

# End-to-End Machine Learning for Speech Processing: Speech-to-Speech, Speech-to-Text and Text-to-Speech

G rard Chollet (CNRS-SAMOVAR)  
and Colleagues from IV, Zaion, CNRS/IMT and SMI

[gerard.chollet@telecom-sudparis.eu](mailto:gerard.chollet@telecom-sudparis.eu)



# What is End-to-End Machine Learning?

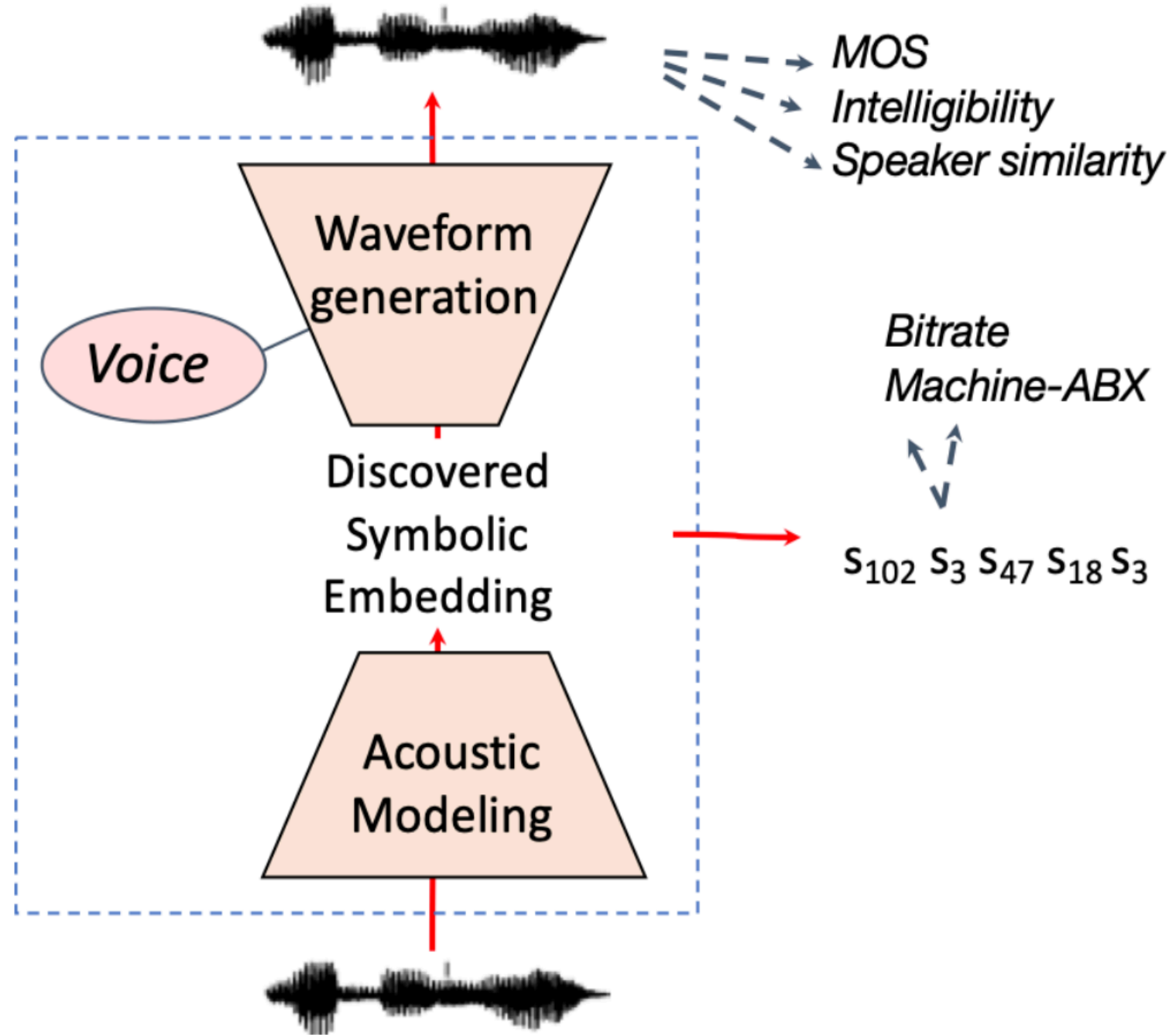
Deep Neural Networks (DNN, CNN, RNN,...) are designed (hand crafted/AutoML) and trained end-to-end to optimize a desired input-output association.

Large data sets (natural/augmented) are necessary to cover all possible (unknown) variations.

An example:

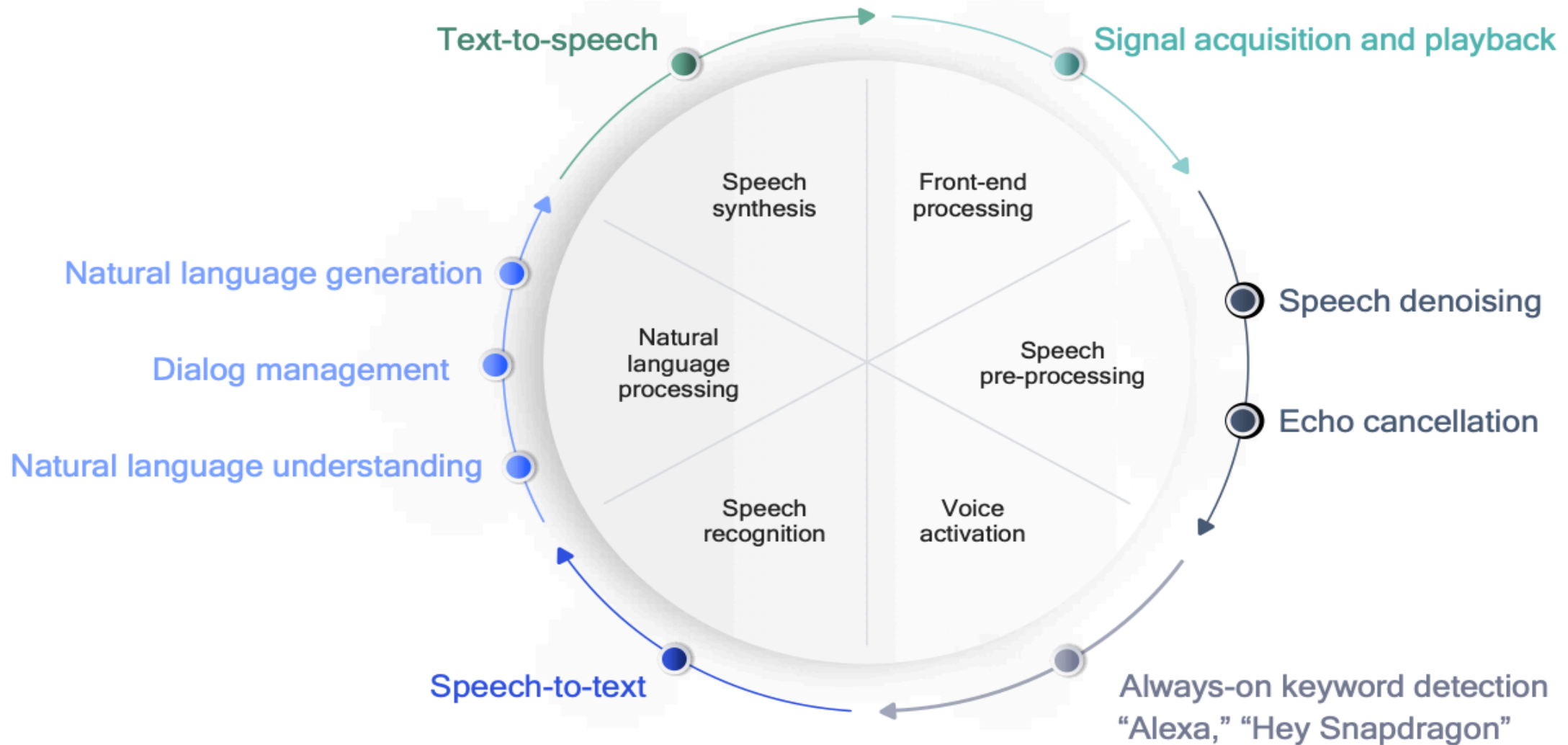
# ZeroSpeech 2019: TTS without T

<https://zerospeech.com/2019/>



# Voice UI components required for an end-to-end solution

Machine speech chain: listener and speaker

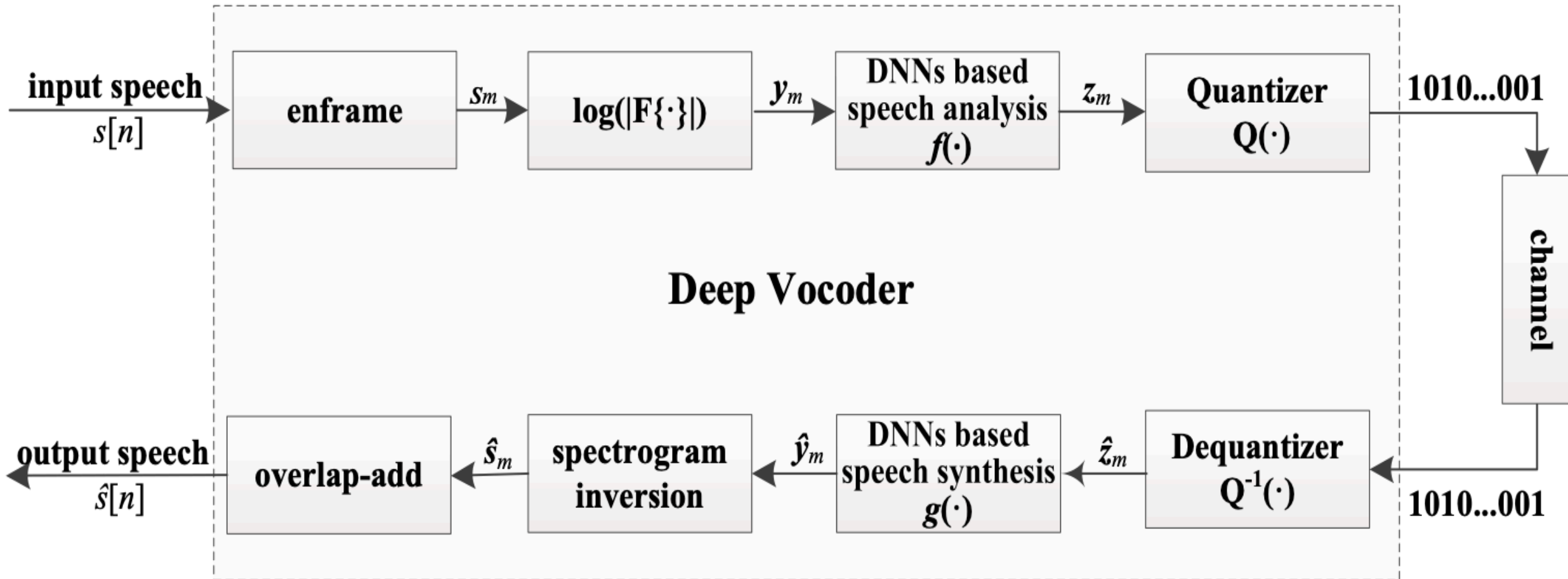


Qualcomm Snapdragon is a product of Qualcomm Technologies, Inc. and/or its subsidiaries.

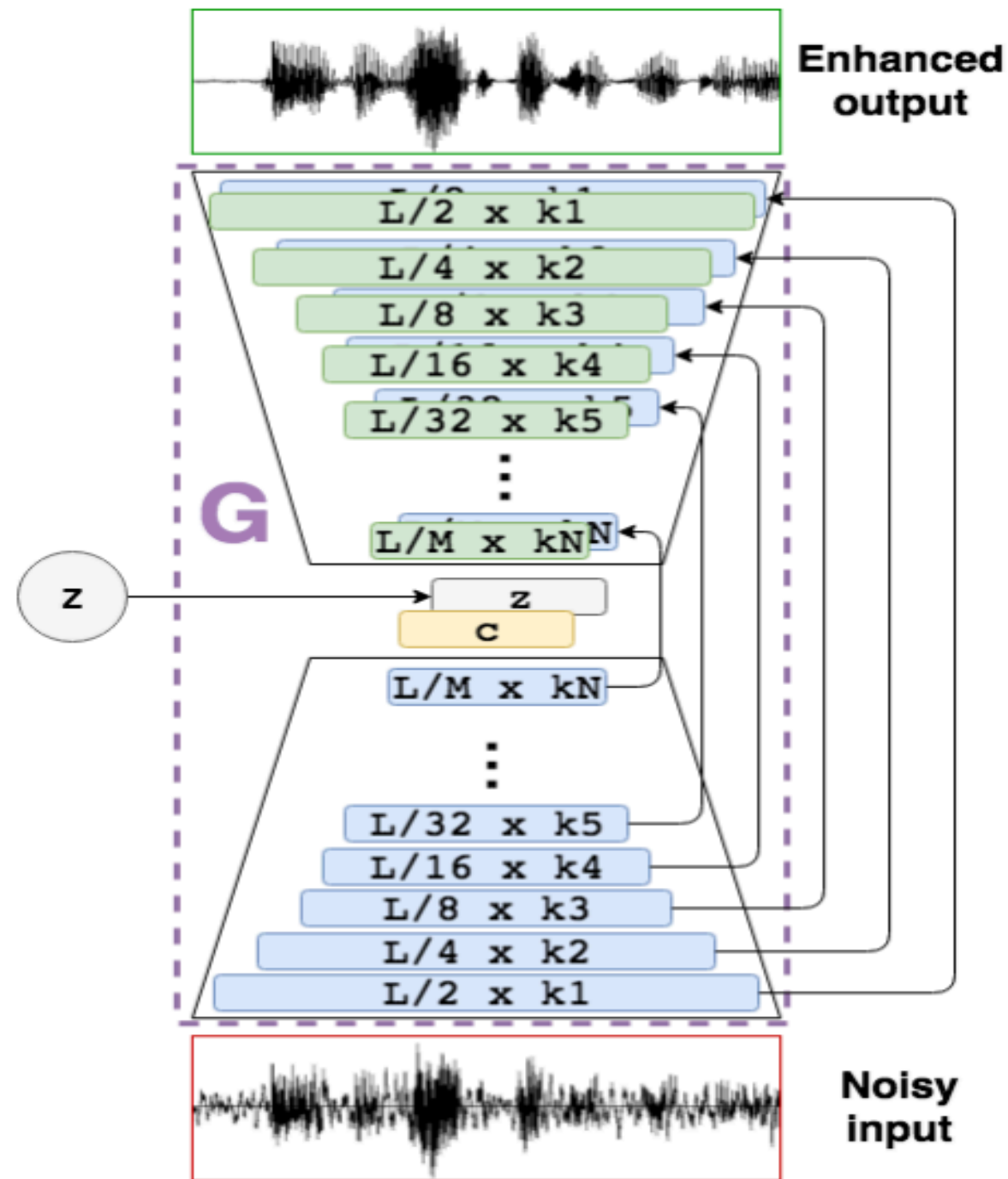
# Speech-to-Speech processing

- Speech coding, compression,
- Speech enhancement, denoising,
- Voice conversion,
- Speech-to-Speech Translation,
- Spoken Dialogue System,
- ...

# DEEP VOCODER: LOW BIT RATE COMPRESSION OF SPEECH WITH DEEP AUTOENCODER

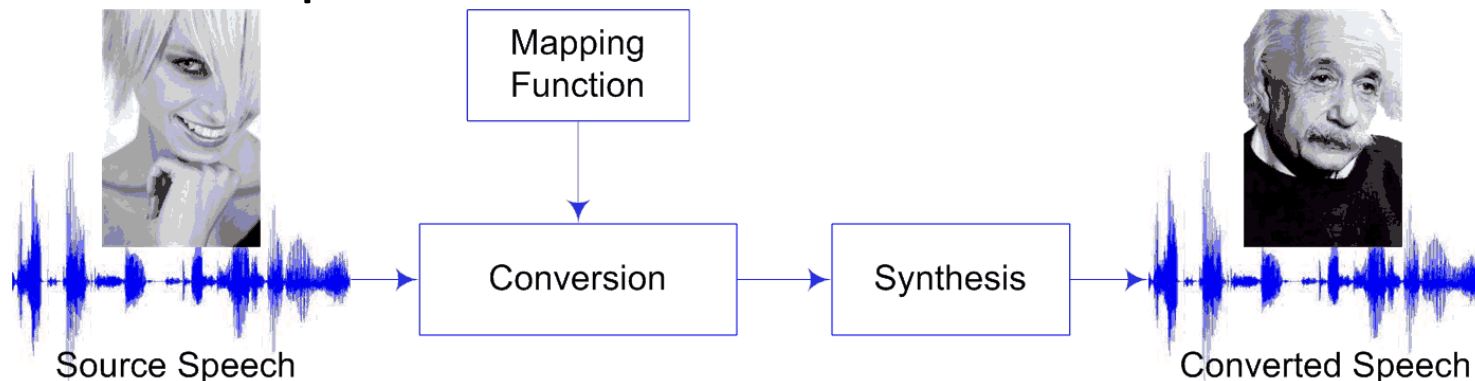


# SEGAN: Speech Enhancement Generative Adversarial Network



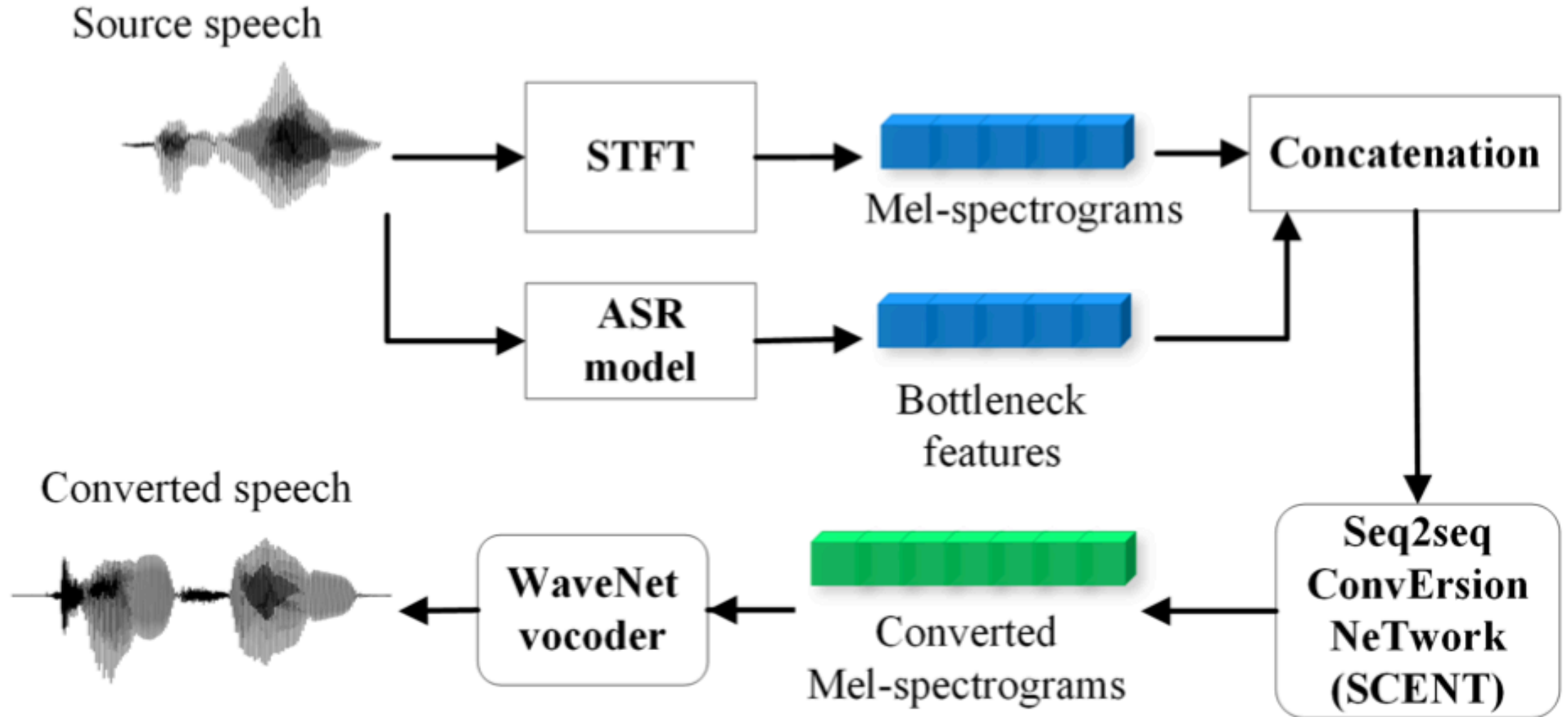
# Speaker Transformation (3/3)

- The estimated transformation rule is applied to an original speech pronounced by the source speaker
- The new utterance sounds like the same speech pronounced by the target speaker
- The last step is the re-synthesis of the signal to reconstruct the “forged” source speech voice

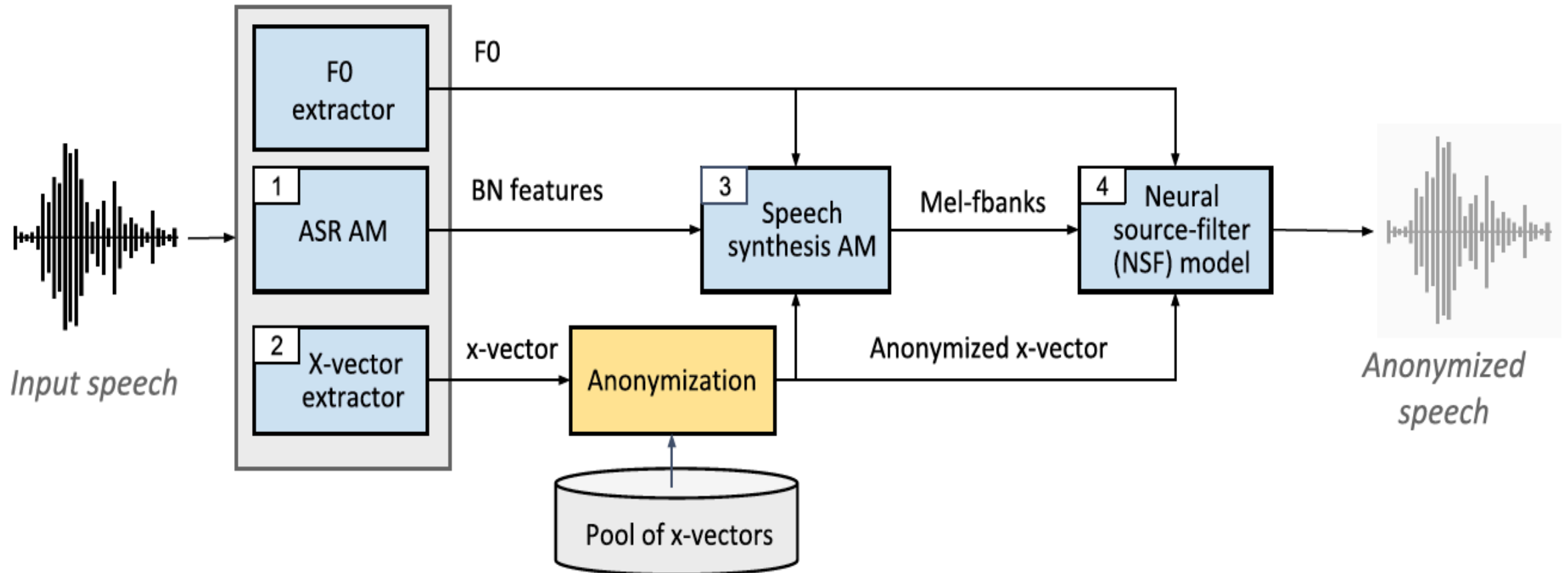




# Sequence-to-Sequence Acoustic Modeling for Voice Conversion

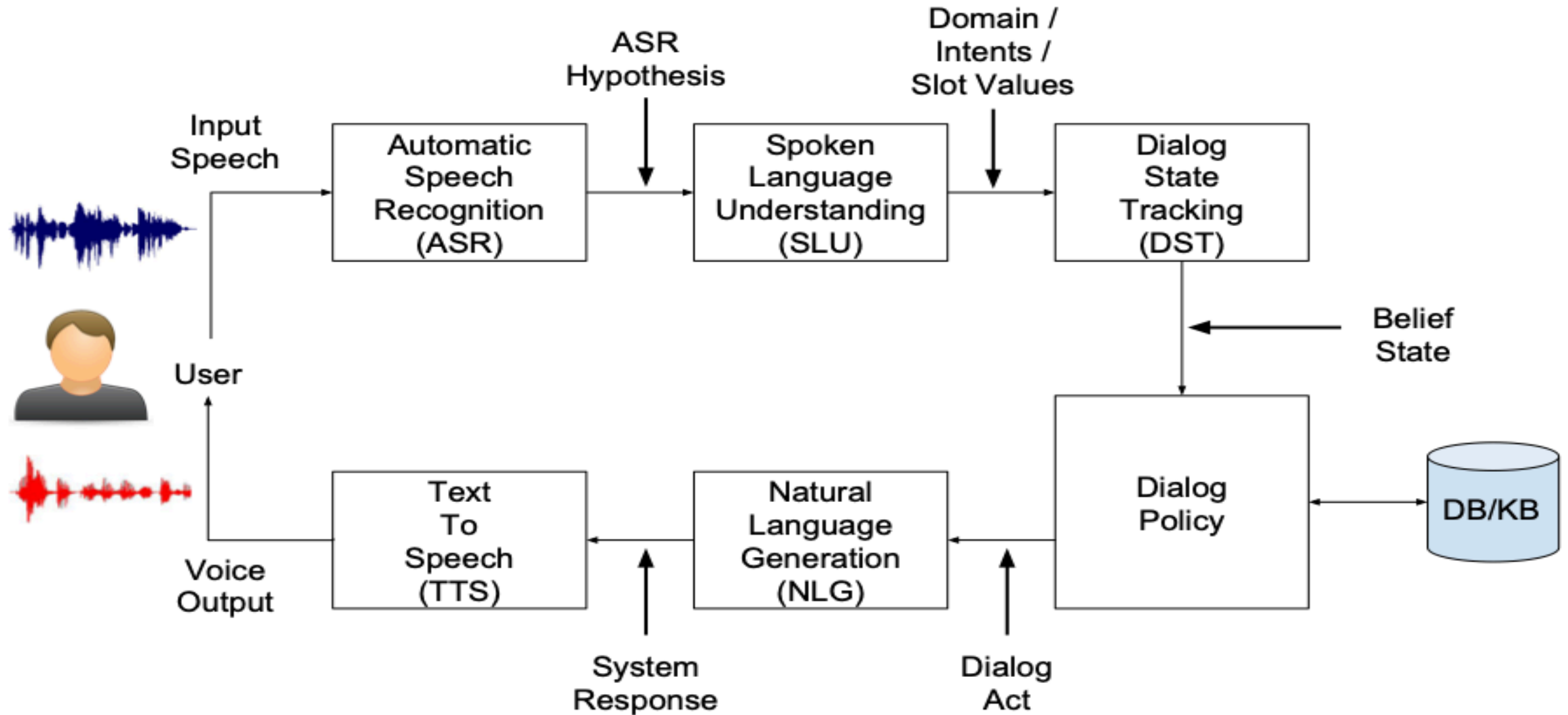


# The VoicePrivacy 2020 Challenge

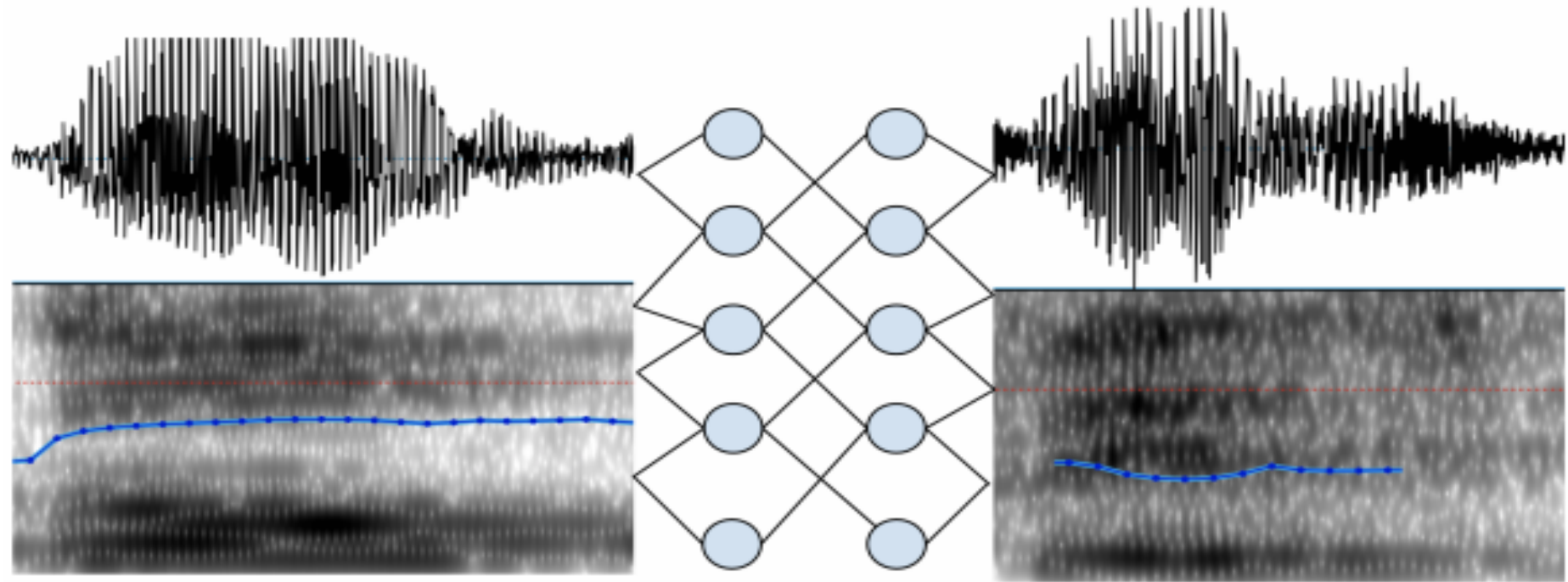


<https://github.com/Voice-Privacy-Challenge/Voice-Privacy-Challenge-2020>

# End-to-End Learning of Task-Oriented Dialogs



# Speech-To-Speech Translation



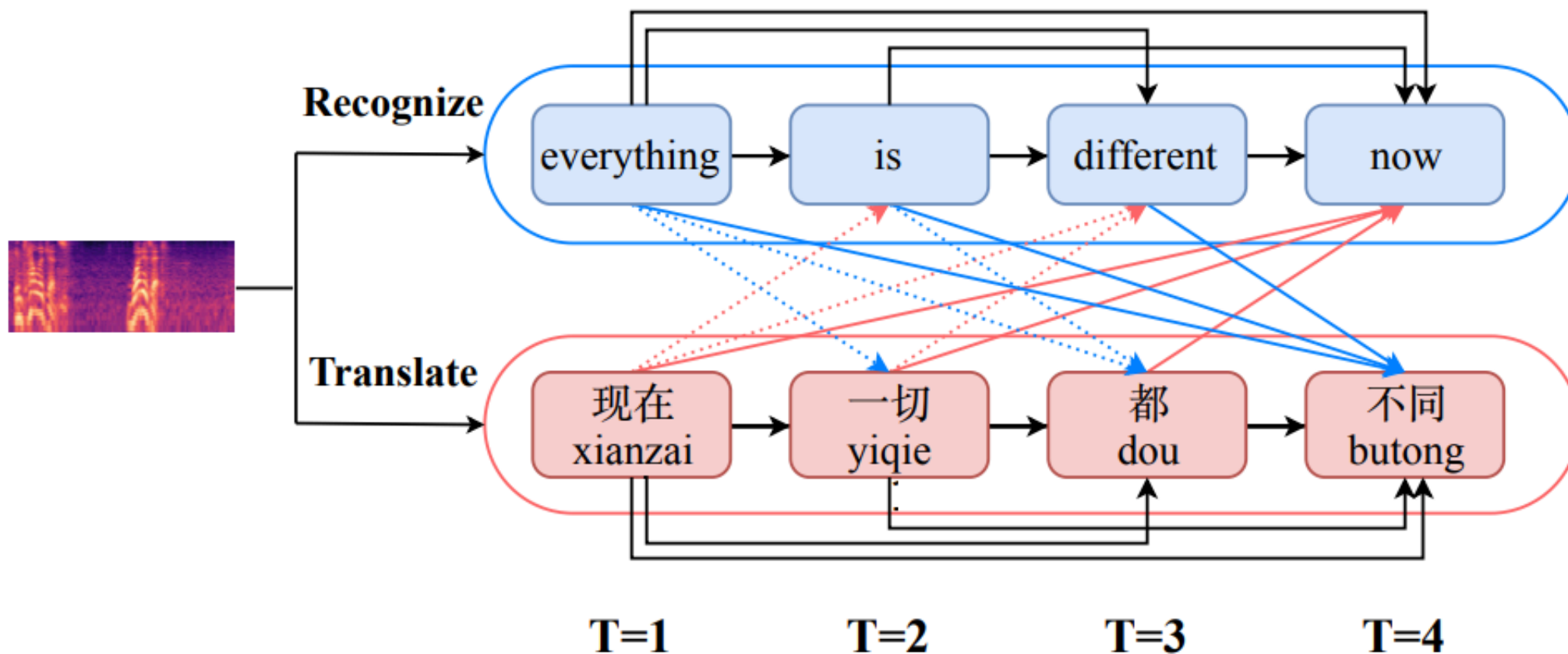
Where is the bus?

公車在哪裡？

[http://web.stanford.edu/class/cs224s/reports/Michelle\\_Guo.pdf](http://web.stanford.edu/class/cs224s/reports/Michelle_Guo.pdf)

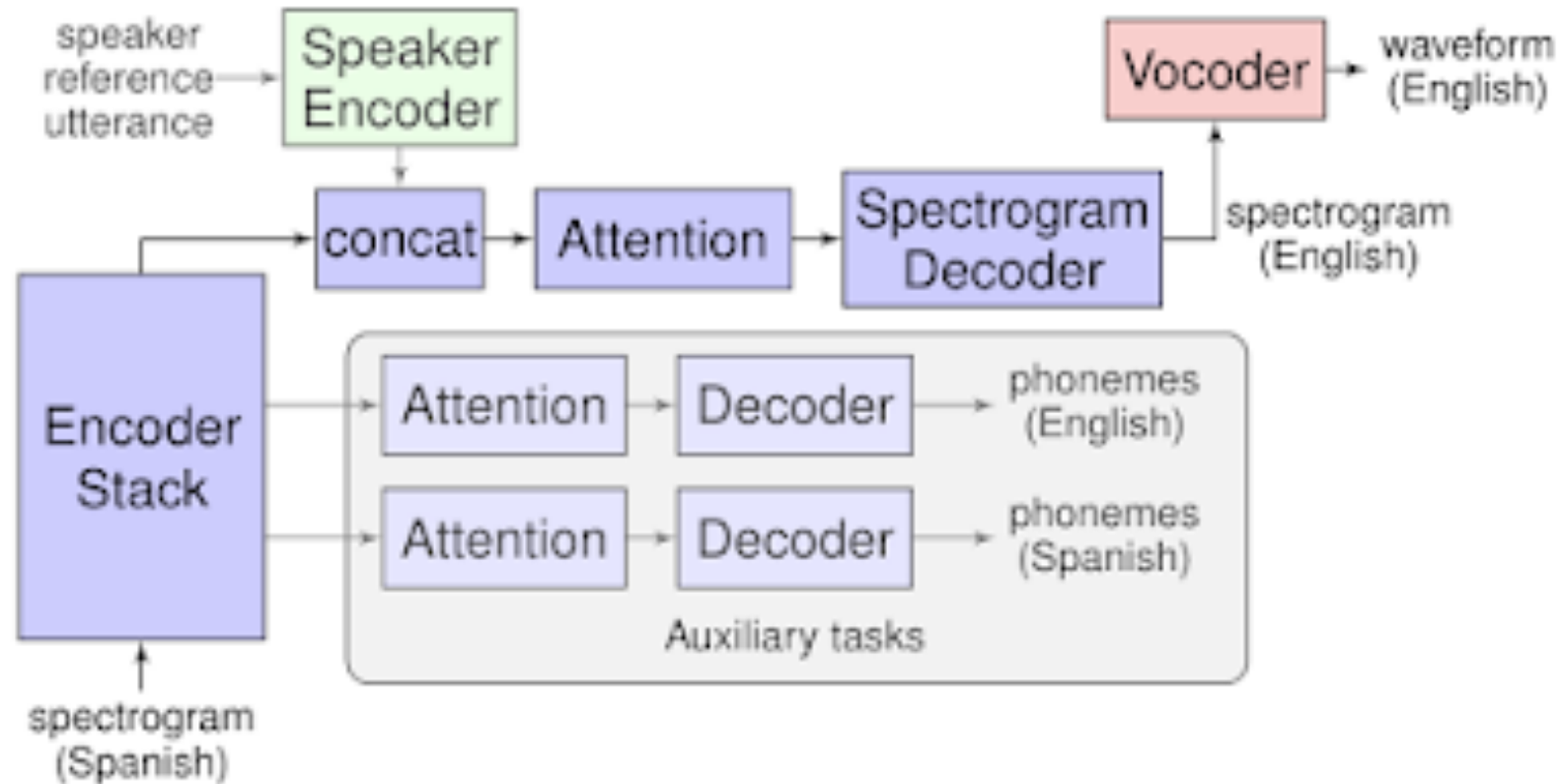
# Synchronous Speech Recognition and Speech-to-Text Translation with Interactive Decoding

<https://arxiv.org/pdf/1912.07240.pdf>



# Translatotron: An End-to-End Speech-to-Speech Translation Model

<https://ai.googleblog.com/2019/05/introducing-translatotron-end-to-end.html>

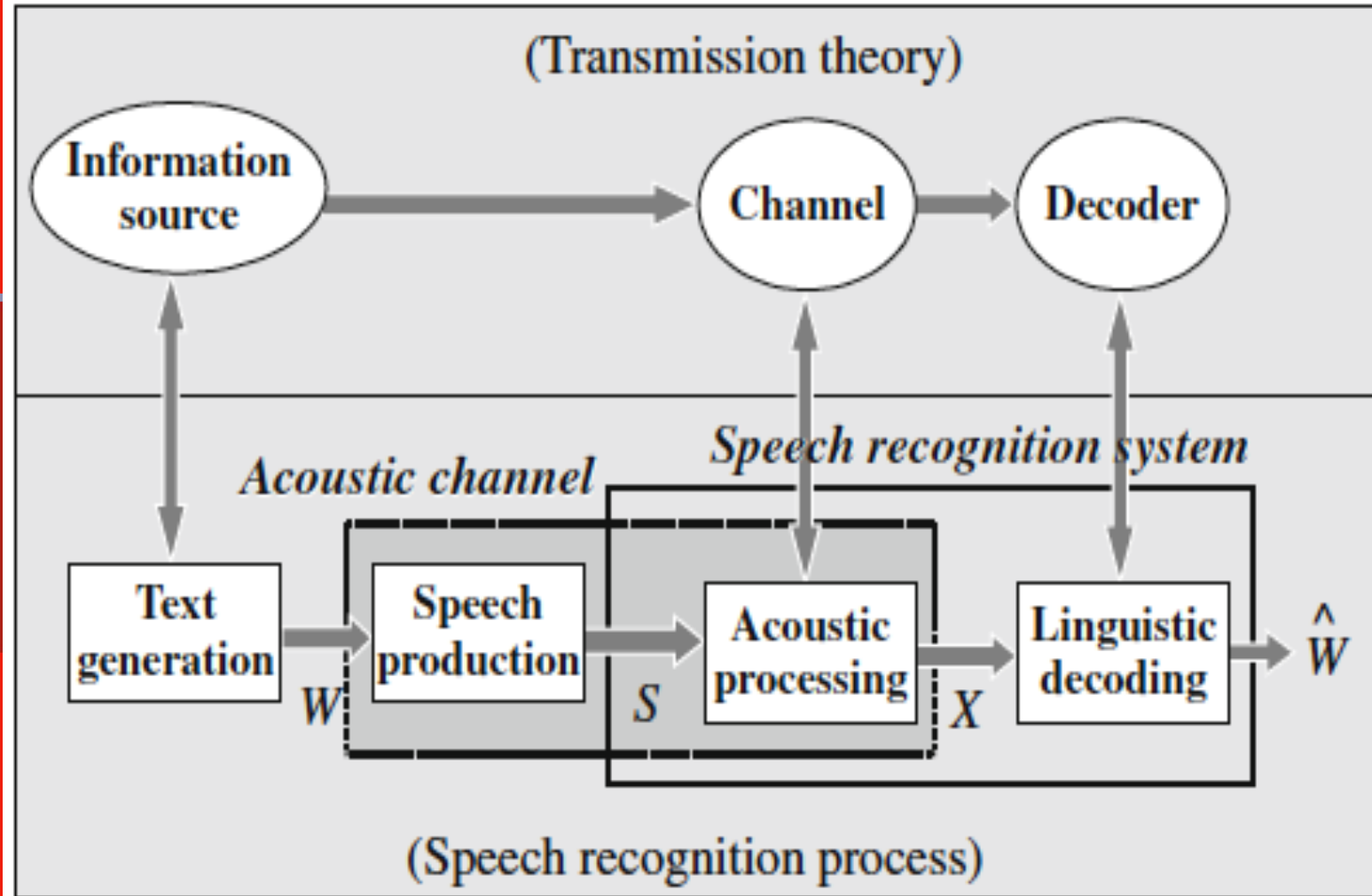


# Speech-to-Text/Symbols

- Automatic Speech Recognition, Wav2Letters, phone alignment,
- Embeddings (acoustic, word, semantic,...)
- Key word spotting,
- Paralinguistic Pattern Recognition:
  - Gender recognition,
  - Age estimation,
  - Emotional state detection,
  - ...
- Speaker identification/verification/diarization,
- Language, dialect, accent recognition,
- ...

# Speech Technology

Theory and Applications



$$\hat{W} = \arg \max_w P(W|X) = \arg \max_w \frac{P(X|W) P(W)}{P(X)}$$



# Towards End-to-End Speech Recognition with Recurrent Neural Networks

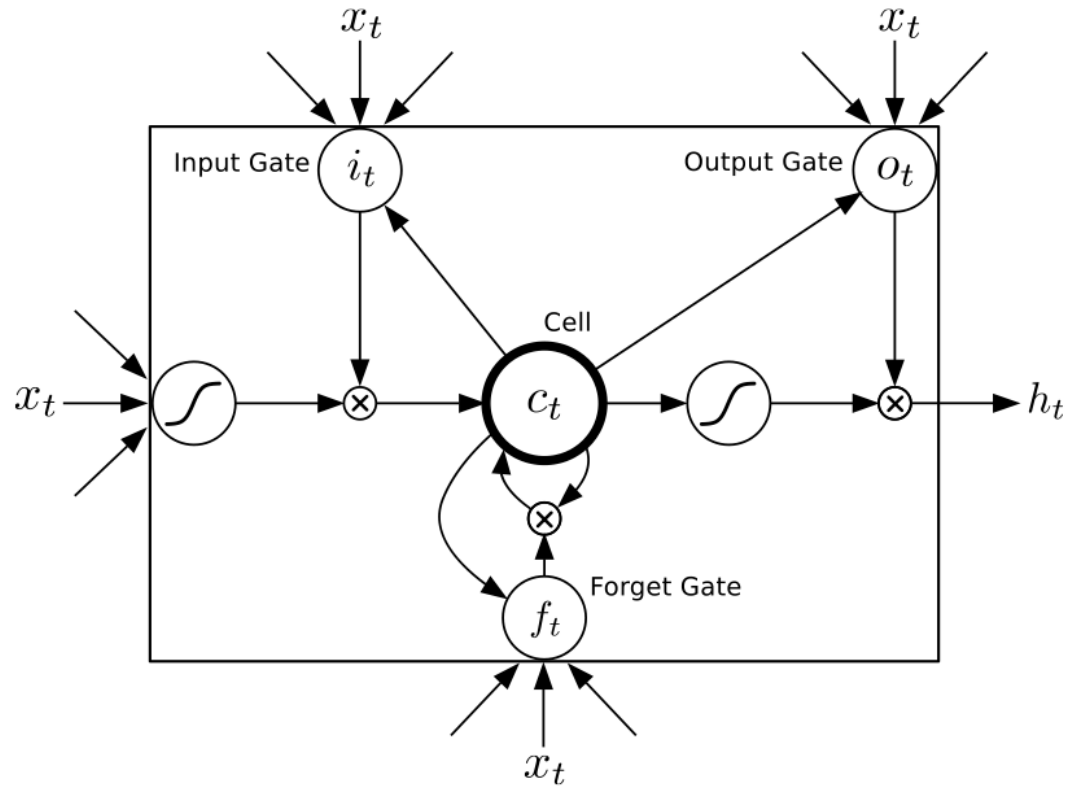


Figure 1. Long Short-term Memory Cell.

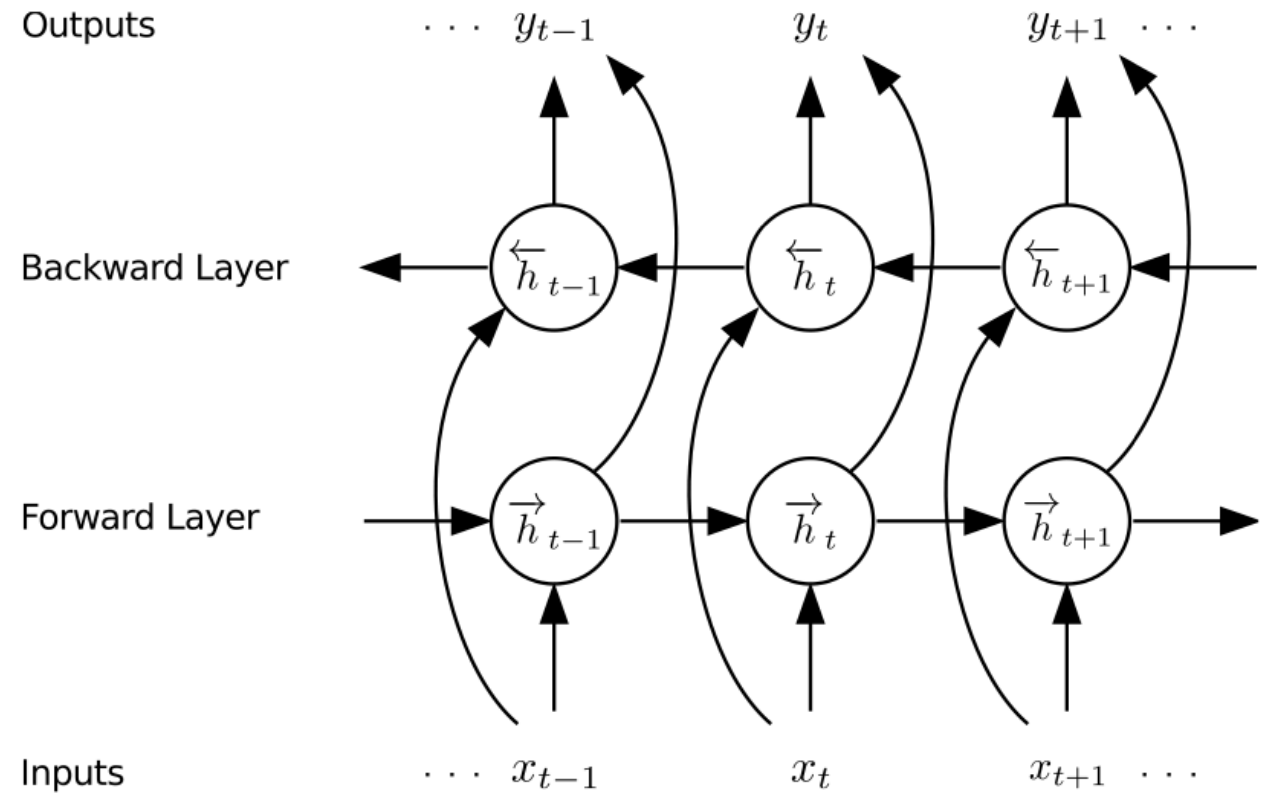
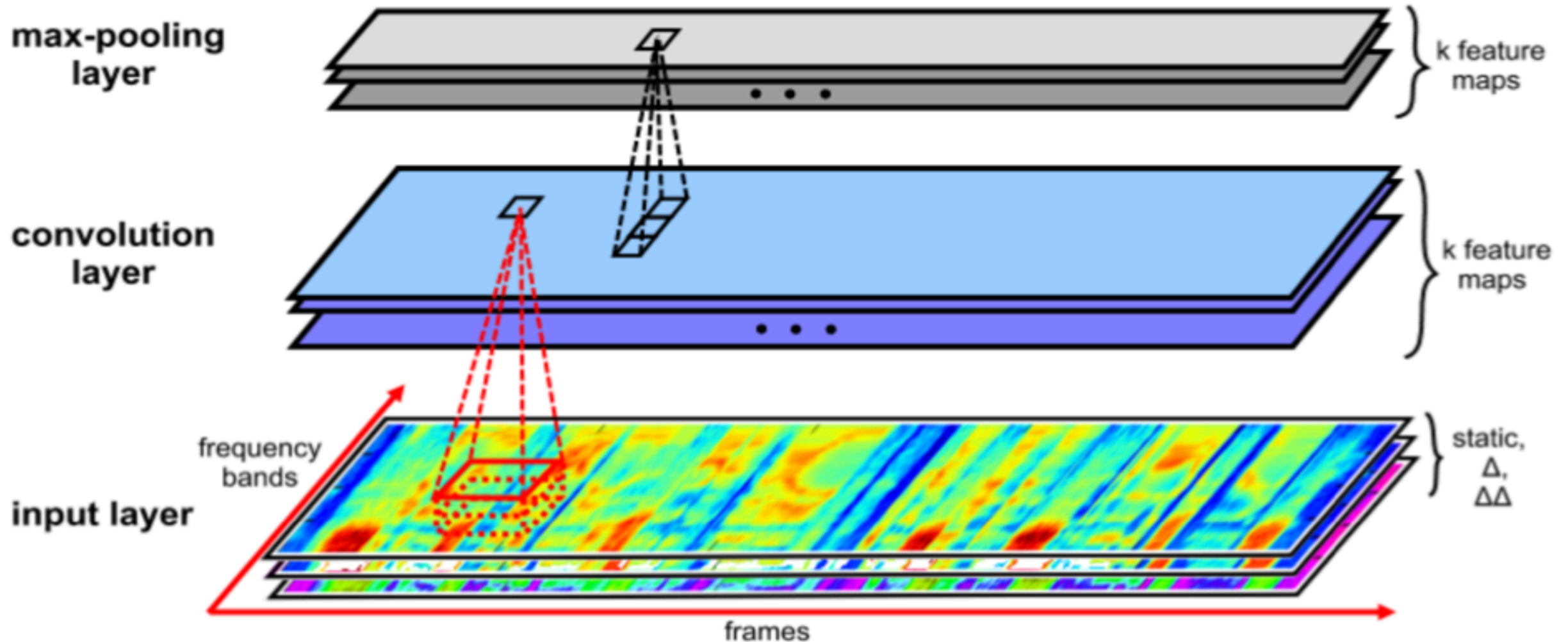
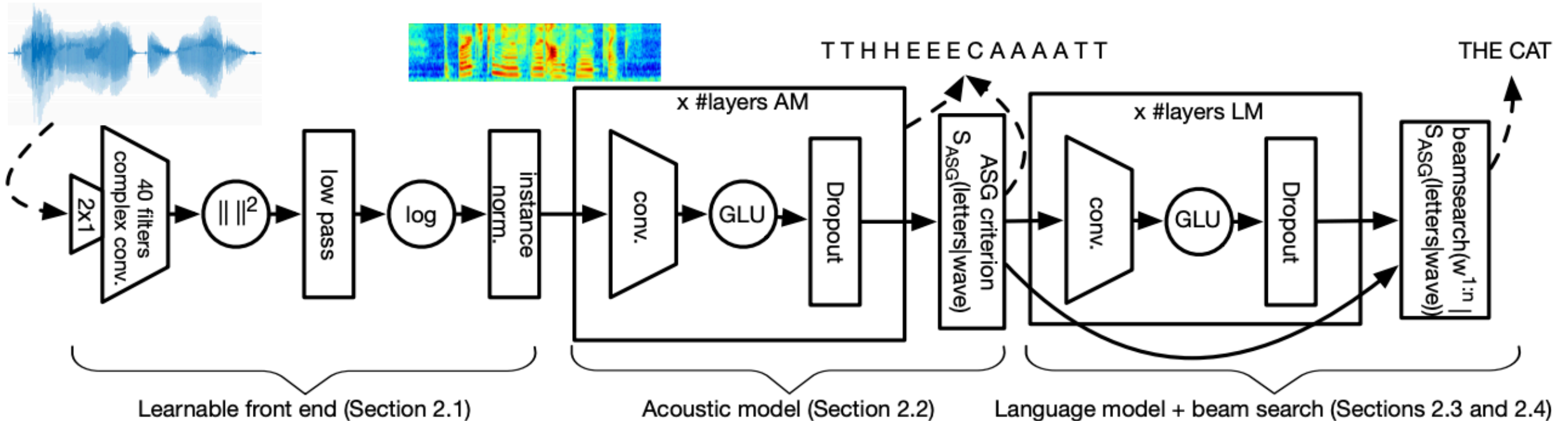


Figure 2. Bidirectional Recurrent Neural Network.

# Towards End-to-End Speech Recognition with Deep Convolutional Neural Networks

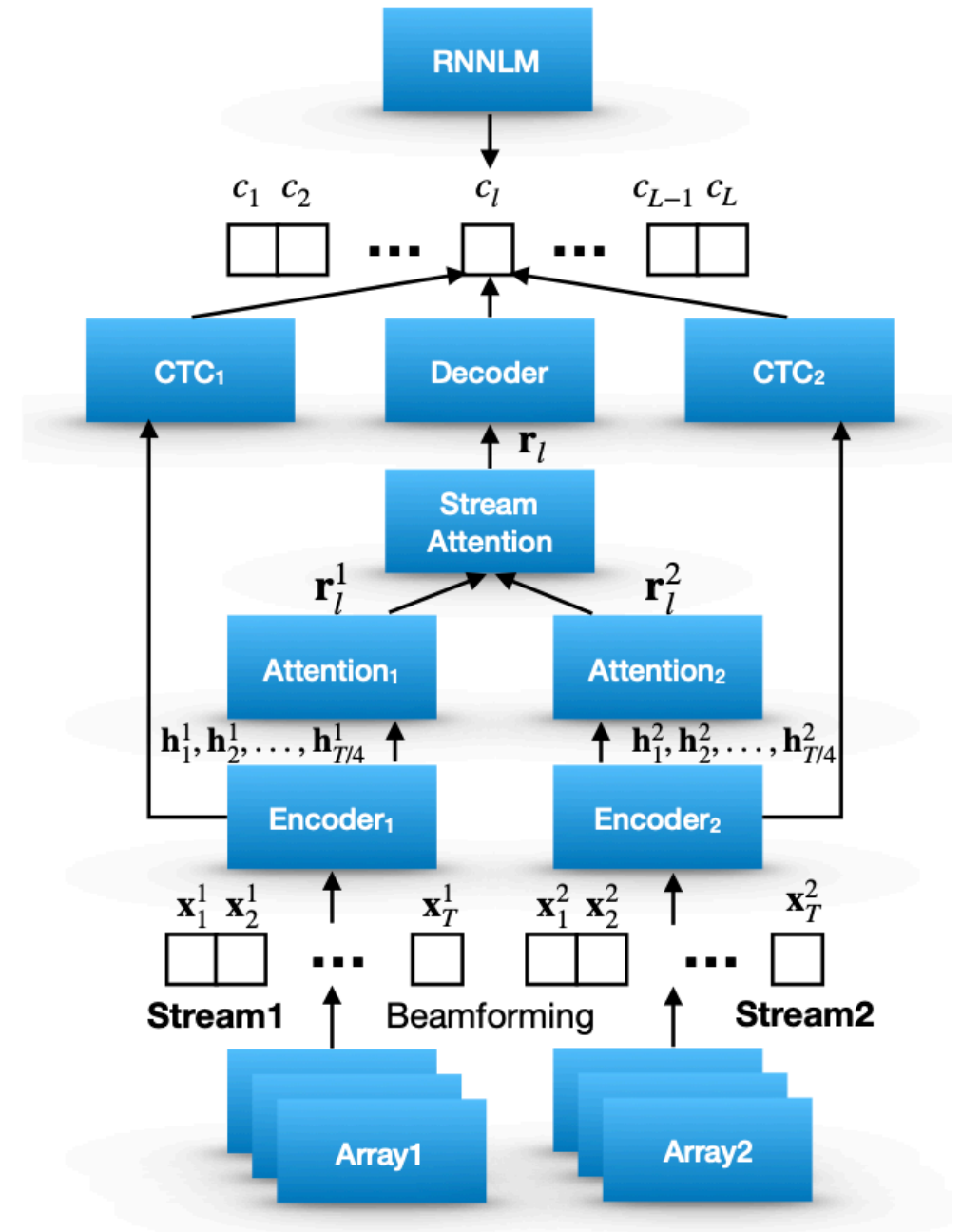


# Facebook Fully Convolutional Speech Recognition



<https://arxiv.org/pdf/1812.06864.pdf>

# STREAM ATTENTION- BASED MULTI-ARRAY END-TO-END SPEECH RECOGNITION

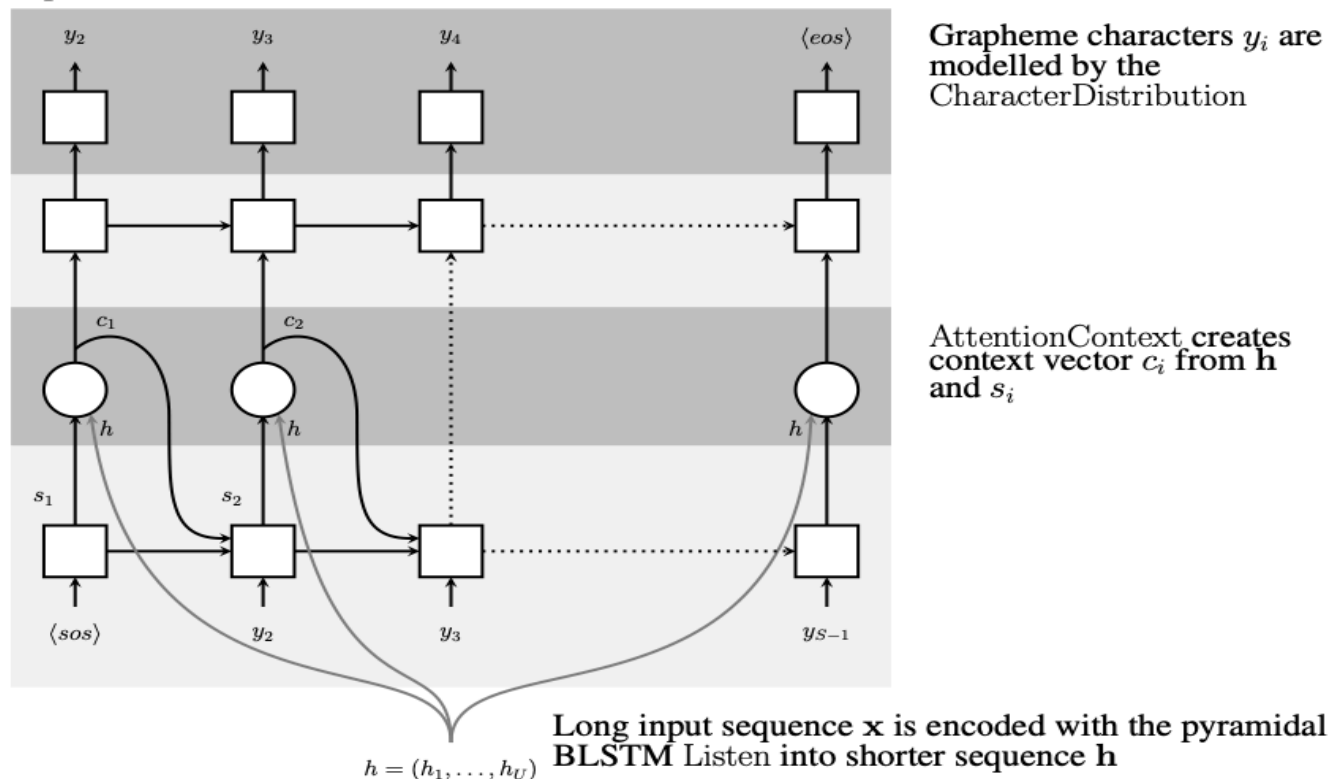


# Listen, Attend and Spell

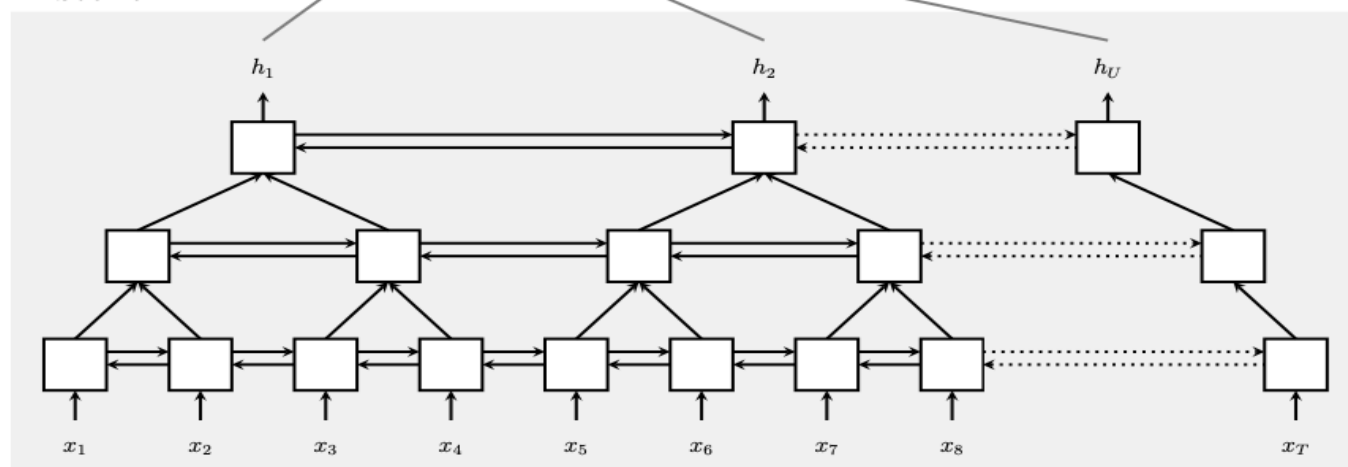
Google Brain

<https://arxiv.org/pdf/1508.01211.pdf>

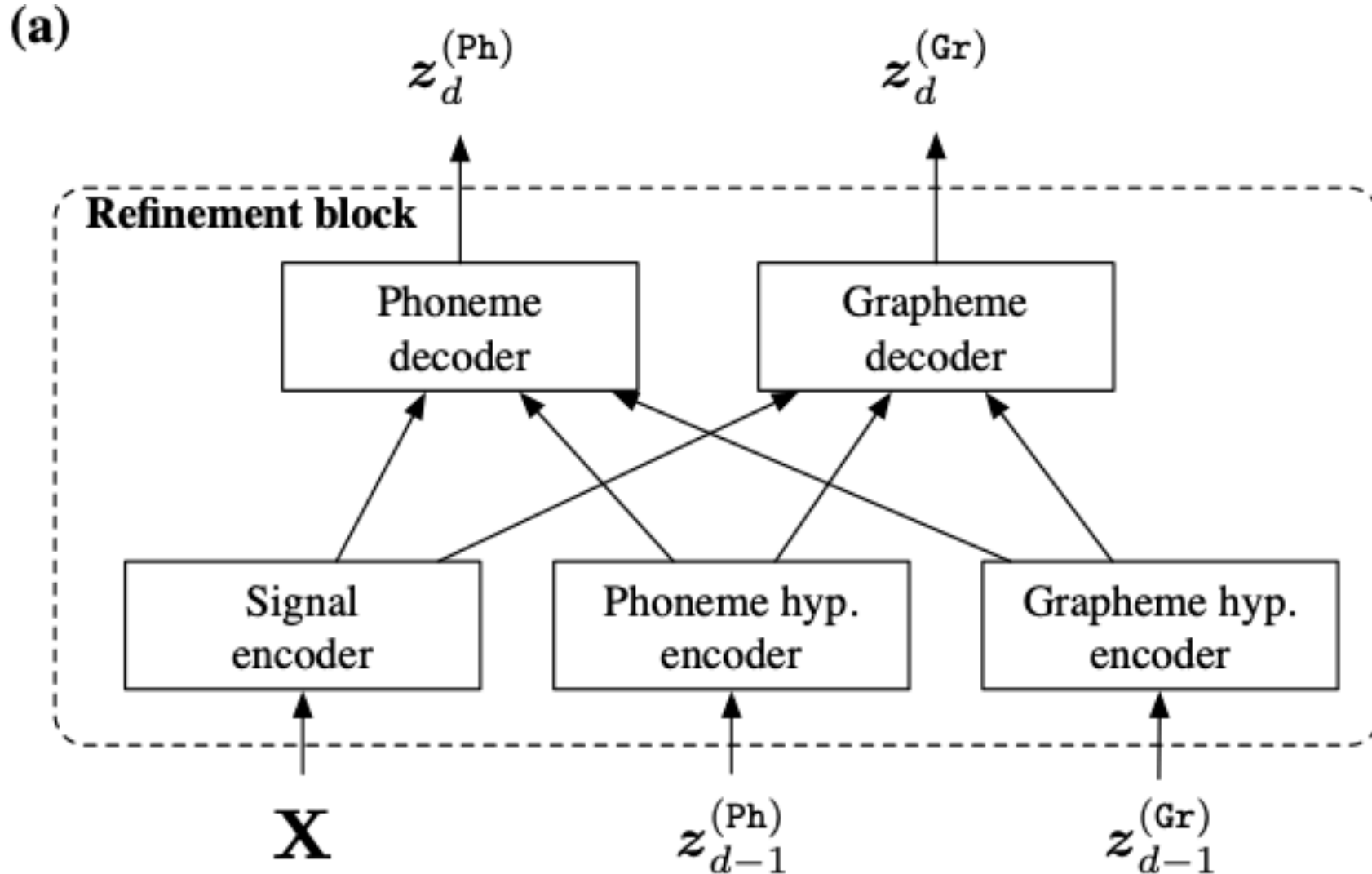
## Speller



## Listener



# JOINT PHONEME-GRAPHEME MODEL FOR END-TO-END SPEECH RECOGNITION

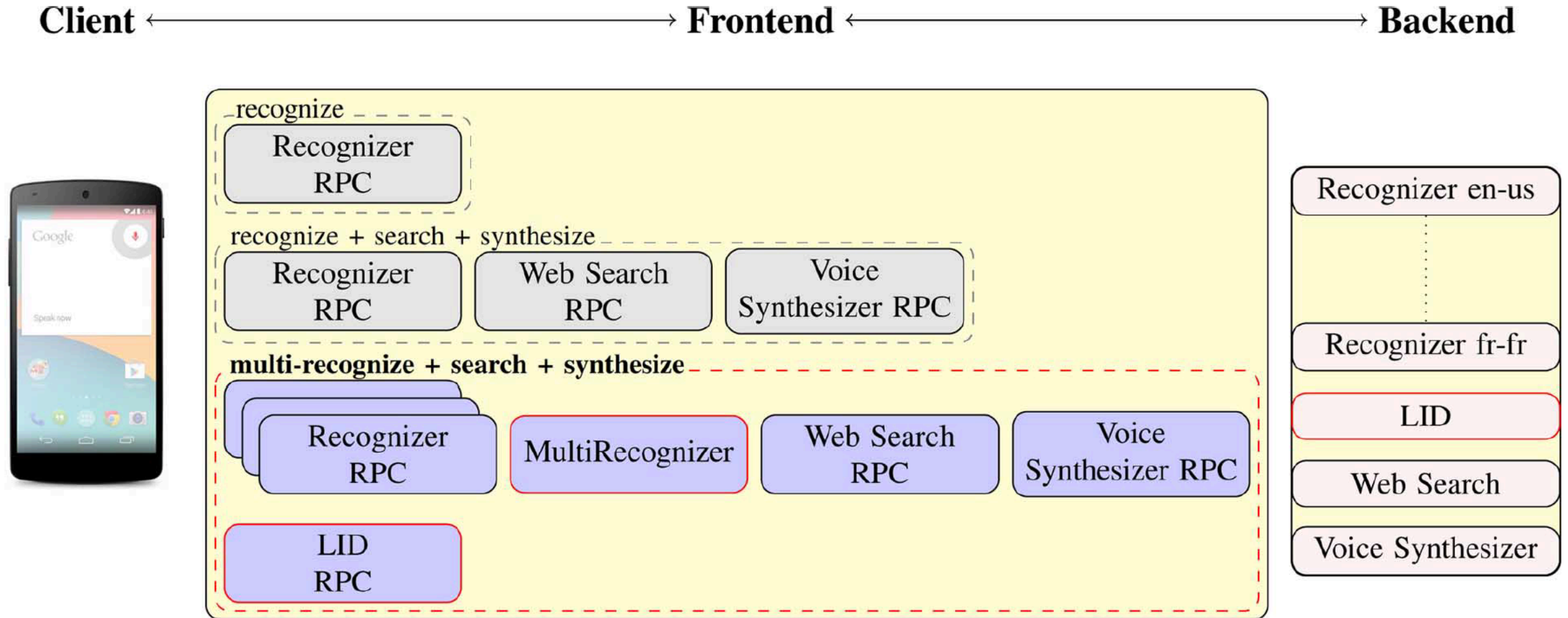


# End-to-End Multilingual Multi- Speaker Speech Recognition





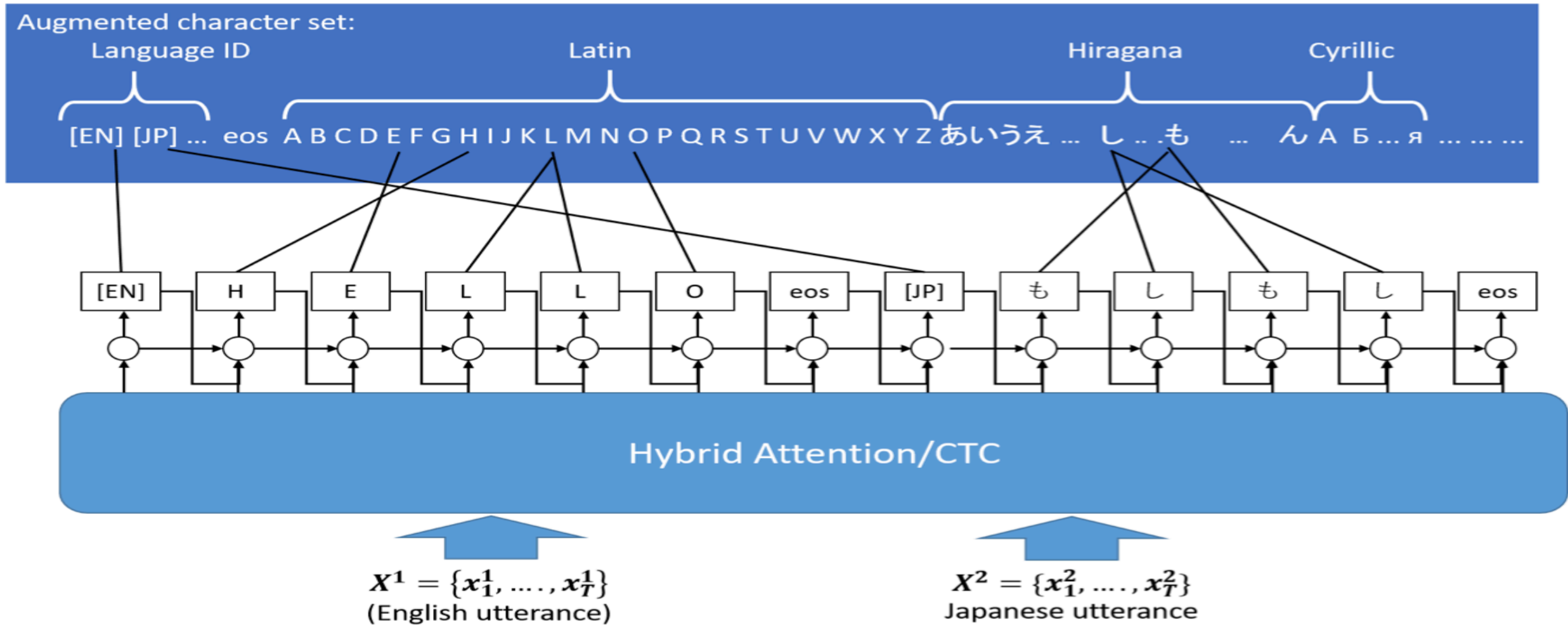
# A Real-Time End-to-End Multilingual Speech Recognition Architecture



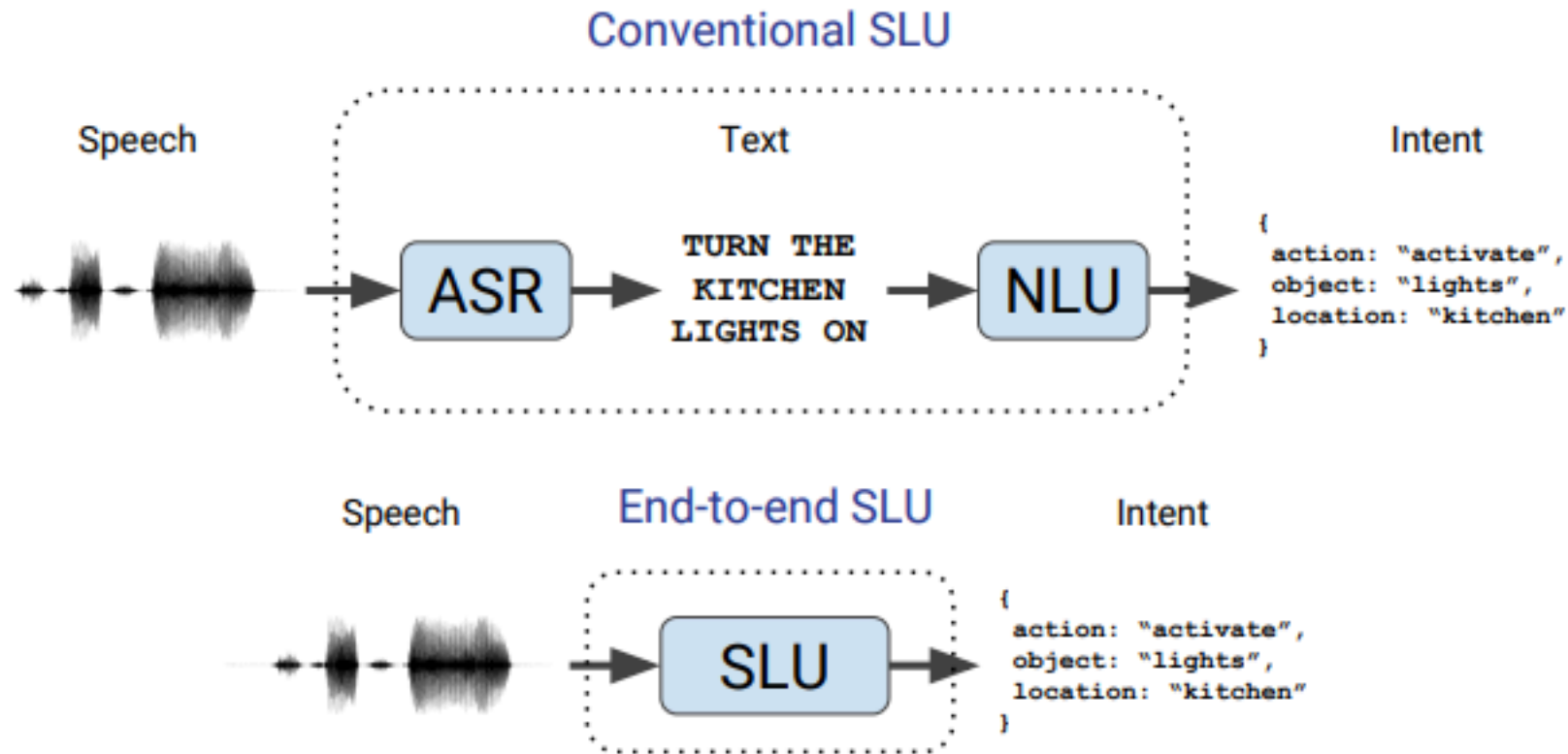
[https://www.researchgate.net/publication/276298472\\_A\\_Real-Time\\_End-to-End\\_Multilingual\\_Speech\\_Recognition\\_Architecture](https://www.researchgate.net/publication/276298472_A_Real-Time_End-to-End_Multilingual_Speech_Recognition_Architecture)



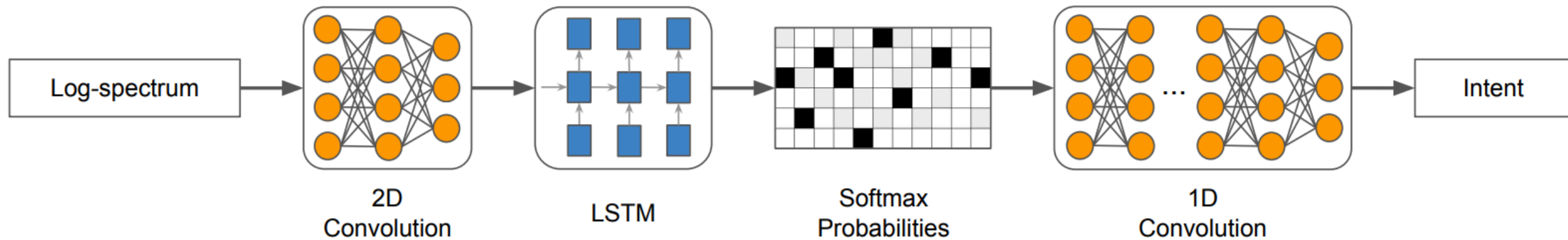
# Language Independent End-to-End Architecture For Joint Language and Speech Recognition



# Speech Model Pre-training for End-to-End Spoken Language Understanding



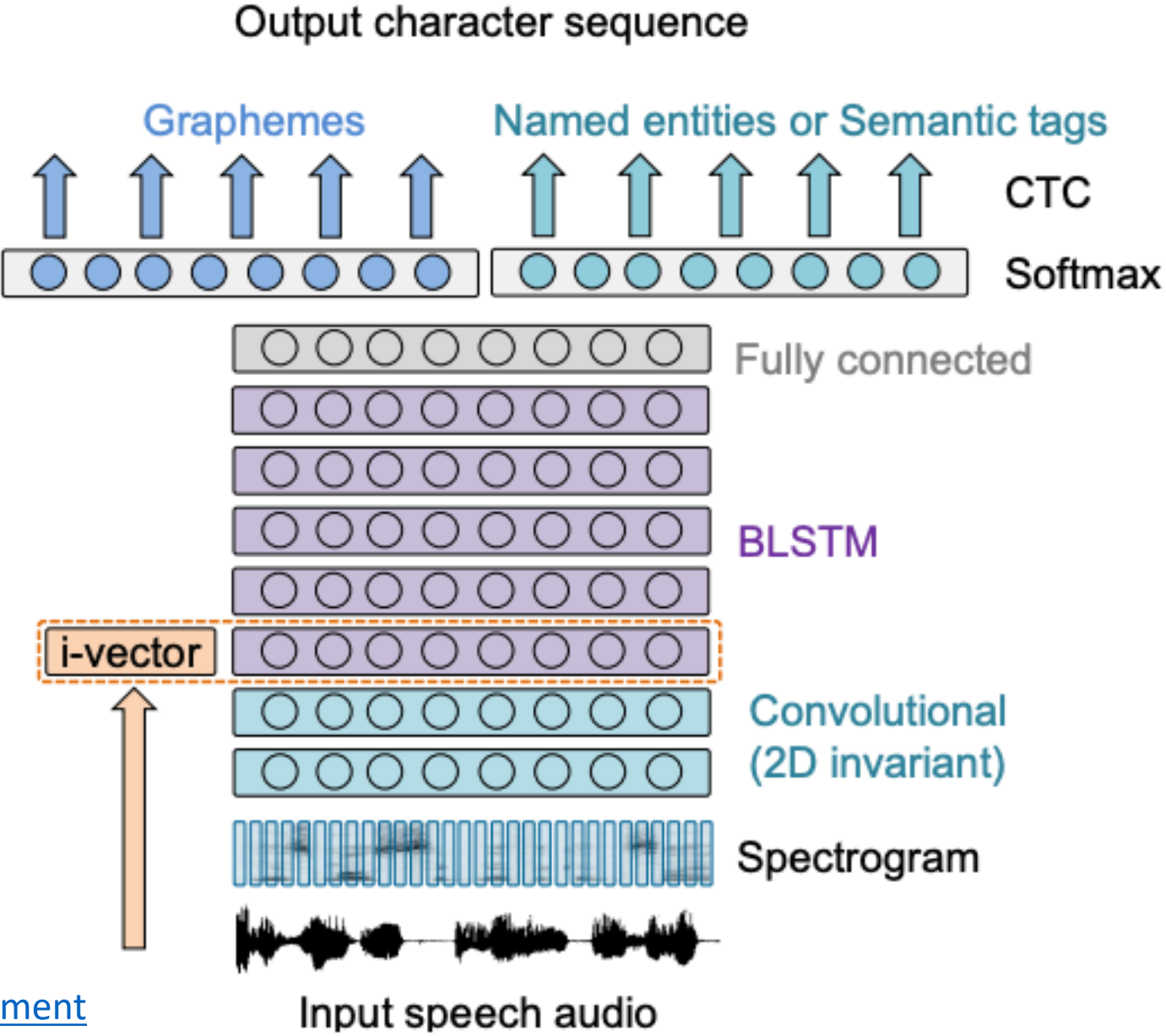
# FROM AUDIO TO SEMANTICS: APPROACHES TO END-TO-END SPOKEN LANGUAGE UNDERSTANDING



<https://arxiv.org/pdf/1809.09190.pdf>

[http://yp-chen.com/static/YPChen18ICASSP\\_Spoken-Language-Understanding.pdf](http://yp-chen.com/static/YPChen18ICASSP_Spoken-Language-Understanding.pdf)

# Recent Advances in End-to-End Spoken Language Understanding



<https://arxiv.org/pdf/1909.13332.pdf>

<https://hal.archives-ouvertes.fr/hal-02307811/document>

Generic concepts

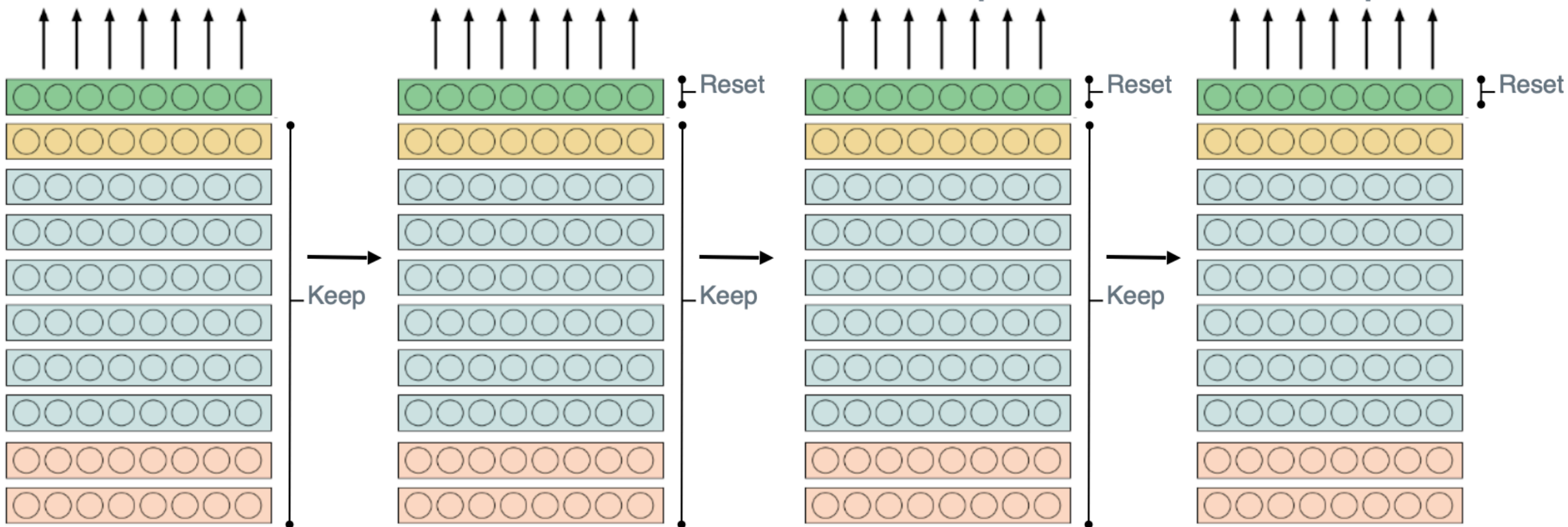
Specific concepts

(ASR)  
Character sequence

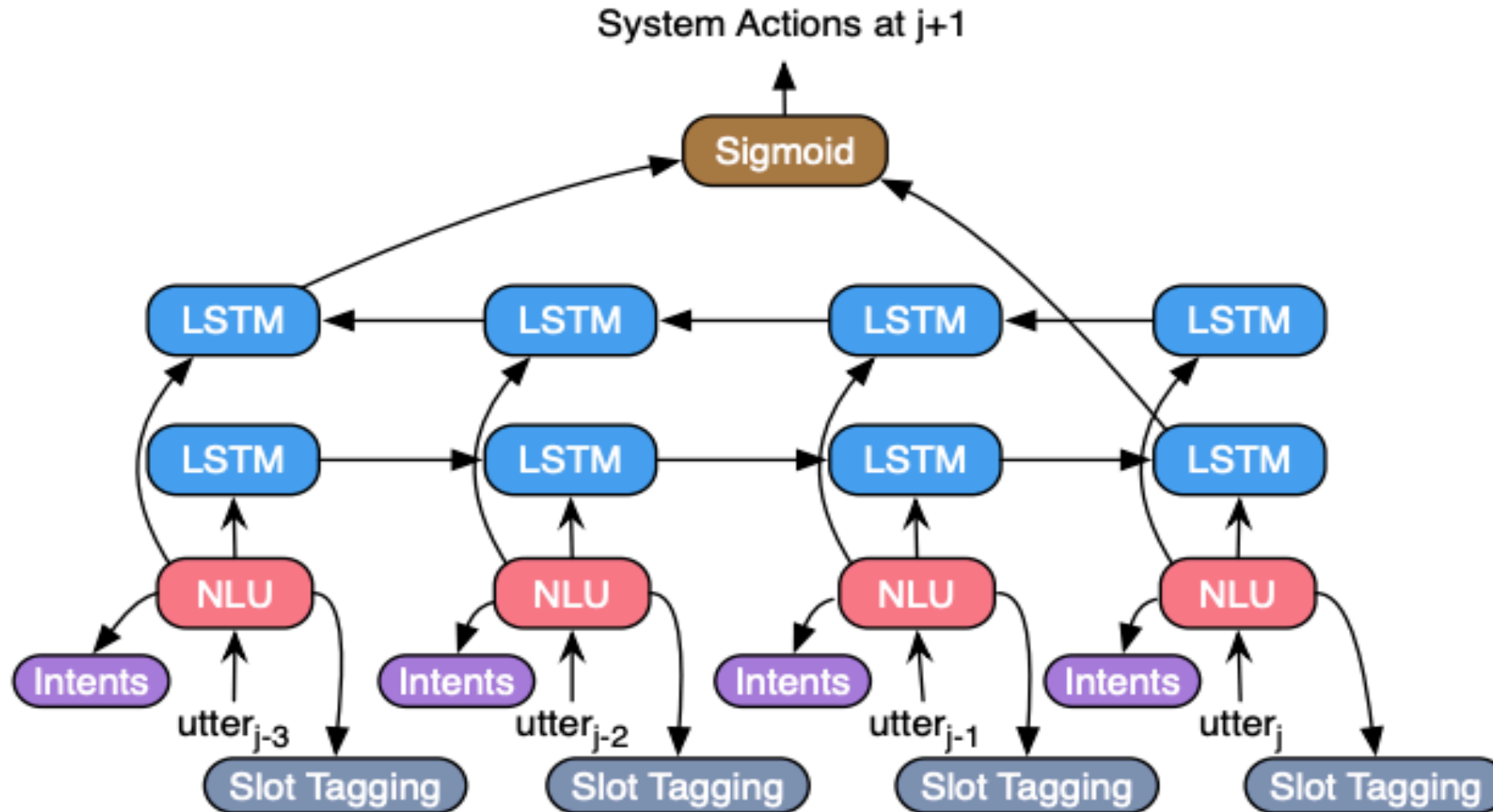
(NER)  
Character sequence  
& named entity

(SC\_mer)  
Character sequence  
& merged semantic  
concepts

(M)  
Character sequence  
& target semantic  
concepts



# END-TO-END JOINT LEARNING OF NATURAL LANGUAGE UNDERSTANDING AND DIALOGUE MANAGER

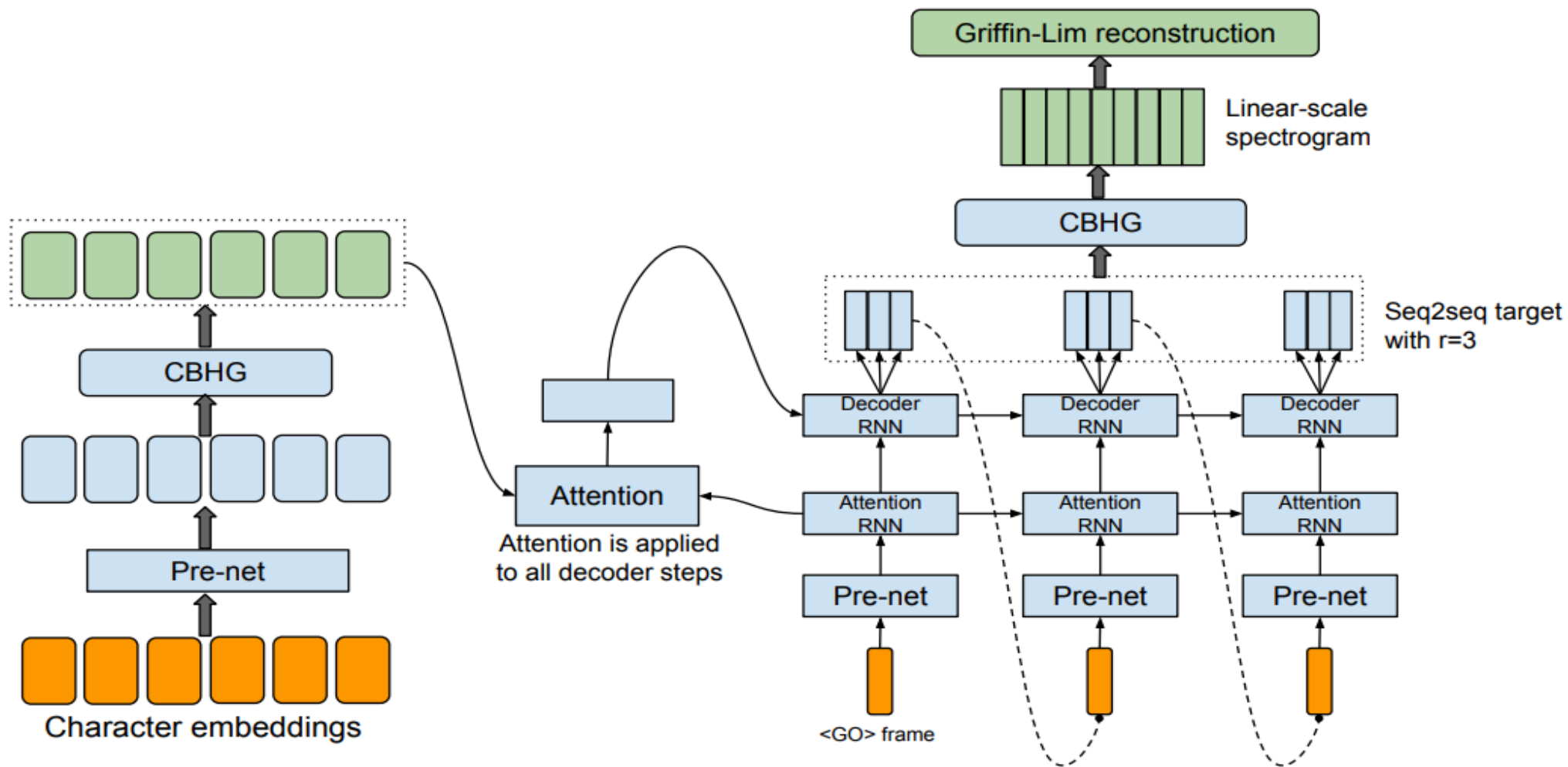


# Text-to-Speech

- End-to-End TTS,
- Text normalisation,
- Lexical ambiguities/Heterophone Homographs,
- Graphemes-to-Phonemes,
- Prosody prediction,
- Emotional Speech,
- Personalisation,
- ...

# <https://google.github.io/tacotron/>

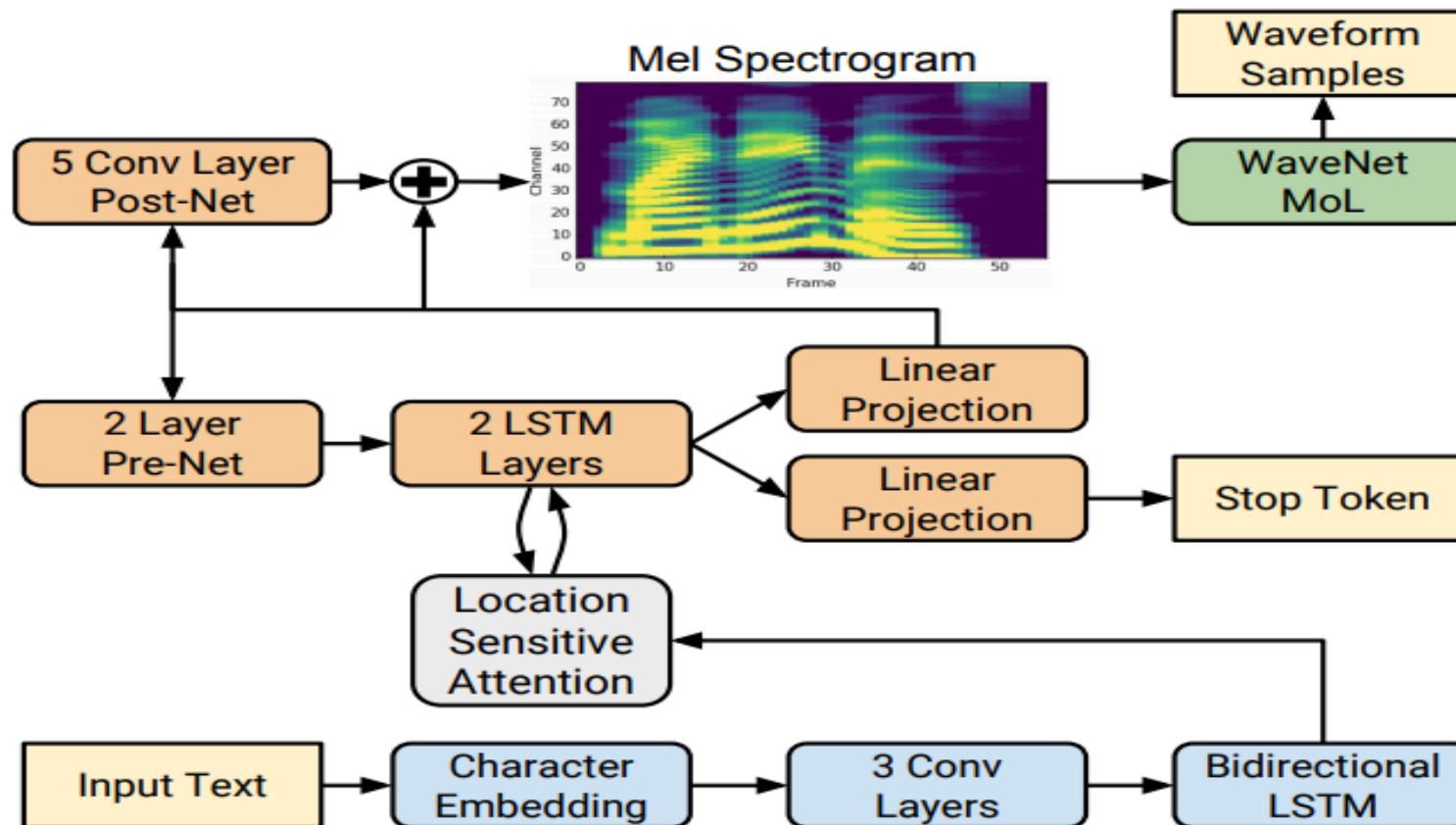
<https://arxiv.org/pdf/1703.10135.pdf>





# NATURAL TTS SYNTHESIS BY CONDITIONING WAVENET ON MEL SPECTROGRAM PREDICTIONS

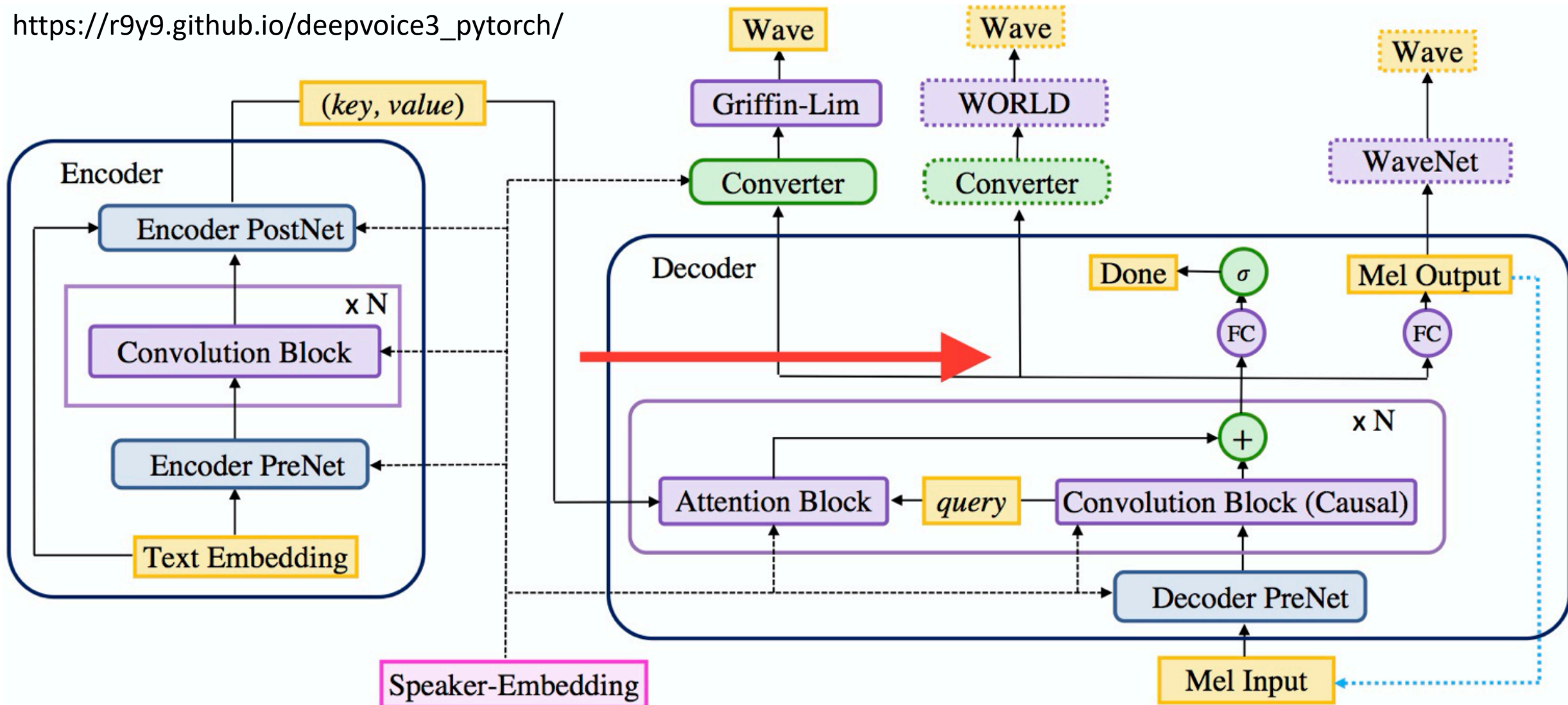
<https://arxiv.org/pdf/1712.05884.pdf>



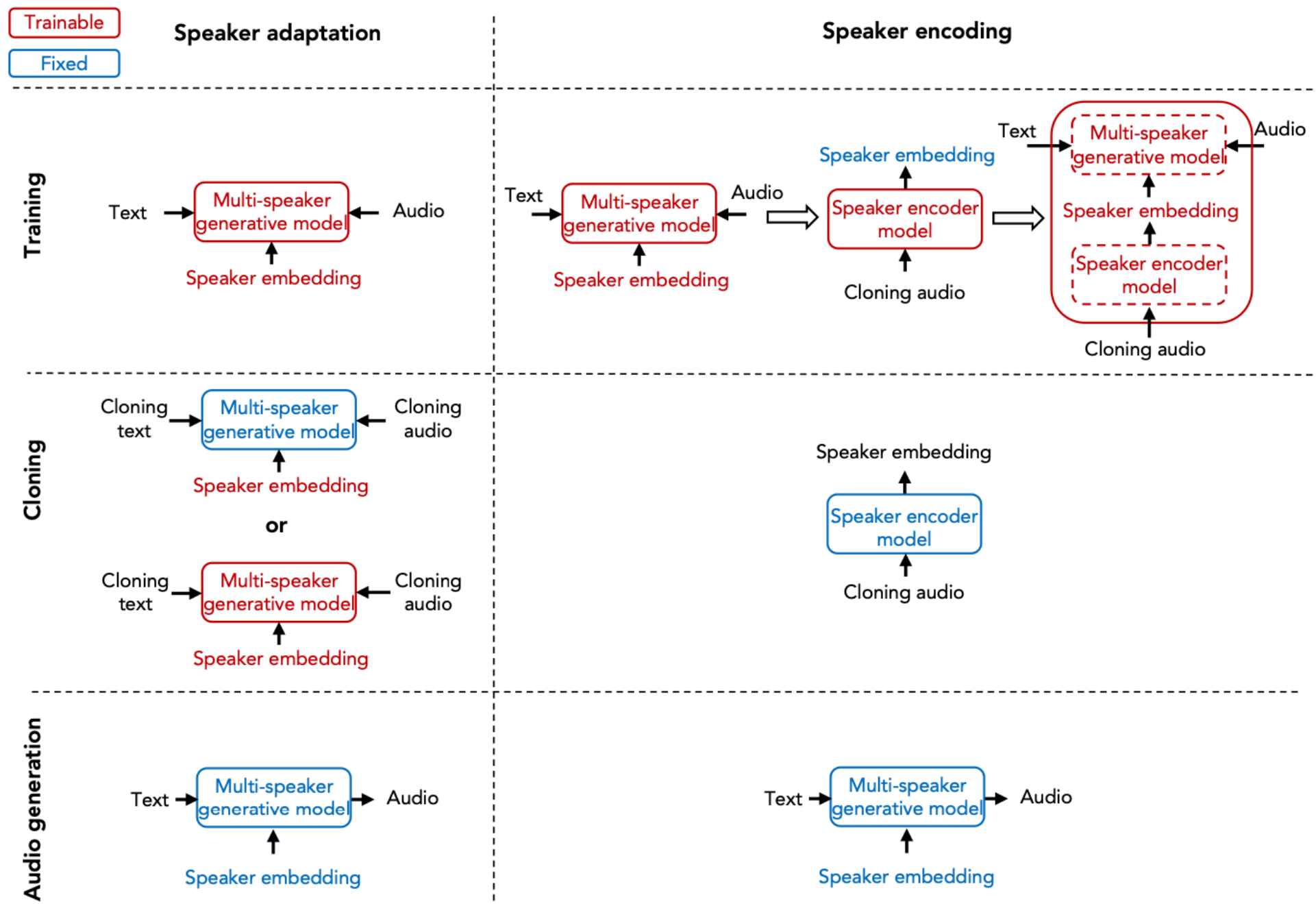


# Deep Voice 3: 2000-Speaker Neural Text-to-Speech

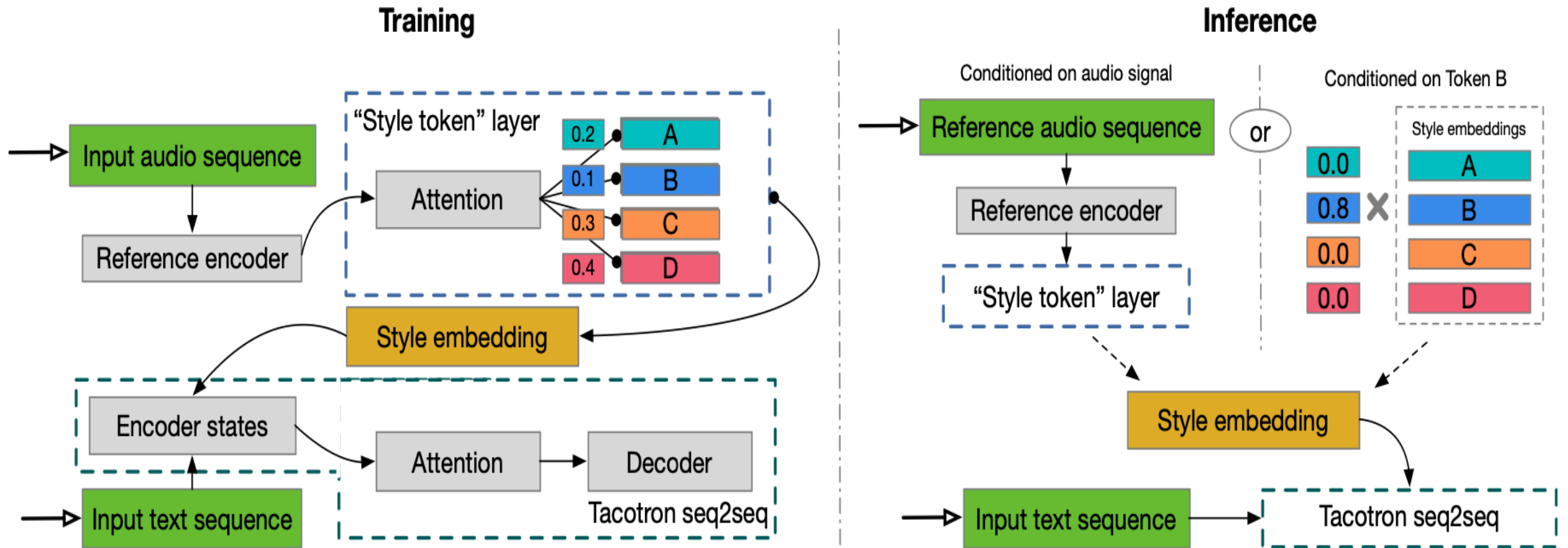
[https://r9y9.github.io/deepvoice3\\_pytorch/](https://r9y9.github.io/deepvoice3_pytorch/)



# Neural Voice Cloning with a Few Samples



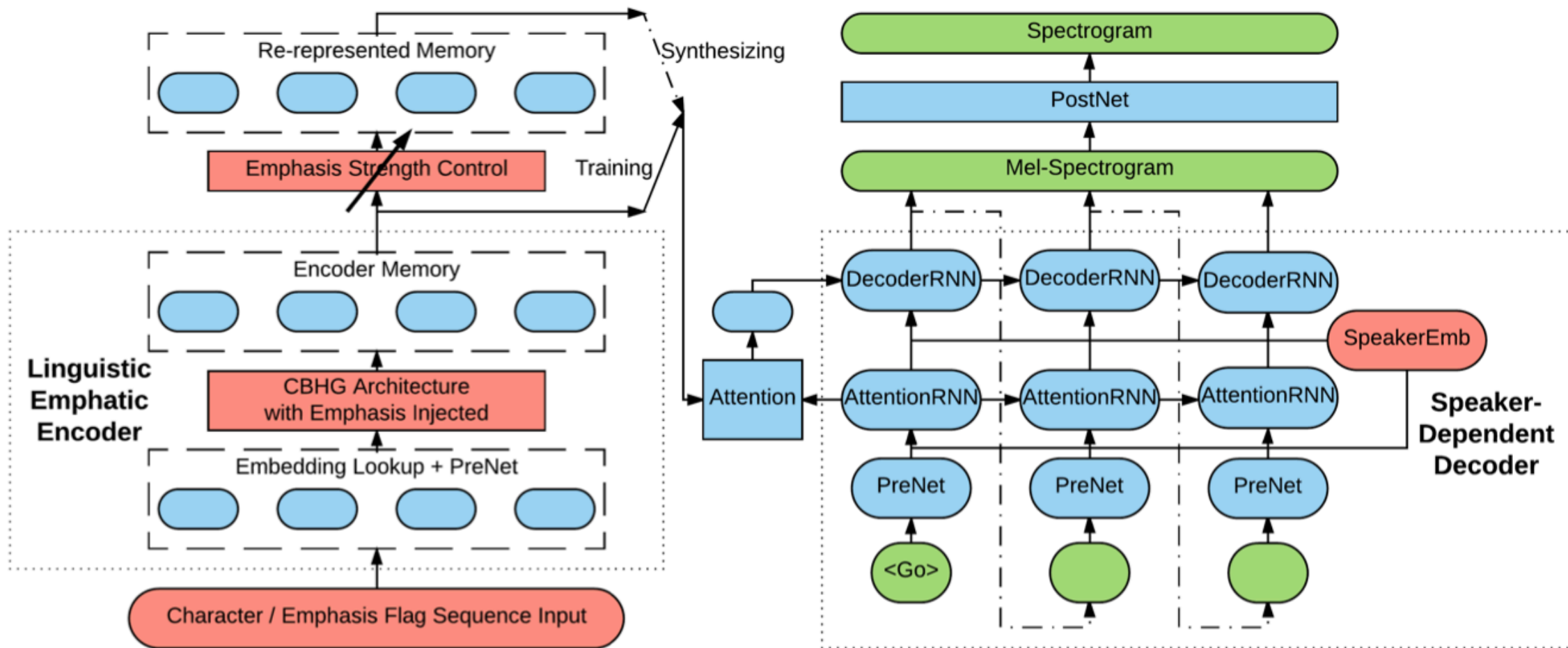
# Style Tokens: Unsupervised Style Modeling, Control and Transfer in End-to-End Speech Synthesis



<https://arxiv.org/pdf/1803.09017.pdf>

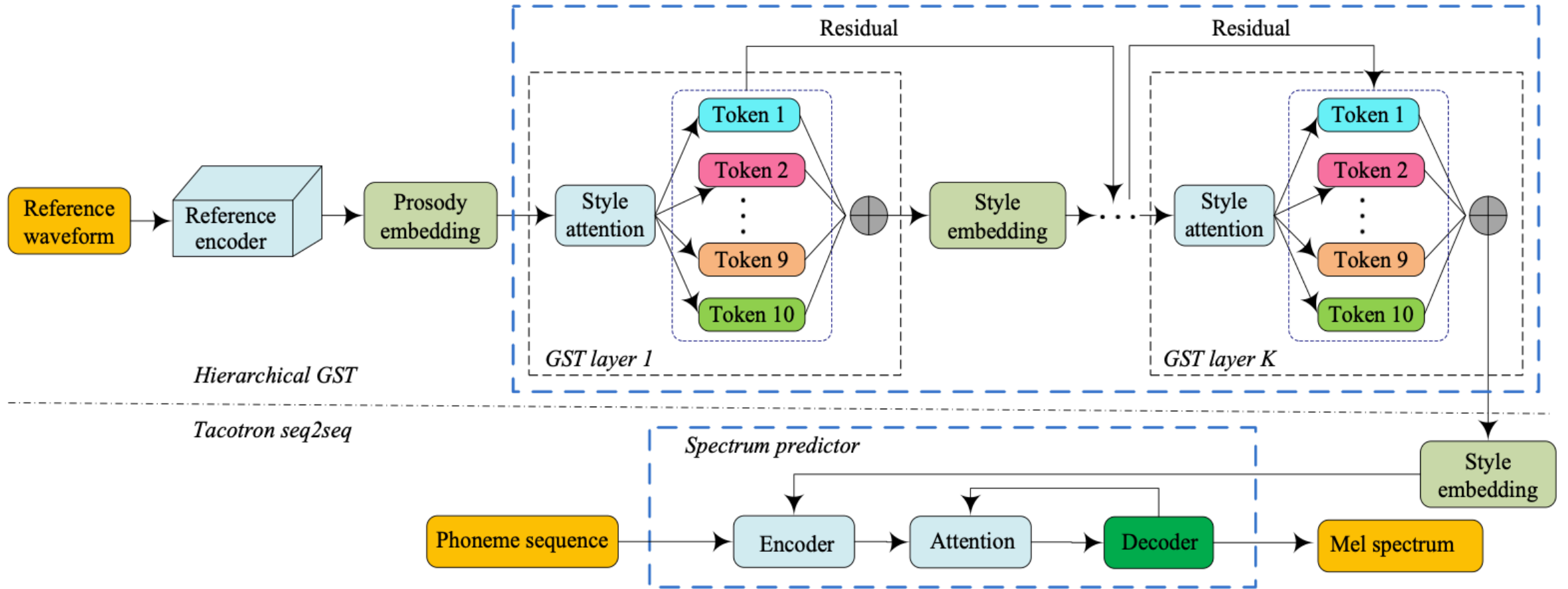
[https://google.github.io/tacotron/publications/global\\_style\\_tokens](https://google.github.io/tacotron/publications/global_style_tokens)

# Emphatic and Control Based on Characteristic Transferring in End-to-End Speech Synthesis



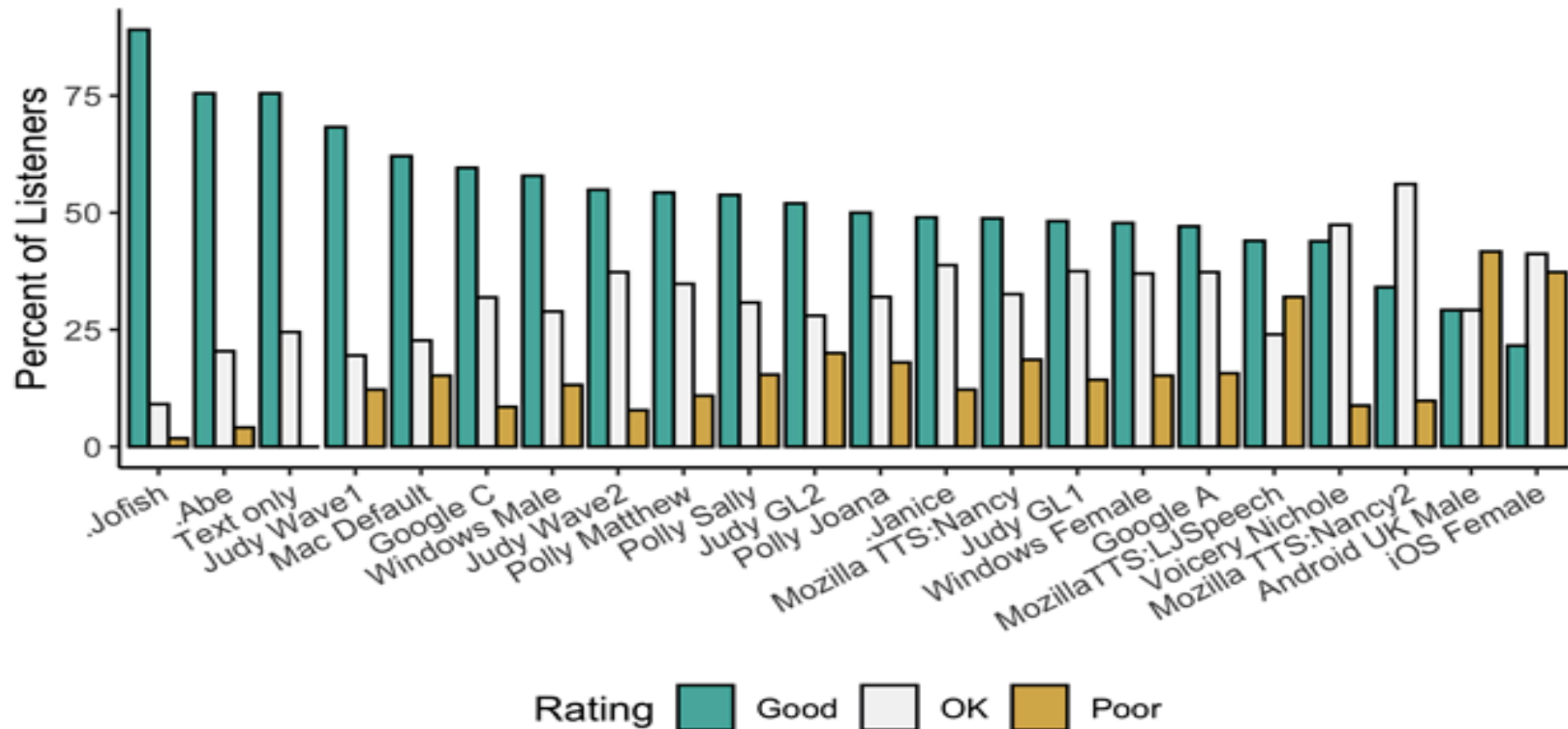


# LEARNING HIERARCHICAL REPRESENTATIONS FOR EXPRESSIVE SPEAKING STYLE IN END-TO-END SPEECH SYNTHESIS



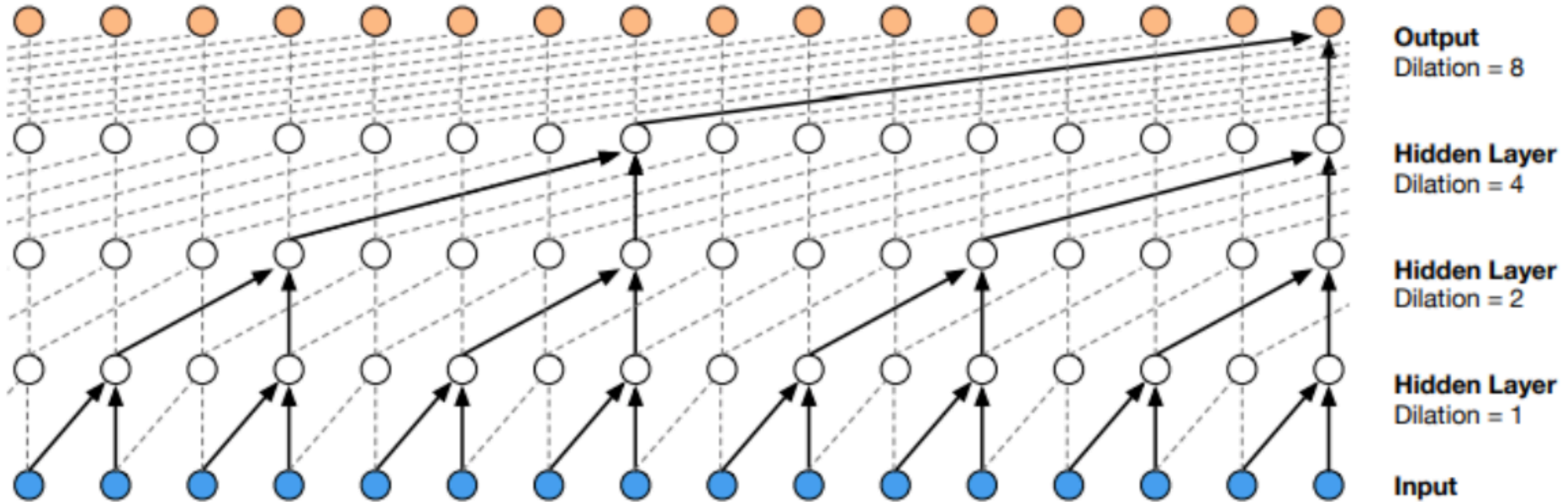
# Mozilla Deep Learning based Text2Speech model

<https://github.com/mozilla/TTS>





# WAVENET: A GENERATIVE MODEL FOR RAW AUDIO



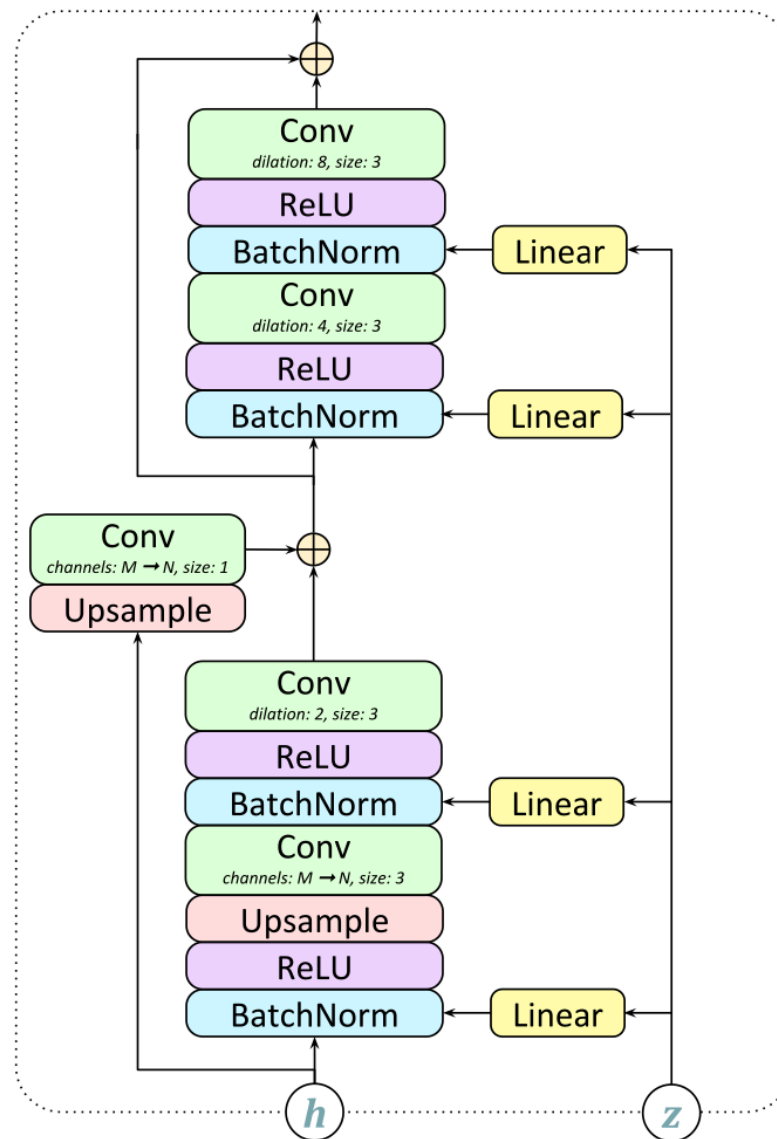
<https://lh3.googleusercontent.com/XMJIFDtIFU7WtZYnPDMDhT6jTFEXvTpY7c7sY70PtTkzEIDm6DvUv5nvHHwbOrUrDckS34alu5fiZx0615p1-nf3rFxxhza387ud=w2048-rw-v1>

<https://www.deepmind.com/blog/article/wavenet-generative-model-raw-audio>

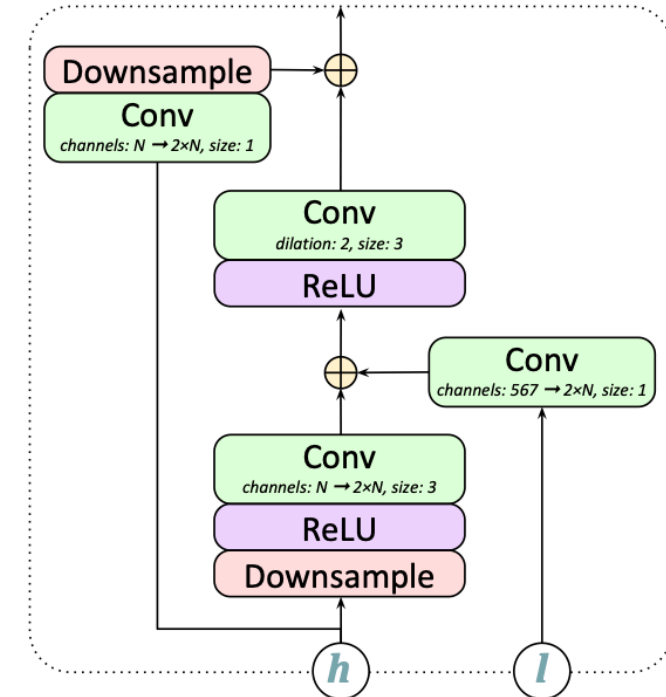
<https://arxiv.org/pdf/1609.03499.pdf>

# HIGH FIDELITY SPEECH SYNTHESIS WITH ADVERSARIAL NETWORKS

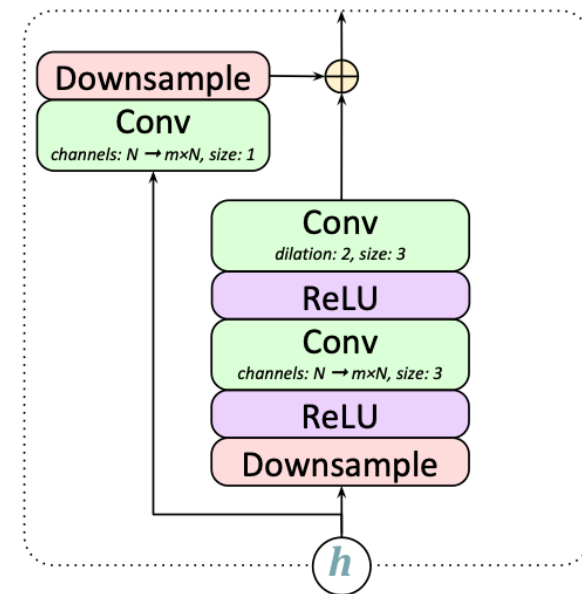
DeepMind



(a) GBlock



(b) Conditional DBlock

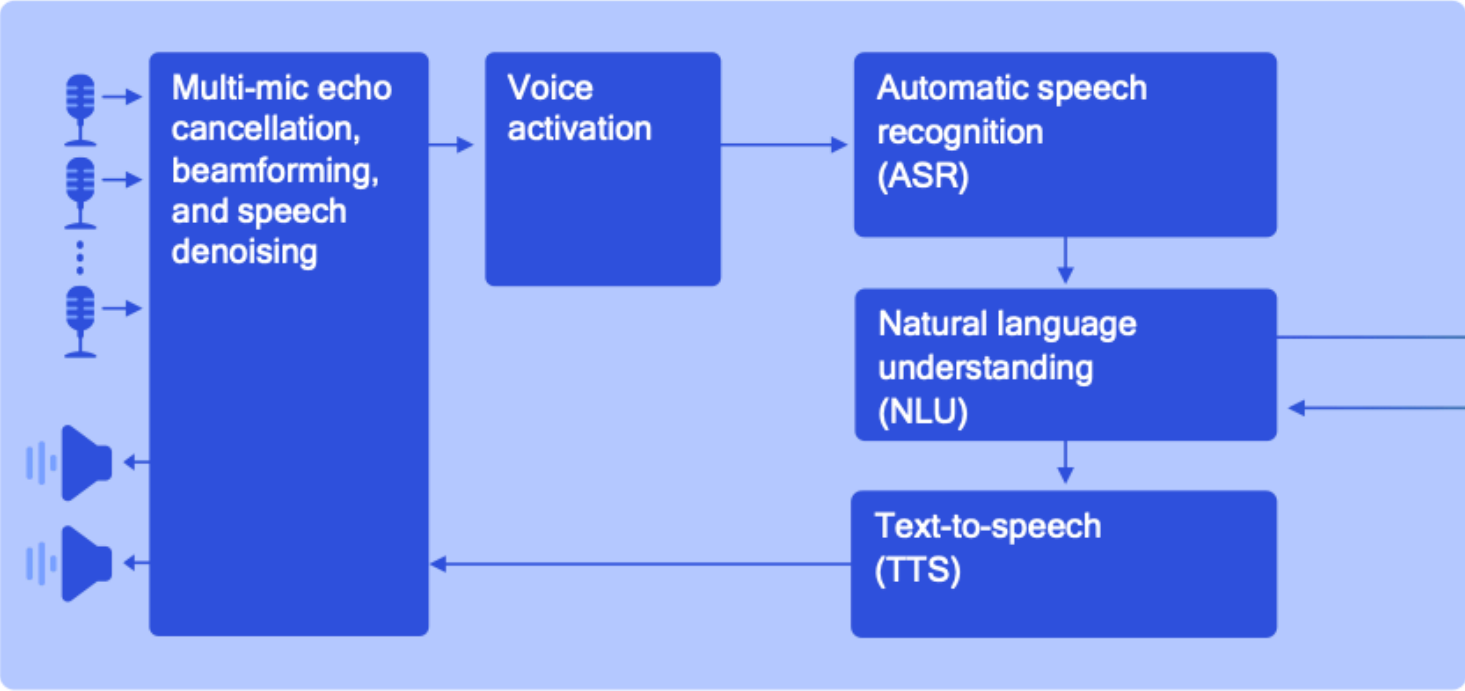


(c) DBlock

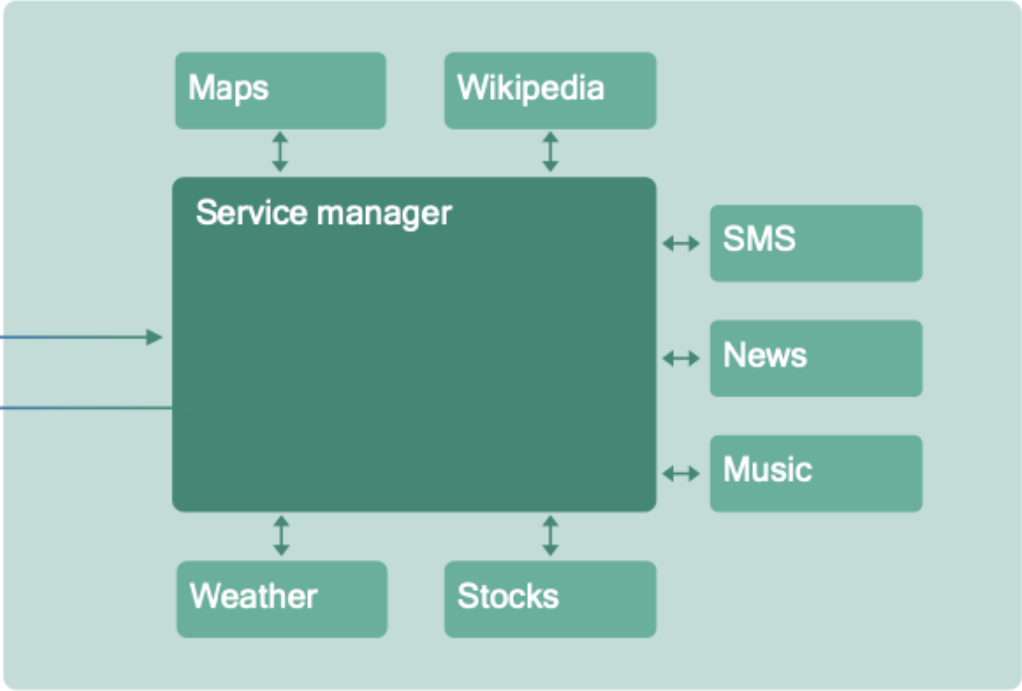
# Moving voice UI functionality to the end device

An end-to-end solution powered by machine learning

## On-device processing (always-on and real-time)



## Cloud processing (services)



## On-device centric (future)

# Some conclusions

- End-To-End Machine Learning (EEML) is a simple solution to train models which are able to solve problems for which one is lacking (or do not trust) a-priori knowledge.
- EEML is particularly suitable in (un-supervised/semi-supervised) (CNRS-SAMOVAR) situations for which large amounts of non-annotated data is available.
- EEML should be bundled (fused) in hybrid solutions with existing (traditional, modular) and symbolic AI techniques.

# End-to-End Machine Learning for Speech Processing: Speech-to-Speech, Speech-to-Text and Text-to-Speech

G rard Chollet (CNRS-SAMOVAR)  
and Colleagues from IV, Zaion, CNRS/IMT and SMI

gerard.chollet@telecom-sudparis.eu

