

# Speech Synthesis, External Memory, Meta-Learning

Milan Straka

 May 20, 2024



EUROPEAN UNION  
European Structural and Investment Fund  
Operational Programme Research,  
Development and Education

Charles University in Prague  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



unless otherwise stated

Our goal is to model speech, using a convolutional auto-regressive model

$$P(\mathbf{x}) = \prod_t P(x_t | x_{t-1}, \dots, x_1).$$

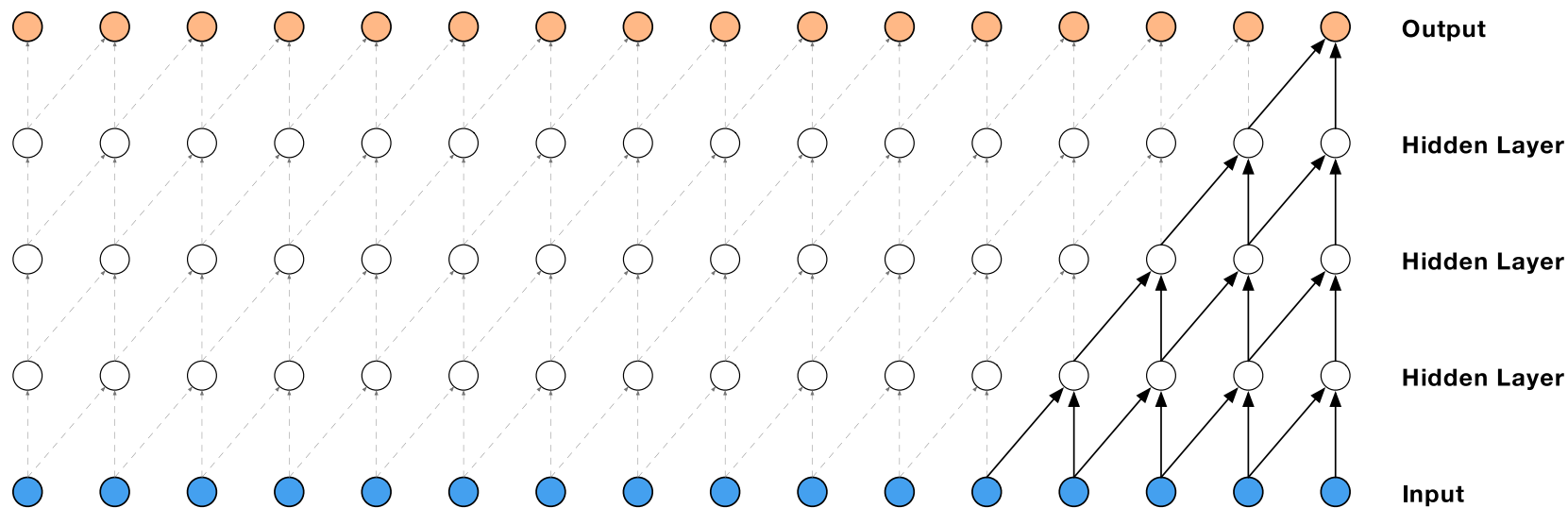


Figure 2: Visualization of a stack of causal convolutional layers.

Figure 2 of "WaveNet: A Generative Model for Raw Audio", <https://arxiv.org/abs/1609.03499>

However, to achieve larger receptive field, we utilize **dilated** (or **atrous**) convolutions:

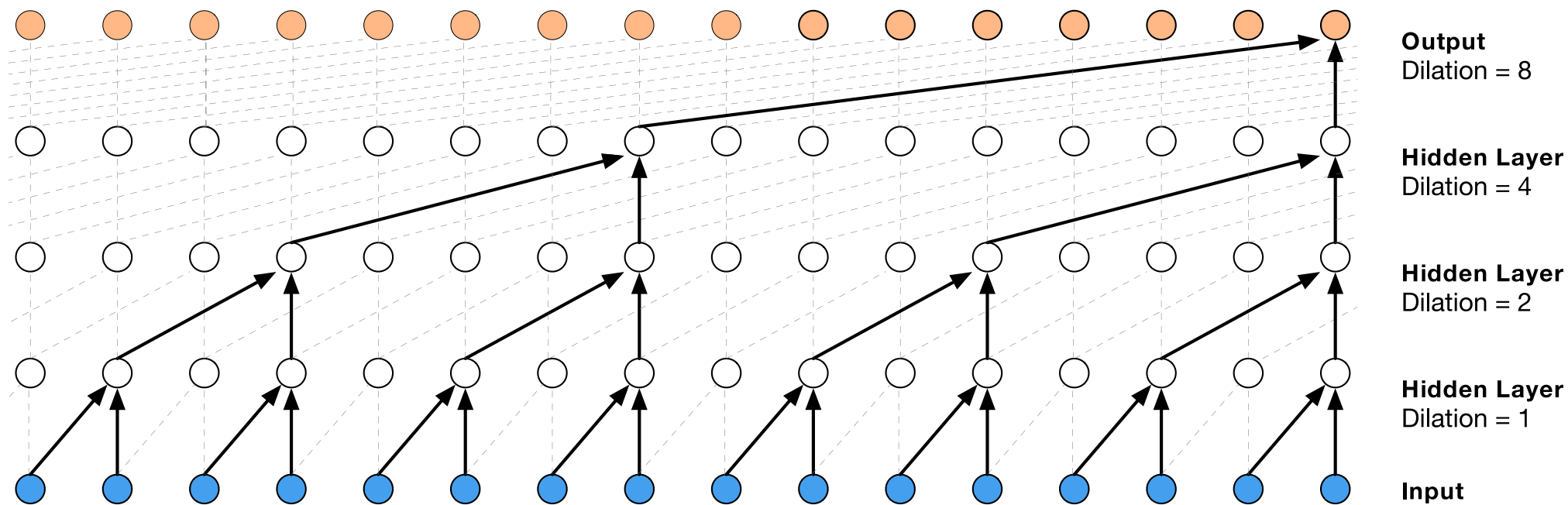
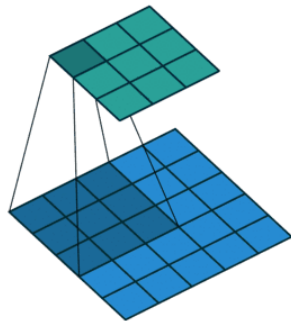


Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

Figure 3 of "WaveNet: A Generative Model for Raw Audio", <https://arxiv.org/abs/1609.03499>

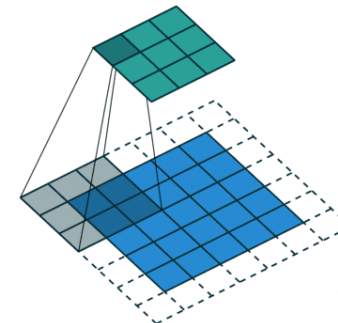
# Dilated Versus Regular Versus Strided Convolutions

## Regular Convolution



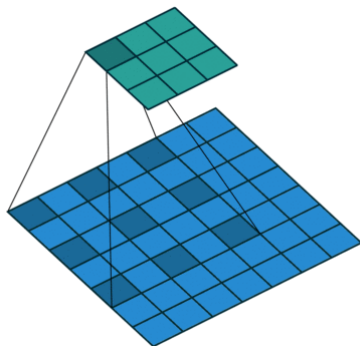
[https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

## Strided Convolution



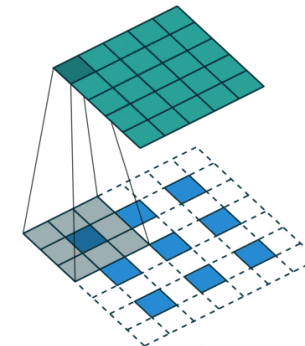
[https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

## Dilated Convolution



[https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

## Transposed Strided Convolution



[https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

## Output Distribution

WaveNet generates audio with 16kHz frequency and 16-bit samples.

However, classification into 65 536 classes would not be efficient. Instead, WaveNet adopts the  $\mu$ -law transformation, which passes the input samples in  $[-1, 1]$  range through the  $\mu$ -law encoding

$$\text{sign}(x) \frac{\log(1 + 255|x|)}{\log(1 + 255)},$$

and the resulting  $[-1, 1]$  range is linearly quantized into 256 buckets.

The model therefore predicts each samples using classification into 256 classes, and then uses the inverse of the above transformation on the model predictions.



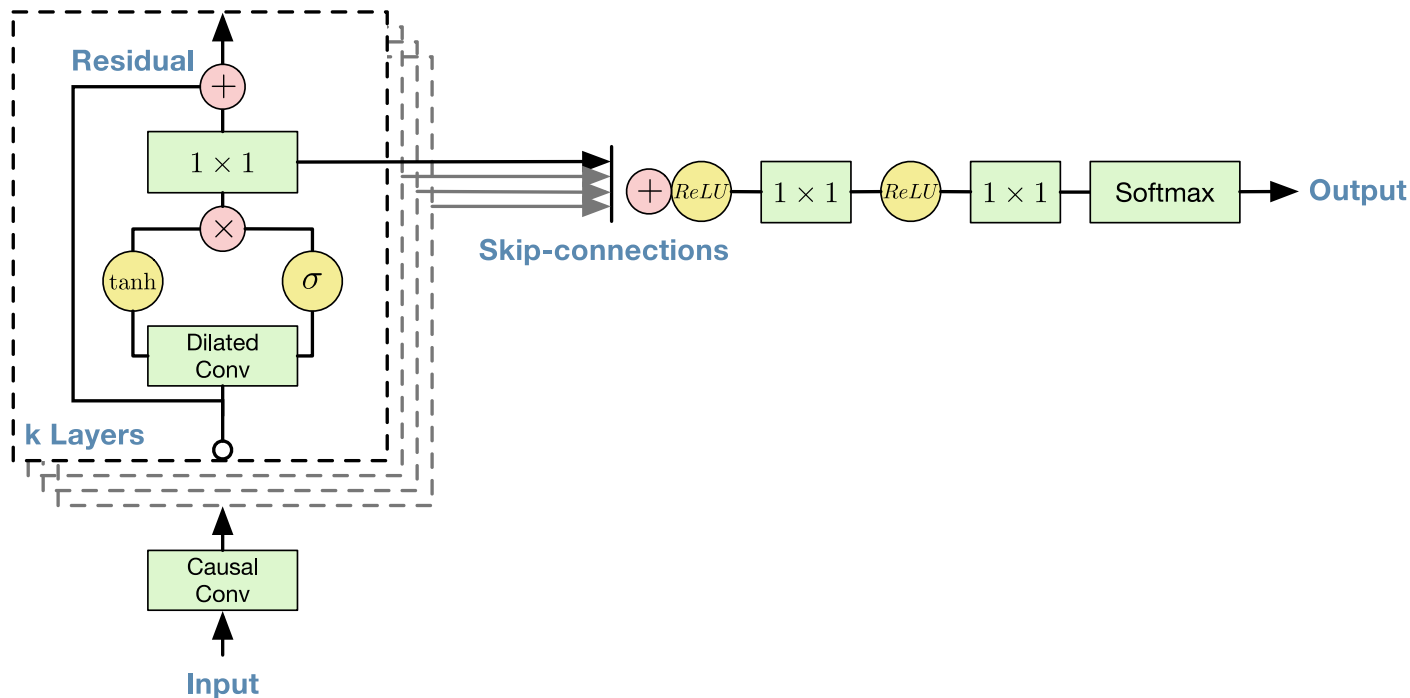


Figure 4: Overview of the residual block and the entire architecture.

Figure 4 of "WaveNet: A Generative Model for Raw Audio", <https://arxiv.org/abs/1609.03499>

The outputs of the dilated convolutions are passed through the *gated activation unit*:

$$\mathbf{z} = \tanh(\mathbf{W}_f * \mathbf{x}) \odot \sigma(\mathbf{W}_g * \mathbf{x}).$$

## Global Conditioning

Global conditioning is performed by a single latent representation  $\mathbf{h}$ , changing the gated activation function to

$$\mathbf{z} = \tanh(\mathbf{W}_f * \mathbf{x} + \mathbf{V}_f \mathbf{h}) \odot \sigma(\mathbf{W}_g * \mathbf{x} + \mathbf{V}_g \mathbf{h}).$$

## Local Conditioning

For local conditioning, we are given a time series  $\mathbf{h}$ , possibly with a lower sampling frequency. We first use transposed convolutions  $\mathbf{y} = f(\mathbf{h})$  to match resolution and then compute analogously to global conditioning

$$\mathbf{z} = \tanh(\mathbf{W}_f * \mathbf{x} + \mathbf{V}_f * \mathbf{y}) \odot \sigma(\mathbf{W}_g * \mathbf{x} + \mathbf{V}_g * \mathbf{y}).$$

The original paper did not mention hyperparameters, but later it was revealed that:

- 30 layers were used
  - grouped into 3 dilation stacks with 10 layers each
  - in a dilation stack, dilation rate increases by a factor of 2, starting with rate 1 and reaching maximum dilation of 512
- kernel size of a dilated convolution is 2 (and increased to 3 in Parallel WaveNet)
- residual connection has dimension 512
- gating layer uses 256+256 hidden units
- the  $1 \times 1$  convolutions in the output step produce 256 filters
- trained for 1 000 000 steps using Adam with a fixed learning rate of  $2e-4$



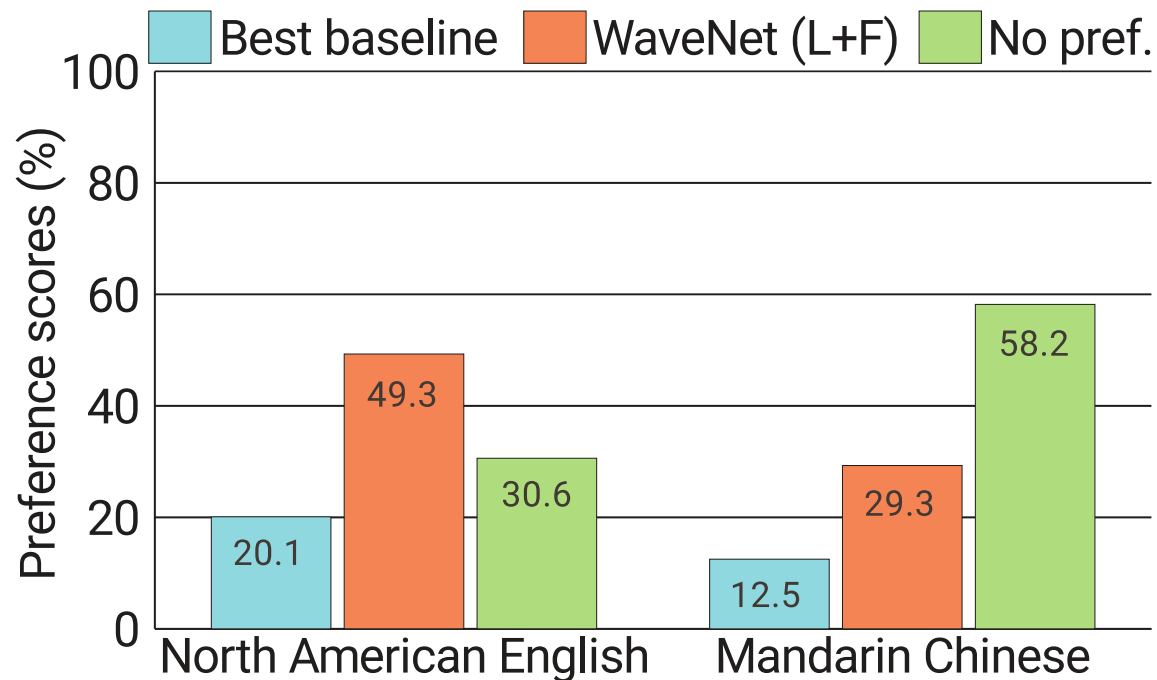


Figure 5: Subjective preference scores (%) of speech samples between (top) two baselines, (middle) two WaveNets, and (bottom) the best baseline and WaveNet. Note that LSTM and Concat correspond to LSTM-RNN-based statistical parametric and HMM-driven unit selection concatenative baseline synthesizers, and WaveNet (L) and WaveNet (L+F) correspond to the WaveNet conditioned on linguistic features only and that conditioned on both linguistic features and  $\log F_0$  values.

Figure 5 of "WaveNet: A Generative Model for Raw Audio", <https://arxiv.org/abs/1609.03499>

# Gated Activations in Transformers

Similar gated activations seem to work the best in Transformers, in the FFN module.

Activation Name	Formula	$\text{FFN}(x; \mathbf{W}_1, \mathbf{W}_2)$
ReLU	$\max(0, x)$	$\max(0, \mathbf{xW}_1)\mathbf{W}_2$
GELU	$x\Phi(x)$	$\text{GELU}(\mathbf{xW}_1)\mathbf{W}_2$
Swish	$x\sigma(x)$	$\text{Swish}(\mathbf{xW}_1)\mathbf{W}_2$

There are several variants of the new gated activations:

Activation Name	Formula	$\text{FFN}(x; \mathbf{W}, \mathbf{V}, \mathbf{W}_2)$
GLU (Gated Linear Unit)	$\sigma(\mathbf{xW} + \mathbf{b}) \odot (\mathbf{xV} + \mathbf{c})$	$(\sigma(\mathbf{xW}) \odot \mathbf{xV})\mathbf{W}_2$
ReGLU	$\max(0, \mathbf{xW} + \mathbf{b}) \odot (\mathbf{xV} + \mathbf{c})$	$(\max(0, \mathbf{xW}) \odot \mathbf{xV})\mathbf{W}_2$
GEGLU	$\text{GELU}(\mathbf{xW} + \mathbf{b}) \odot (\mathbf{xV} + \mathbf{c})$	$(\text{GELU}(\mathbf{xW}) \odot \mathbf{xV})\mathbf{W}_2$
SwiGLU	$\text{Swish}(\mathbf{xW} + \mathbf{b}) \odot (\mathbf{xV} + \mathbf{c})$	$(\text{Swish}(\mathbf{xW}) \odot \mathbf{xV})\mathbf{W}_2$

# Gated Activations in Transformers

	Score Average	CoLA MCC	SST-2 Acc	MRPC F1	MRPC Acc	STSB PCC	STSB SCC	QQP F1	QQP Acc	MNLIm Acc	MNLImm Acc	QNLI Acc	RTE Acc	EM	F1
FFN <sub>ReLU</sub>	83.80	51.32	94.04	<b>93.08</b>	<b>90.20</b>	89.64	89.42	89.01	91.75	85.83	86.42	92.81	80.14	83.18	90.87
FFN <sub>GELU</sub>	83.86	53.48	94.04	92.81	<b>90.20</b>	89.69	89.49	88.63	91.62	85.89	86.13	92.39	80.51	83.09	90.79
FFN <sub>Swish</sub>	83.60	49.79	93.69	92.31	89.46	89.20	88.98	88.84	91.67	85.22	85.02	92.33	81.23	83.25	90.76
FFN <sub>GLU</sub>	84.20	49.16	94.27	92.39	89.46	89.46	89.35	88.79	91.62	86.36	86.18	92.92	<b>84.12</b>	82.88	90.69
FFN <sub>GEGLU</sub>	84.12	53.65	93.92	92.68	89.71	90.26	90.13	89.11	91.85	86.15	86.17	92.81	79.42	83.55	91.12
FFN <sub>Bilinear</sub>	83.79	51.02	<b>94.38</b>	92.28	89.46	90.06	89.84	88.95	91.69	<b>86.90</b>	<b>87.08</b>	92.92	81.95	<b>83.82</b>	91.06
FFN <sub>SwiGLU</sub>	84.36	51.59	93.92	92.23	88.97	<b>90.32</b>	<b>90.13</b>	<b>89.14</b>	<b>91.87</b>	86.45	86.47	<b>92.93</b>	83.39	83.42	91.03
FFN <sub>ReLU</sub>	<b>84.67</b>	<b>56.16</b>	<b>94.38</b>	92.06	89.22	89.97	89.85	88.86	91.72	86.20	86.40	92.68	81.59	83.53	<b>91.18</b>

Table 2 of "GLU Variants Improve Transformer", <https://arxiv.org/abs/2002.05202> Table 4 of "GLU Variants Improve Transformer", <https://arxiv.org/abs/2002.05202>

Model	Params	Ops	Step/s	Early loss	Final loss	SGLUE	XSum	WebQ	WMT EnDe
Vanilla Transformer	223M	11.1T	3.50	2.182 ± 0.005	1.838	71.66	17.78	23.02	26.62
GeLU	223M	11.1T	3.58	2.179 ± 0.003	1.838	<b>75.79</b>	<b>17.86</b>	<b>25.13</b>	26.47
Swish	223M	11.1T	3.62	2.186 ± 0.003	1.847	<b>73.77</b>	17.74	<b>24.34</b>	<b>26.75</b>
ELU	223M	11.1T	3.56	2.270 ± 0.007	1.932	67.83	16.73	23.02	26.08
GLU	223M	11.1T	3.59	2.174 ± 0.003	<b>1.814</b>	<b>74.20</b>	<b>17.42</b>	24.34	<b>27.12</b>
GeGLU	223M	11.1T	3.55	2.130 ± 0.006	<b>1.792</b>	<b>75.96</b>	<b>18.27</b>	<b>24.87</b>	<b>26.87</b>
ReGLU	223M	11.1T	3.57	2.145 ± 0.004	<b>1.803</b>	<b>76.17</b>	<b>18.36</b>	<b>24.87</b>	<b>27.02</b>
SeLU	223M	11.1T	3.55	2.315 ± 0.004	1.948	68.76	16.76	22.75	25.99
SwiGLU	223M	11.1T	3.53	2.127 ± 0.003	<b>1.789</b>	<b>76.00</b>	<b>18.20</b>	<b>24.34</b>	<b>27.02</b>
LiGLU	223M	11.1T	3.59	2.149 ± 0.005	<b>1.798</b>	<b>75.34</b>	<b>17.97</b>	<b>24.34</b>	26.53
Sigmoid	223M	11.1T	3.63	2.291 ± 0.019	1.867	<b>74.31</b>	17.51	23.02	26.30
Softplus	223M	11.1T	3.47	2.207 ± 0.011	1.850	<b>72.45</b>	17.65	<b>24.34</b>	<b>26.89</b>

Table 1 of "Do Transformer Modifications Transfer Across Implementations and Applications?", <https://arxiv.org/abs/2102.11972>

Parallel WaveNet is an improvement of the original WaveNet by the same authors.

First, the output distribution was changed from 256  $\mu$ -law values to a Mixture of Logistic (suggested in another paper – PixelCNN++, but reused in other architectures since):

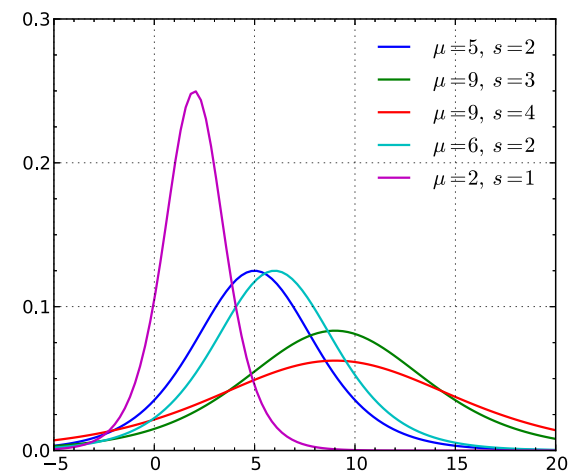
$$x \sim \sum_i \pi_i \text{Logistic}(\mu_i, s_i).$$

The logistic distribution is a distribution with a  $\sigma$  as cumulative density function (where the mean and scale is parametrized by  $\mu$  and  $s$ ). Therefore, we can write

$$P(x|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{s}) = \sum_i \pi_i \left[ \sigma\left(\frac{x + 0.5 - \mu_i}{s_i}\right) - \sigma\left(\frac{x - 0.5 - \mu_i}{s_i}\right) \right],$$

where we replace  $-0.5$  and  $0.5$  in the edge cases by  $-\infty$  and  $\infty$ .

In Parallel WaveNet teacher, 10 mixture components are used.



<https://commons.wikimedia.org/wiki/File:Logisticpdffunction.svg>

Auto-regressive (sequential) inference is extremely slow in WaveNet.

Instead, we model  $P(x_t)$  as  $P(x_t | \mathbf{z}_{<t}) = \text{Logistic}(x_t; \mu^1(\mathbf{z}_{<t}), s^1(\mathbf{z}_{<t}))$  for a *random*  $\mathbf{z}$  drawn from a logistic distribution  $\text{Logistic}(\mathbf{0}, \mathbf{1})$ . Therefore, using the reparametrization trick,

$$\mathbf{x}_t^1 = \mu^1(\mathbf{z}_{<t}) + z_t \cdot s^1(\mathbf{z}_{<t}).$$

Usually, one iteration of the algorithm does not produce good enough results – consequently, 4 iterations were used by the authors. In further iterations,

$$\mathbf{x}_t^i = \mu^i(\mathbf{x}_{<t}^{i-1}) + x_t^{i-1} \cdot s^i(\mathbf{x}_{<t}^{i-1}).$$

After  $N$  iterations,  $P(\mathbf{x}_t^N | \mathbf{z}_{<t})$  is a logistic distribution with location  $\mu^{\text{tot}}$  and scale  $s^{\text{tot}}$ :

$$\mu_t^{\text{tot}} = \sum_{i=1}^N \mu^i(\mathbf{x}_{<t}^{i-1}) \cdot \left( \prod_{j>i}^N s^j(\mathbf{x}_{<t}^{j-1}) \right) \quad \text{and} \quad s_t^{\text{tot}} = \prod_{i=1}^N s^i(\mathbf{x}_{<t}^{i-1}),$$

where we have denoted  $\mathbf{z}$  as  $\mathbf{x}^0$  for convenience.

The consequences of changing the model to

$$\begin{aligned} \mathbf{x}_t^1 &= \mu^1(\mathbf{z}_{<t}) + z_t \cdot \mathbf{s}^1(\mathbf{z}_{<t}) \\ \mathbf{x}_t^i &= \mu^i(\mathbf{x}_{<t}^{i-1}) + x_t^{i-1} \cdot \mathbf{s}^i(\mathbf{x}_{<t}^{i-1}) \end{aligned}$$

are:

- During inference, the prediction can be computed in parallel, because  $x_t^i$  depends only on  $\mathbf{x}_{<t}^{i-1}$ , not on  $\mathbf{x}_{<t}^i$ .
- However, we cannot perform training in parallel. If we try maximizing the log-likelihood of an input sequence  $\mathbf{x}^1$ , we need to find out which  $\mathbf{z}$  sequence generates it.
  - The  $z_1$  can be computed using  $x_1^1$ .
  - However,  $z_2$  depends not only on  $x_1^1$  and  $x_2^1$ , but also on  $z_1$ ; generally,  $z_t$  depends on  $\mathbf{x}^1$  and also on all  $\mathbf{z}_{<t}$ , and can be computed only sequentially.

Therefore, WaveNet can perform parallel training and sequential inference, while the proposed model can perform parallel inference but sequential training.

# Probability Density Distillation

The authors propose to train the network by a **probability density distillation** using a teacher WaveNet (producing a mixture of logistic with 10 components) with KL-divergence as a loss.

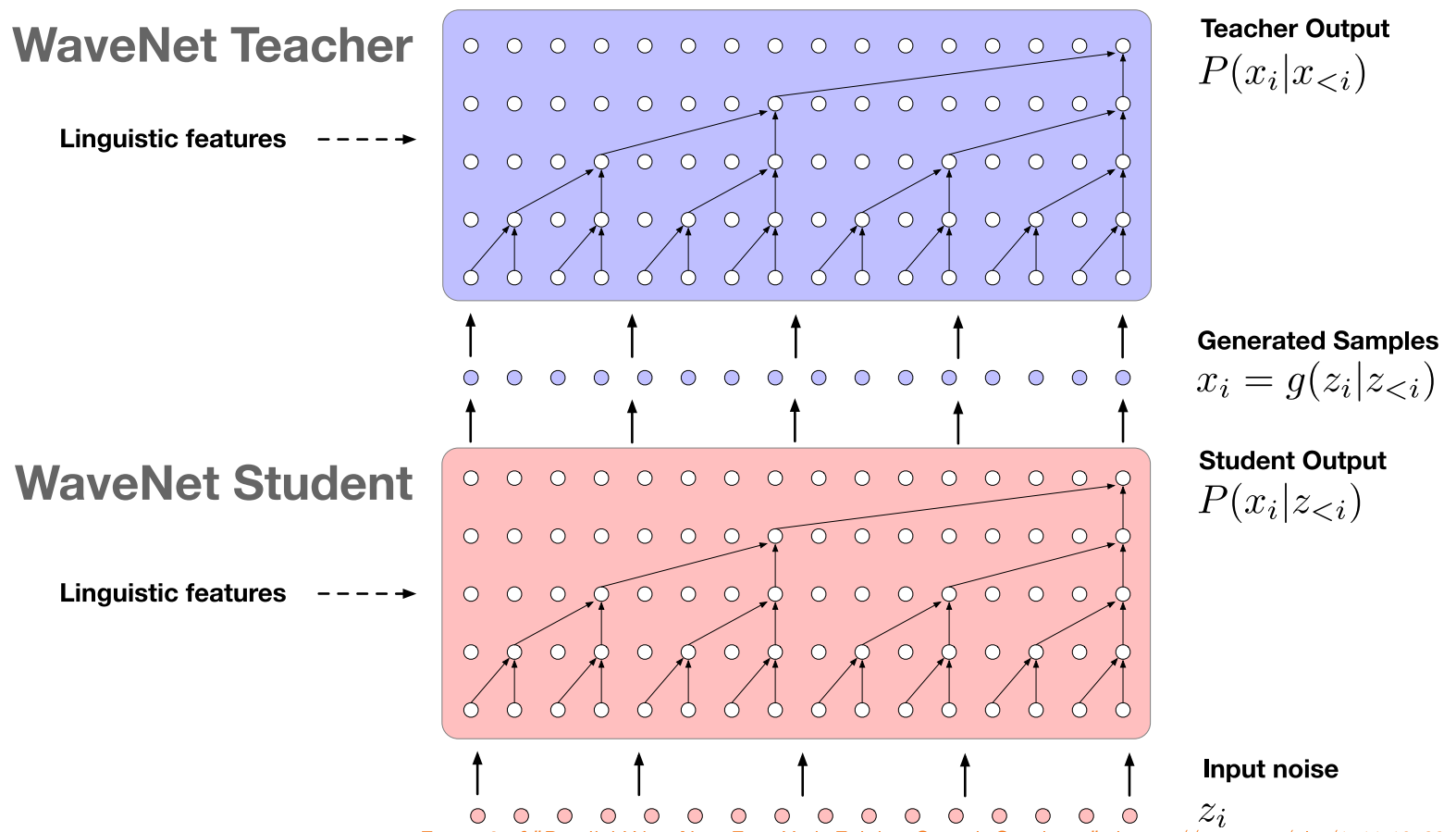


Figure 2 of "Parallel WaveNet: Fast High-Fidelity Speech Synthesis", <https://arxiv.org/abs/1711.10433>

Therefore, instead of computing  $\mathbf{z}$  from some gold  $\mathbf{x}_g$ , we

- sample a random  $\mathbf{z}$ ;
- generate the output  $\mathbf{x}$ ;
- use the teacher WaveNet model to estimate the log-likelihood of  $\mathbf{x}$ ;
- update the student to match the log-likelihood of the teacher.

Denoting the teacher distribution as  $P_T$  and the student distribution as  $P_S$ , the loss is

$$D_{\text{KL}}(P_S || P_T) = H(P_S, P_T) - H(P_S).$$

Therefore, we do not only minimize cross-entropy, but we also try to keep the entropy of the student as high as possible – it is indeed crucial not to match just the mode of the teacher.

- Consider a teacher generating white noise, where every sample comes from  $\mathcal{N}(0, 1)$  – in this case, the cross-entropy loss of a constant  $\mathbf{0}$ , complete silence, would be maximal.

In a sense, probability density distillation is similar to GANs. However, the teacher is kept fixed, and the student does not attempt to fool it but to match its distribution instead.



# Probability Density Distillation Details

Because the entropy of a logistic distribution  $\text{Logistic}(\mu, s)$  is  $\log s + 2$ , the entropy term  $H(P_S)$  can be rewritten as follows:

$$\begin{aligned} H(P_S) &= \mathbb{E}_{z \sim \text{Logistic}(0,1)} \left[ \sum_{t=1}^T -\log p_S(x_t | \mathbf{z}_{<t}) \right] \\ &= \mathbb{E}_{z \sim \text{Logistic}(0,1)} \left[ \sum_{t=1}^T \log s(\mathbf{z}_{<t}, \boldsymbol{\theta}) \right] + 2T. \end{aligned}$$

Therefore, this term can be computed without having to generate  $\mathbf{x}$ .

# Probability Density Distillation Details

However, the cross-entropy term  $H(P_S, P_T)$  requires sampling from  $P_S$  to estimate:

$$\begin{aligned}
 H(P_S, P_T) &= \int_{\mathbf{x}} -P_S(\mathbf{x}) \log P_T(\mathbf{x}) \\
 &= \sum_{t=1}^T \int_{\mathbf{x}} -P_S(\mathbf{x}) \log P_T(x_t | \mathbf{x}_{<t}) \\
 &= \sum_{t=1}^T \int_{\mathbf{x}} -P_S(\mathbf{x}_{<t}) P_S(x_t | \mathbf{x}_{<t}) P_S(\mathbf{x}_{>t} | \mathbf{x}_{\leq t}) \log P_T(x_t | \mathbf{x}_{<t}) \\
 &= \sum_{t=1}^T \mathbb{E}_{P_S(\mathbf{x}_{<t})} \left[ \int_{x_t} -P_S(x_t | \mathbf{x}_{<t}) \log P_T(x_t | \mathbf{x}_{<t}) \underbrace{\int_{\mathbf{x}_{>t}} P_S(\mathbf{x}_{>t} | \mathbf{x}_{\leq t})}_{1} \right] \\
 &= \sum_{t=1}^T \mathbb{E}_{P_S(\mathbf{x}_{<t})} H\left(P_S(x_t | \mathbf{x}_{<t}), P_T(x_t | \mathbf{x}_{<t})\right).
 \end{aligned}$$

$$H(P_S, P_T) = \sum_{t=1}^T \mathbb{E}_{P_S(\mathbf{x}_{<t})} H\left(P_S(x_t | \mathbf{x}_{<t}), P_T(x_t | \mathbf{x}_{<t})\right)$$

We can therefore estimate  $H(P_S, P_T)$  by:

- drawing a single sample  $\mathbf{x}$  from the student  $P_S$  [a  $Logistic(\boldsymbol{\mu}^{\text{tot}}, \mathbf{s}^{\text{tot}})$ ],
- compute all  $P_T(x_t | \mathbf{x}_{<t})$  from the teacher in parallel [*mixture of logistic distributions*],
- and finally evaluate  $H(P_S(x_t | \mathbf{x}_{<t}), P_T(x_t | \mathbf{x}_{<t}))$  by sampling multiple different  $x_t$  from the  $P_S(x_t | \mathbf{x}_{<t})$ .

The authors state that this unbiased estimator has a much lower variance than naively evaluating a single sequence sample under the teacher using the original formulation.

Finally, analogously to the normal distribution, the logistic distribution offers the **reparametrization trick**. Therefore, we can differentiate  $\log P_T(x_t | \mathbf{x}_{<t})$  with respect to both  $x_t$  and  $\mathbf{x}_{<t}$  (while the categorical distribution is differentiable only with respect to  $\mathbf{x}_{<t}$ ).

The Parallel WaveNet model consists of 4 iterations with 10, 10, 10, 30 layers, respectively. The dimension of the residuals and the gating units is 64 (compared to 512 in WaveNet).

The Parallel WaveNet generates over 500k samples per second, compared to ~170 samples per second of a regular WaveNet – more than a 1000 times speedup.

Method	Subjective 5-scale MOS
<b>16kHz, 8-bit <math>\mu</math>-law, 25h data:</b>	
LSTM-RNN parametric [27]	$3.67 \pm 0.098$
HMM-driven concatenative [27]	$3.86 \pm 0.137$
WaveNet [27]	$4.21 \pm 0.081$
<b>24kHz, 16-bit linear PCM, 65h data:</b>	
HMM-driven concatenative	$4.19 \pm 0.097$
Autoregressive WaveNet	$4.41 \pm 0.069$
Distilled WaveNet	$4.41 \pm 0.078$

Table 1 of "Parallel WaveNet: Fast High-Fidelity Speech Synthesis", <https://arxiv.org/abs/1711.10433>

For comparison, using a single iteration with 30 layers achieve MOS of 4.21.

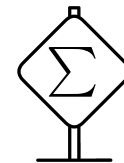
The Parallel WaveNet can be trained to generate speech of multiple speakers (using the global conditioning). Because such a model needs larger capacity, it used 30 layers in every iteration (instead of 10, 10, 10, 30).

	Parametric	Concatenative	Distilled WaveNet
<b>English speaker 1</b> (female - 65h data)	3.88	4.19	4.41
<b>English speaker 2</b> (male - 21h data)	3.96	4.09	4.34
<b>English speaker 3</b> (male - 10h data)	3.77	3.65	4.47
<b>English speaker 4</b> (female - 9h data)	3.42	3.40	3.97
<b>Japanese speaker</b> (female - 28h data)	4.07	3.47	4.23

Table 2: Comparison of MOS scores on English and Japanese with multi-speaker distilled WaveNets. Note that some speakers sounded less appealing to people and always get lower MOS, however distilled parallel WaveNet always achieved significantly better results.

*Table 2 of "Parallel WaveNet: Fast High-Fidelity Speech Synthesis", <https://arxiv.org/abs/1711.10433>*

To generate high-quality audio, the probability density distillation is not entirely sufficient. The authors therefore introduce additional losses:

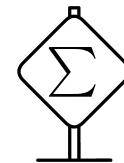


- **power loss**: ensures the power in different frequency bands is on average similar between the generated speech and human speech. For a conditioned training data  $(\mathbf{x}, \mathbf{c})$  and WaveNet student  $g$ , the loss is

$$\| \text{STFT}(g(\mathbf{z}, \mathbf{c})) - \text{STFT}(\mathbf{x}) \|^2.$$

- **perceptual loss**: apart from the power in frequency bands, we can use a pre-trained classifier to extract features from generated and human speech and add a loss measuring their difference. The authors propose the loss as squared Frobenius norm of differences between Gram matrices (uncentered covariance matrices) of features of a WaveNet-like classifier predicting phones from raw audio.
- **contrastive loss**: to make the model respect the conditioning instead of generating outputs with high likelihood independent on the conditioning, the authors propose a contrastive distillation loss ( $\gamma = 0.3$  is used in the paper):

$$D_{\text{KL}}(P_S(\mathbf{c}_1) || P_T(\mathbf{c}_1)) - \gamma D_{\text{KL}}(P_S(\mathbf{c}_1) || P_T(\mathbf{c}_2)).$$



Method	Preference Scores versus baseline concatenative system Win - Lose - Neutral
<b>Losses used</b>	
KL + Power	60% - 15% - 25%
KL + Power + Perceptual	66% - 10% - 24%
KL + Power + Perceptual + Contrastive (= default)	65% - 9% - 26%

Table 3: Performance with respect to different combinations of loss terms. We report preference comparison scores since their mean opinion scores tend to be very close and inconclusive.

*Table 3 of "Parallel WaveNet: Fast High-Fidelity Speech Synthesis", <https://arxiv.org/abs/1711.10433>*

Tacotron 2 model presents end-to-end speech synthesis directly from text. It consists of two components trained separately:

- a seq2seq model processing input characters and generating mel spectrograms;
- a Parallel WaveNet generating the speech from Mel spectrograms.

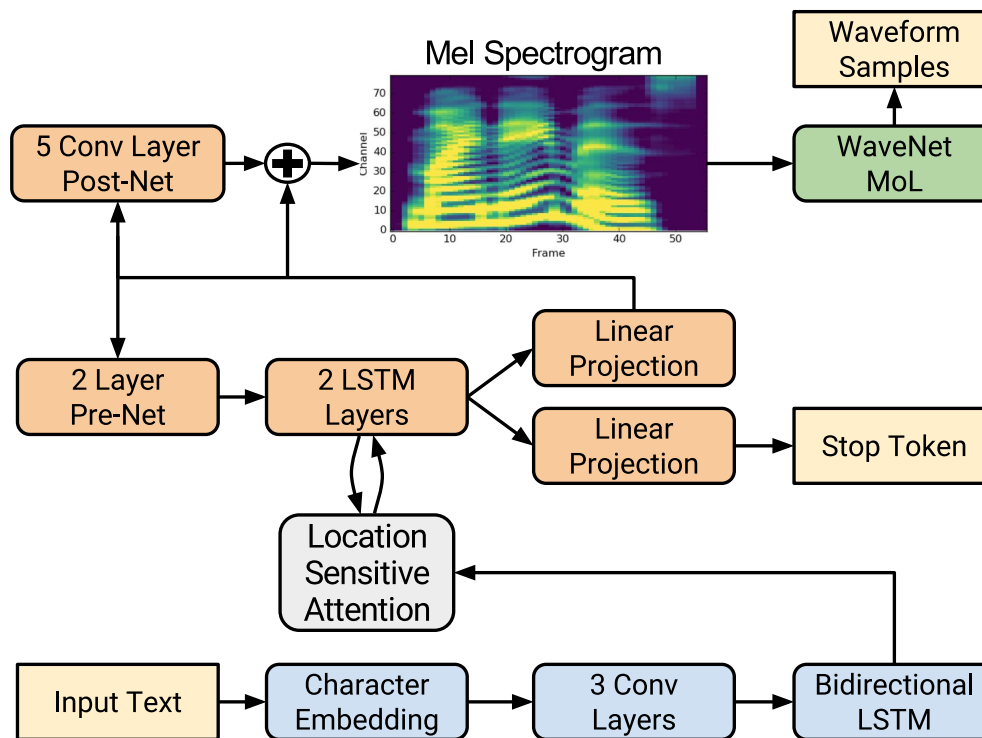


Figure 1 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>



The Mel spectrograms are computed using STFT (short-time Fourier transform).

- The authors propose a frame size of 50ms, 12.5ms frame hop, and a Hann window.
- STFT magnitudes are transformed into 80-channel Mel scale spanning 175Hz to 7.6kHz, followed by a log dynamic range compression (clipping input values to at least 0.01).

To make sequential processing of input characters easier, Tacotron 2 utilizes *location-sensitive attention*, which is an extension of the additive attention. While the additive (Bahdanau) attention computes

$$\alpha_i = \text{Attend}(\mathbf{s}_{i-1}, \mathbf{h}), \quad \alpha_{ij} = \text{softmax}(\mathbf{v}^\top \tanh(\mathbf{V}\mathbf{h}_j + \mathbf{W}\mathbf{s}_{i-1} + \mathbf{b})),$$

the location-sensitive attention also inputs the previous time step attention weights into the current attention computation:

$$\alpha_i = \text{Attend}(\mathbf{s}_{i-1}, \mathbf{h}, \alpha_{i-1}).$$

In detail, the previous attention weights are processed by a 1-D convolution with kernel  $\mathbf{F}$ :

$$\alpha_{ij} = \text{softmax}(\mathbf{v}^\top \tanh(\mathbf{V}\mathbf{h}_j + \mathbf{W}\mathbf{s}_{i-1} + (\mathbf{F} * \alpha_{i-1})_j + \mathbf{b})).$$

System	MOS
Parametric	3.492 $\pm$ 0.096
Tacotron (Griffin-Lim)	4.001 $\pm$ 0.087
Concatenative	4.166 $\pm$ 0.091
WaveNet (Linguistic)	4.341 $\pm$ 0.051
Ground truth	4.582 $\pm$ 0.053
<b>Tacotron 2 (this paper)</b>	<b>4.526 <math>\pm</math> 0.066</b>

Table 1 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

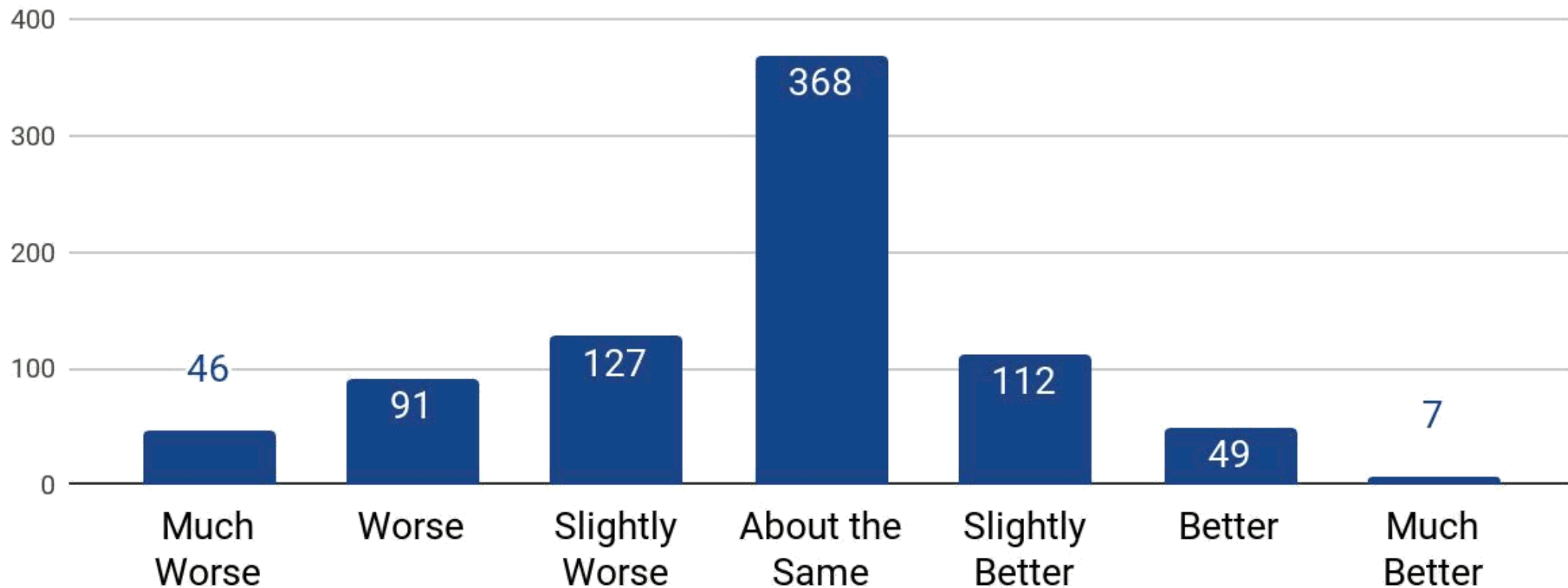


Figure 2 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

You can listen to samples at <https://google.github.io/tacotron/publications/tacotron2/>

Training	Synthesis	
	Predicted	Ground truth
Predicted	$4.526 \pm 0.066$	$4.449 \pm 0.060$
Ground truth	$4.362 \pm 0.066$	$4.522 \pm 0.055$

**Table 2.** Comparison of evaluated MOS for our system when WaveNet trained on predicted/ground truth mel spectrograms are made to synthesize from predicted/ground truth mel spectrograms.

*Table 2 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>*

System	MOS
Tacotron 2 (Linear + G-L)	$3.944 \pm 0.091$
Tacotron 2 (Linear + WaveNet)	$4.510 \pm 0.054$
Tacotron 2 (Mel + WaveNet)	<b><math>4.526 \pm 0.066</math></b>

**Table 3.** Comparison of evaluated MOS for Griffin-Lim vs. WaveNet as a vocoder, and using 1,025-dimensional linear spectrograms vs. 80-dimensional mel spectrograms as conditioning inputs to WaveNet.

*Table 3 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>*

Total layers	Num cycles	Dilation cycle size	Receptive field (samples / ms)	MOS
30	3	10	6,139 / 255.8	$4.526 \pm 0.066$
24	4	6	505 / 21.0	$4.547 \pm 0.056$
12	2	6	253 / 10.5	$4.481 \pm 0.059$
30	30	1	61 / 2.5	$3.930 \pm 0.076$

**Table 4.** WaveNet with various layer and receptive field sizes.

*Table 4 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>*

So far, all input information was stored either directly in network weights, or in a state of a recurrent network.

However, mammal brains seem to operate with a **working memory** – a capacity for short-term storage of information and its rule-based manipulation.

We can therefore try to introduce an external memory to a neural network. The memory  $M$  will be a matrix, where rows correspond to memory cells.

# Neural Turing Machines

The network will control the memory using a controller which reads from the memory and writes to it. Although the original paper also considered a feed-forward (non-recurrent) controller, usually the controller is a recurrent LSTM network.

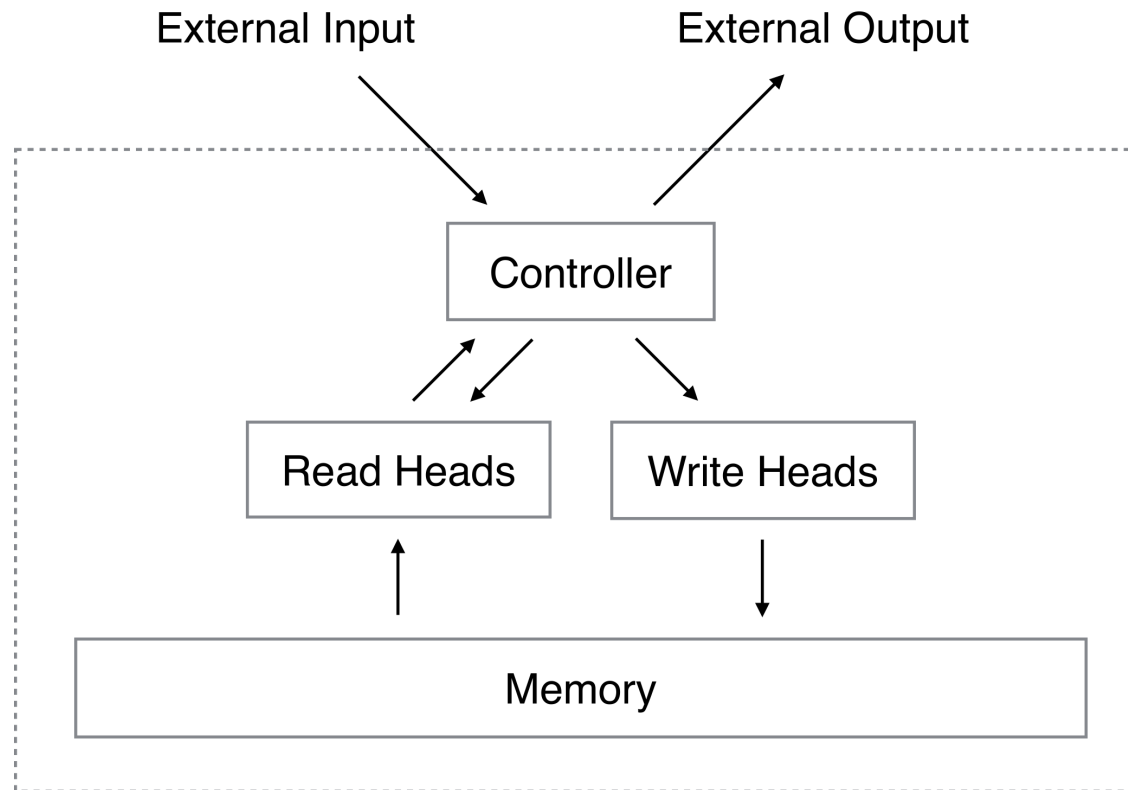


Figure 1 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>

## Reading

To read the memory in a differentiable way, the controller at time  $t$  emits a read distribution  $w_t$  over memory locations, and the returned read vector  $r_t$  is then

$$r_t = \sum_i w_t(i) \cdot M_t(i).$$

## Writing

Writing is performed in two steps – an **erase** followed by an **add**: the controller at time  $t$  emits a write distribution  $w_t$  over memory locations, together with an *erase vector*  $e_t$  and an *add vector*  $a_t$ . The memory is then updated as

$$M_t(i) = M_{t-1}(i) [1 - w_t(i)e_t] + w_t(i)a_t.$$

# Neural Turing Machine

The addressing mechanism is designed to allow both

- content addressing, and
- location addressing.

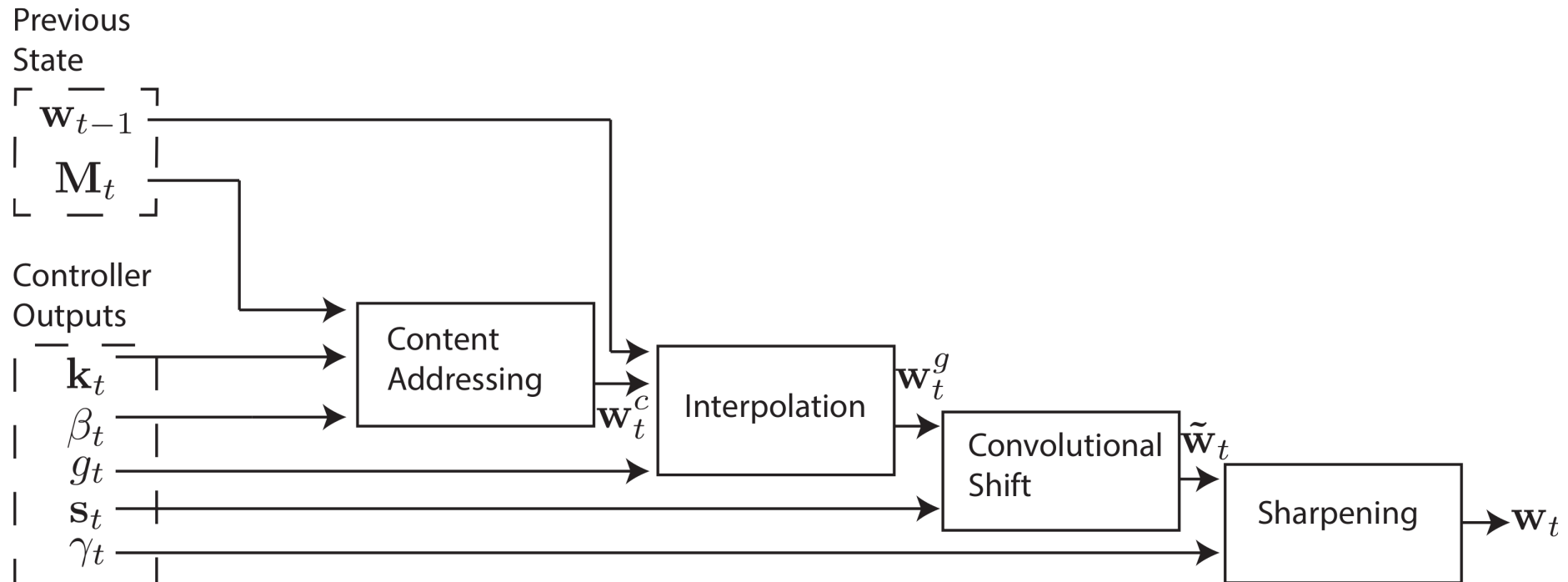


Figure 2 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>



## Content Addressing

Content addressing starts by the controller emitting the *key vector*  $\mathbf{k}_t$ , which is compared to all memory locations  $\mathbf{M}_t(i)$ , generating a distribution using a softmax with temperature  $\beta_t$ .

$$w_t^c(i) = \frac{\exp(\beta_t \cdot \text{distance}(\mathbf{k}_t, \mathbf{M}_t(i)))}{\sum_j \exp(\beta_t \cdot \text{distance}(\mathbf{k}_t, \mathbf{M}_t(j)))}$$

The distance measure is usually the cosine similarity

$$\text{distance}(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a}^T \mathbf{b}}{\|\mathbf{a}\| \cdot \|\mathbf{b}\|}.$$

## Location-Based Addressing

To allow iterative access to memory, the controller might decide to reuse the memory location from the previous timestep. Specifically, it emits an *interpolation gate*  $g_t \in (0, 1)$  and sets

$$\mathbf{w}_t^g = g_t \mathbf{w}_t^c + (1 - g_t) \mathbf{w}_{t-1}.$$

Then, the current weighting may be shifted, i.e., the controller might decide to “rotate” the weights by a small integer. For a given range (the simplest case are only shifts  $\{-1, 0, 1\}$ ), the network emits a softmax distribution over the shifts, and the weights are then defined using a circular convolution

$$\tilde{w}_t(i) = \sum_j w_t^g(j) s_t(i - j).$$

Finally, not to lose precision over time, the controller emits a *sharpening factor*  $\gamma_t \geq 1$ , and the final memory location weights are  $w_t(i) = \tilde{w}_t(i)^{\gamma_t} / \sum_j \tilde{w}_t(j)^{\gamma_t}$ .

## Overall Execution

Even if not specified in the original paper, following the DNC paper, the LSTM controller can be implemented as a (potentially deep) LSTM. Assuming  $R$  read heads and one write head, the input is  $\mathbf{x}_t$  and  $R$  read vectors  $\mathbf{r}_{t-1}^1, \dots, \mathbf{r}_{t-1}^R$  from the previous time step, the output of the controller are vectors  $(\boldsymbol{\nu}_t, \boldsymbol{\xi}_t)$ , and the final output is  $\mathbf{y}_t = \boldsymbol{\nu}_t + \mathbf{W}_r [\mathbf{r}_t^1, \dots, \mathbf{r}_t^R]$ . The  $\boldsymbol{\xi}_t$  is a concatenation of

$$\mathbf{k}_t^1, \beta_t^1, \mathbf{g}_t^1, \mathbf{s}_t^1, \gamma_t^1, \mathbf{k}_t^2, \beta_t^2, \mathbf{g}_t^2, \mathbf{s}_t^2, \gamma_t^2, \dots, \mathbf{k}_t^w, \beta_t^w, \mathbf{g}_t^w, \mathbf{s}_t^w, \gamma_t^w, \mathbf{e}_t^w, \mathbf{a}_t^w.$$

## Copy Task

Repeat the same sequence as given on input. Trained with sequences of length up to 20.

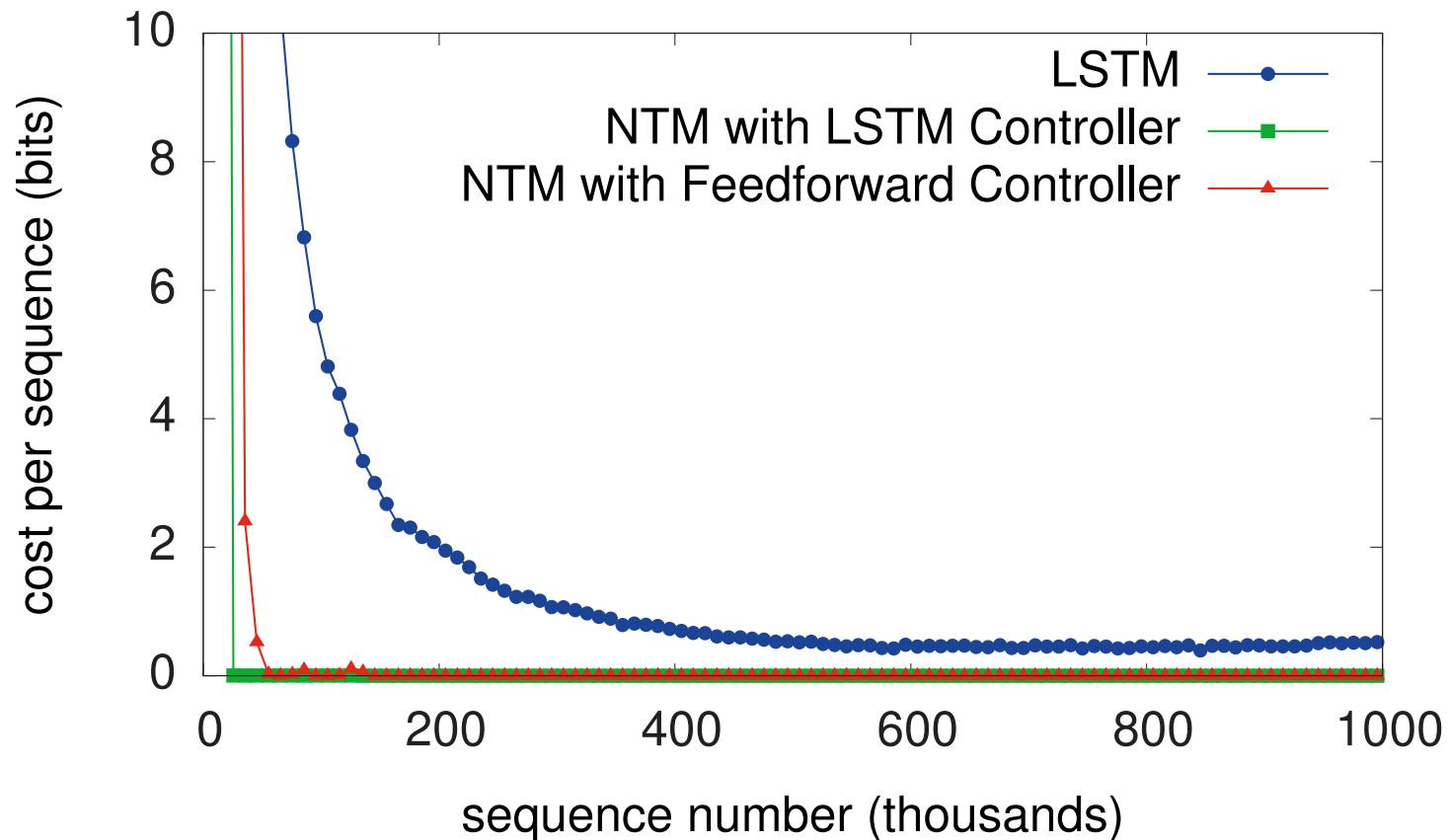
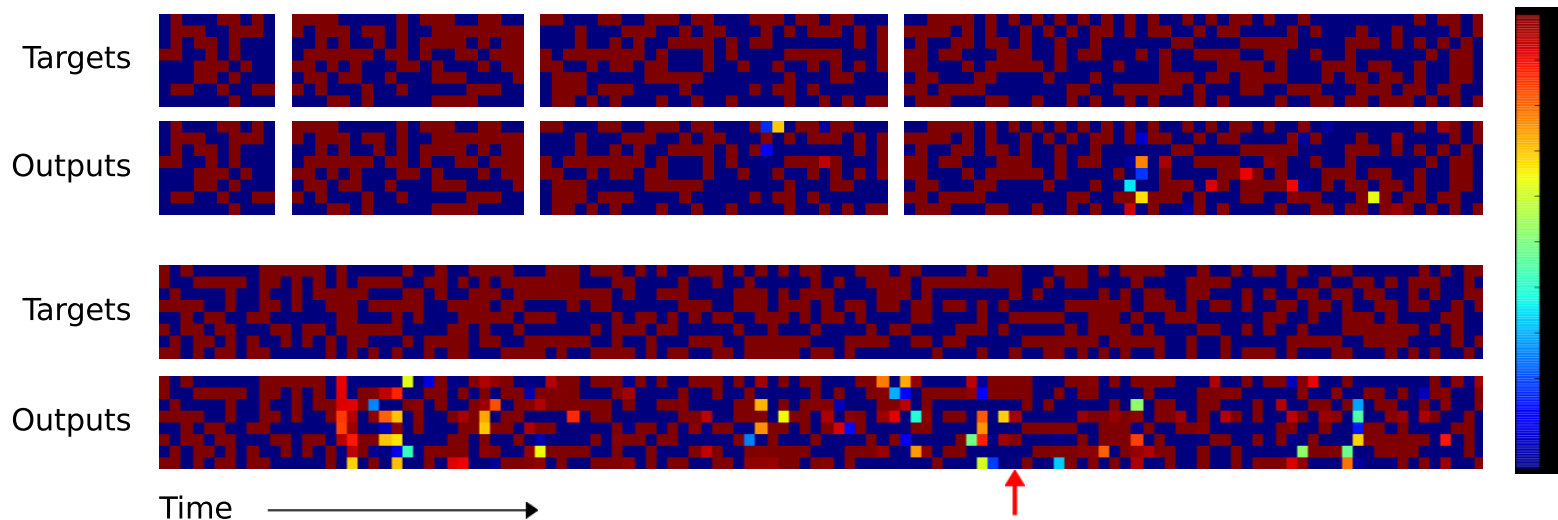
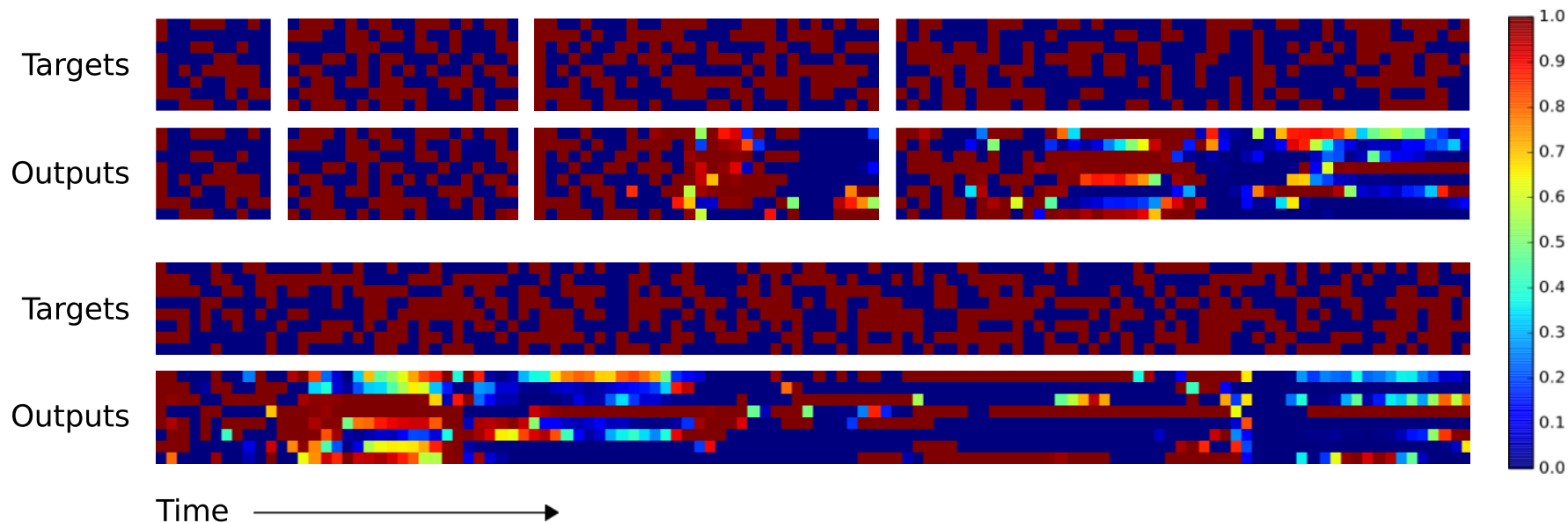


Figure 3 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>



**Figure 4: NTM Generalisation on the Copy Task.** The four pairs of plots in the top row depict network outputs and corresponding copy targets for test sequences of length 10, 20, 30, and 50, respectively. The plots in the bottom row are for a length 120 sequence. The network was only trained on sequences of up to length 20. The first four sequences are reproduced with high confidence and very few mistakes. The longest one has a few more local errors and one global error: at the point indicated by the red arrow at the bottom, a single vector is duplicated, pushing all subsequent vectors one step back. Despite being subjectively close to a correct copy, this leads to a high loss.

Figure 4 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>



**Figure 5: LSTM Generalisation on the Copy Task.** The plots show inputs and outputs for the same sequence lengths as Figure 4. Like NTM, LSTM learns to reproduce sequences of up to length 20 almost perfectly. However it clearly fails to generalise to longer sequences. Also note that the length of the accurate prefix decreases as the sequence length increases, suggesting that the network has trouble retaining information for long periods.

*Figure 5 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>*

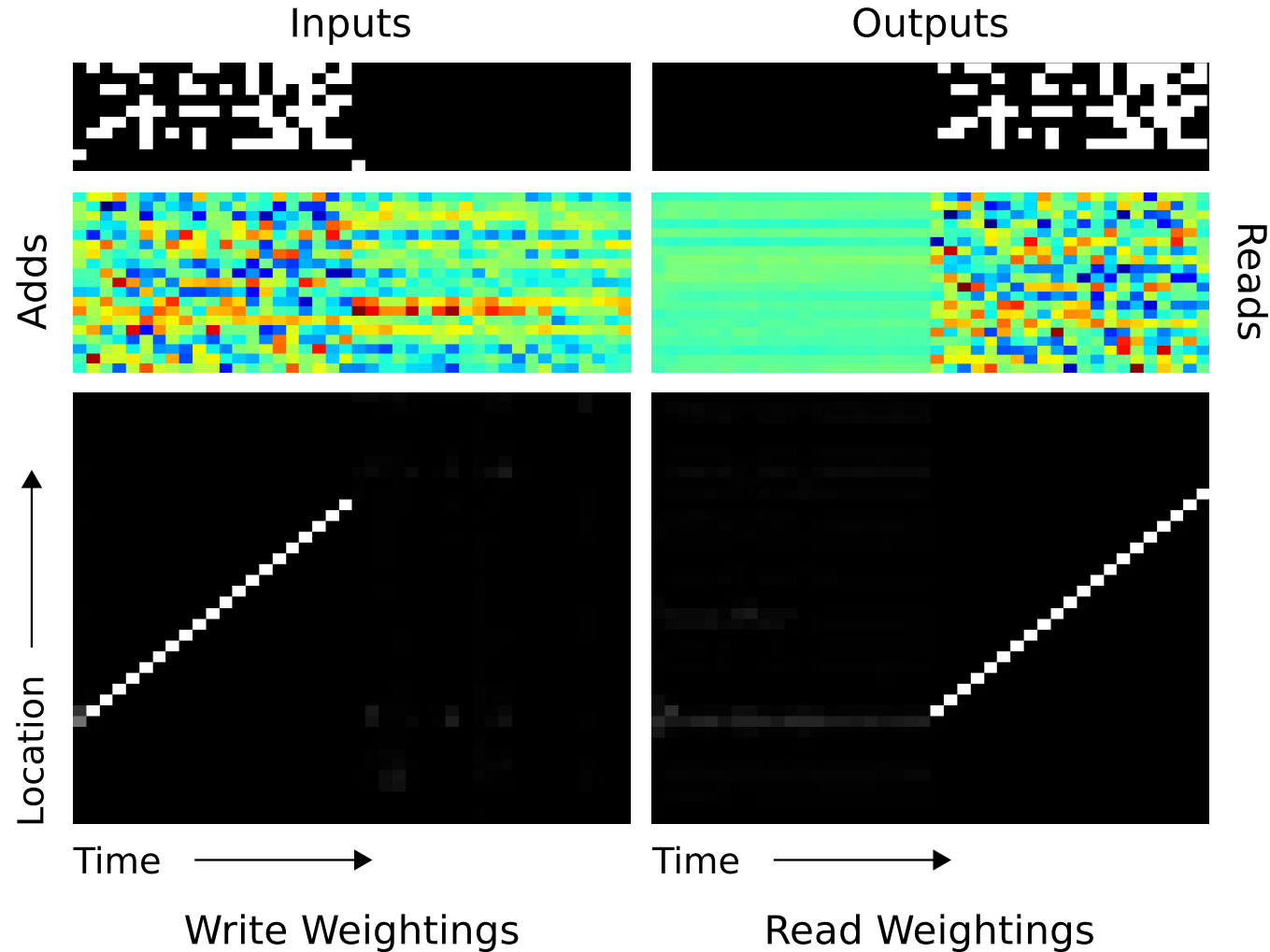


Figure 6 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>

## Associative Recall

In associative recall, a sequence is given on input, consisting of subsequences of length 3. Then a randomly chosen subsequence is presented on input and the goal is to produce the following subsequence.



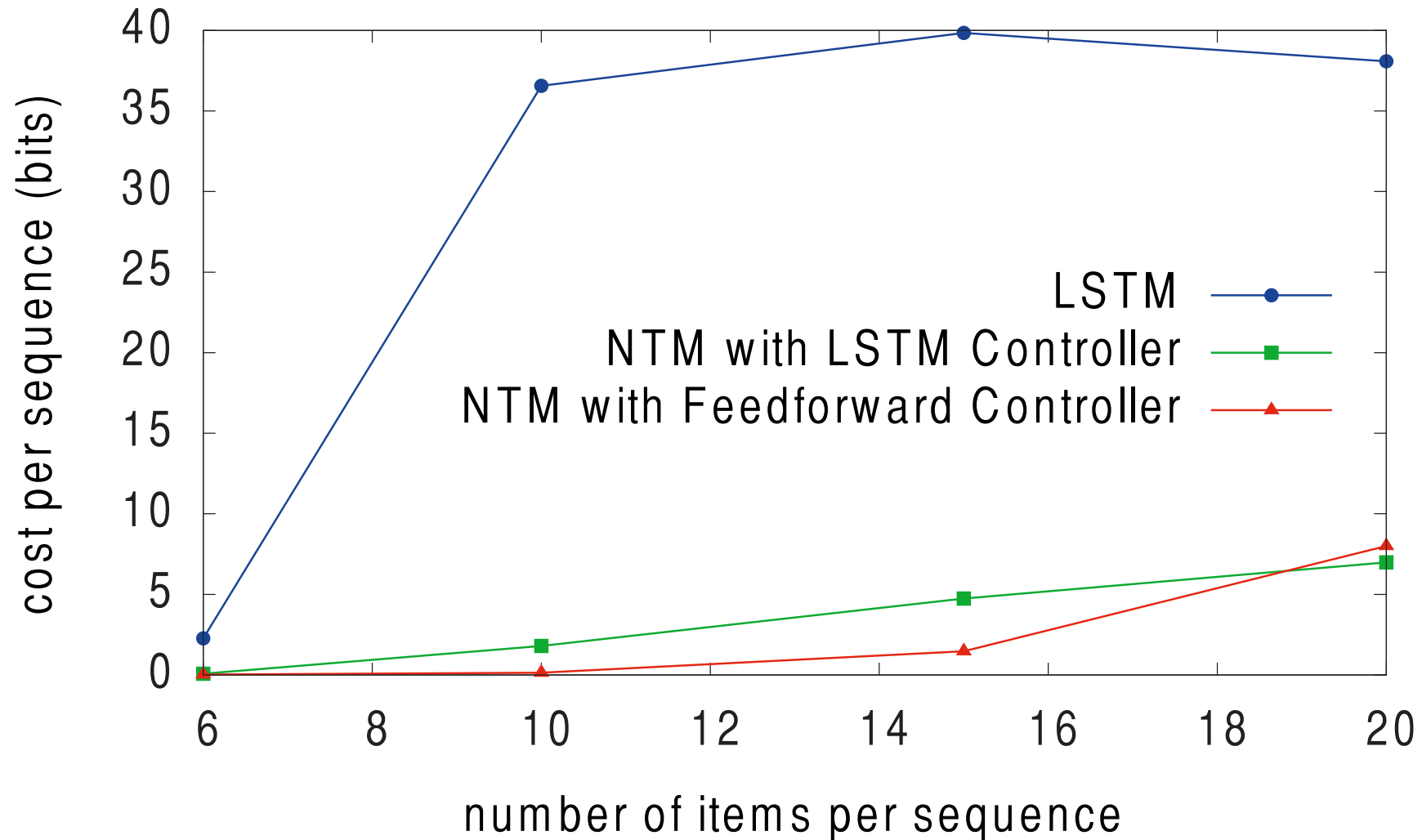


Figure 11 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>

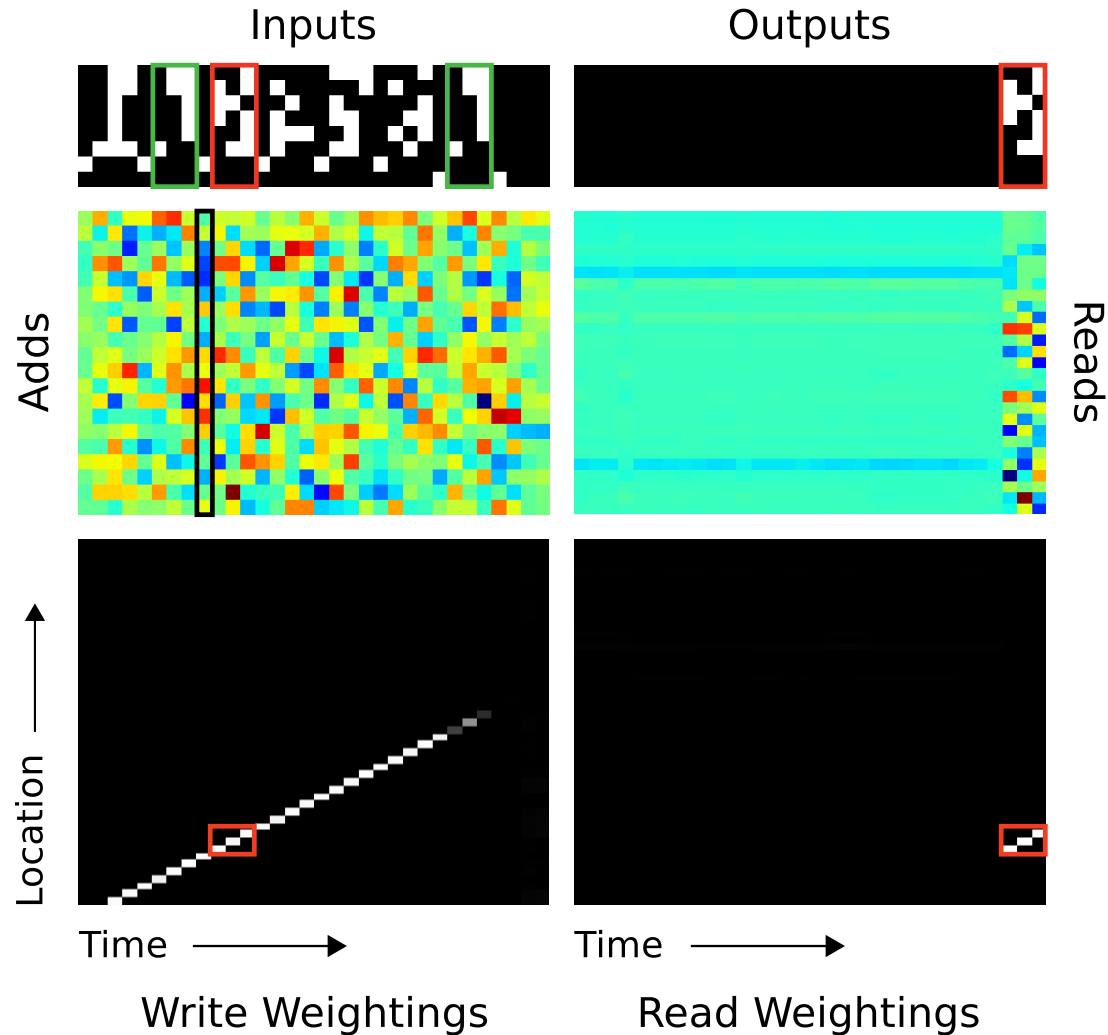


Figure 12 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>

# Memory-augmented Neural Networks

We now focus on a **learning to learn** task – consider a network, which should learn classification into a user-defined hierarchy by observing ideally a small number of samples.

Apart from finetuning the model and storing the information in the *weights*, an alternative is to store the samples in **external memory**. Therefore, the model learns how to store the data and access it efficiently, which allows it to learn without changing its weights.

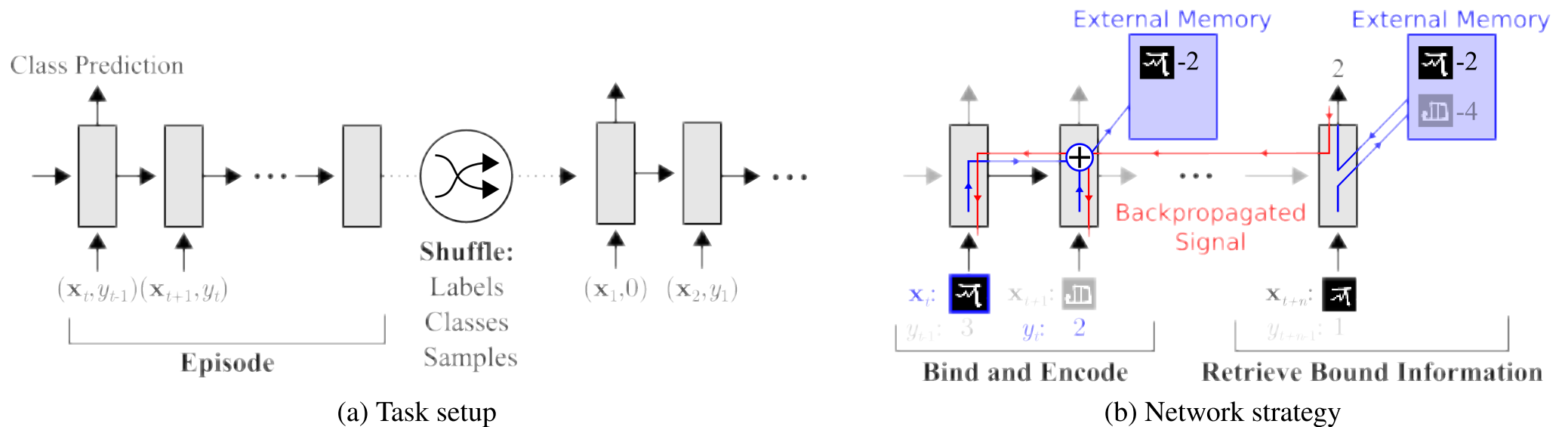


Figure 1 of "One-shot learning with Memory-Augmented Neural Networks", <https://arxiv.org/abs/1605.06065>

$$K(\mathbf{k}_t, \mathbf{M}_t(i)) = \frac{\mathbf{k}_t \cdot \mathbf{M}_t(i)}{\|\mathbf{k}_t\| \|\mathbf{M}_t(i)\|}, \quad (2)$$

which is used to produce a read-weight vector,  $\mathbf{w}_t^r$ , with elements computed according to a softmax:

$$w_t^r(i) \leftarrow \frac{\exp(K(\mathbf{k}_t, \mathbf{M}_t(i)))}{\sum_j \exp(K(\mathbf{k}_t, \mathbf{M}_t(j)))}. \quad (3)$$

A memory,  $\mathbf{r}_t$ , is retrieved using this weight vector:

$$\mathbf{r}_t \leftarrow \sum_i w_t^r(i) \mathbf{M}_t(i). \quad (4)$$

Page 3 of "One-shot learning with Memory-Augmented Neural Networks", <https://arxiv.org/abs/1605.06065>

$$\mathbf{w}_t^u \leftarrow \gamma \mathbf{w}_{t-1}^u + \mathbf{w}_t^r + \mathbf{w}_t^w. \quad (5)$$

Here,  $\gamma$  is a decay parameter and  $\mathbf{w}_t^r$  is computed as in (3). The *least-used* weights,  $\mathbf{w}_t^{lu}$ , for a given time-step can then be computed using  $\mathbf{w}_t^u$ . First, we introduce the notation  $m(\mathbf{v}, n)$  to denote the  $n^{th}$  smallest element of the vector  $\mathbf{v}$ . Elements of  $\mathbf{w}_t^{lu}$  are set accordingly:

$$w_t^{lu}(i) = \begin{cases} 0 & \text{if } w_t^u(i) > m(\mathbf{w}_t^u, n) \\ 1 & \text{if } w_t^u(i) \leq m(\mathbf{w}_t^u, n) \end{cases}, \quad (6)$$

where  $n$  is set to equal the number of reads to memory. To obtain the write weights  $\mathbf{w}_t^w$ , a learnable sigmoid gate parameter is used to compute a convex combination of the previous read weights and previous least-used weights:

$$\mathbf{w}_t^w \leftarrow \sigma(\alpha) \mathbf{w}_{t-1}^r + (1 - \sigma(\alpha)) \mathbf{w}_{t-1}^{lu}. \quad (7)$$

Here,  $\sigma(\cdot)$  is a sigmoid function,  $\frac{1}{1+e^{-x}}$ , and  $\alpha$  is a scalar gate parameter to interpolate between the weights. Prior to writing to memory, the least used memory location is computed from  $\mathbf{w}_{t-1}^u$  and is set to zero. Writing to memory then occurs in accordance with the computed vector of write weights:

$$\mathbf{M}_t(i) \leftarrow \mathbf{M}_{t-1}(i) + w_t^w(i) \mathbf{k}_t, \forall i \quad (8)$$

Page 4 of "One-shot learning with Memory-Augmented Neural Networks", <https://arxiv.org/abs/1605.06065>

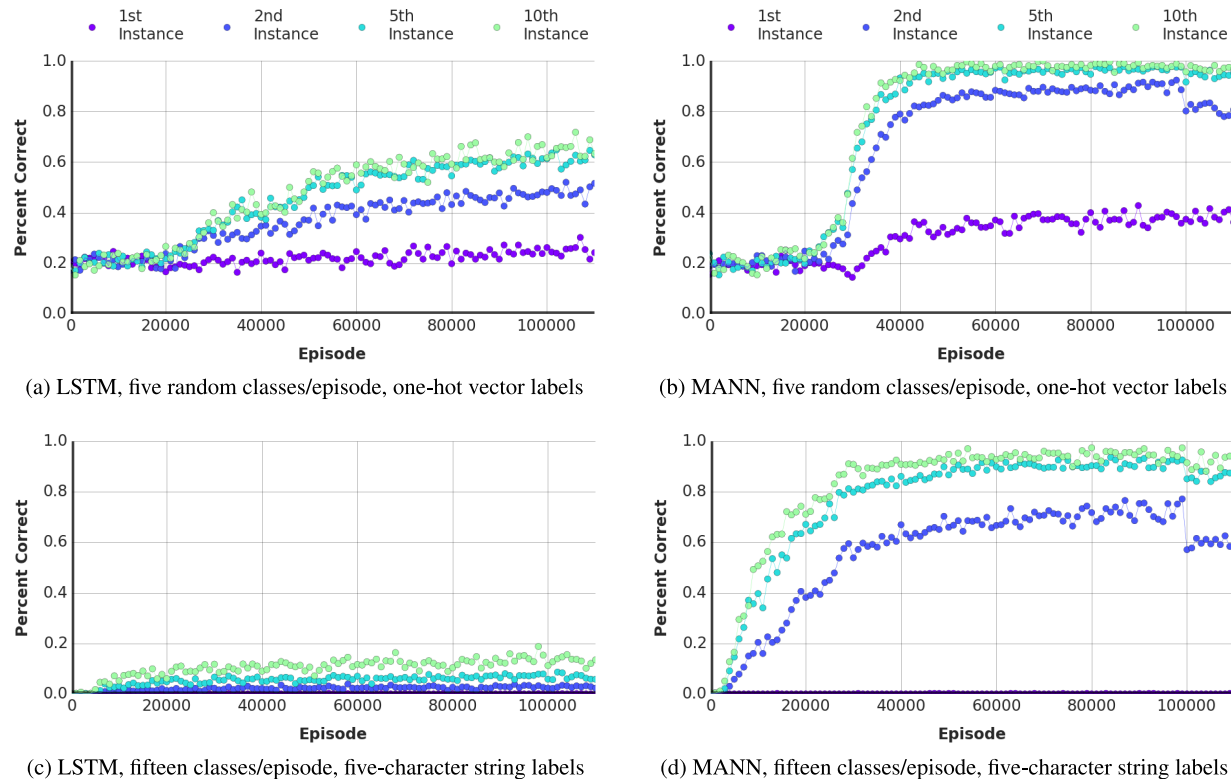


Figure 2. Omniglot classification. The network was given either five (a-b) or up to fifteen (c-d) random classes per episode, which were of length 50 or 100 respectively. Labels were one-hot vectors in (a-b), and five-character strings in (c-d). In (b), first instance accuracy is above chance, indicating that the MANN is performing “educated guesses” for new classes based on the classes it has already seen and stored in memory. In (c-d), first instance accuracy is poor, as is expected, since it must make a guess from 3125 random strings. Second instance accuracy, however, approaches 80% during training for the MANN (d). At the 100,000 episode mark the network was tested, without further learning, on distinct classes withheld from the training set, and exhibited comparable performance.

Figure 2 of "One-shot learning with Memory-Augmented Neural Networks", <https://arxiv.org/abs/1605.06065>

Table 1. Test-set classification accuracies for humans compared to machine algorithms trained on the Omniglot dataset, using one-hot encodings of labels and five classes presented per episode.

MODEL	INSTANCE (% CORRECT)					
	1 <sup>ST</sup>	2 <sup>ND</sup>	3 <sup>RD</sup>	4 <sup>TH</sup>	5 <sup>TH</sup>	10 <sup>TH</sup>
HUMAN	34.5	57.3	70.1	71.8	81.4	92.4
FEEDFORWARD	24.4	19.6	21.1	19.9	22.8	19.5
LSTM	24.4	49.5	55.3	61.0	63.6	62.5
MANN	<b>36.4</b>	<b>82.8</b>	<b>91.0</b>	<b>92.6</b>	<b>94.9</b>	<b>98.1</b>

Table 1 of "One-shot learning with Memory-Augmented Neural Networks", <https://arxiv.org/abs/1605.06065>

Table 2. Test-set classification accuracies for various architectures on the Omniglot dataset after 100000 episodes of training, using five-character-long strings as labels. See the supplemental information for an explanation of 1st instance accuracies for the kNN classifier.

MODEL	CONTROLLER	# OF CLASSES	INSTANCE (% CORRECT)					
			1 <sup>ST</sup>	2 <sup>ND</sup>	3 <sup>RD</sup>	4 <sup>TH</sup>	5 <sup>TH</sup>	10 <sup>TH</sup>
KNN (RAW PIXELS)	–	15	0.5	18.7	23.3	26.5	29.1	37.0
KNN (DEEP FEATURES)	–	15	0.4	32.7	41.2	47.1	50.6	60.0
FEEDFORWARD	–	15	0.0	0.1	0.0	0.0	0.0	0.0
LSTM	–	15	0.0	2.2	2.9	4.3	5.6	12.7
MANN (LRUA)	FEEDFORWARD	15	0.1	12.8	22.3	28.8	32.2	43.4
MANN (LRUA)	LSTM	15	0.1	<b>62.6</b>	<b>79.3</b>	<b>86.6</b>	<b>88.7</b>	<b>95.3</b>
MANN (NTM)	LSTM	15	0.0	35.4	61.2	71.7	77.7	88.4

Table 2 of "One-shot learning with Memory-Augmented Neural Networks", <https://arxiv.org/abs/1605.06065>

NTM was later extended to a Differentiable Neural Computer.

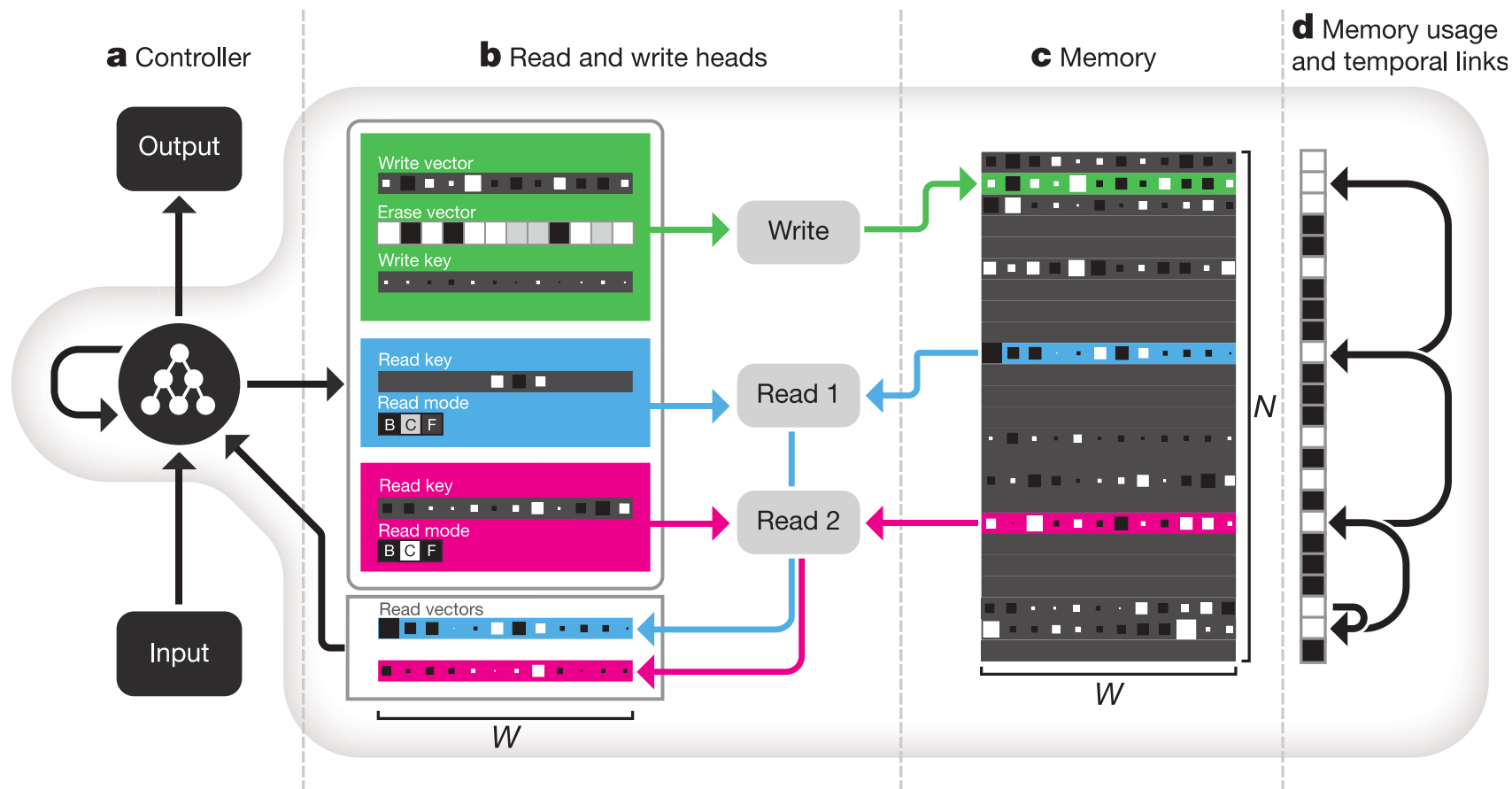


Figure 1 of "Hybrid computing using a neural network with dynamic external memory", <https://www.nature.com/articles/nature20101>

The DNC contains multiple read heads and one write head.

The controller is a deep LSTM network, with input at time  $t$  being the current input  $\mathbf{x}_t$  and  $R$  read vectors  $\mathbf{r}_{t-1}^1, \dots, \mathbf{r}_{t-1}^R$  from previous time step. The output of the controller are vectors  $(\boldsymbol{\nu}_t, \boldsymbol{\xi}_t)$ , and the final output is  $\mathbf{y}_t = \boldsymbol{\nu}_t + W_r [\mathbf{r}_t^1, \dots, \mathbf{r}_t^R]$ . The  $\boldsymbol{\xi}_t$  is a concatenation of parameters for read and write heads (keys, gates, sharpening parameters, ...).

In DNC, the usage of every memory location is tracked, which enables performing dynamic allocation – at each time step, a cell with least usage can be allocated.

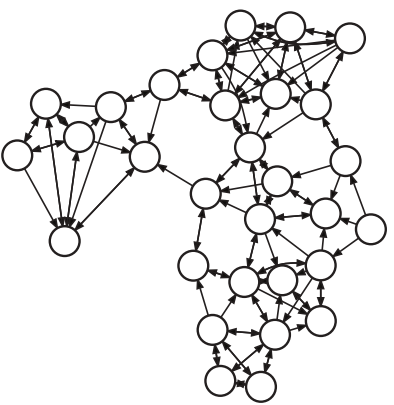
Furthermore, for every memory location, we track which memory location was written to previously and subsequently, allowing to recover sequences in the order in which it was written, independently on the real indices used.

The write weighting is defined as a weighted combination of the allocation weighting and write content weighting, and read weighting is computed as a weighted combination of read content weighting, previous write weighting, and subsequent write weighting.

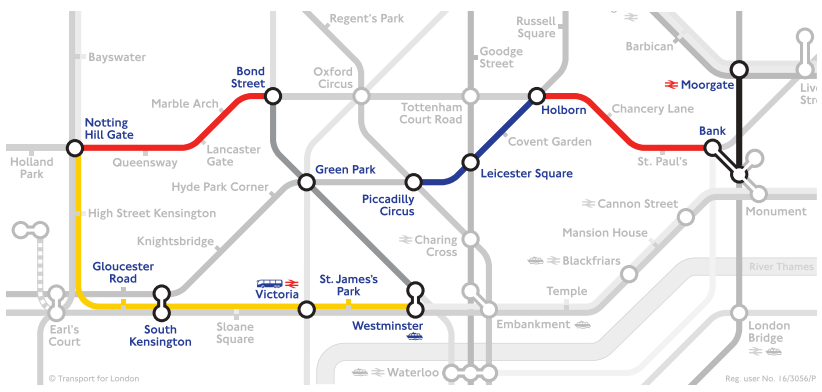


# Differentiable Neural Computer

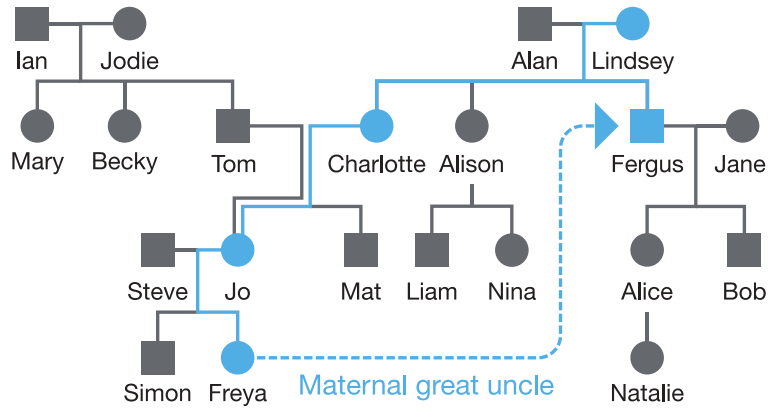
**a** Random graph



**b** London Underground



**c** Family tree



<p><b>Underground input:</b>  (OxfordCircus, TottenhamCtRd, Central)  (TottenhamCtRd, OxfordCircus, Central)  (BakerSt, Marylebone, Circle)  (BakerSt, Marylebone, Bakerloo)  (BakerSt, OxfordCircus, Bakerloo)  ⋮  (LeicesterSq, CharingCross, Northern)  (TottenhamCtRd, LeicesterSq, Northern)  (OxfordCircus, PiccadillyCircus, Bakerloo)  (OxfordCircus, NottingHillGate, Central)  (OxfordCircus, Euston, Victoria)</p> <p>84 edges in total</p>	<p><b>Traversal question:</b>  (BondSt, _, Central),  (., _, Circle), (., _, Circle),  (., _, Circle), (., _, Circle),  (., _, Jubilee), (., _, Jubilee),</p> <p><b>Answer:</b>  (BondSt, NottingHillGate, Central)  (NottingHillGate, GloucesterRd, Circle)  ⋮  (Westminster, GreenPark, Jubilee)  (GreenPark, BondSt, Jubilee)</p>	<p><b>Shortest-path question:</b>  (Moorgate, PiccadillyCircus, _)</p> <p><b>Answer:</b>  (Moorgate, Bank, Northern)  (Bank, Holborn, Central)  (Holborn, LeicesterSq, Piccadilly)  (LeicesterSq, PiccadillyCircus, Piccadilly)</p>	<p><b>Family tree input:</b>  (Charlotte, Alan, Father)  (Simon, Steve, Father)  (Steve, Simon, Son1)  (Nina, Alison, Mother)  (Lindsey, Fergus, Son1)  ⋮  (Bob, Jane, Mother)  (Natalie, Alice, Mother)  (Mary, Ian, Father)  (Jane, Alice, Daughter1)  (Mat, Charlotte, Mother)</p> <p>54 edges in total</p>	<p><b>Inference question:</b>  (Freya, _, MaternalGreatUncle)</p> <p><b>Answer:</b>  (Freya, Fergus, MaternalGreatUncle)</p>
--	--	---	--	--

Figure 2 of "Hybrid computing using a neural network with dynamic external memory", <https://www.nature.com/articles/nature20101>

# Differentiable Neural Computer

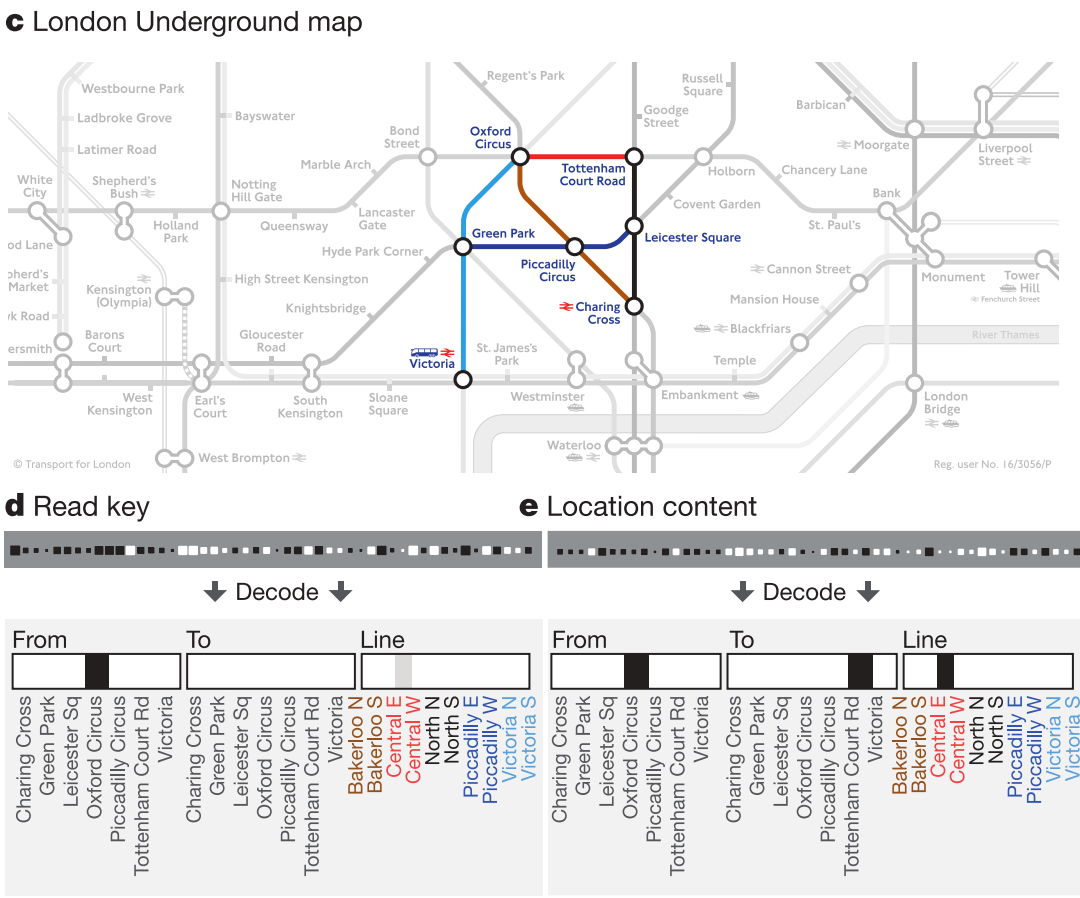
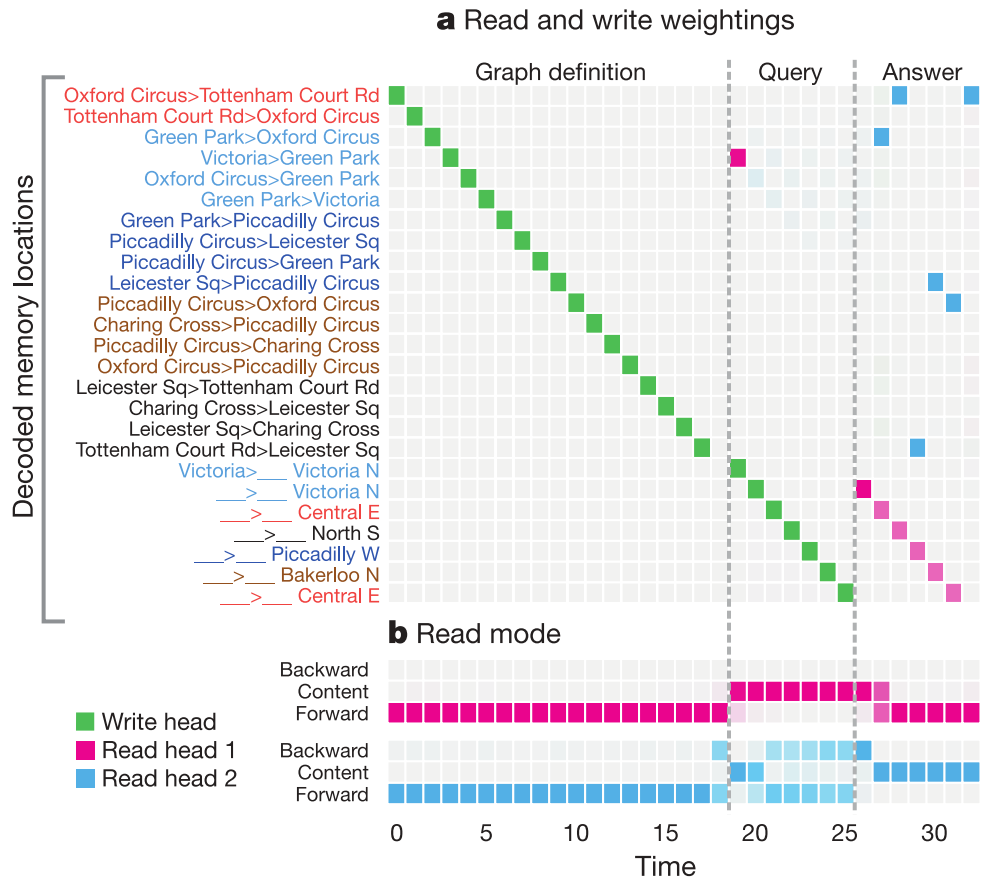


Figure 3 of "Hybrid computing using a neural network with dynamic external memory", <https://www.nature.com/articles/nature20101>

Token Turing Machines are a sequential, autoregressive Transformer model with memory for sequential visual understanding.

In every time step  $t$ , we get

- a fixed-size memory  $\mathbf{M}^t \in \mathbb{R}^{m \times d}$ ,
- a fixed-size input  $\mathbf{I}^t \in \mathbb{R}^{n \times d}$ .

In every step, we perform:

- $\mathbf{Z}^t \leftarrow \text{Read}(\mathbf{M}^t, \mathbf{I}^t)$ ,
- $\mathbf{O}^t \leftarrow \text{Process}(\mathbf{Z}^t)$ ,
- $\mathbf{M}^{t+1} \leftarrow \text{Write}(\mathbf{M}^t, \mathbf{I}^t, \mathbf{O}^t)$ ,
- $\mathbf{Y}^t \leftarrow \text{Output}(\mathbf{O}^t)$ .

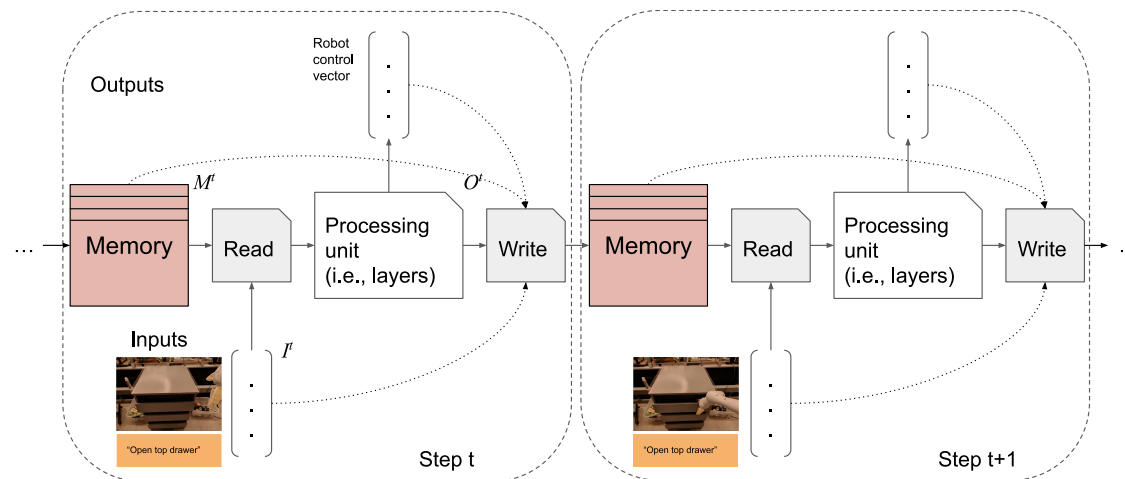


Figure 1. TTM overview with robot examples. Each dotted rectangle shows TTM at each step.

Figure 1 of "Token Turing Machines", <https://arxiv.org/abs/2211.09119>

Token **summarization**  $S_k$  summarizes a sequence of  $\mathbf{V} \in \mathbb{R}^{p \times d}$  tokens into  $k \ll p$  tokens  $\mathbf{Z} \in \mathbb{R}^{k \times d}$ .

First, weight vector  $\mathbf{w}_i$  is computed as

$$\mathbf{w}_i \leftarrow \text{softmax}(\mathbf{q}_i^T \mathbf{V}^T / \sqrt{d}),$$

and then

$$\mathbf{z}_i \leftarrow \mathbf{w}_i^T \mathbf{V}.$$

Finally, read operation is implemented as a summarization to  $r$  tokens

$$\mathbf{Z}^t \leftarrow \text{Read}(\mathbf{I}^t, \mathbf{M}^t) \stackrel{\text{def}}{=} S_k(\mathbf{M}^t \parallel \mathbf{I}^t).$$

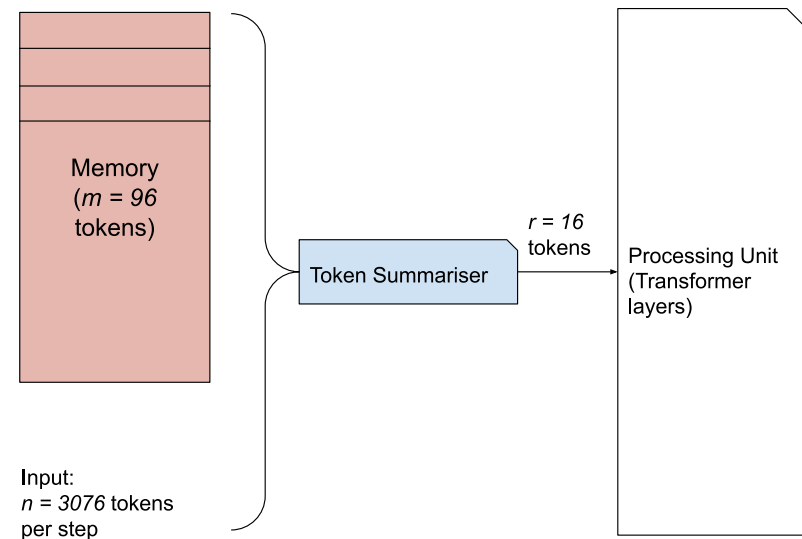


Figure 2. TTM Read. Note how it greatly reduces the computation of the subsequent processing module by summarising the input sequence as well.

Figure 2 of "Token Turing Machines", <https://arxiv.org/abs/2211.09119>

The processing unit computes  $r$  tokens  $\mathbf{O}^t$  from the  $r$  read tokens  $\mathbf{Z}^t$ , by using for example the Transformer architecture:

$$\mathbf{O}^t \leftarrow \text{Process}(\mathbf{Z}^t) \stackrel{\text{def}}{=} \text{Transformer}(\mathbf{Z}^t).$$

The output is computed using a linear output head on the  $\mathbf{O}^t$  tokens:

$$\mathbf{Y}^t \leftarrow \text{Output}(\mathbf{O}^t) \stackrel{\text{def}}{=} \text{flatten}(\mathbf{O}^t) \mathbf{W}_o.$$

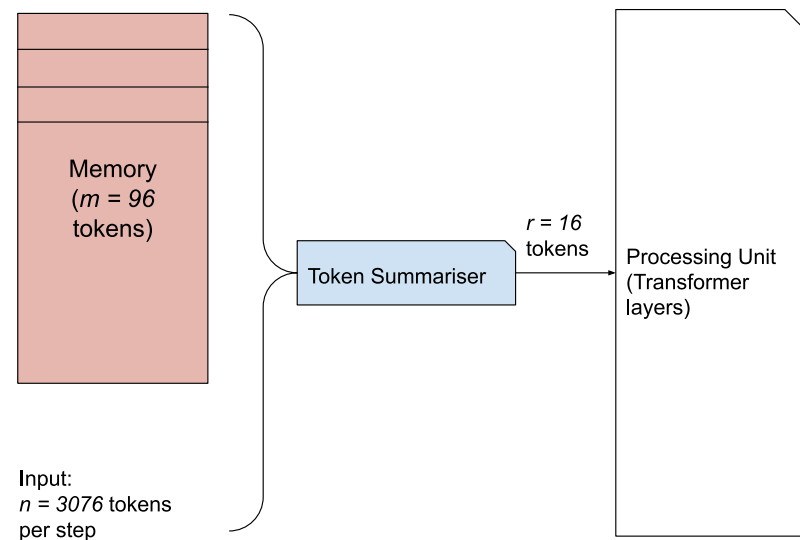


Figure 2. TTM Read. Note how it greatly reduces the computation of the subsequent processing module by summarising the input sequence as well.

Figure 2 of "Token Turing Machines", <https://arxiv.org/abs/2211.09119>

# Write Operation

The write operation is another summarization operation summarizing  $n$  input cells,  $m$  input memory cells, and  $r$  output cells into a new memory of  $n$  cells:

$$\mathbf{M}^{t+1} \leftarrow \text{Write}(\mathbf{M}^t, \mathbf{I}^t, \mathbf{O}^t) \stackrel{\text{def}}{=} \stackrel{\text{def}}{=} S_m(\mathbf{M}^t \parallel \mathbf{I}^t \parallel \mathbf{O}^t).$$

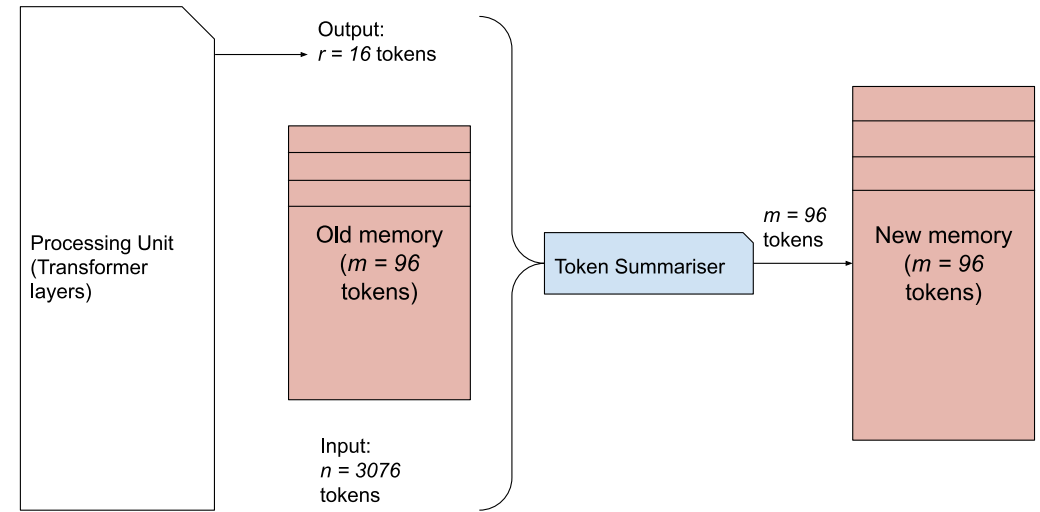


Figure 3. TTM Write, formulated as the token summarisation.

Figure 3 of "Token Turing Machines", <https://arxiv.org/abs/2211.09119>

Method	Setting	modality	mAP
I3D + super-events [52]	offline	RGB + Flow	19.41
I3D + super-events + TGM [53]	offline	RGB + Flow	22.30
I3D + STGCN [28]	offline	RGB + Flow	19.09
I3D + biGRU + VS-ST-MPNN [48]	offline	RGB + Object	23.7
Coarse-Fine (w/ X3D) [40]	offline	RGB	25.1
I3D + CTRN [16]	offline	RGB	25.3
I3D + MS-TCT [17]	offline	RGB	25.4
I3D + PDAN [18]	offline	RGB + Flow	26.5
I3D + CTRN [16]	offline	RGB + Flow	27.8
I3D [10]	online	RGB + Flow	17.22
X3D [26]	online	RGB	18.87
ViViT-B [3]	online	RGB	23.18
ViViT-B + TTM (ours)	online	RGB	26.34
ViViT-L [3]	online	RGB	26.01
ViViT-L + TTM (ours)	online	RGB	28.79

Table 1. Comparison with the state-of-the-art methods on Charades temporal activity detection.

Table 1 of "Token Turing Machines", <https://arxiv.org/abs/2211.09119>



Method	mAP	GFLOPS
ViViT only	23.18	-
<i>Alternative temporal models</i>		
Temporal MLP Mixer (tokens=96)	24.41	0.382
Causal Transformer (tokens=96)	25.85	0.523
Temporal Transformer (tokens=96)	25.61	1.269
Temporal MLP Mixer (tokens=3360)	24.26	13.317
Causal Transformer (tokens=3360)	25.88	29.695
Temporal Transformer (tokens=3360)	25.53	112.836
<i>Alternative recurrent networks</i>		
LSTM	23.96	0.107
Recurrent Transformer (tokens=16+16)	25.97	0.410
Recurrent Transformer (tokens=3136+16)	25.97	17.10
<i>Token Turing Machines</i>		
TTM-Mixer ( $n = 16$ )	25.83	0.089
TTM-Transformer ( $n = 16$ )	26.24	0.228
TTM-Mixer ( $n = 3136$ )	26.14	0.704
TTM-Transformer ( $n = 3136$ )	26.34	0.842

Table 2. TTM vs. different sequence modeling methods. ViViT-B was used as the backbone. TTM-Transformer means we use Transformer as the processing unit, and TTM-Mixer means we use MLP Mixer as the processing unit. FLOP measure is for computation in addition to the backbone.

Architecture	mAP	GFLOPS
MLP	23.34	0.689
MLP Mixer	26.14	0.704
Transformer	26.34	0.842

Table 3. Using different processing unit architectures in TTMs.

Method	mAP	GFLOPS
Pooling	25.75	0.206
MLP	26.34	0.842
Latent query	26.75	8.537

Table 4. Different Token Summarisation

Method	mAP	GFLOPS
Concatenate (Memorizing Transformer-style)	20.97	0.920
Erase and Add (NTM-style write)	25.86	0.423
TTM without memory	22.65	0.842
TTM	26.34	0.842

Table 5. TTM vs. different history/memory update. They all use Transformer processing units, and MLP-based token summarisations. The number of input tokens per step,  $n = 3176$ .

Tables 2-5 of "Token Turing Machines", <https://arxiv.org/abs/2211.09119>