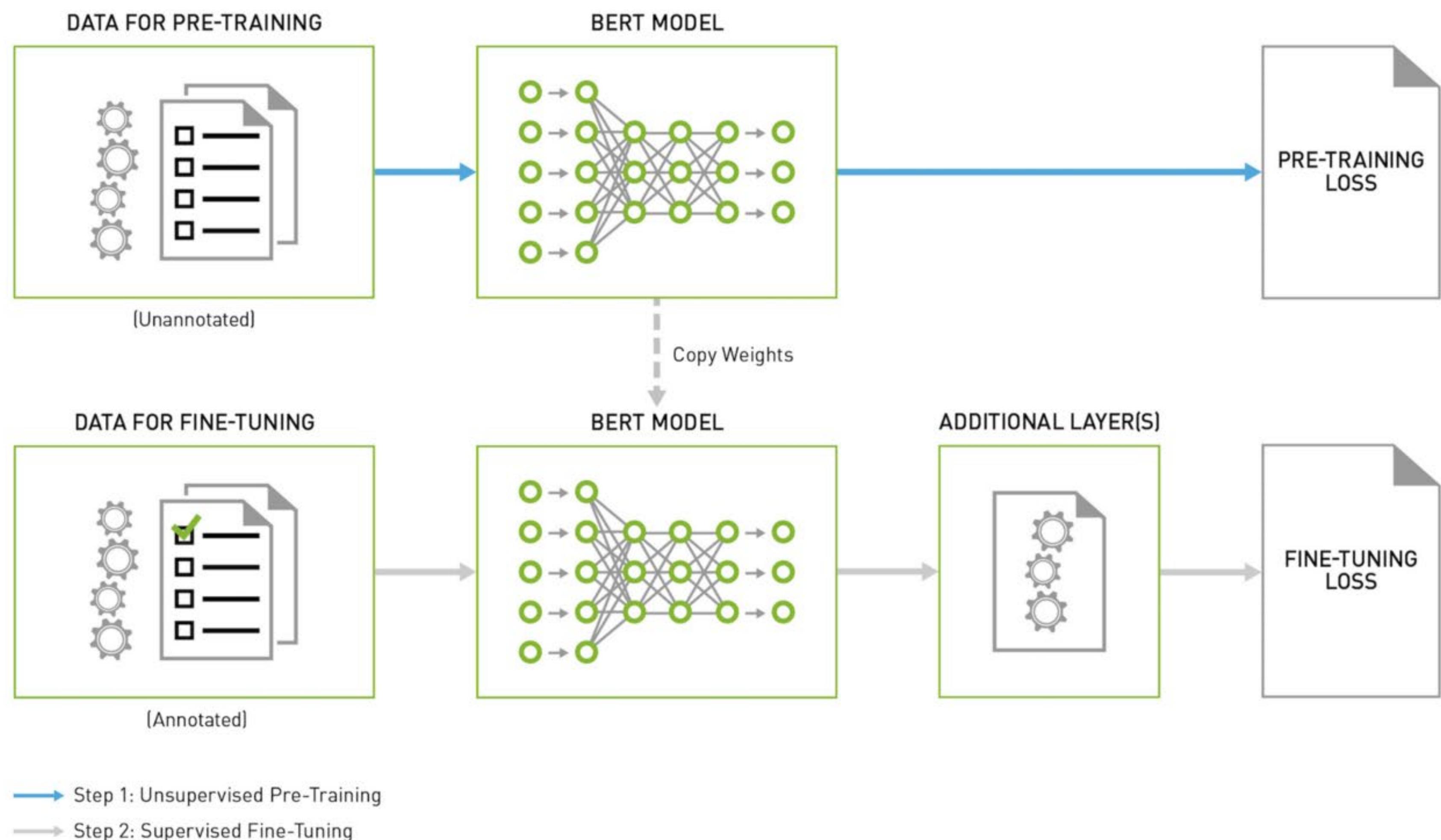
 *hi my name feels data | score 0.66*  
 *hi my name is dana | score = 0.79*

# ASR Pipeline



# Post Processing

## Punctuation, Capitalization



when is the next flight to new York  
OU 00 00 00 00 00 OU ?U  
When is the next flight to New York?

# Post Processing

## Normalization: Grammar Rules

### Weighted finite-state transducer (WFST) grammars

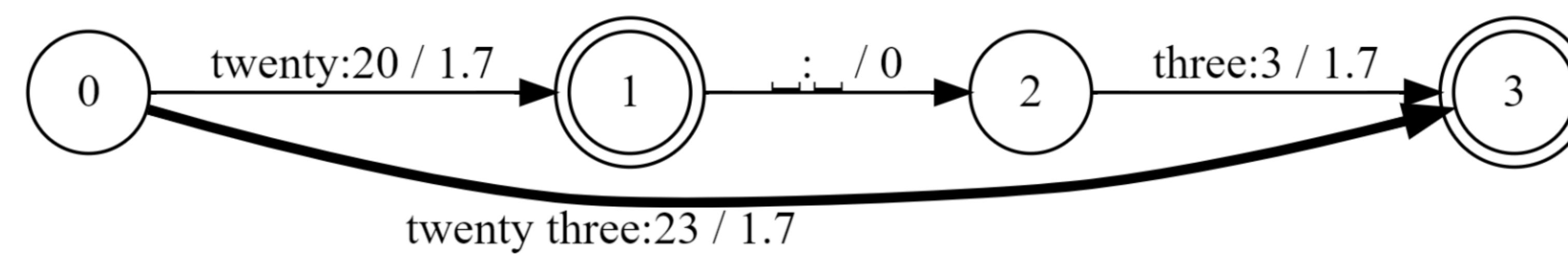
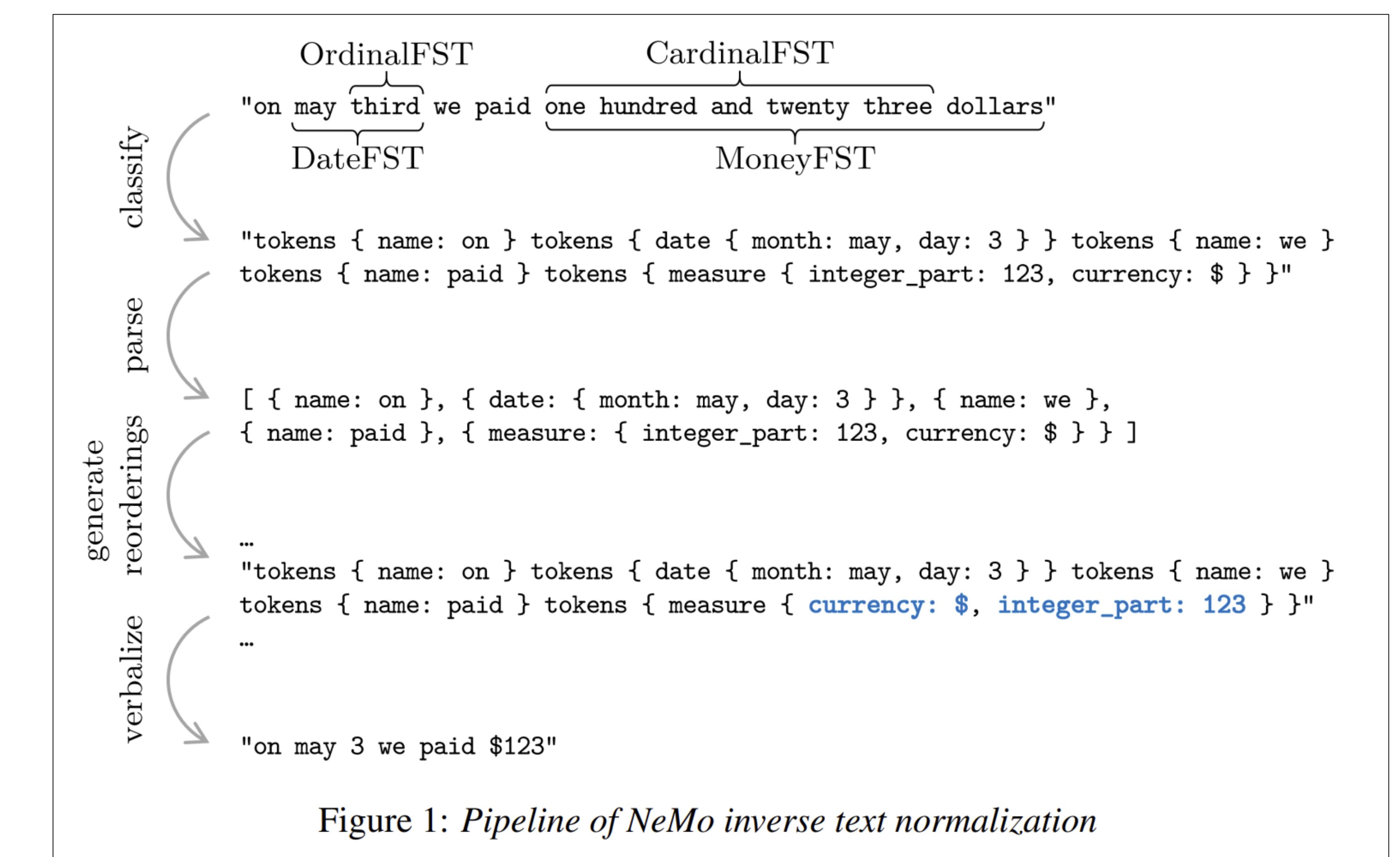


Figure 5: In case of ambiguity the path with the smallest sum of weights will be chosen. Here “twenty three” is transduced to “23” instead of “20 3”.



# Other Speech Tasks



# Speech Tasks

## Speech Command Recognition

The task of classifying an input audio pattern into a discrete set of classes.

E.g.: Yes, No, "Go", "Stop", "Left", "Down"

## Voice Activity Detection (VAD)

The task of predicting which parts of input audio contain speech versus background noise.

## Spoken Language Identification (LangID)

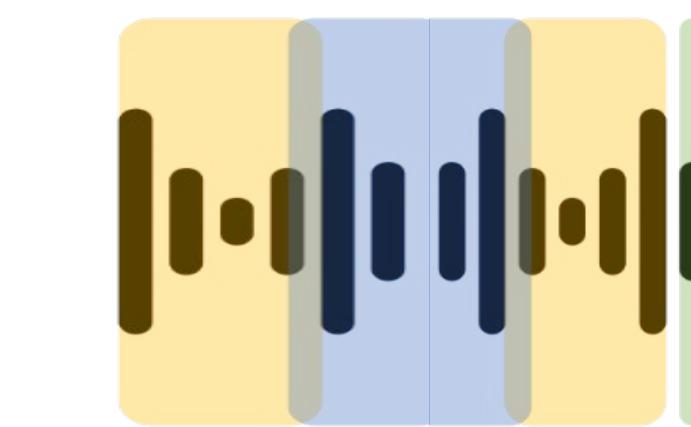
The task of recognizing the language of the spoken utterance automatically.

## Audio Sentiment Classification

Extends the conventional text-based sentiment analysis to depend on the acoustic features extracted from speech.

## Speaker Diarization

Group same speakers





# ASR (Part 1)

## Lecture

- What is Conversational AI? Complexity of Deploying Conversational AI Applications
- Deep Learning Fundamentals
- Speech and Text Processing
- Automatic Speech Recognition
- ASR Tools: Riva, NGC, and NeMo
- Lab Overview

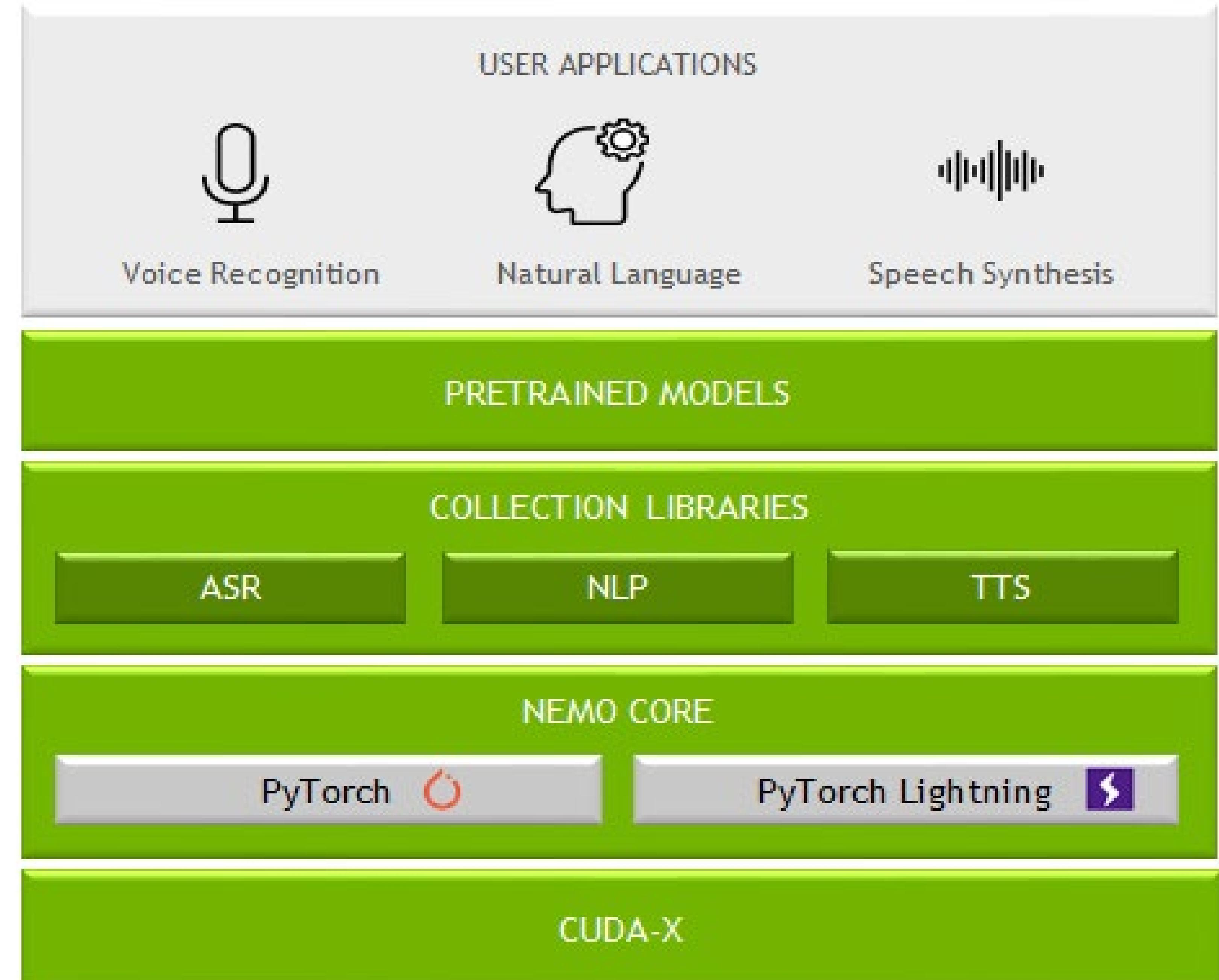
## Lab

- Investigate ASR Pipeline with NeMo and Riva

# What is NeMo?

Toolkit to build state-of-the-art conversational AI models

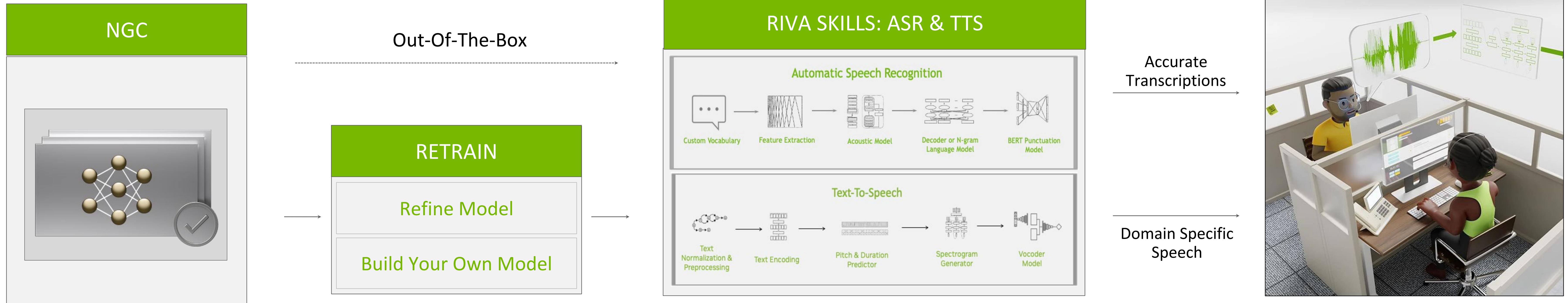
- DL-Based Speech & Language Understanding Models
- Include Semantic Checking for Correct-by-Construction Architectures
- Open-Sourced
- Integrated with PyTorch & PyTorch Lightning
- Easy-to-Use APIs
- Optimized Training Performance
- 100+ Pre-Trained GPU-Optimized Checkpoints
- Scale to 1000s of NVIDIA GPUs



<https://ngc.nvidia.com/catalog/containers/nvidia:nemo>

# NVIDIA Riva

Fully customizable, GPU-accelerated SDK for real-time Speech AI



SOTA Pre-Trained Models

Train for Any  
Domain and Language

- Fully Customizable for the Best Possible Accuracy
- GPU-Accelerated, Real-Time
- Scale to Hundreds of Thousands of Users
- On-Prem, in Any Cloud, at the Edge, or Embedded

SOTA – State-Of-The-Art

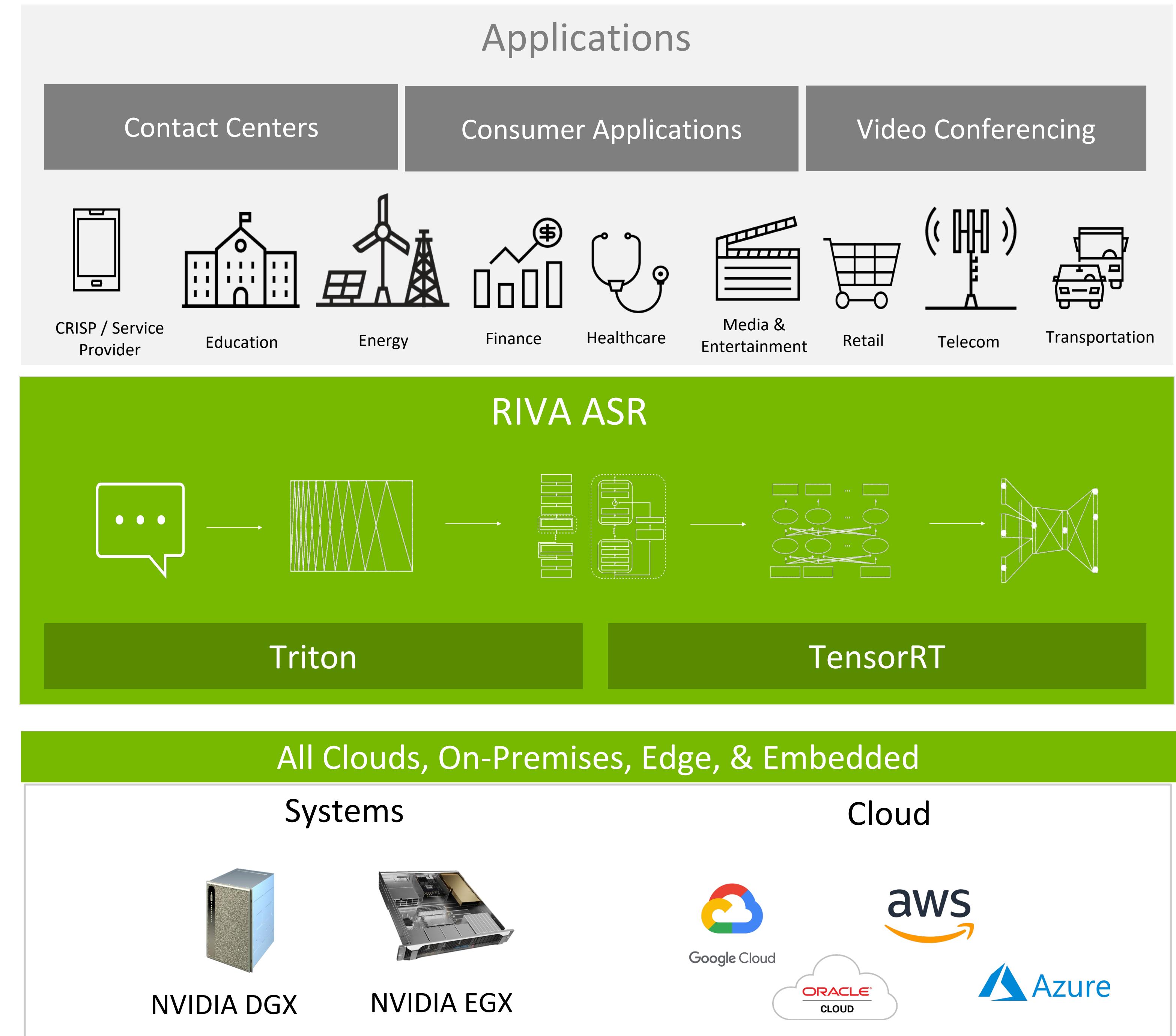
| ASR – Automatic Speech Recognition

| TTS – Text-To-Speech

# NVIDIA Riva Automatic Speech Recognition

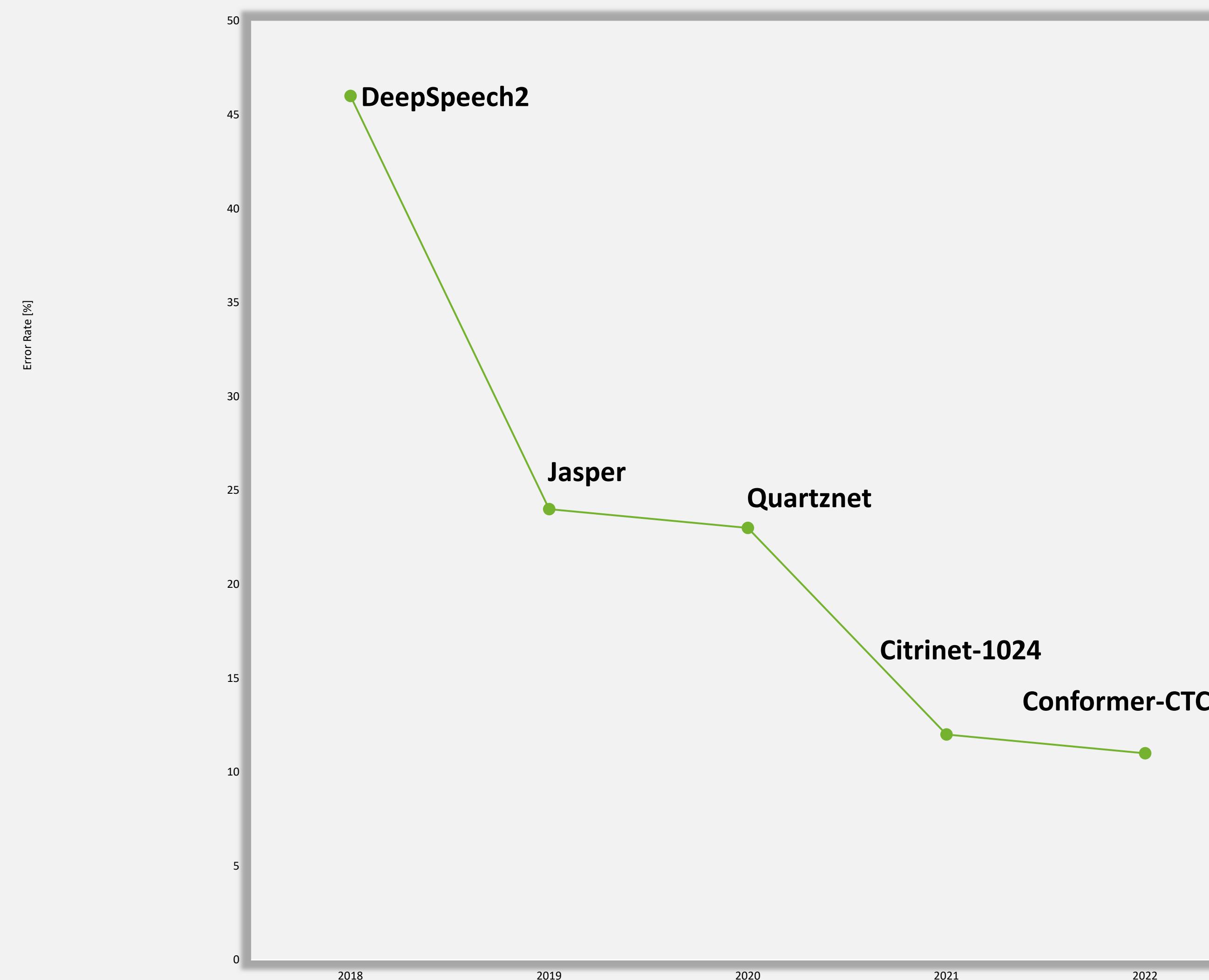
The best possible ASR with SOTA deep learning

- **Accurate ASR with SOTA\* models** trained for 1M+ hrs on 70K hrs of speech
- Support for:
  - **Languages:** English, Spanish, Mandarin, Hindi, Russian, Korean, German, French, Arabic, Italian, Japanese, Brazilian-Portuguese, European-Spanish
- **2X accuracy improvement** with customizations for:
  - Industry specific jargon
  - Accents & dialects
  - Noisy environments
- **Real-time** performance far below 300ms for interactive speech apps
- **High scale** of 100s thousands of concurrent streams
- **Runs anywhere:** all clouds, on-premises, at the edge, embedded

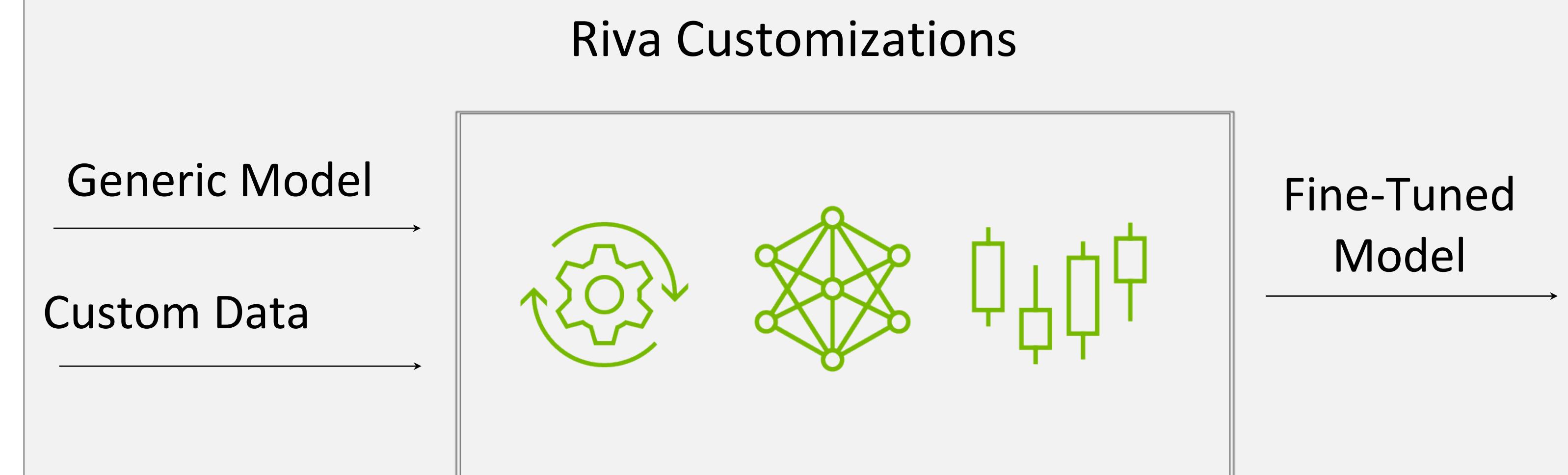


# NVIDIA Riva Delivers World-Class ASR

4X Accuracy with SOTA Model Architectures

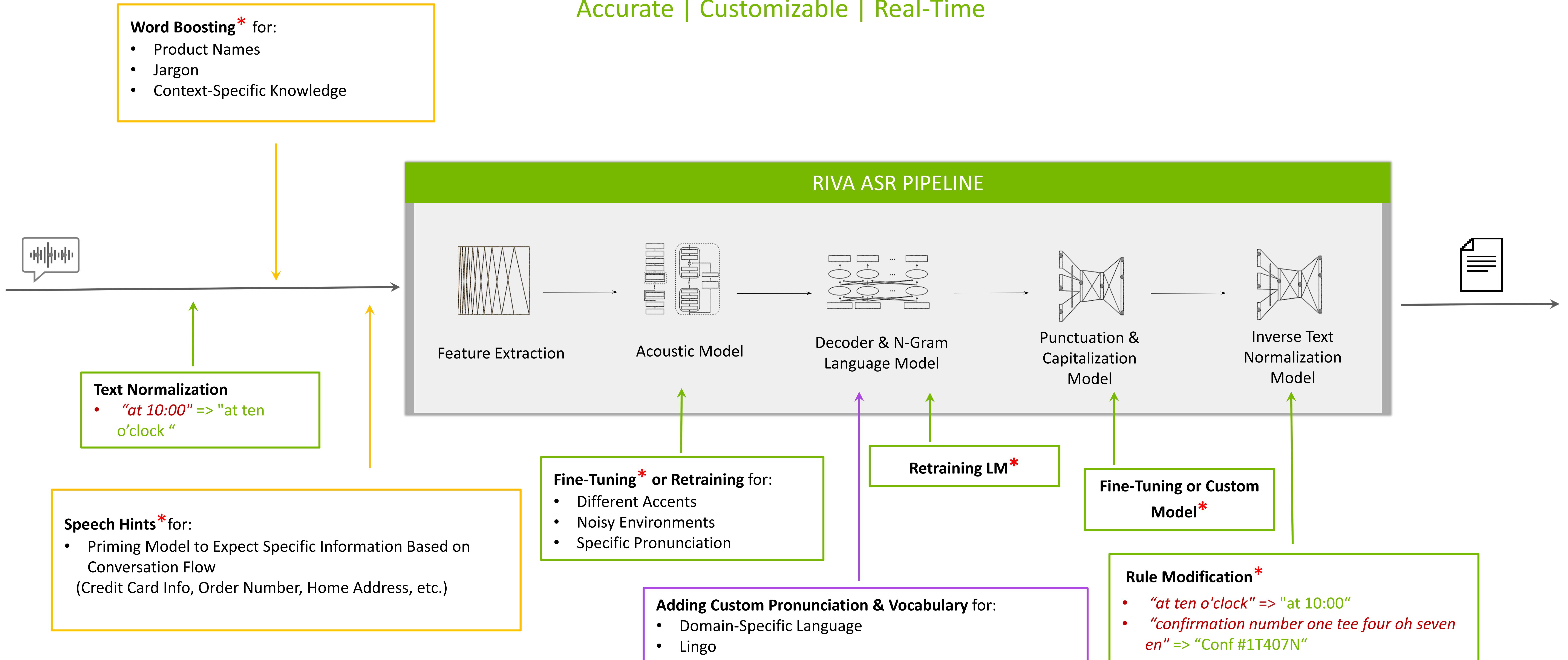


Additional 2X Accuracy with Customizations



# Customizations Across Riva ASR Pipeline

Accurate | Customizable | Real-Time



→ Training (server-side) – Need to rebuild/restart server

→ Deployment

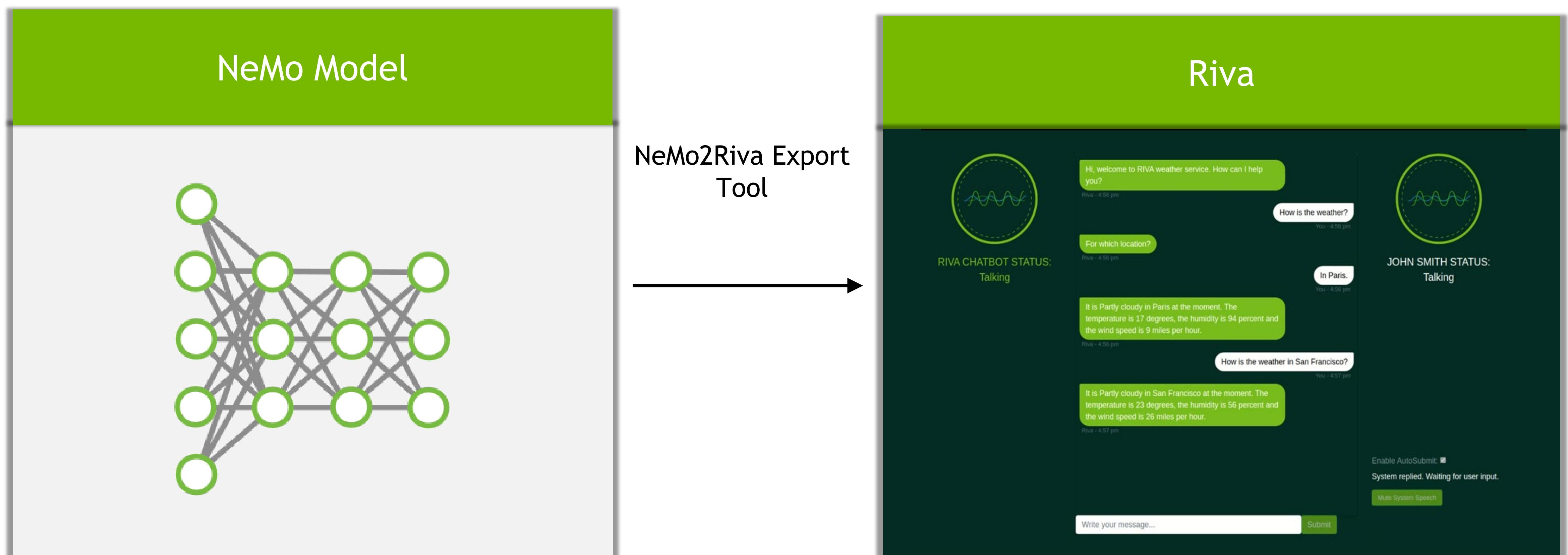
→ Inference (client-side) – No need to rebuild/start Riva server

\* Generic ASR Services offer black box with rigid & selectively limited customization

# Deploy to Production

Generate high-performance inference

- Quickly export NeMo models to Riva
- Support for speech and language models across multiple languages
- Step-by-step deployment instructions in documentation



Exporting NeMo Models to Riva



# ASR (Part 1)

## Lecture

- What is Conversational AI? Complexity of Deploying Conversational AI Applications
- Deep Learning Fundamentals
- Speech and Text Processing
- Automatic Speech Recognition
- ASR Tools: Riva, NGC, and NeMo
- Lab Overview

## Lab

- Investigate ASR Pipeline with NeMo and Riva

# Lab Overview

## ASR Pipeline

1. NGC Credential (Required for Riva deployment)
2. Investigate AI models for Automatic Speech Recognition
  - Gain an understanding about the ASR pipeline and the various acoustic models
  - Use transcript decoders to select the best transcriptions
  - Use ITN, capitalization, and punctuation to improve written transcriptions
  - Use a language identification model to identify what language was spoken
3. ASR deployment with Riva in several languages
  - Launch Riva ASR service for English and Spanish
  - Request the ASR service using a Python client API



# Nvidia RIVA

Where to start?

The screenshot shows the NVIDIA NGC Catalog interface. At the top, there's a navigation bar with the NVIDIA logo, the text "NVIDIA NGC | CATALOG", and "Welcome Guest". Below the navigation bar, the URL "Catalog > Resources > Riva Skills Quick Start" is visible. The main content area has a title "Riva Skills Quick Start" and a "Download" button. A navigation menu at the top of the content area includes "Overview" (which is underlined), "Version History", "File Browser", "Release Notes", "Related Collections", and "More". The "Overview" section contains the following content:

- Quick Start Guide for Data Center Platforms**: A brief description stating that NVIDIA Riva supports two architectures, Linux x86\_64 and Linux ARM64, referred to as **data center** (x86\_64) and **embedded** (ARM64). It notes that instructions are applicable to data center users.
- A link to the [NVIDIA Riva Developer Forum](#) for more information and questions.
- Prerequisites**: A list of requirements:
  - You have access and are [logged into](#) NVIDIA NGC. Refer to the [NGC Getting Started Guide](#).
  - You have access to an NVIDIA Volta, NVIDIA Turing, or an NVIDIA Ampere Architecture-based A100 GPU. Refer to the [Support Matrix](#).
  - You have Docker installed with support for NVIDIA GPUs. Refer to the [Support Matrix](#).
  - For Riva Enterprise customers, ensure you have the NGC\_API\_Key and NGC\_Organization details to enable the Riva Enterprise metrics collection. Refer to the [Configuration](#) section.
- Getting Started with Riva for Data Center Platforms**: A note stating that Riva includes Quick Start scripts to help get started with Riva Speech AI Skills, intended for deploying services locally, testing, and running example applications.

On the left side of the content area, there's a sidebar with the following sections and their values:

- Description**: Scripts and utilities for getting started with Riva Speech Skills
- Publisher**: NVIDIA
- Use Case**: Other
- Framework**: Other
- Latest Version**: 2.7.0
- Modified**: October 28, 2022
- Compressed Size**: 2.7.0

1 Configure Riva

2 Initialize Riva

3 Start Riva

[Download Riva from NGC](#)



# Nvidia RIVA

## Launch Riva ASR services

1

### Service Configuration

*Modify the config.sh file*

#### Config.sh

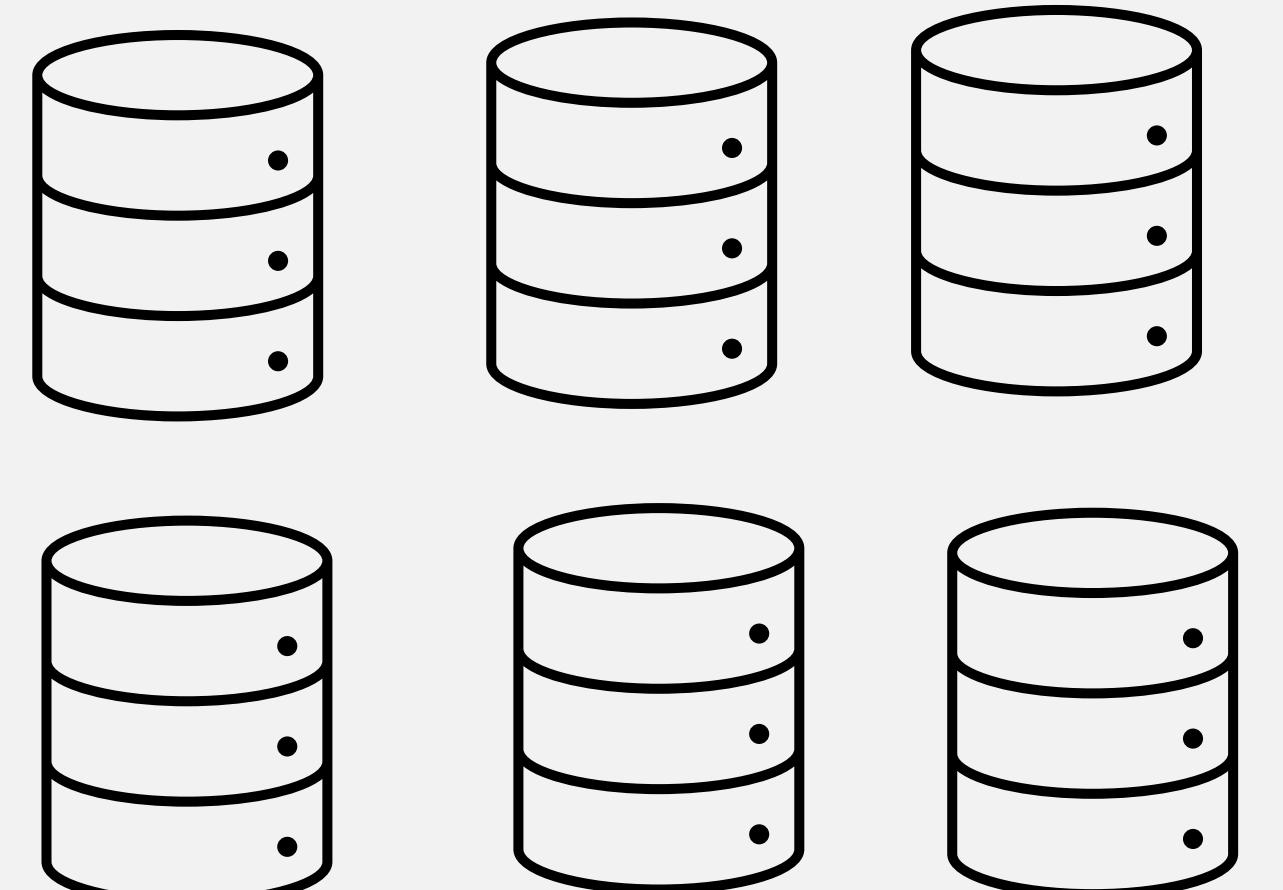
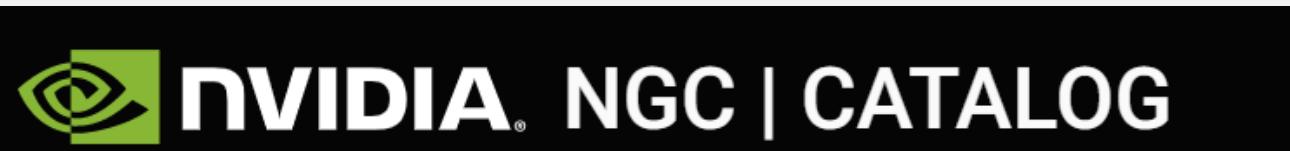
```
...
# Enable or Disable Riva Services
service_enabled_asr=true
service_enabled_nlp=false
service_enabled_tts=false
...
# Language code to fetch models of a specify language
# Currently only ASR supports languages other than English
# Supported language codes: en-US, en-GB, de-DE, es-US, ru-RU,
zh-CN, hi-IN, fr-FR, ko-KR, pt-BR
language_code=("en-US")
...
# Models ($riva_model_loc/models)
# During the riva_init process, the RMIR files in
$riva_model_loc/rmir
# are inspected and optimized for deployment. The optimized
versions are
# stored in $riva_model_loc/models. The riva server
exclusively uses these
# optimized versions.
riva_model_loc="/path/to/riva-model-repo"
```

2

### Service Initialization

*Run: bash riva\_init.sh*

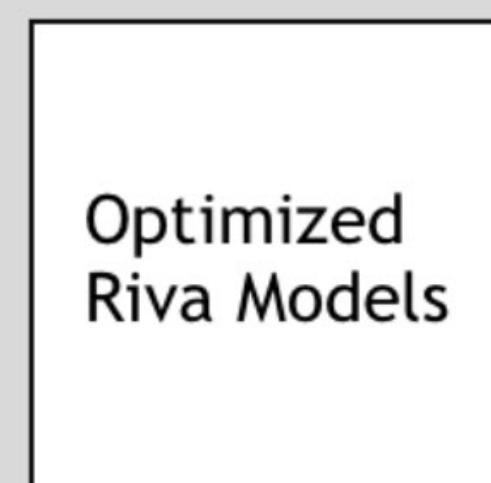
#### Download Models



Download containers and the prebuild RMIR  
files  
Riva Model Intermediate Representation

#### riva-deploy

Riva Model Repository

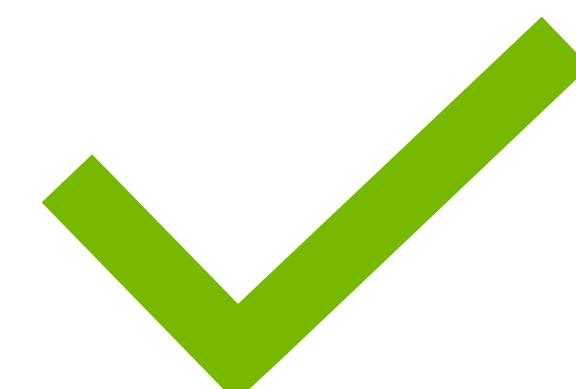


Convert RMIR into the Riva model repository.  
Neural networks are exported and optimized  
to run on the target platform

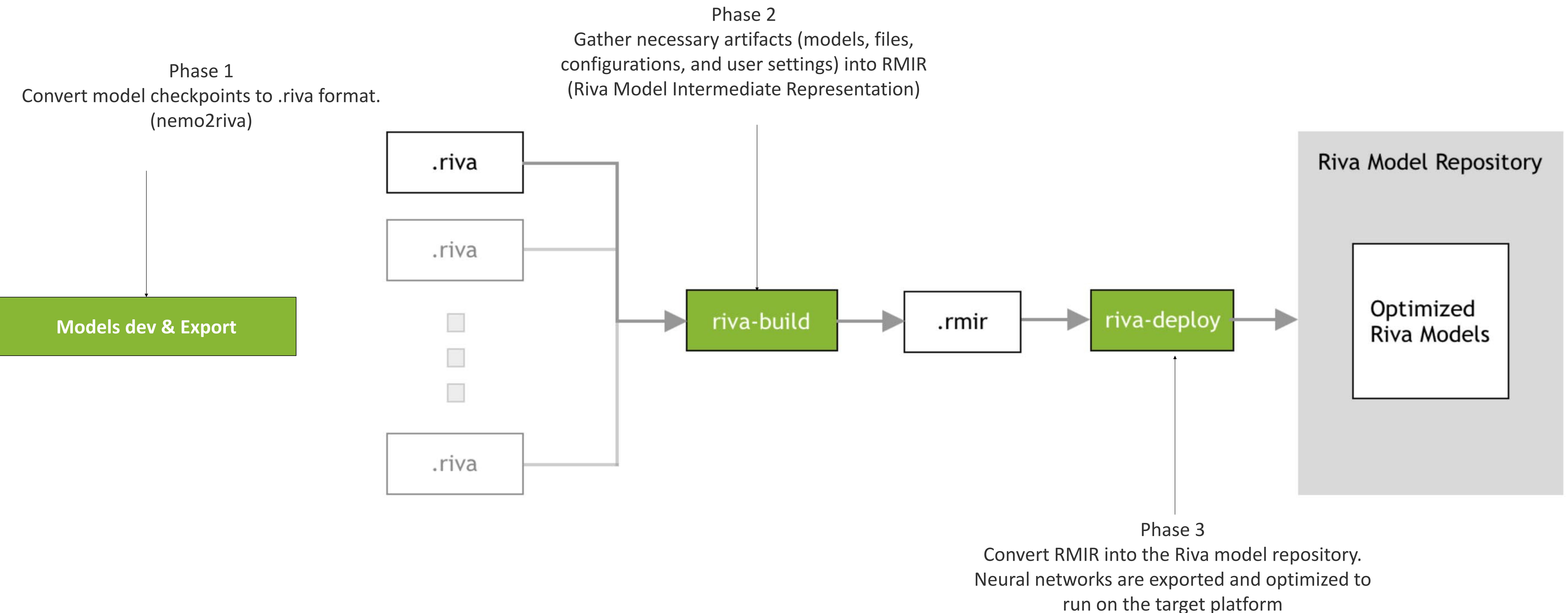
3

### Service Start

*Run: bash riva\_start.sh*

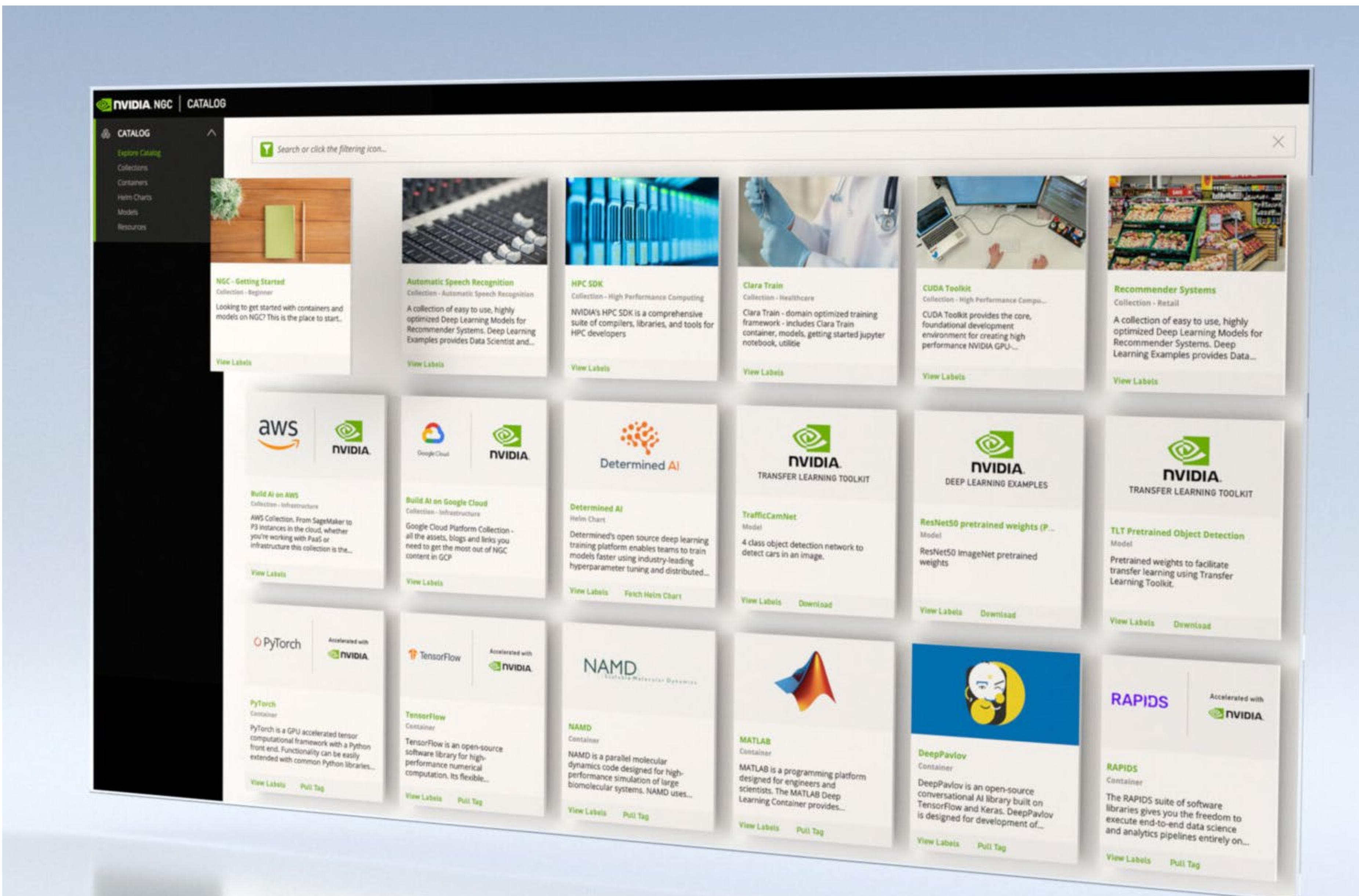


# NVIDIA Riva



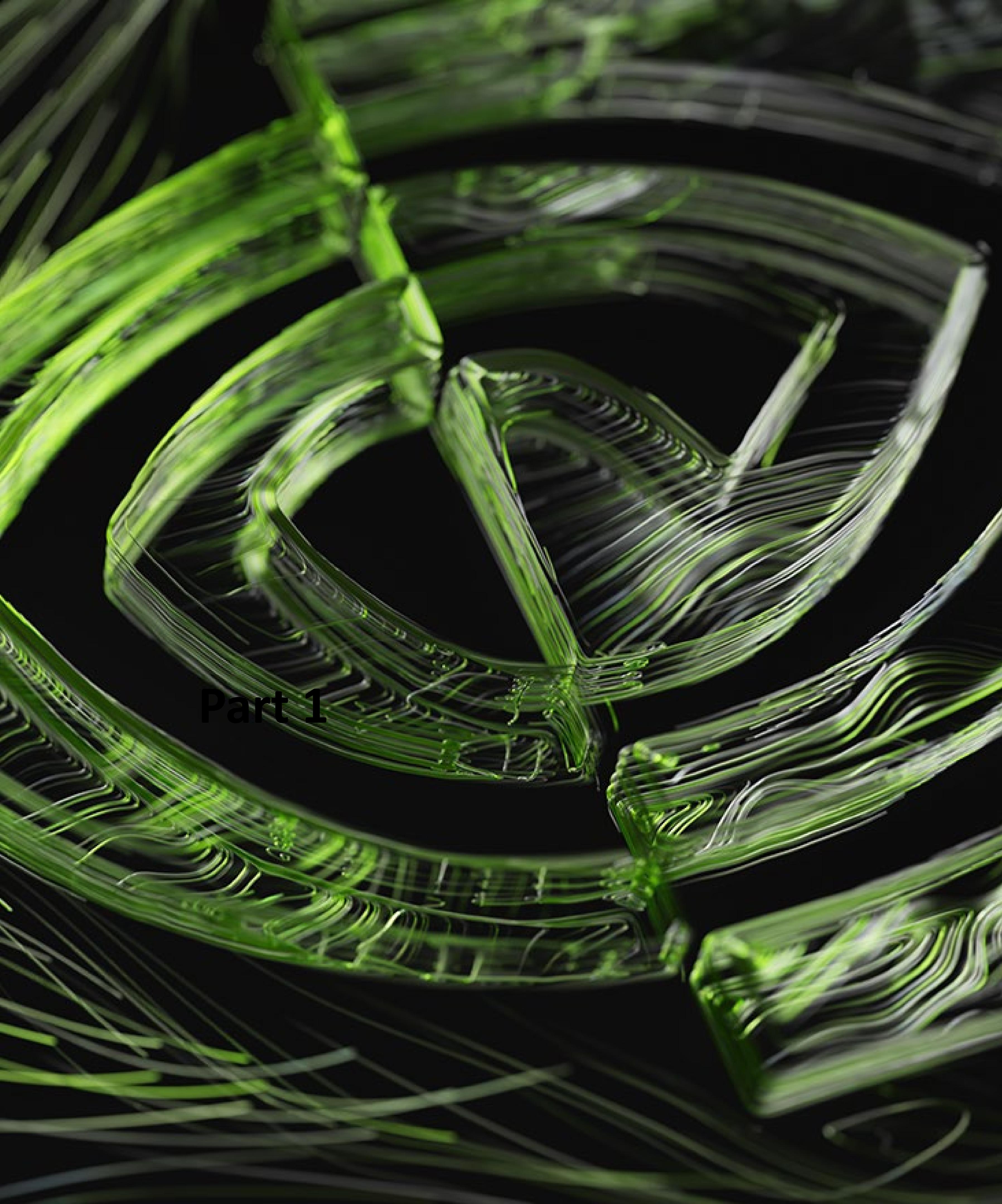
# NVIDIA Riva

## NVIDIA GPU Cloud (NGC) setup



A screenshot of the NVIDIA NGC Setup interface. The left sidebar has sections for CATALOG, PRIVATE REGISTRY, and BASE COMMAND. The main area is titled 'Setup'. It features a 'Generate API Key' section with a key icon and the text: 'Generate your own API key in order to use the NGC service through the Docker client or through NGC CLI.' A green 'Get API Key' button is highlighted with a red oval. At the bottom right are 'Documentation' and 'Downloads' buttons.

A screenshot of the NGC CLI usage instructions. It shows a terminal window with the command '\$ nvc config set'. A red arrow points from the 'Get API Key' button in the previous screenshot to this command. Below it, another terminal window shows '\$ docker login nvcr.io' followed by 'Username: \$oauthtoken' and 'Password: &lt;Your Key&gt;'. A second red arrow points from the 'Get API Key' button to the '\$oauthtoken' placeholder in the password field.



# ASR (Part 1)

## Lecture

- What is Conversational AI? Complexity of Deploying Conversational AI Applications
- Deep Learning Fundamentals
- Speech and Text Processing
- Automatic Speech Recognition
- ASR Tools: Riva, NGC, and NeMo
- Lab Overview

## Lab

- Investigate ASR Pipeline with NeMo and Riva

!

Lab

