

# Neural Computation

Monday Session 27 November

# Attention Is All You Need

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Lukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

- Sequence-to-Sequence Architecture



- Hugely influential
- Basis for ChatGPT (**G**enerative **P**re-trained **T**ransformer)
- Also for vision

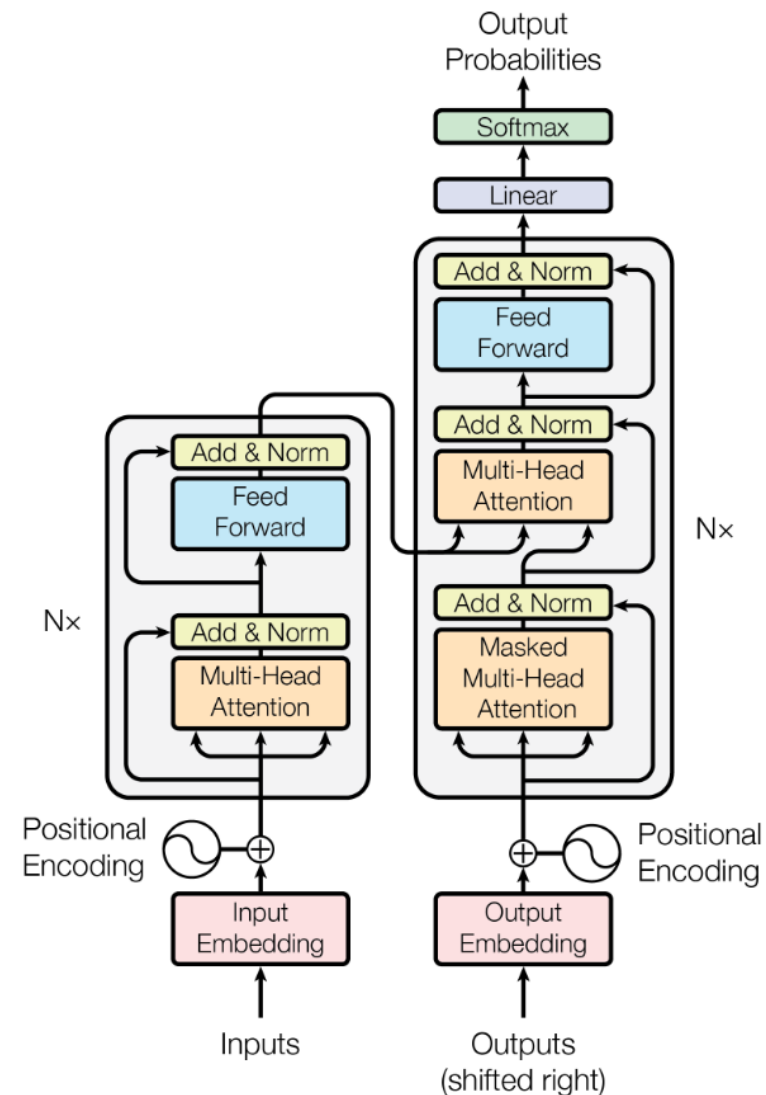
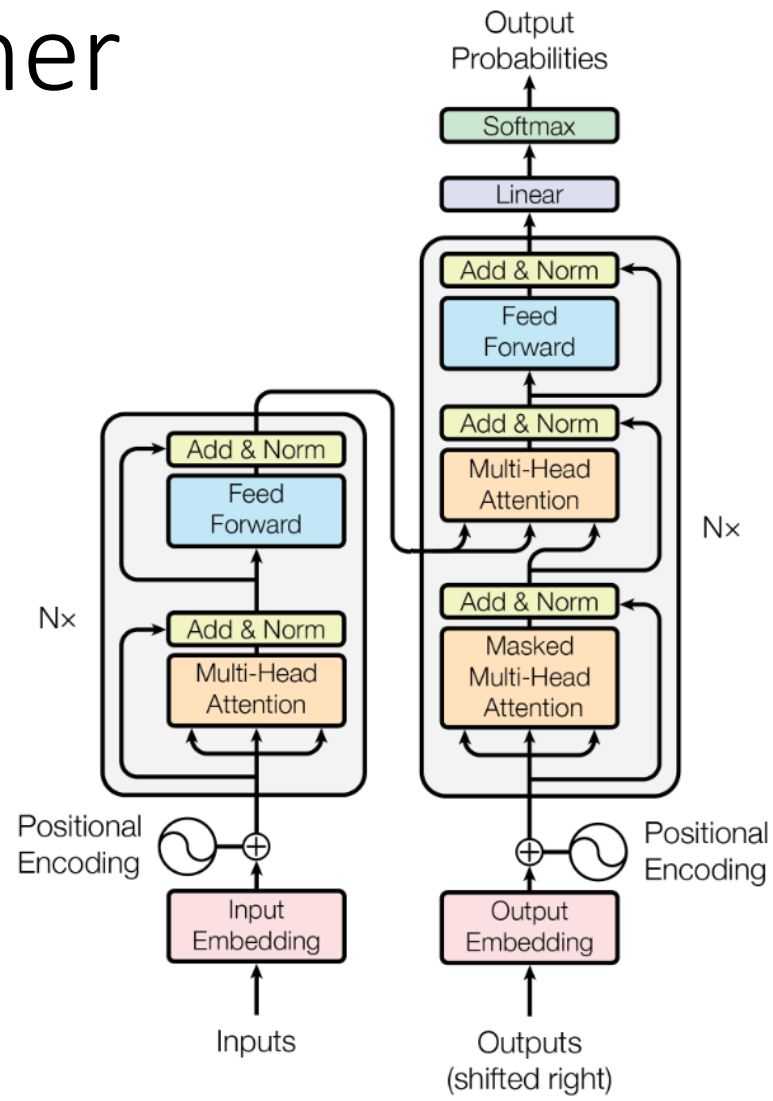
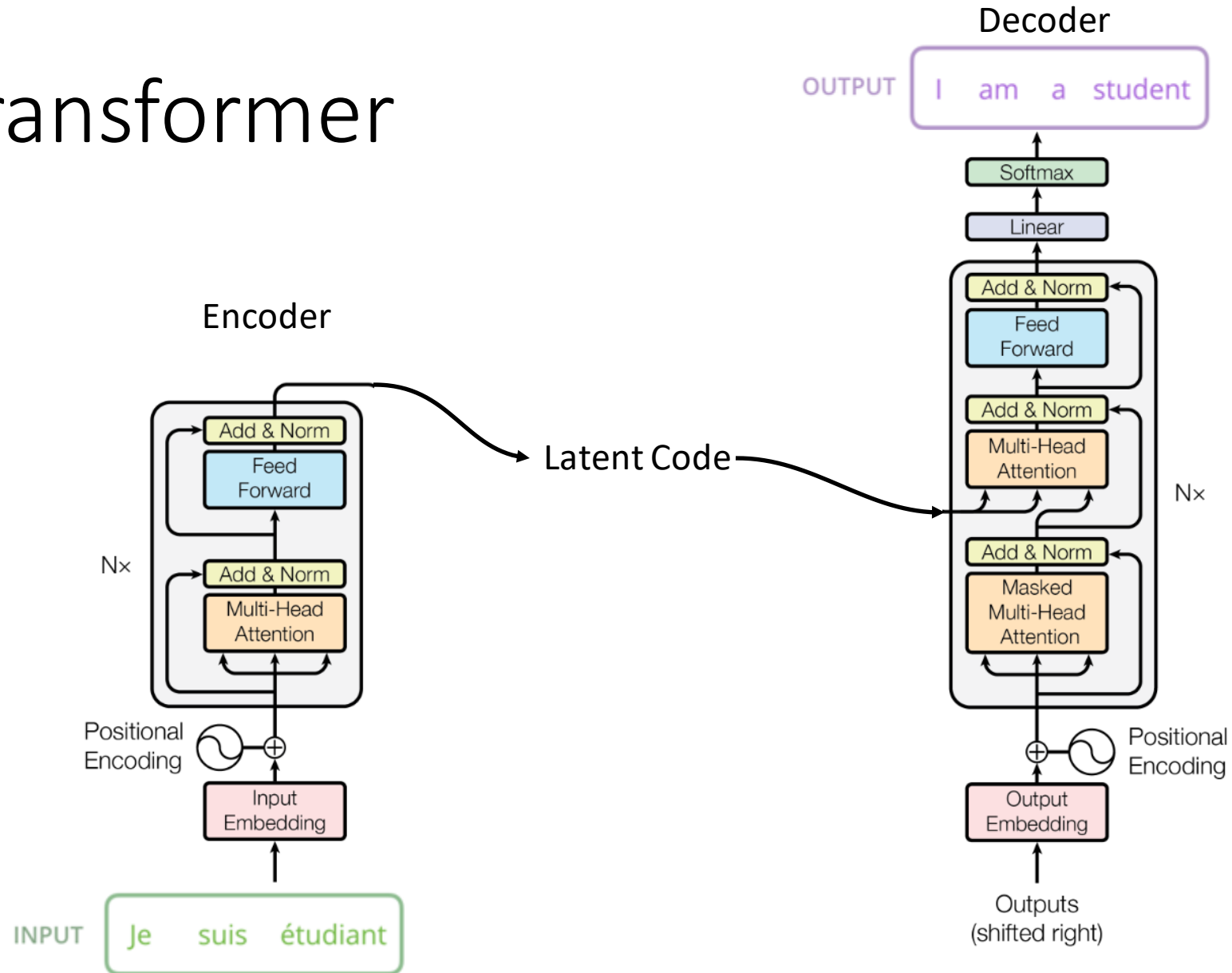


Figure 1: The Transformer - model architecture.

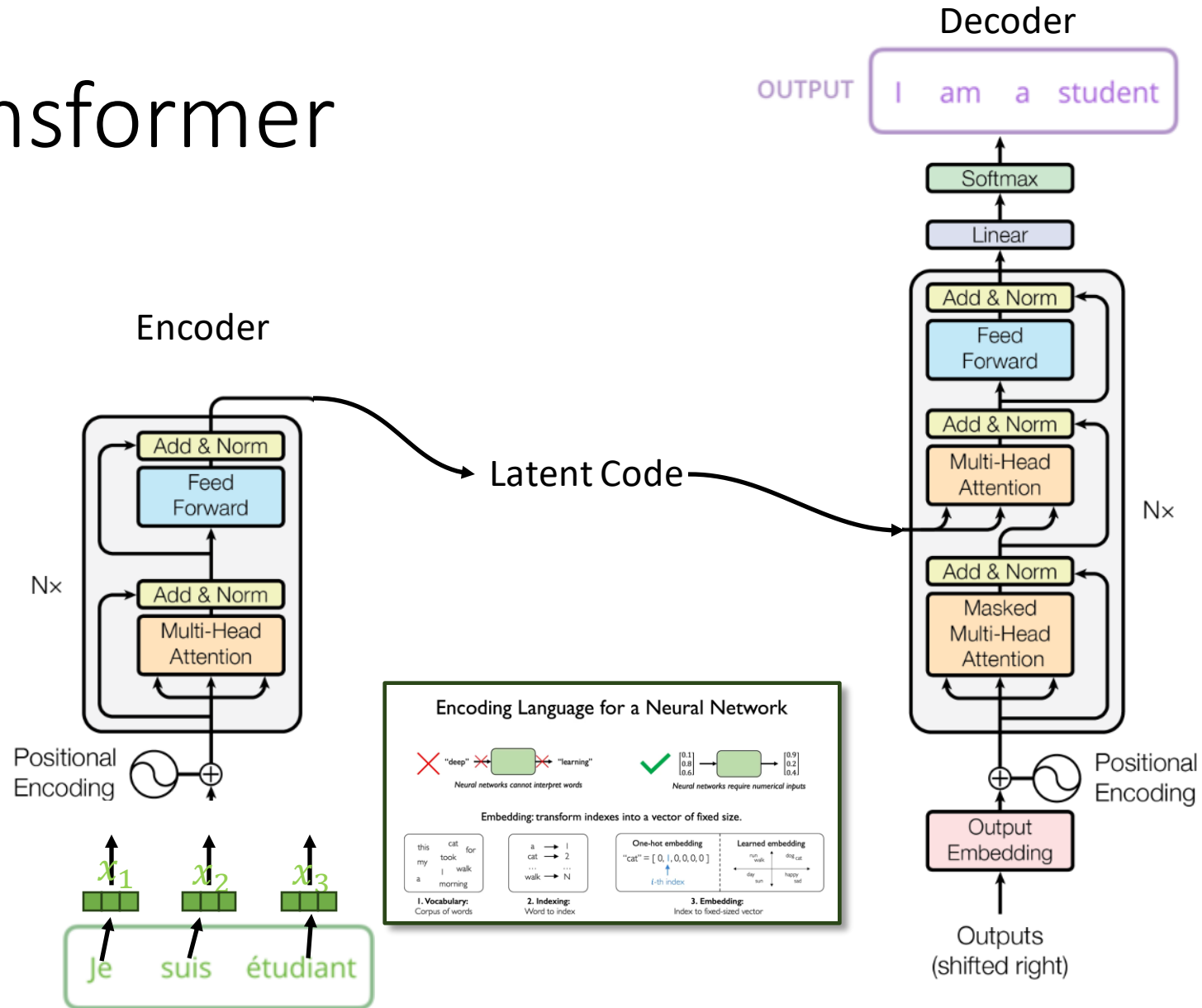
# The Transformer



# The Transformer



# The Transformer



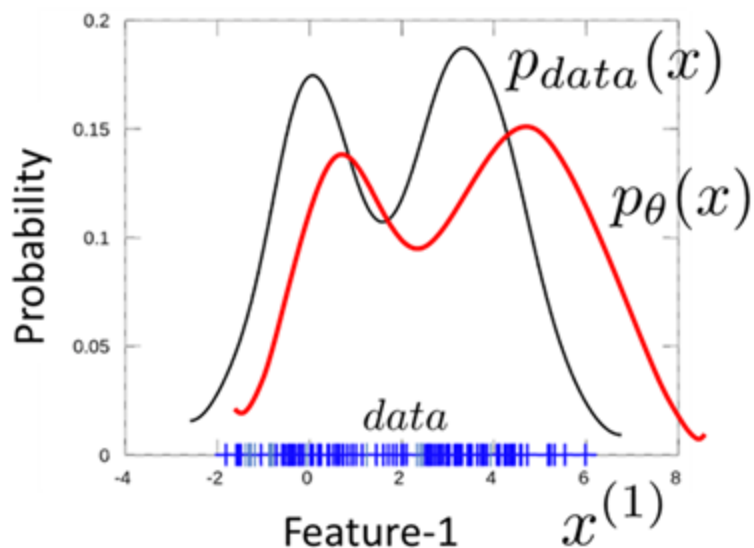
# Probability Density Estimation

One of the main aims of unsupervised approaches and Generative Modelling.

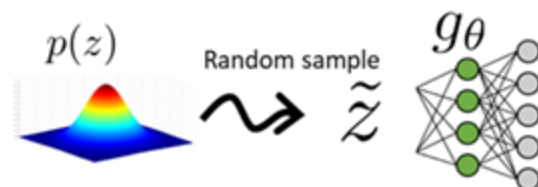
## Goal of Density Estimation:

We could try to fit a probabilistic model  $p_{\theta}(x)$  to the data, to learn their underlying distribution  $p_{data}(x)$ .

How? By learning its parameters  $\theta$  so that:  $p_{\theta}(x) \approx p_{data}(x)$



But we cannot always do that directly.  
Perhaps we cannot compute  $p_{data}(x)$  or  $p_{\theta}(x)$ .  
Instead, we could do PDE *indirectly*:  
*Enforce samples from model to be similar to real data* instead:



$$x \sim p_{data}(x)$$

7	3	4	6	1	8	1	0
9	8	0	3	1	3	7	0
2	9	6	0	1	6	7	1
9	7	6	5	5	8	8	3
4	4	8	7	3	6	4	6
6	3	6	8	9	9	4	4
0	7	8	1	0	0	1	8
5	7	1	7	5	5	9	9

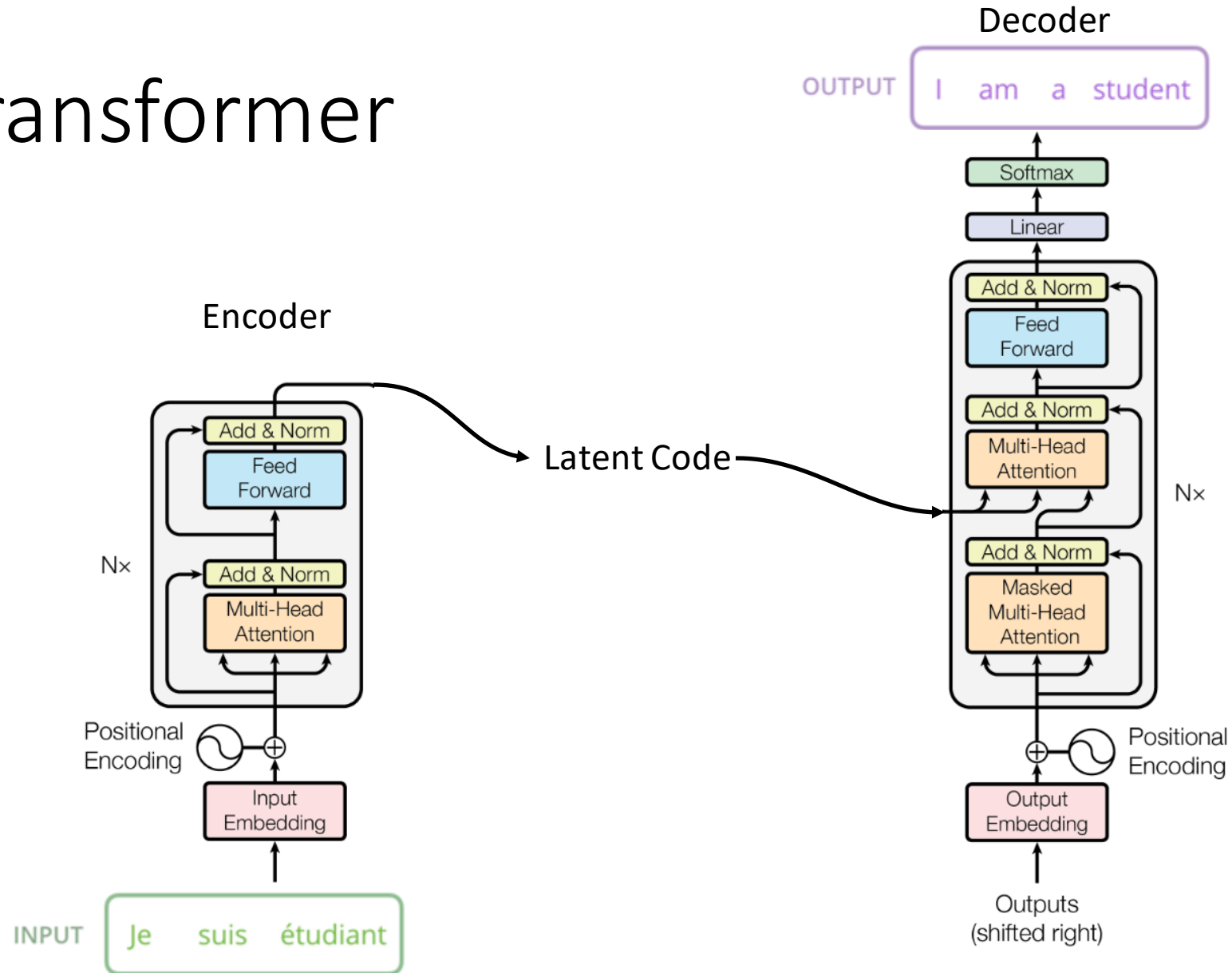


$$\tilde{x} \sim p_{\theta}(x)$$

1	1	1	2	1	1	7	
5	7	6	7	7	1	7	
7	6	9	1	3	3	1	
8	7	5	9	3	1	7	
7	7	5	0	7	7	7	
5	0	4	5	7	8	7	
5	5	1	1	5	1	1	
6	1	2	1	2	4	1	

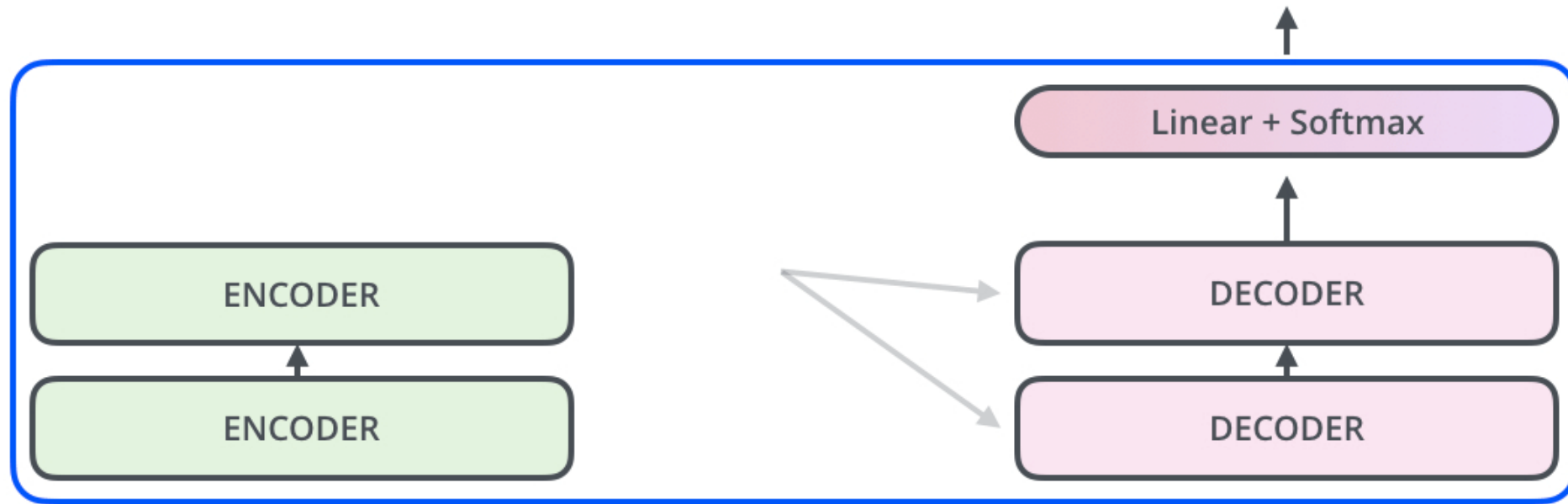
**Both VAEs and GANs can be seen as following this approach.**

# The Transformer

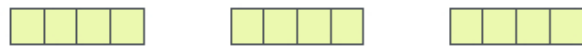


Decoding time step: 1 2 3 4 5 6

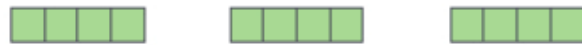
OUTPUT



EMBEDDING  
WITH TIME  
SIGNAL



EMBEDDINGS

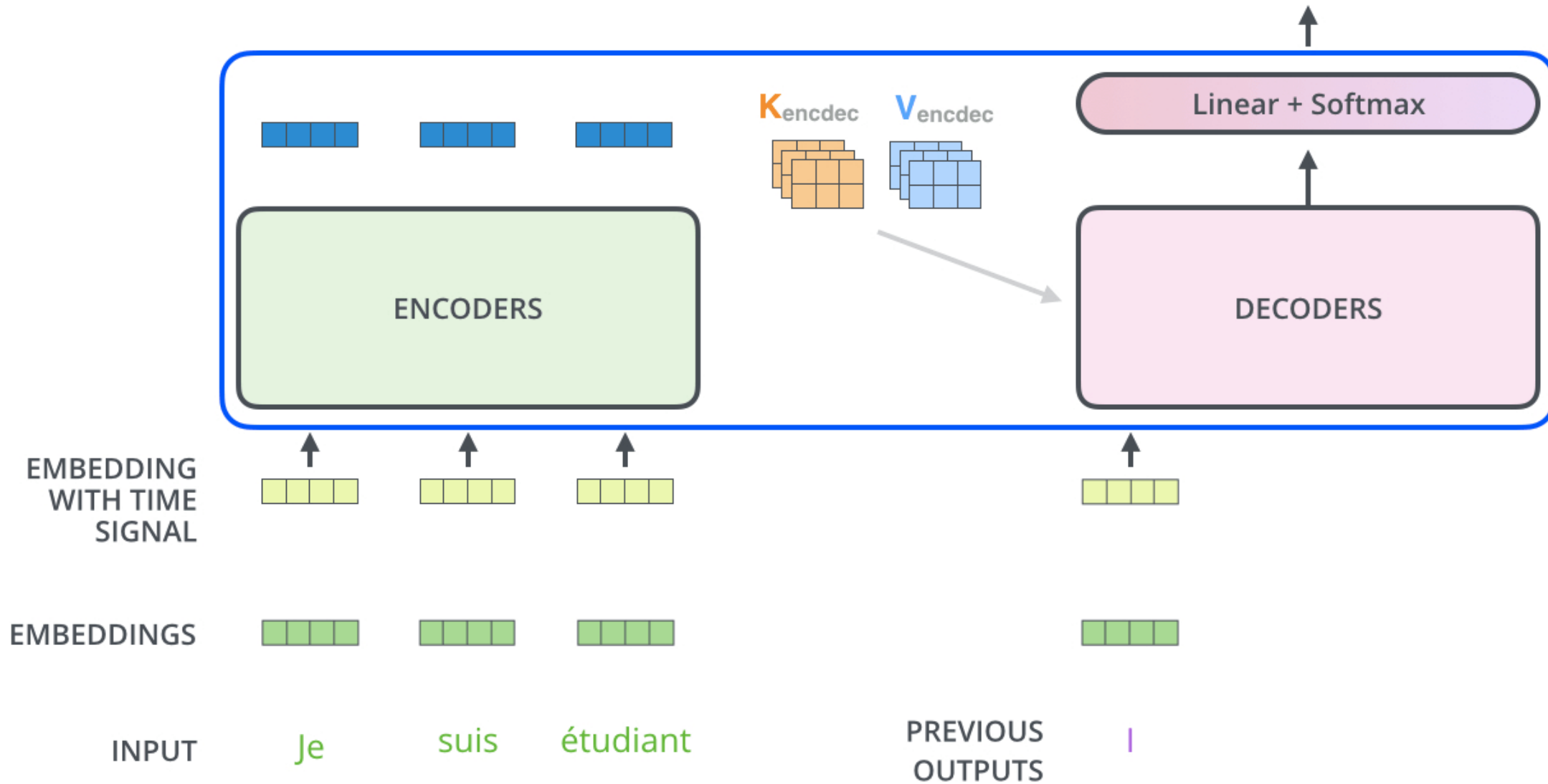


INPUT Je suis étudiant



Decoding time step: 1 2 3 4 5 6

OUTPUT |

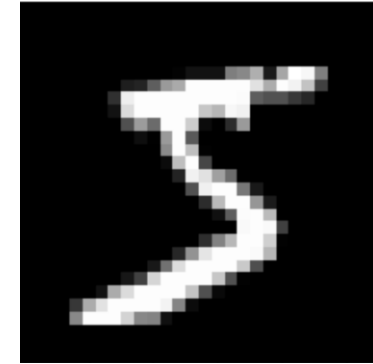


# Autoregressive Models

## Problem:

- Model high dimensional difficult distribution

$$p_{\theta}(\mathbf{x}) = p_{\text{data}}(\mathbf{x}), \text{ with } \mathbf{x} = (x_1, \dots, x_n)$$



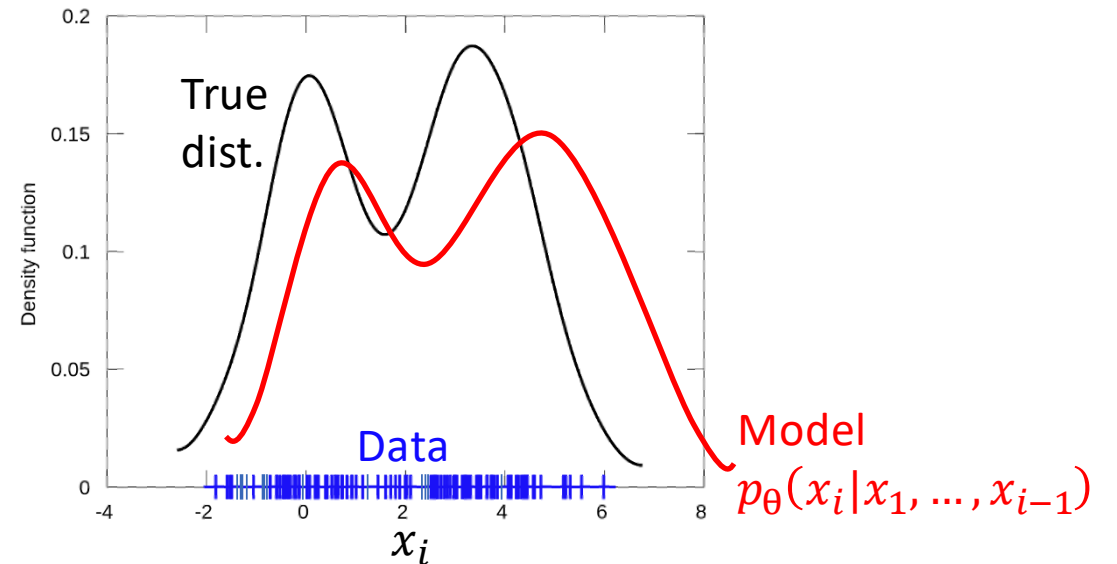
## Idea:

- Factorise distribution

$$p_{\theta}(\mathbf{x}) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

Neural network:

- Parameters  $\theta$
- Input  $x_1, \dots, x_{i-1}$
- Output dist. over  $x_i$



# WAVENET: A GENERATIVE MODEL FOR RAW AUDIO

**Aäron van den Oord**

**Sander Dieleman**

**Heiga Zen<sup>†</sup>**

**Karen Simonyan**

**Oriol Vinyals**

**Alex Graves**

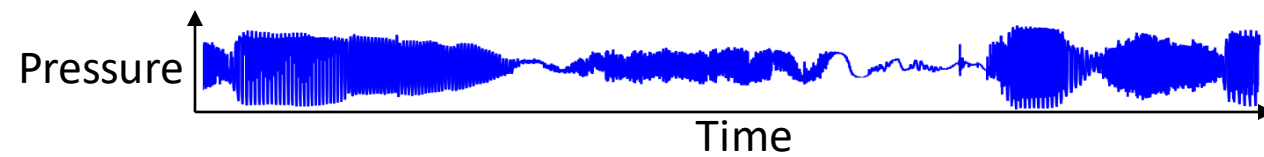
**Nal Kalchbrenner**

**Andrew Senior**

**Koray Kavukcuoglu**

{avdnoord, sedielem, heigazen, simonyan, vinyals, graves, nalk, andrewsenior, korayk}@google.com  
Google DeepMind, London, UK

<sup>†</sup> Google, London, UK



# WAVENET: A GENERATIVE MODEL FOR RAW AUDIO

Aäron van den Oord

Sander Dieleman

Heiga Zen<sup>†</sup>

Karen Simonyan

Oriol Vinyals

Alex Graves

Nal Kalchbrenner

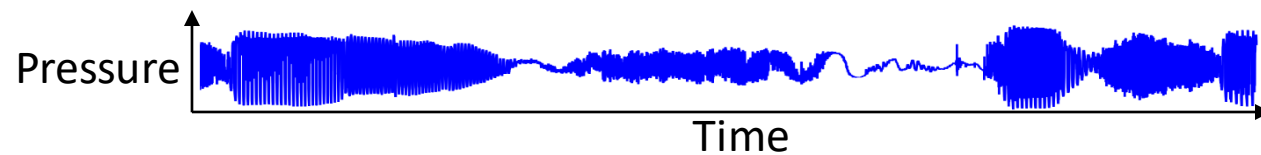
Andrew Senior

Koray Kavukcuoglu

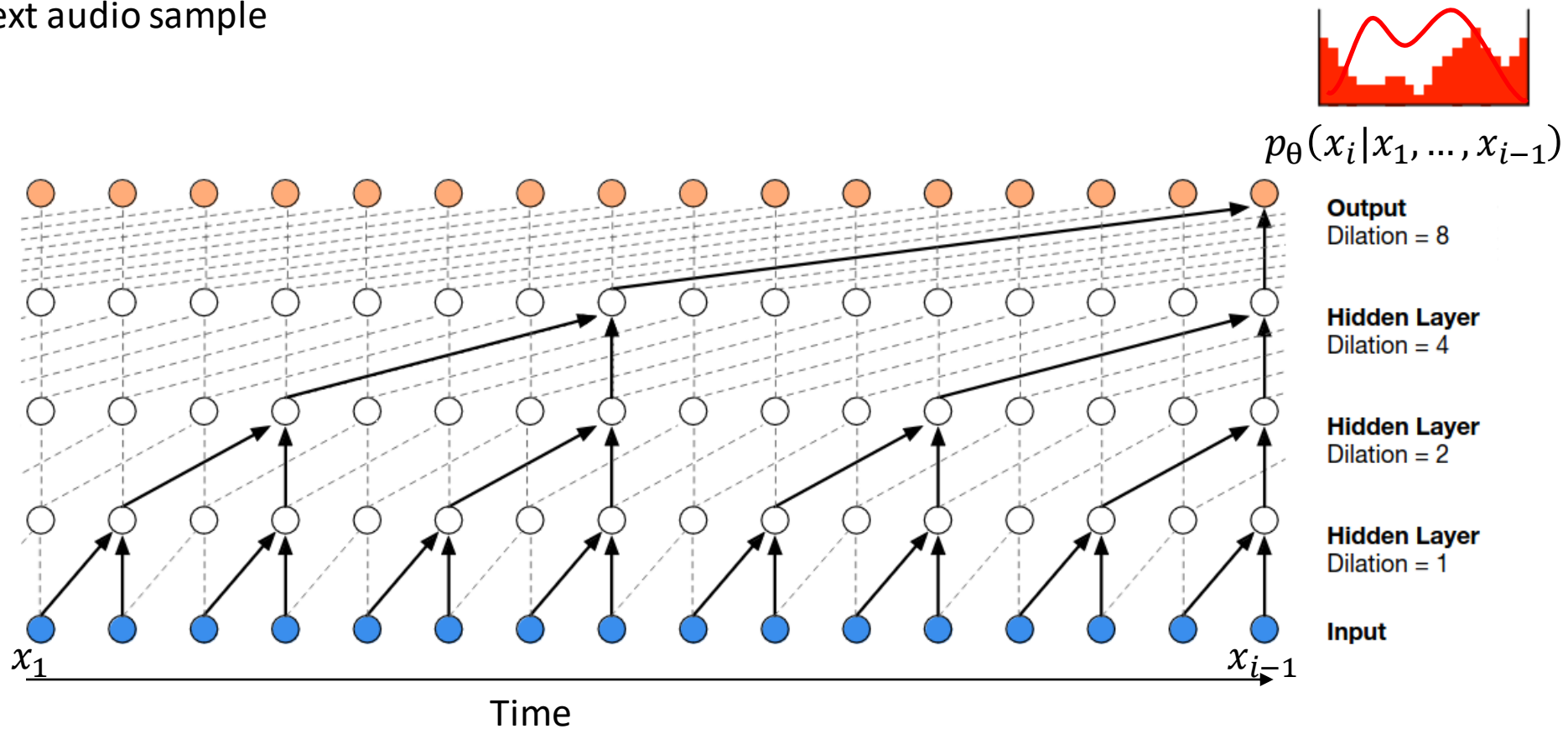
{avdnoord, sedielem, heigazen, simonyan, vinyals, graves, nalk, andrewsenior, korayk}@google.com

Google DeepMind, London, UK

<sup>†</sup> Google, London, UK



- Predict dist. for next audio sample



# WAVENET: A GENERATIVE MODEL FOR RAW AUDIO

Aäron van den Oord

Sander Dieleman

Heiga Zen<sup>†</sup>

Karen Simonyan

Oriol Vinyals

Alex Graves

Nal Kalchbrenner

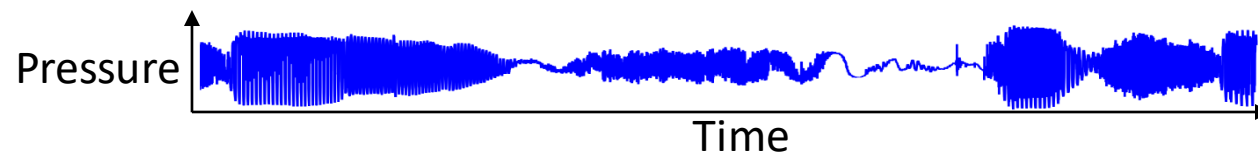
Andrew Senior

Koray Kavukcuoglu

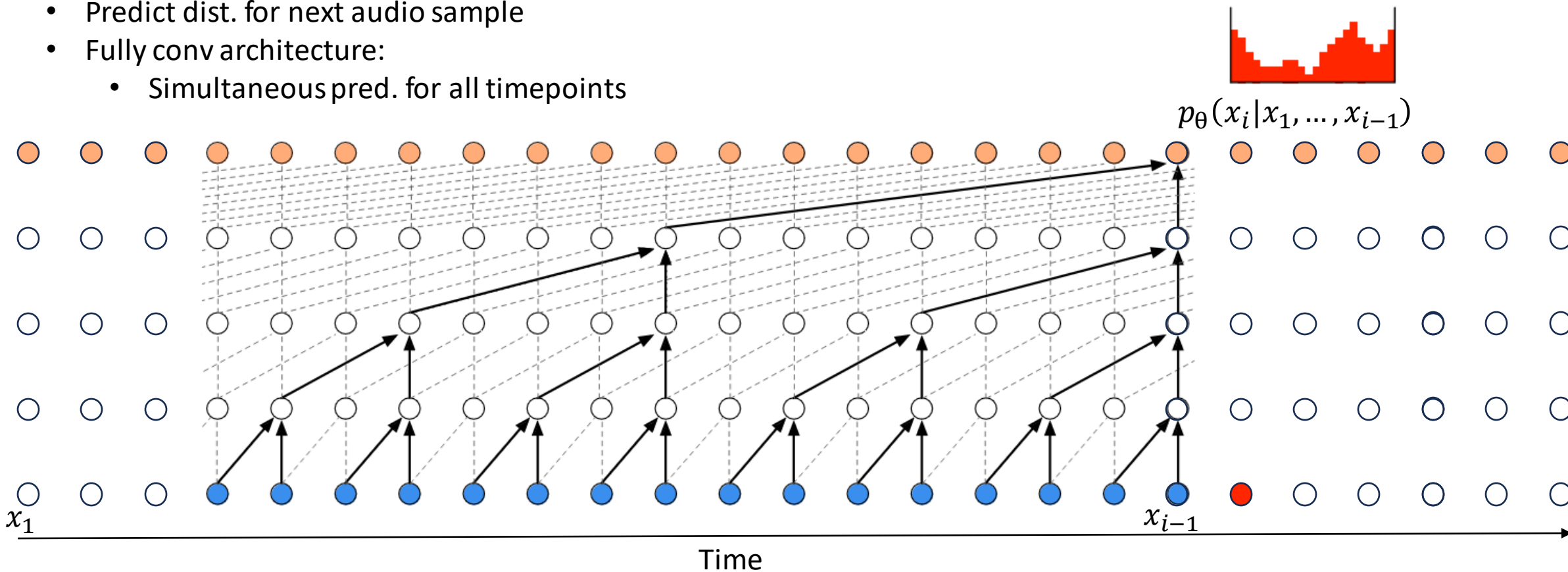
{avdnoord, sedielem, heigazen, simonyan, vinyals, graves, nalk, andrewsenior, korayk}@google.com

Google DeepMind, London, UK

<sup>†</sup> Google, London, UK



- Predict dist. for next audio sample
- Fully conv architecture:
  - Simultaneous pred. for all timepoints



# WAVENET: A GENERATIVE MODEL FOR RAW AUDIO

Aäron van den Oord

Sander Dieleman

Heiga Zen<sup>†</sup>

Karen Simonyan

Oriol Vinyals

Alex Graves

Nal Kalchbrenner

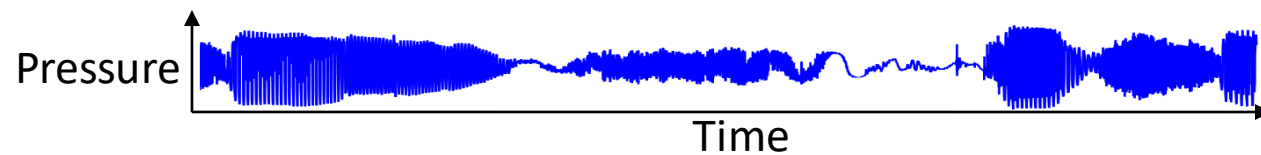
Andrew Senior

Koray Kavukcuoglu

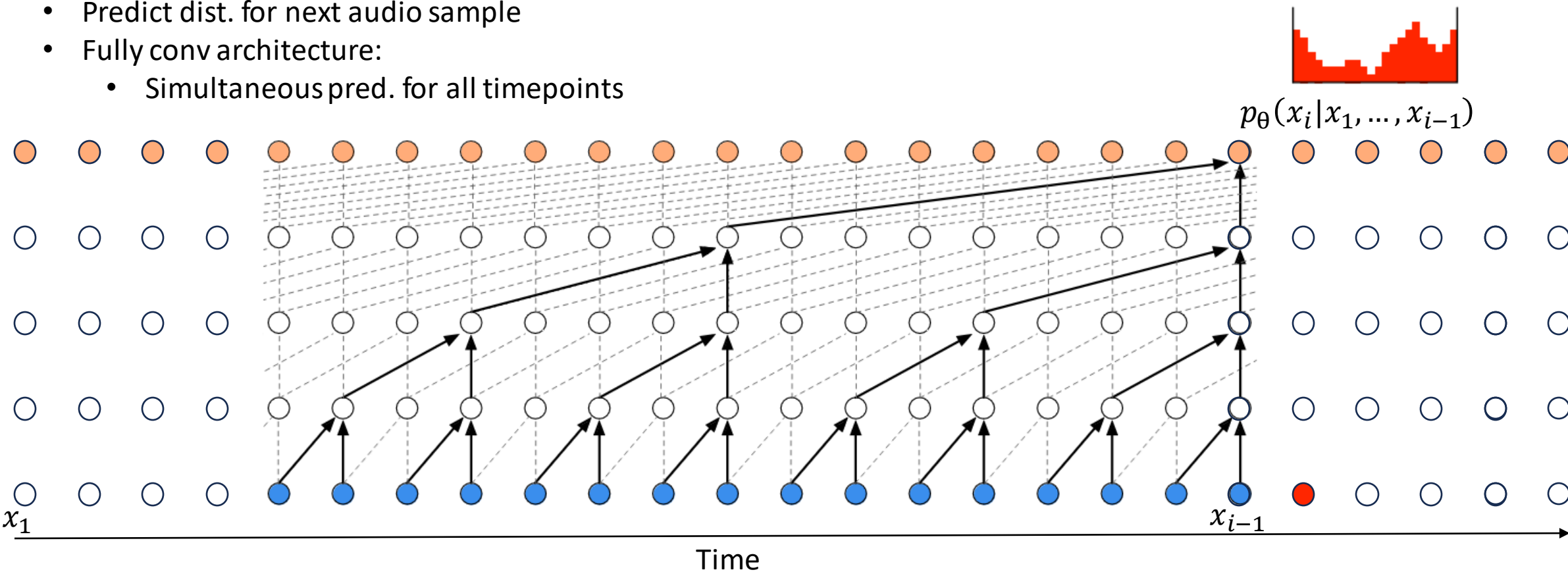
{avdnoord, sedielem, heigazen, simonyan, vinyals, graves, nalk, andrewsenior, korayk}@google.com

Google DeepMind, London, UK

<sup>†</sup> Google, London, UK



- Predict dist. for next audio sample
- Fully conv architecture:
  - Simultaneous pred. for all timepoints



# WAVENET: A GENERATIVE MODEL FOR RAW AUDIO

Aäron van den Oord

Sander Dieleman

Heiga Zen<sup>†</sup>

Karen Simonyan

Oriol Vinyals

Alex Graves

Nal Kalchbrenner

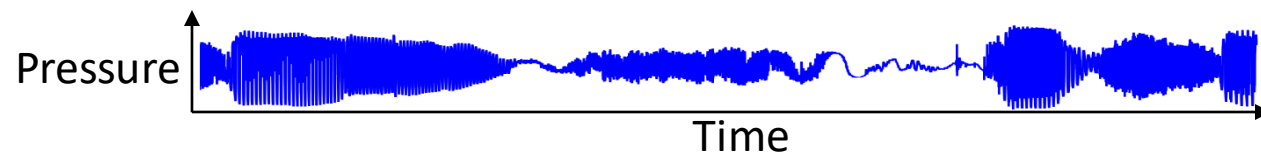
Andrew Senior

Koray Kavukcuoglu

{avdnoord, sedielem, heigazen, simonyan, vinyals, graves, nalk, andrewsenior, korayk}@google.com

Google DeepMind, London, UK

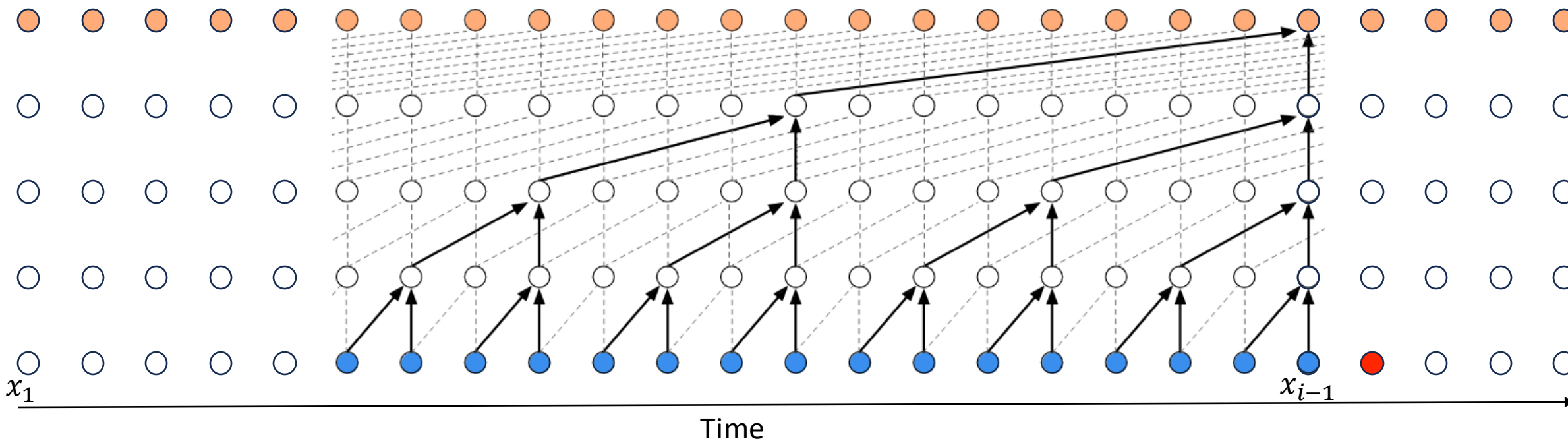
<sup>†</sup> Google, London, UK



- Predict dist. for next audio sample
- Fully conv architecture:
  - Simultaneous pred. for all timepoints

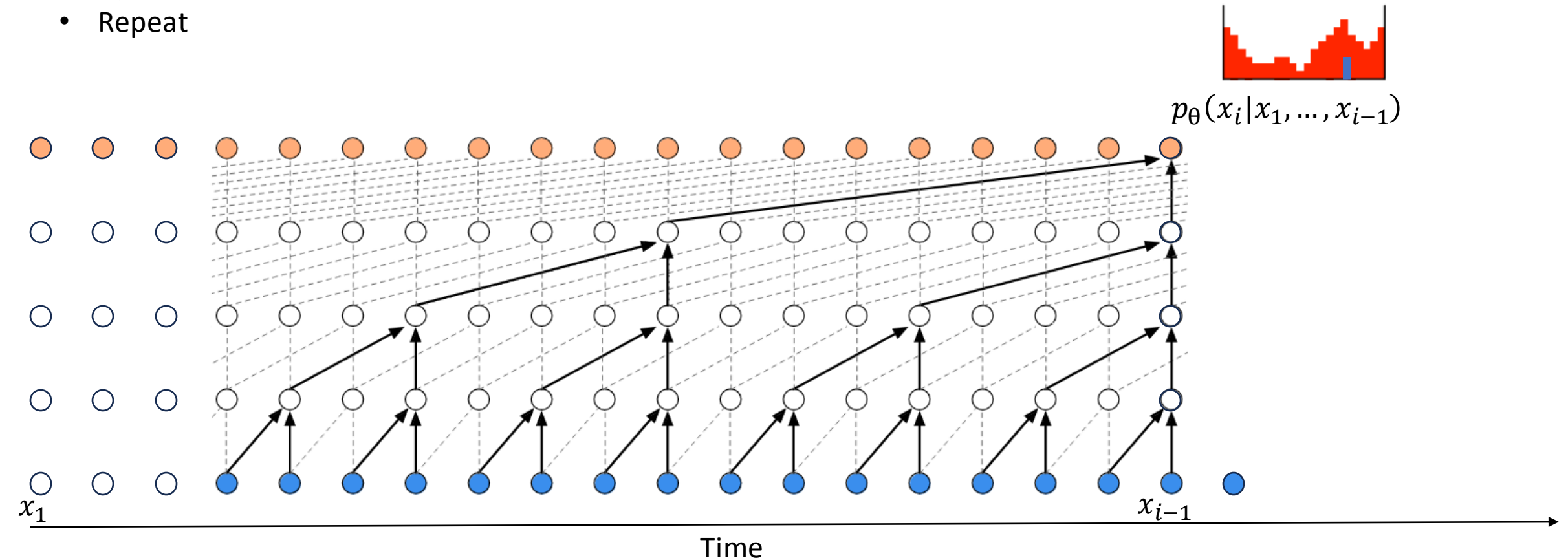


$$p_{\theta}(x_i | x_1, \dots, x_{i-1})$$



# Sampling from the Model

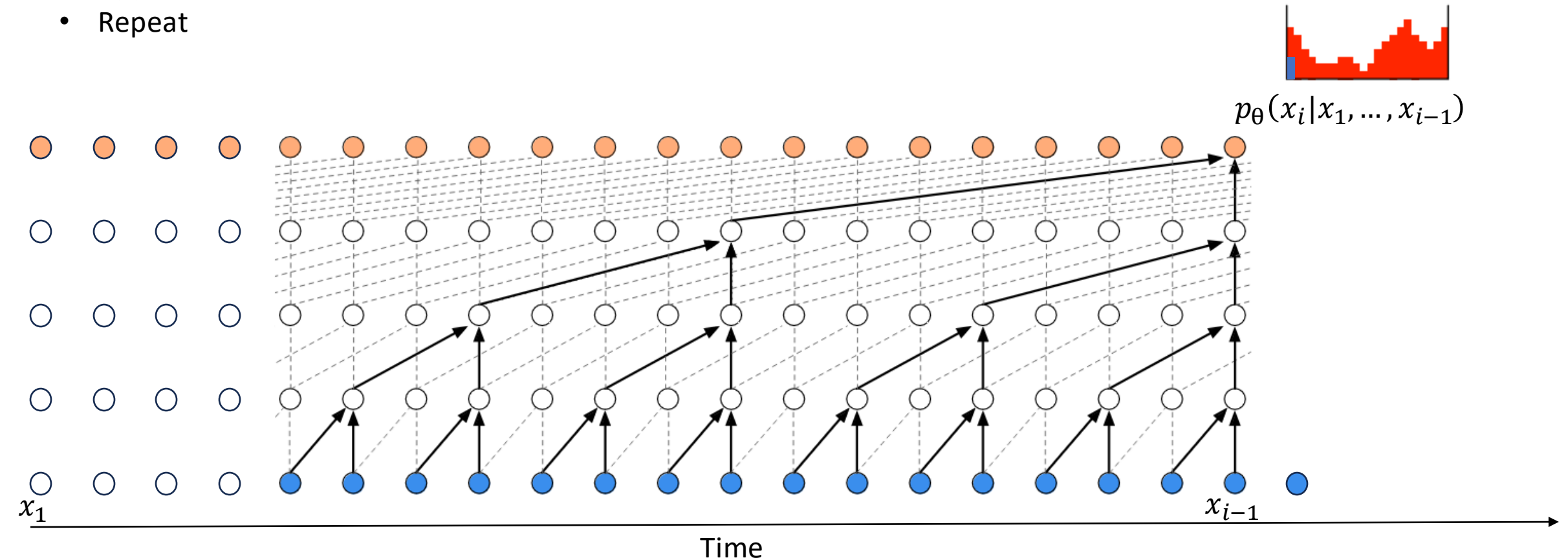
- Predict dist. for next audio sample
- Sample from distribution
- Append new sample
- Repeat





# Sampling from the Model

- Predict dist. for next audio sample
- Sample from distribution
- Append new sample
- Repeat

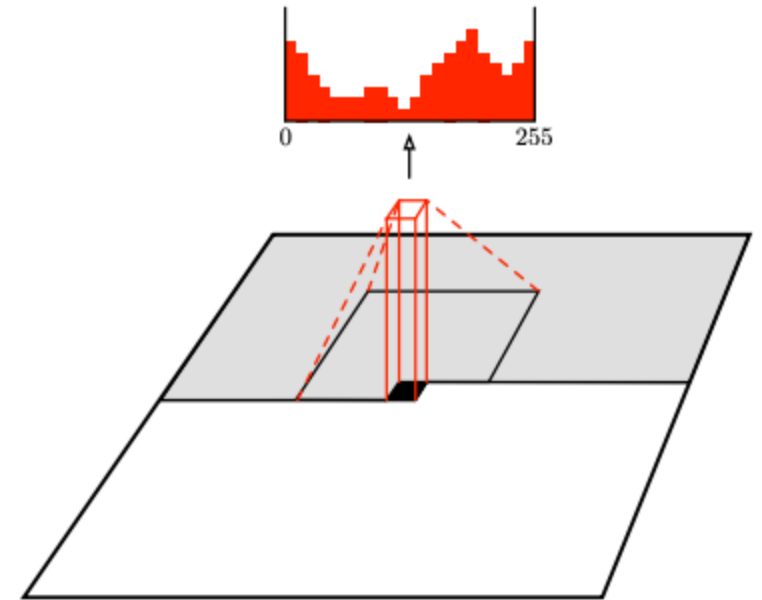
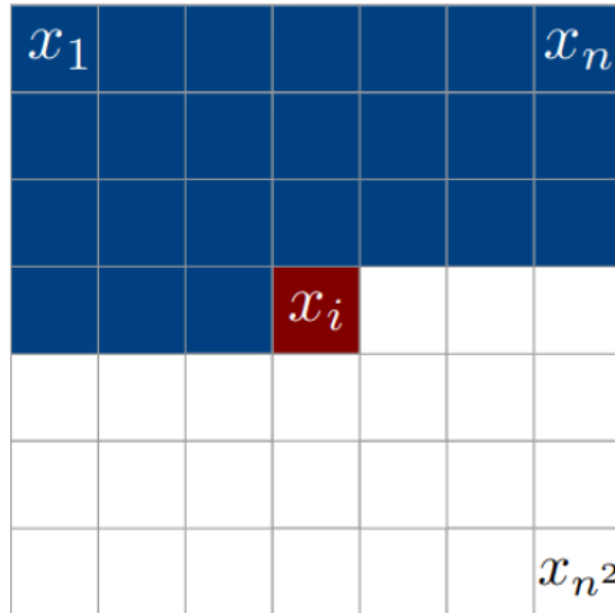


# Pixel Recurrent Neural Networks

**Aäron van den Oord**  
**Nal Kalchbrenner**  
**Koray Kavukcuoglu**

Google DeepMind

AVDNOORD@GOOGLE.COM  
NALK@GOOGLE.COM  
KORAYK@GOOGLE.COM



# Pixel Recurrent Neural Networks

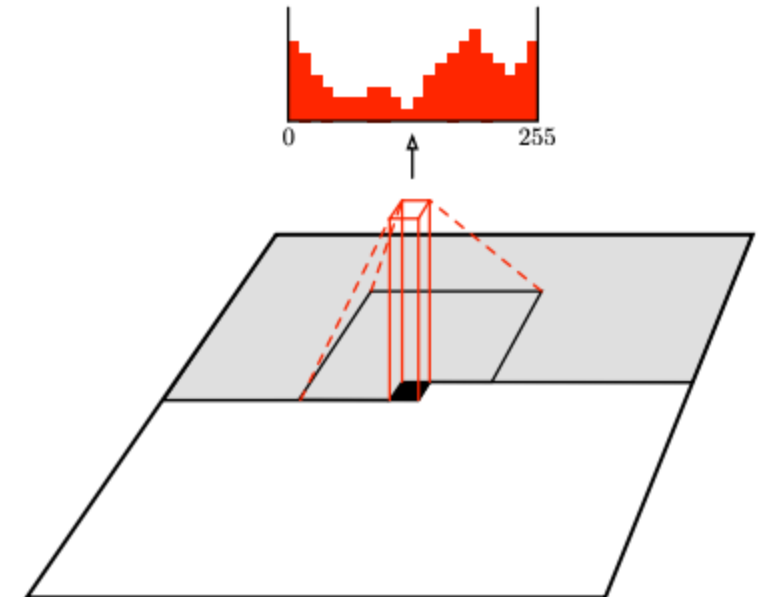
**Aäron van den Oord**  
**Nal Kalchbrenner**  
**Koray Kavukcuoglu**

Google DeepMind

## Image Generation:

- Sample one pixel
- Apply network
- Repeat

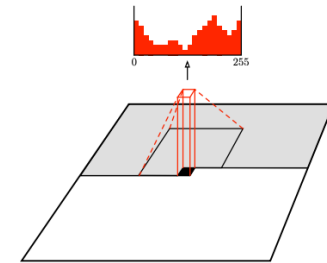
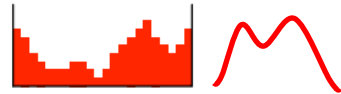
AVDNOORD@GOOGLE.COM  
NALK@GOOGLE.COM  
KORAYK@GOOGLE.COM



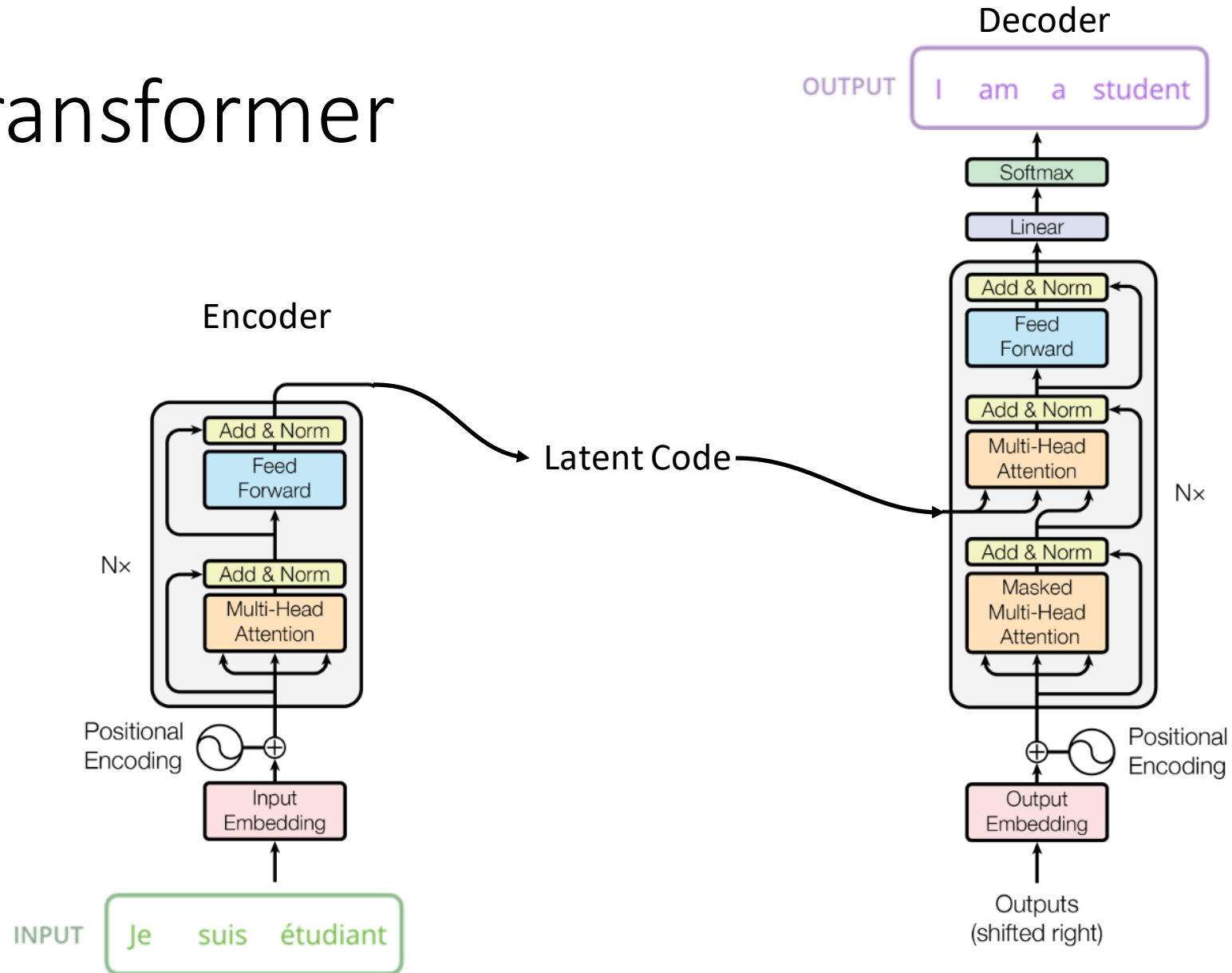
# Summary

- Interpret data as sequence
- Train neural network
  - Input: previous values  $(x_1, \dots, x_{i-1})$
  - Distribution of possible next values  $p_\theta(x_i | x_1, \dots, x_{i-1})$ 
    - E.g. as histogram
    - Or Parametric dist.
  - Ensure correct receptive field, e.g. special convolutions
- Sampling:
  - One sample at a time
  - Slow, involves repeated application of model

$$p_\theta(\mathbf{x}) = \prod_{i=1}^n p_\theta(x_i | x_1, \dots, x_{i-1})$$



# The Transformer



# The Encoder

- Process set of tokens
- Tokens remain separate
  - (except for attention layer)
- Tokens don't have order
  - (except for positional encoding)

