# Speech and Language Processing

## Lecture 4
## Neural network based speech recognition and synthesis

Information and Communications Engineering Course
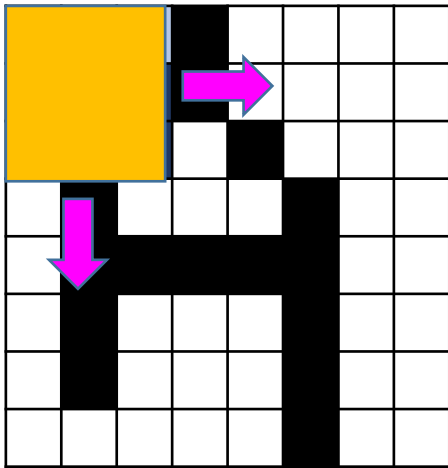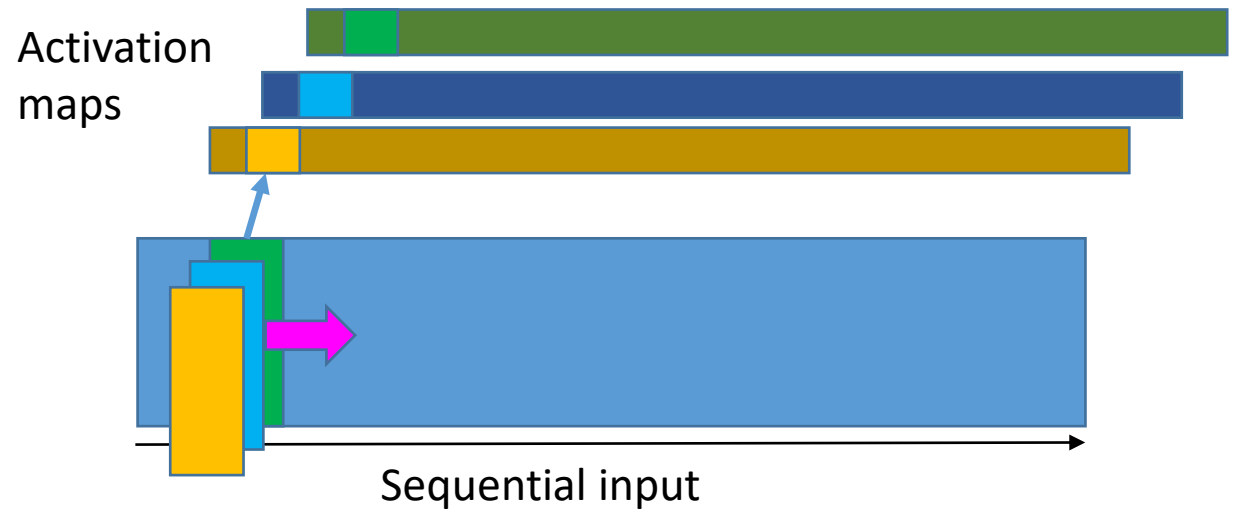
Takahiro Shinozaki

Manabu Okumura

2024/10/1

# Some Basic Functional Elements

# 1D-CNN



Activation maps
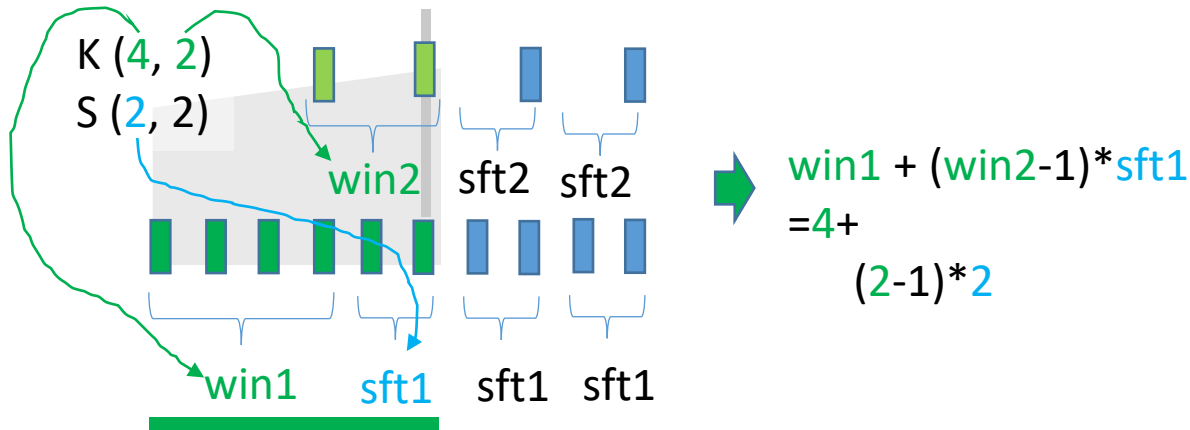
Sequential input

2D-CNN                    1D-CNN

3

# Receptive Field Length of Cascaded 1D Convolution

1st convolution: Kernel (window) Size = 4, Stride (shift) = 2.

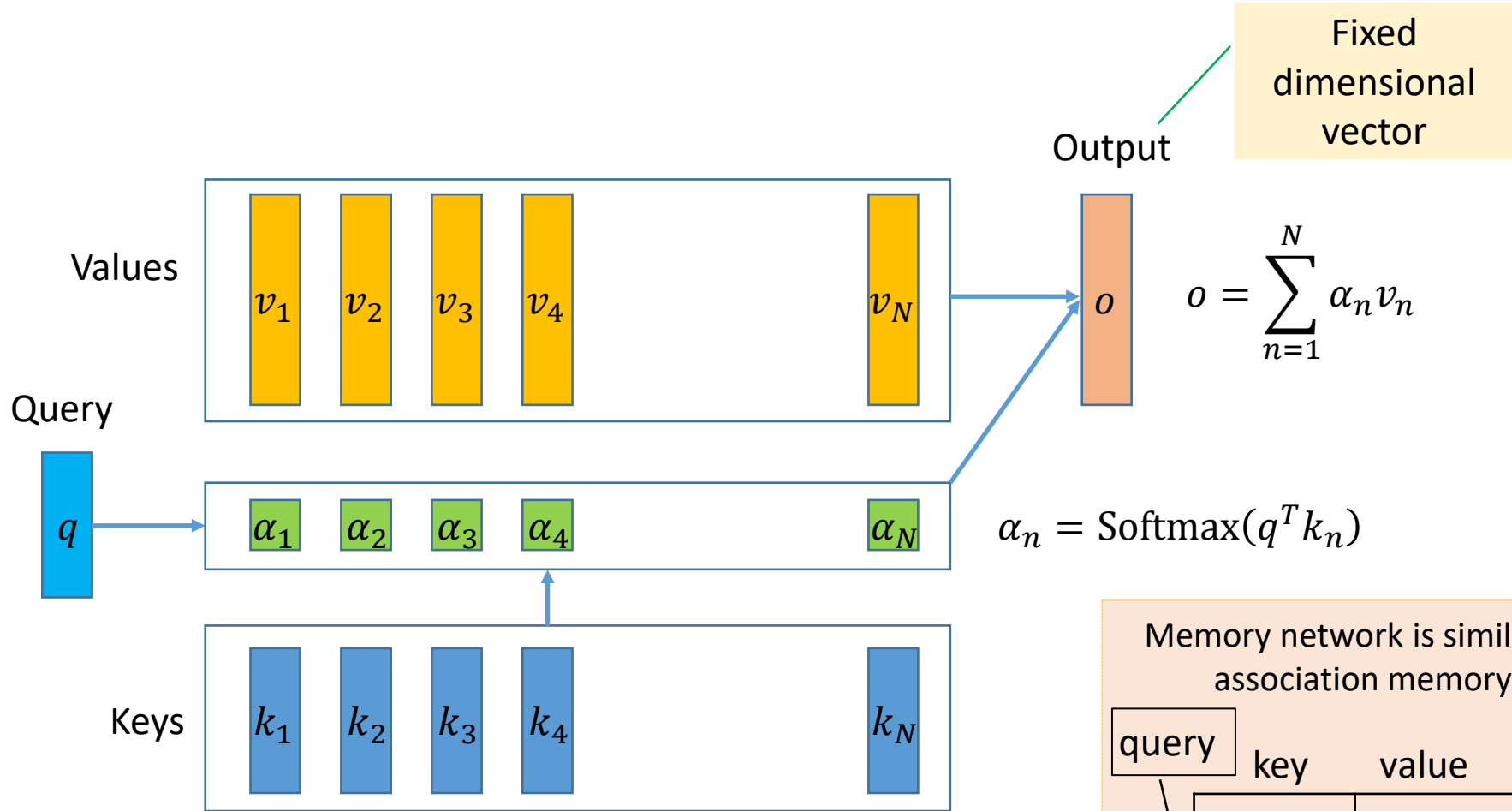2nd convolution: Kernel (window) Size = 2, Stride (shift) = 2.



$win1 + (win2-1)*sft1$

$=4+$

$(2-1)*2$

K (10, 8, 4, 4, 4), S (5, 4, 2, 2, 2)

((((10                         ← Window width of 1st layer @ input sample rate

+(8-1)*5)                    ← Window width of 2st layer @ input sample rate

+(4-1)*(5*4))              ← Window width of 3st layer @ input sample rate

+(4-1)*(5*4*2))         ← Window width of 4st layer @ input sample rate

+(4-1)*(5*4*2*2))    ← Window width of 5st layer @ input sample rate

=465

# Memory Network

Output

Fixed dimensional vector

Values

$$v_1 \quad v_2 \quad v_3 \quad v_4 \qquad v_N$$

$o$

$$o = \sum_{n=1}^{N} \alpha_n v_n$$

Query

$q$

$$\alpha_1 \quad \alpha_2 \quad \alpha_3 \quad \alpha_4 \qquad \alpha_N$$

$$\alpha_n = \text{Softmax}(q^T k_n)$$

Keys

$$k_1 \quad k_2 \quad k_3 \quad k_4 \qquad k_N$$

Memory network is similar to association memory

query

| key | value |
|--------|-------|
| Apple | 120 |
| Banana | 100 |
| Durian | 5500 |

# Fixed-Dimensional Embeddings of Sequences

- ## Use RNN



Sequential input

Fixed dimensional embedding vector

- ## 1D CNN with Global average pooling



Activation maps

Global average pooling

Sequential input

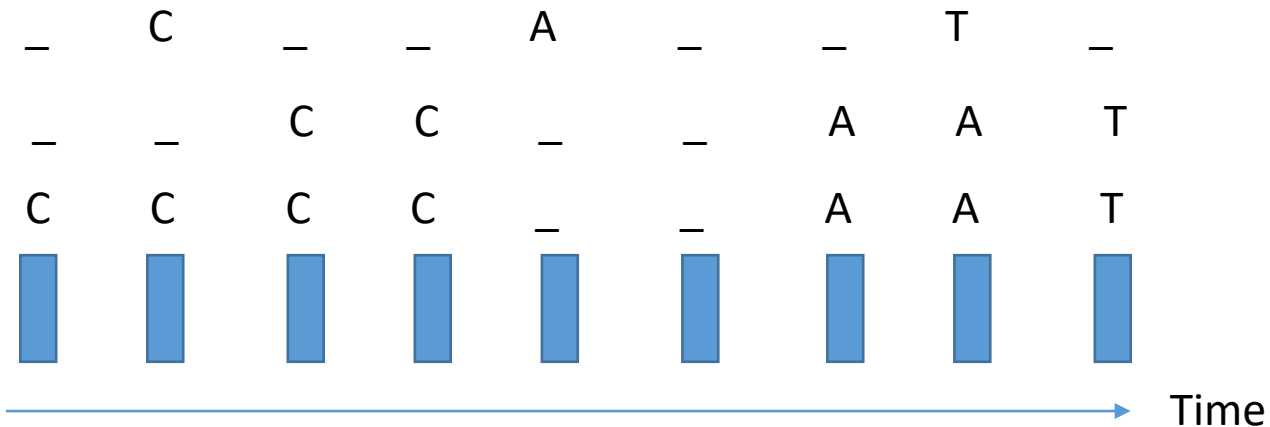# Connectionist Temporal Classification (CTC) Loss

- Assume:
  - We have frame-wise character (or word etc.) prediction for a sequence of time frames
  - A blank label is included as a special character
  - Reference text is a sequence of characters whose length is smaller than the frame sequence
- Matching of the prediction and the reference
  - Form output by collapsing repeated characters and removing blank character from the predicted sequence
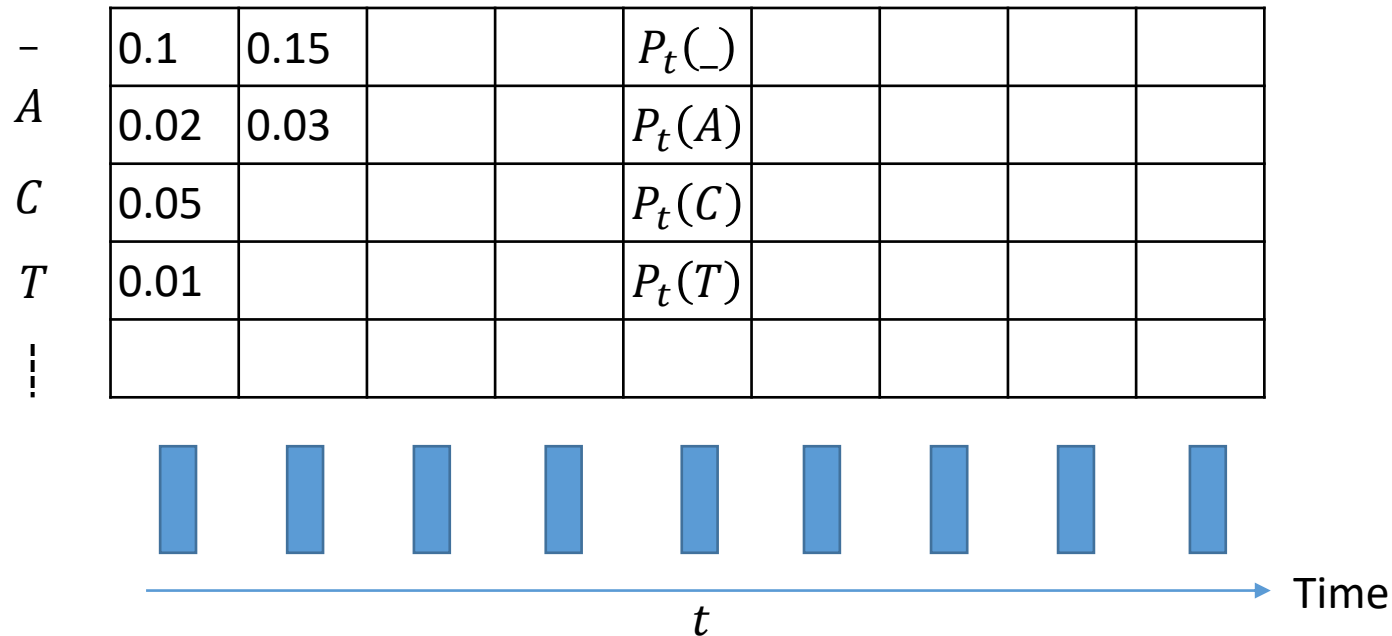
Reference: C A T



A. Graves+, "Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks," ICML, 2006https://www.cs.toronto.edu/~graves/icml_2006.pdf

7

# Probability of Predicted Sequence

- Probability of Predicted Sequence is a product of frame-wise prediction probabilities

$$P(\_\ \_\ C\ C\ \_\ \_\ A\ A\ T) = P_1(\_)P_2(\_)P_3(C)P_4(C)P_5(\_)P_6(\_)P_7(A)P_8(A)P_9(T)$$

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| _ | 0.1 | 0.15 | | | $P_t(\_)$ | | | | |
| $A$ | 0.02 | 0.03 | | | $P_t(A)$ | | | | |
| $C$ | 0.05 | | | | $P_t(C)$ | | | | |
| $T$ | 0.01 | | | | $P_t(T)$ | | | | |
| ⋮ | | | | | | | | | |

Time

$t$

# Probability of Character Sequence

- Probability of character sequence (like the reference) is a sum of probabilities of all the matching predicted sequences

$$P(CAT) = P(\_\ \_\ C\ C\ \_\ \_\ A\ A\ T) + P(C\ \ C\ C\ C\ \_\ \_\ A\ A\ T) \cdots$$

$$= \sum_{\pi \in \mathcal{B}^{-1}(CAT)} P(\pi) = \sum_{\pi \in \mathcal{B}^{-1}(CAT)} \prod_t P_t(\pi_t)$$
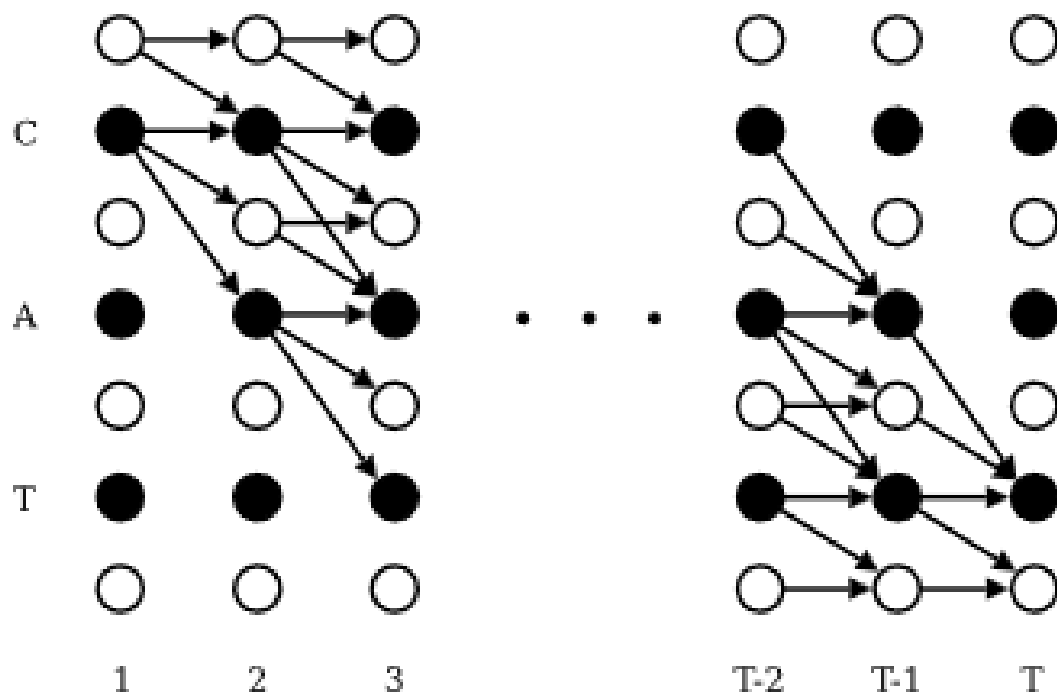
$\mathcal{B}$: The contraction function of CTC.

   e.g. $\mathcal{B}(\_\ \_\ C\ C\ \_\ \_\ A\ A\ T) = \mathcal{B}(C\ \ C\ C\ C\ \_\ \_\ A\ A\ T) = CAT$

$\mathcal{B}^{-1}$: Inverse of the contraction function (one-to-many mapping)

# Efficient Probability Evaluation

The probability of the character sequence is efficiently evaluated by the forward algorithm



Black circles represent labels, and white circles represent blanks

# CTC Loss

CTC loss
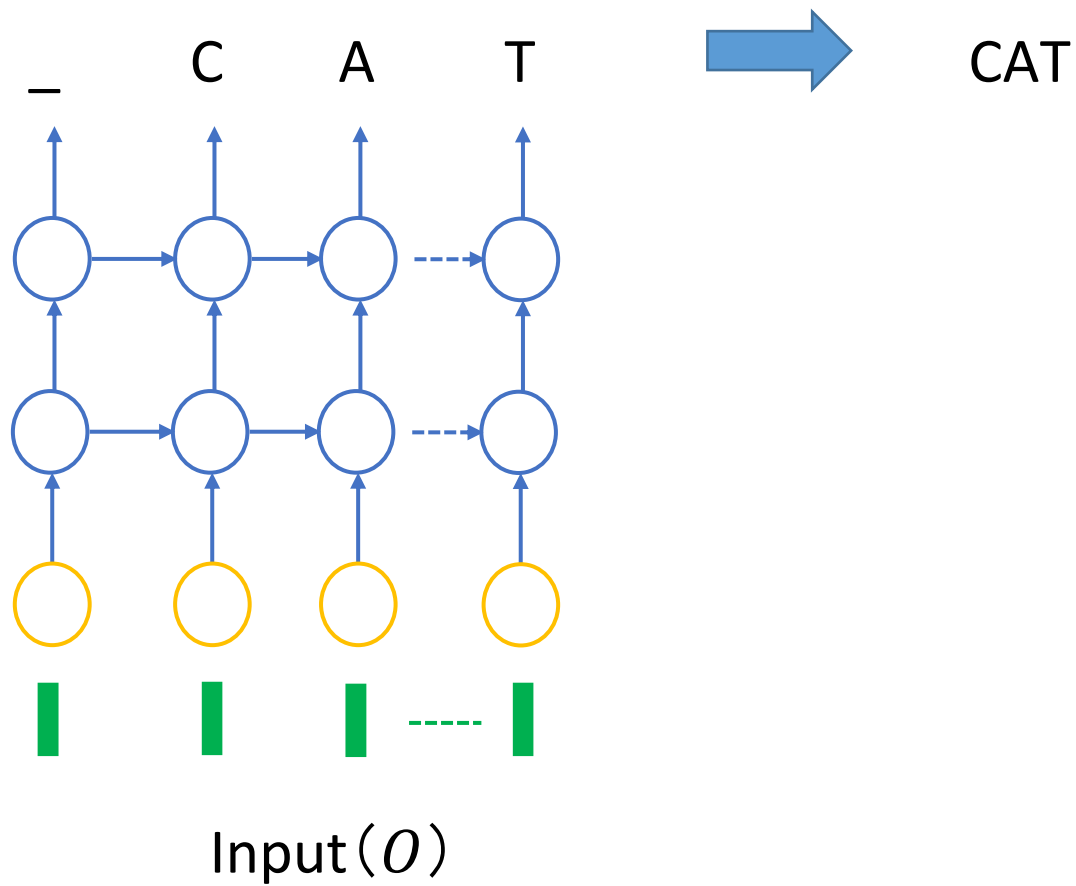
$$\mathcal{L}(S) = -logP(S) = -log \sum_{\pi \in \mathcal{B}^{-1}(S)} \prod_t P_t(\pi_t)$$

It's gradient

$$\frac{\partial}{\partial P_u(c)} \mathcal{L}(S) = -\frac{\sum_{\pi \in \mathcal{B}^{-1}(S), \pi_u = c} \prod_{t \neq u} P_t(\pi_t)}{\sum_{\pi \in \mathcal{B}^{-1}(S)} \prod_t P_t(\pi_t)}$$
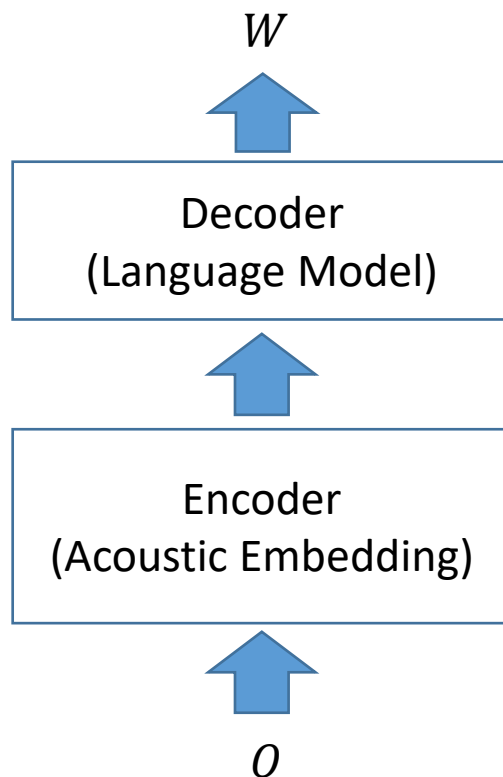
# Neural Network Based Speech Recognition
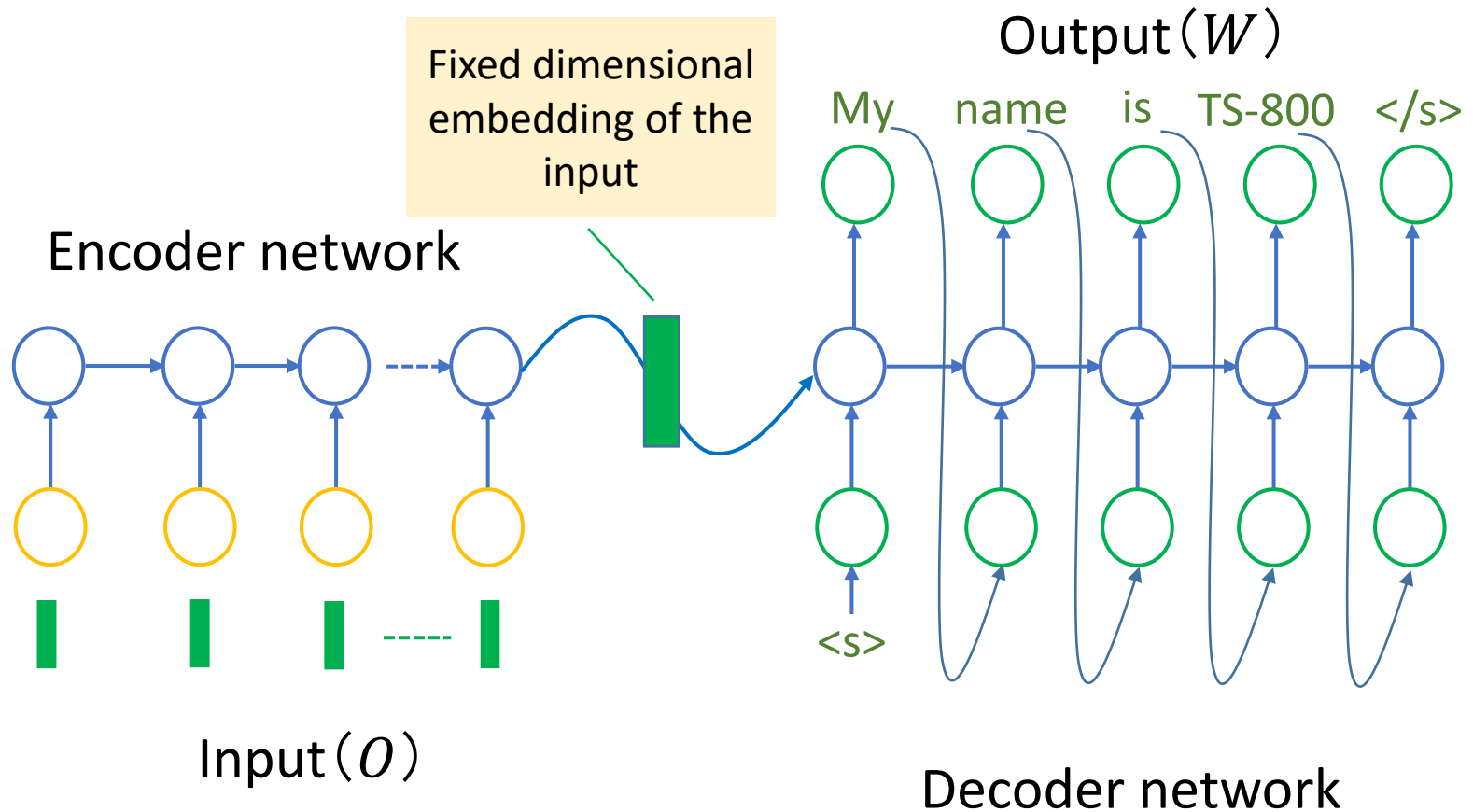
# RNN+CTC

# Encoder-Decoder Networks

Language model models a probability of a sentence $P(W)$.
By conditioning it with an acoustic input $O$, we get the discriminative modeling based speech recognizer $P(W|O)$.

$$W$$

⬆

| Decoder |
| :---: |
| (Language Model) |

⬆

| Encoder |
| :---: |
| (Acoustic Embedding) |

⬆

$$O$$

# Encoder-Decoder Network

Directly models $P(W|O)$

Output $(W)$

My    name    is    TS-800    </s>
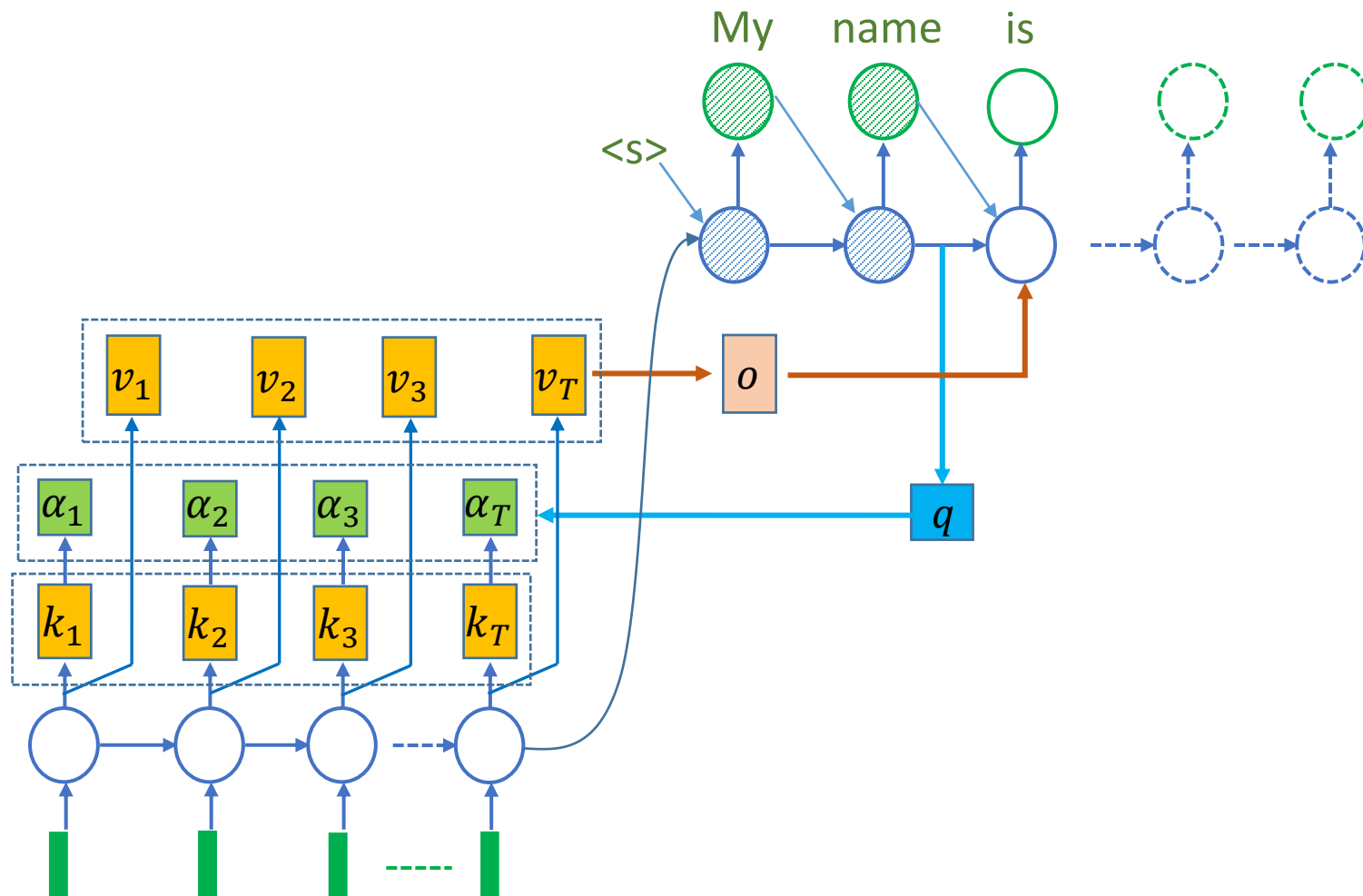
Fixed dimensional embedding of the input

Encoder network
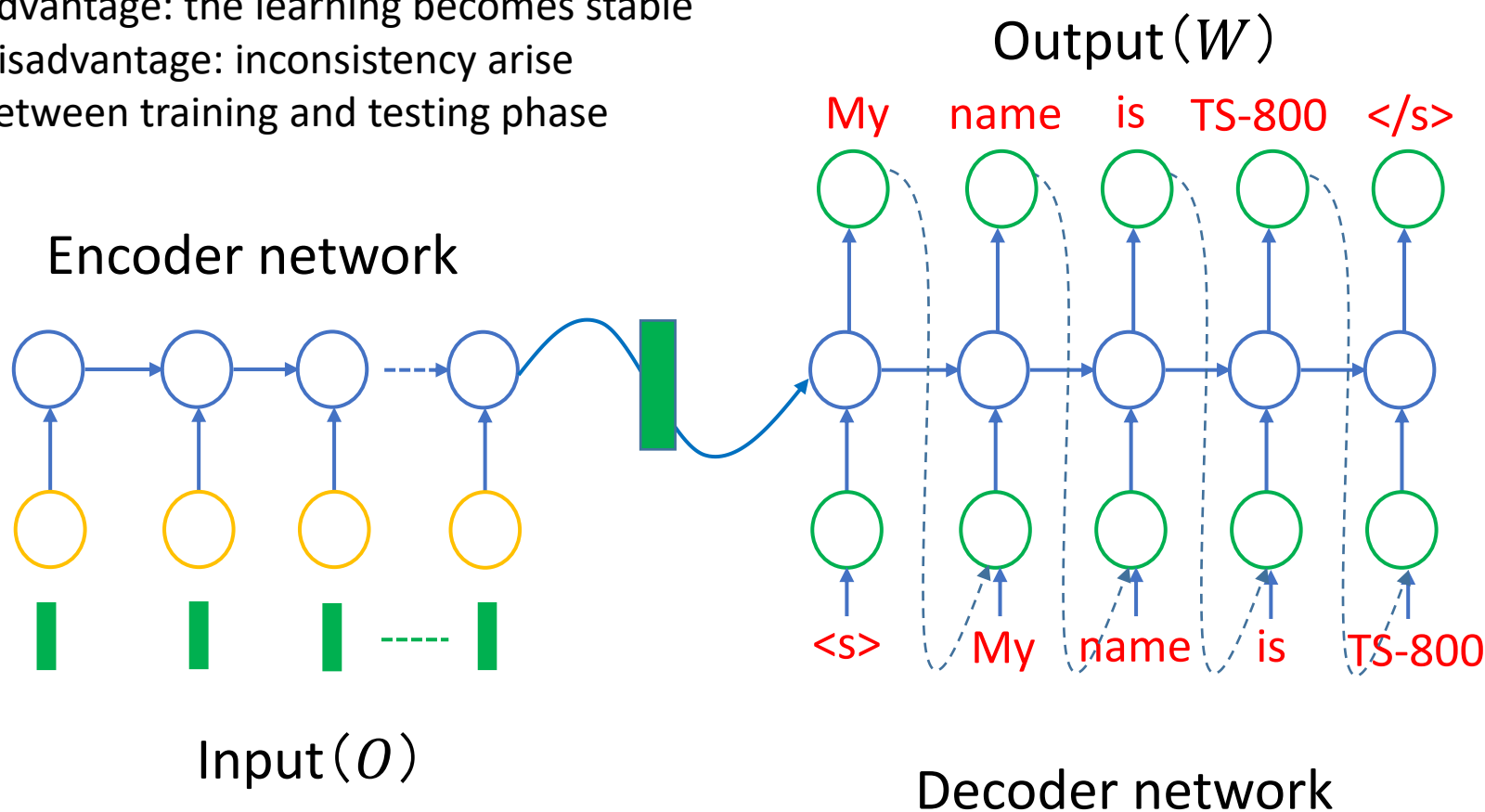
<s>

Input $(O)$

Decoder network

# Attention Encoder-Decoder
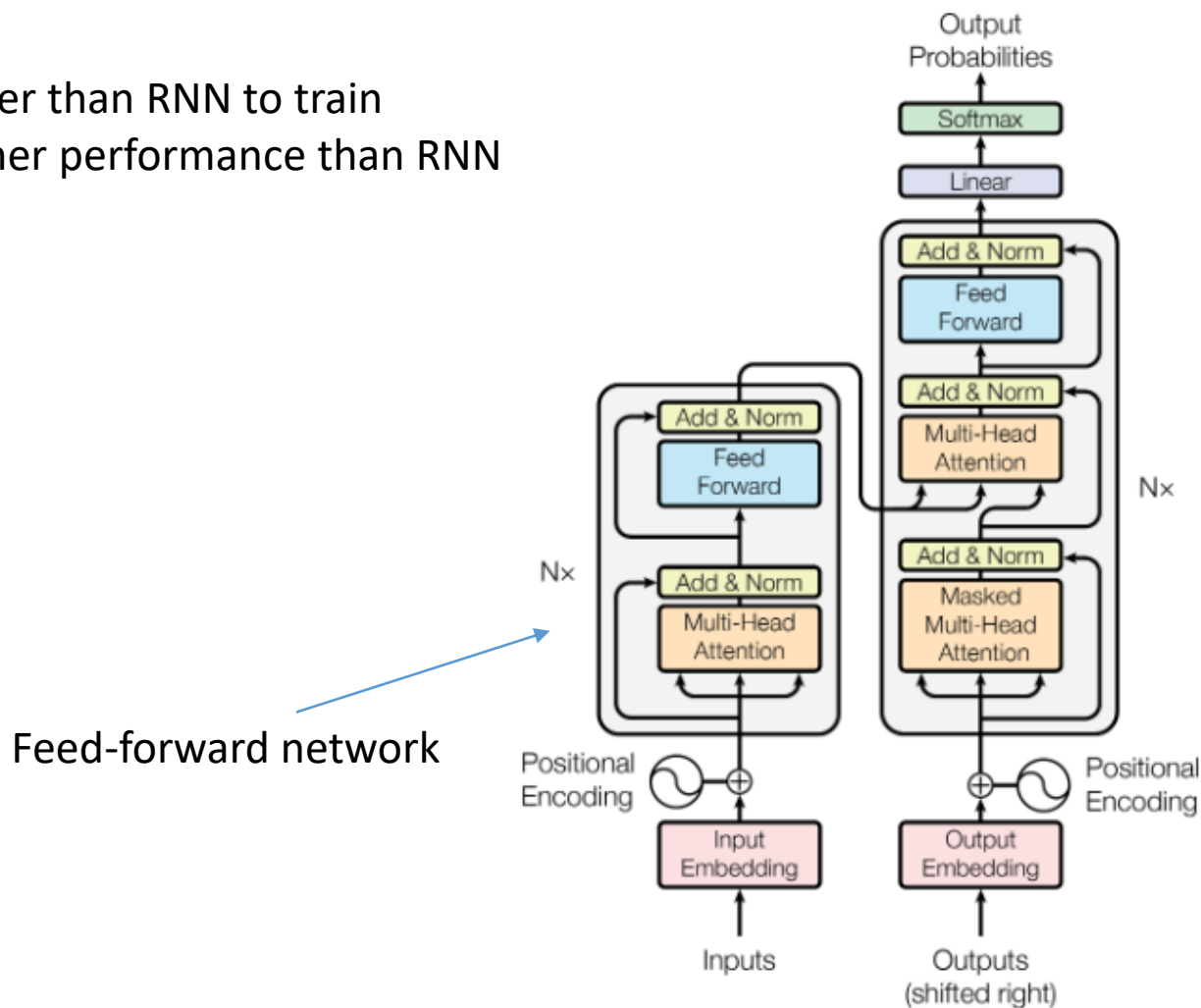
# Teacher Forcing

When training the encoder-decoder network, teacher forcing uses reference words in the decoder input instead of the predicted words

- Advantage: the learning becomes stable
- Disadvantage: inconsistency arise between training and testing phase

Output $(W)$

My    name    is    TS-800    </s>

Encoder network

Input $(O)$

<s>    My    name    is    TS-800

Decoder network

# Transformer

- Faster than RNN to train
- Higher performance than RNN

Feed-forward network



Output
Probabilities

Softmax

Linear

Add & Norm
Feed
Forward

Add & Norm
Multi-Head
Attention

Nx

Add & Norm
Masked
Multi-Head
Attention

Add & Norm
Feed
Forward

Nx

Add & Norm
Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

Positional
Encoding

Output
Embedding

Outputs
(shifted right)

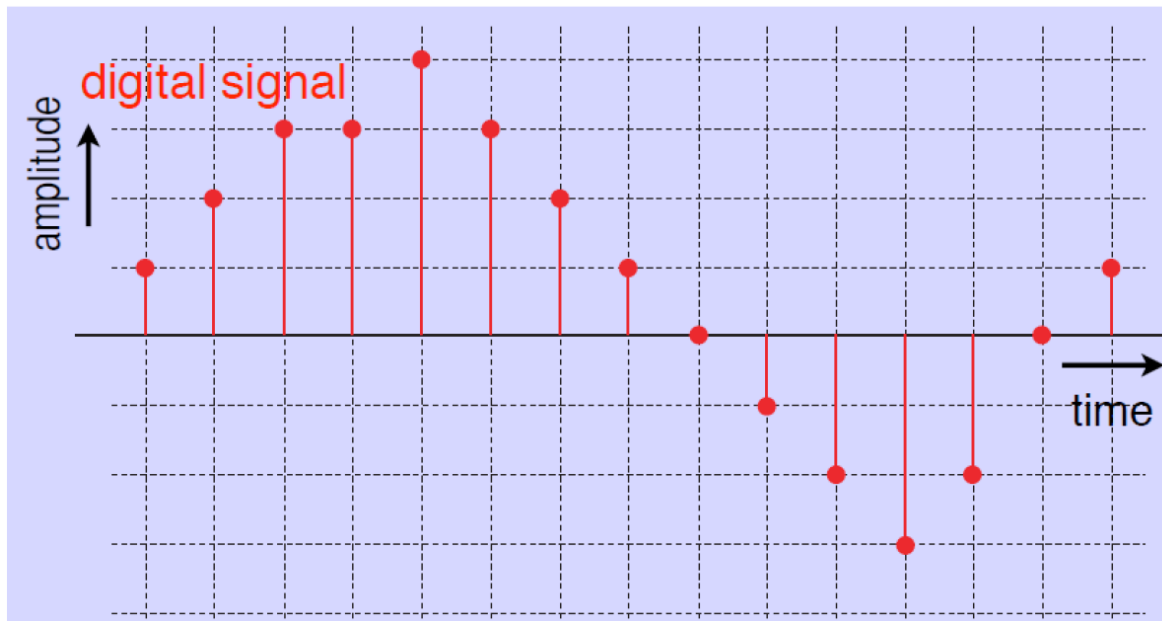A. Vaswani+, Attention Is All You Need, 2017

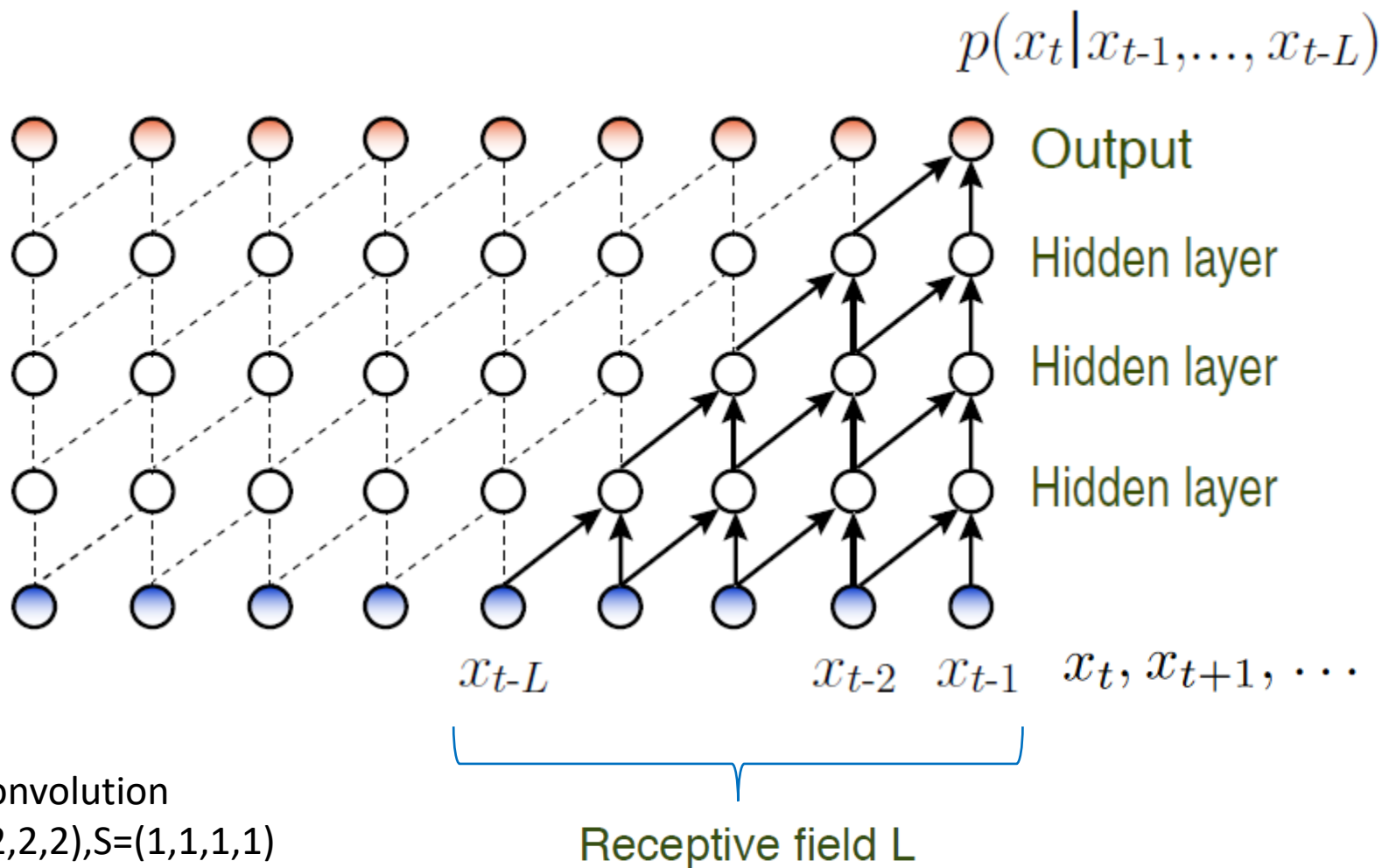# Neural Network Based Speech Synthesis

# WAVENET

- A DNN based generative raw waveform model
  [van den Oord, et al., 2016]

$$p(\mathbf{x}) = \prod_{t=1}^{T} p(x_t \mid x_1, \ldots, x_{t-1})$$

# Causal Convolution

The prediction emitted at time step $t$ is independent of future time steps $t, t+1, \cdots$



$$p(x_t | x_{t-1}, \ldots, x_{t-L})$$

Output

Hidden layer

Hidden layer

Hidden layer

$x_{t-L}$  $x_{t-2}$  $x_{t-1}$  $x_t, x_{t+1}, \cdots$

Receptive field L

1-D convolution
K=(2,2,2,2),S=(1,1,1,1)

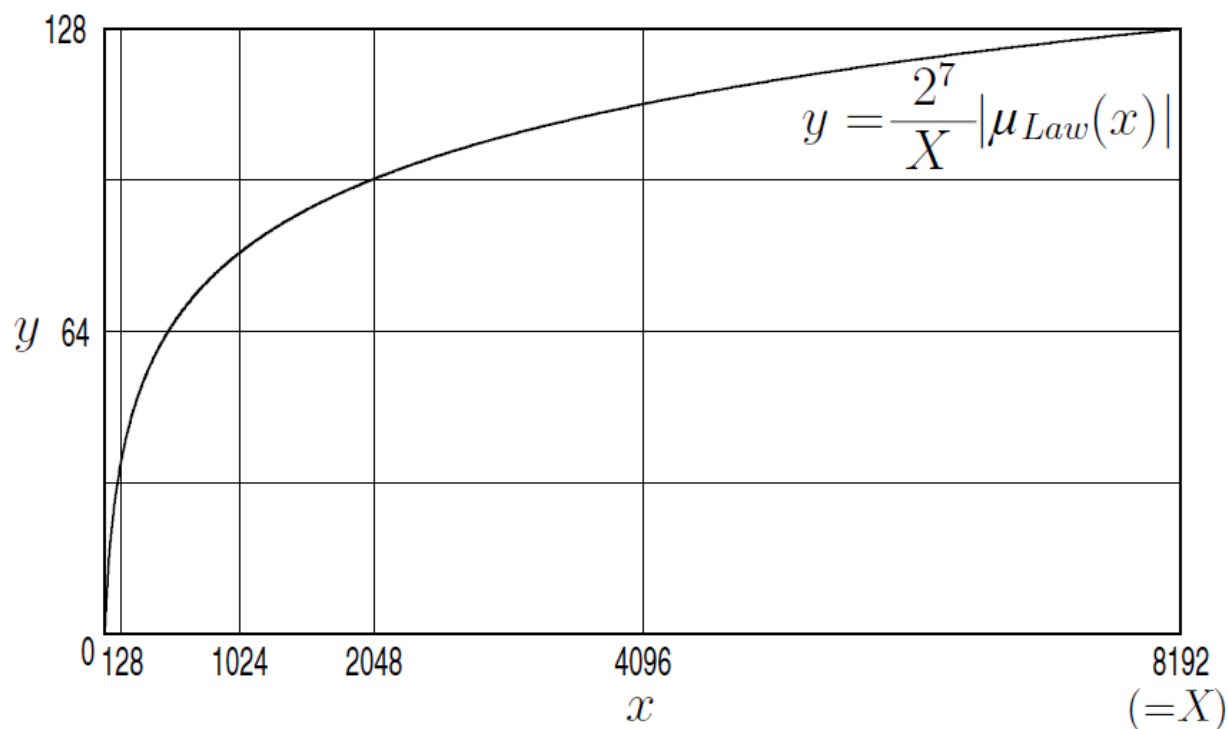# Categorical Prediction of Amplitude

- Discrete prediction by the softmax function is used, as it is found to work better than continuous regression
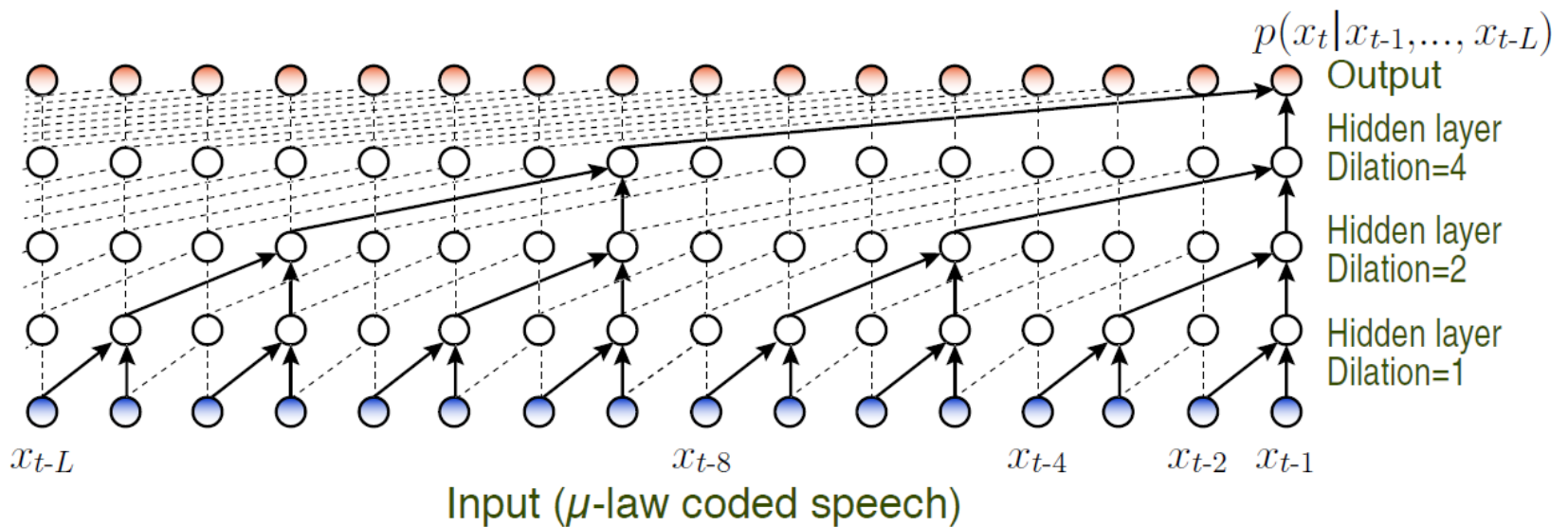


1 millisecond

# μ-Law Coding

$$\mu_{Law}(x) = \text{sign}(x) X \frac{\log\left(1 + \mu |x| / X\right)}{\log\left(1 + \mu\right)} \ , \quad |x| \leq X$$

- $\mu = 255$ for 8bit PCM



$$y = \frac{2^7}{X} \left| \mu_{Law}(x) \right|$$

# Dilated Causal Convolution



$p(x_t|x_{t-1},\ldots,x_{t-L})$

Output

Hidden layer
Dilation=4

Hidden layer
Dilation=2

Hidden layer
Dilation=1

$x_{t-L}$ $x_{t-8}$ $x_{t-4}$ $x_{t-2}$ $x_{t-1}$

Input ($\mu$-law coded speech)

1-D convolution
K=(2,2,2,2),S=(2,2,2,2)

# Signal Generation

- Random sampling from estimated distribution



$$p(x_t|\hat{x}_{t-1},...,\hat{x}_{t-L})$$

$$\hat{x}_{t-L} \quad \cdots \quad \hat{x}_{t-8} \quad\quad\quad \hat{x}_{t-4} \quad\quad \hat{x}_{t-2} \quad \hat{x}_{t-1} \quad \hat{x}_t \sim p(x_t|\hat{x}_{t-1},...,\hat{x}_{t-L})$$

# Conditional WavNet

$$\ln p(\boldsymbol{x} \mid \boldsymbol{h}) \approx \sum_{t=1}^{T} \ln p(x_t \mid x_{t-1},\ldots,x_{t-L},\boldsymbol{h})$$

- Auxiliary input $\boldsymbol{h}$ : F0, mel spectrum, spectrogram, etc.
- Receptive field L: several hundreds milliseconds

Acoustic features

$\boldsymbol{h}$

Conditional WaveNet

$$p(x_t \mid \hat{x}_{t-1},\ldots,\hat{x}_{t-L},\boldsymbol{h})$$

$$\hat{x}_t \sim p(x_t \mid \hat{x}_{t-1},\ldots,\hat{x}_{t-L},\boldsymbol{h})$$

$\hat{x}_{t-L}$ $\qquad$ $\hat{x}_{t-2}$ $\hat{x}_{t-1}$

$z^{-1}$ $\quad$ $z^{-1}$ $\quad$ $z^{-1}$ $\quad$ $z^{-1}$ $\quad$ $z^{-1}$

https://deepmind.com/blog/article/wavenet-generative-model-raw-audio
Last visited 2023/5/23

26

# Tacotron 2

A neural network architecture for speech synthesis directly from text



*The figure is cited from J. Shen et al., "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram predictions", ICASSP 2018.

https://google.github.io/tacotron/publications/tacotron/index.html
Last visited 2023/5/23

# Exercise (Q4.1, Q4.2)

Q4.1

What is the receptive field length of 1-D convolution when K=(2,2,2,2,2),S=(1,1,1,1,1)?

Q4.2

What is the receptive field length of 1-D convolution when K=(2,2,2,2,2),S=(2,2,2,2,2)?