

# NASNet, Speech Synthesis, External Memory Networks

Milan Straka

 May 16, 2022



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

- We can design neural network architectures using reinforcement learning.
- The designed network is encoded as a sequence of elements, and is generated using an **RNN controller**, which is trained using the REINFORCE with baseline algorithm.

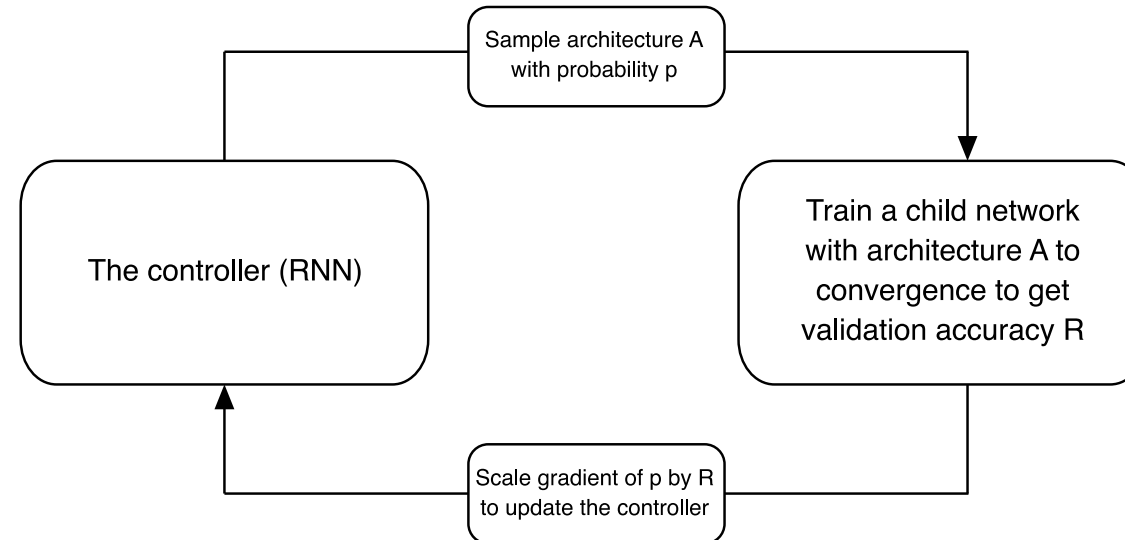


Figure 1 of "Learning Transferable Architectures for Scalable Image Recognition", <https://arxiv.org/abs/1707.07012>

- For every generated sequence, the corresponding network is trained on CIFAR-10 and the development accuracy is used as a return.

The overall architecture of the designed network is fixed and only the Normal Cells and Reduction Cells are generated by the controller.

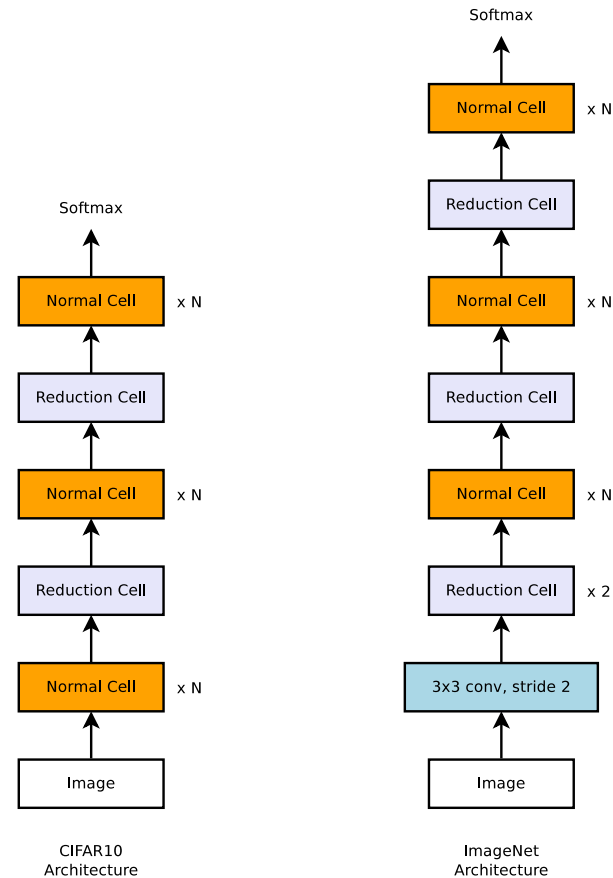


Figure 2 of "Learning Transferable Architectures for Scalable Image Recognition", <https://arxiv.org/abs/1707.07012>

- Each cell is composed of  $B$  blocks ( $B = 5$  is used in NASNet).
- Each block is designed by a RNN controller generating 5 parameters.

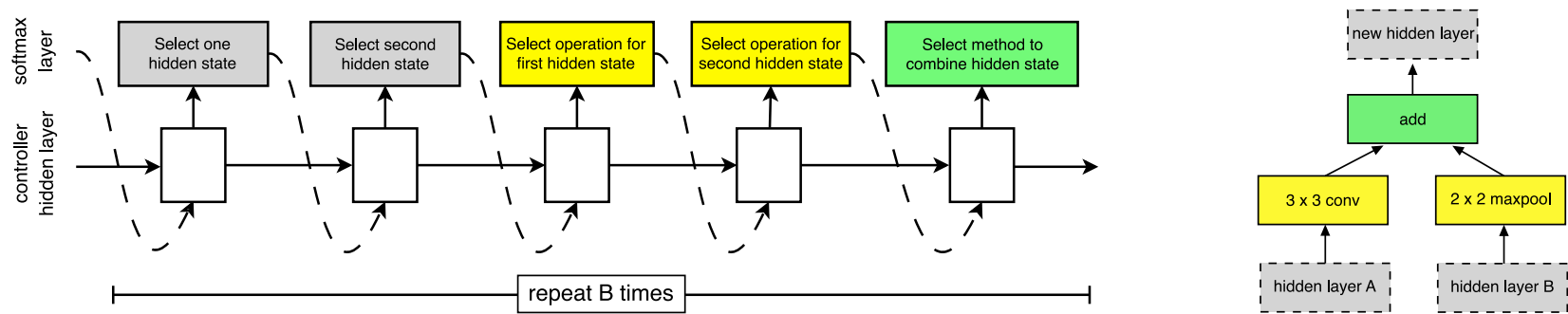


Figure 3. Controller model architecture for recursively constructing one block of a convolutional cell. Each block requires selecting 5 discrete parameters, each of which corresponds to the output of a softmax layer. Example constructed block shown on right. A convolutional cell contains  $B$  blocks, hence the controller contains  $5B$  softmax layers for predicting the architecture of a convolutional cell. In our experiments, the number of blocks  $B$  is 5.

Figure 3 of "Learning Transferable Architectures for Scalable Image Recognition", <https://arxiv.org/abs/1707.07012>

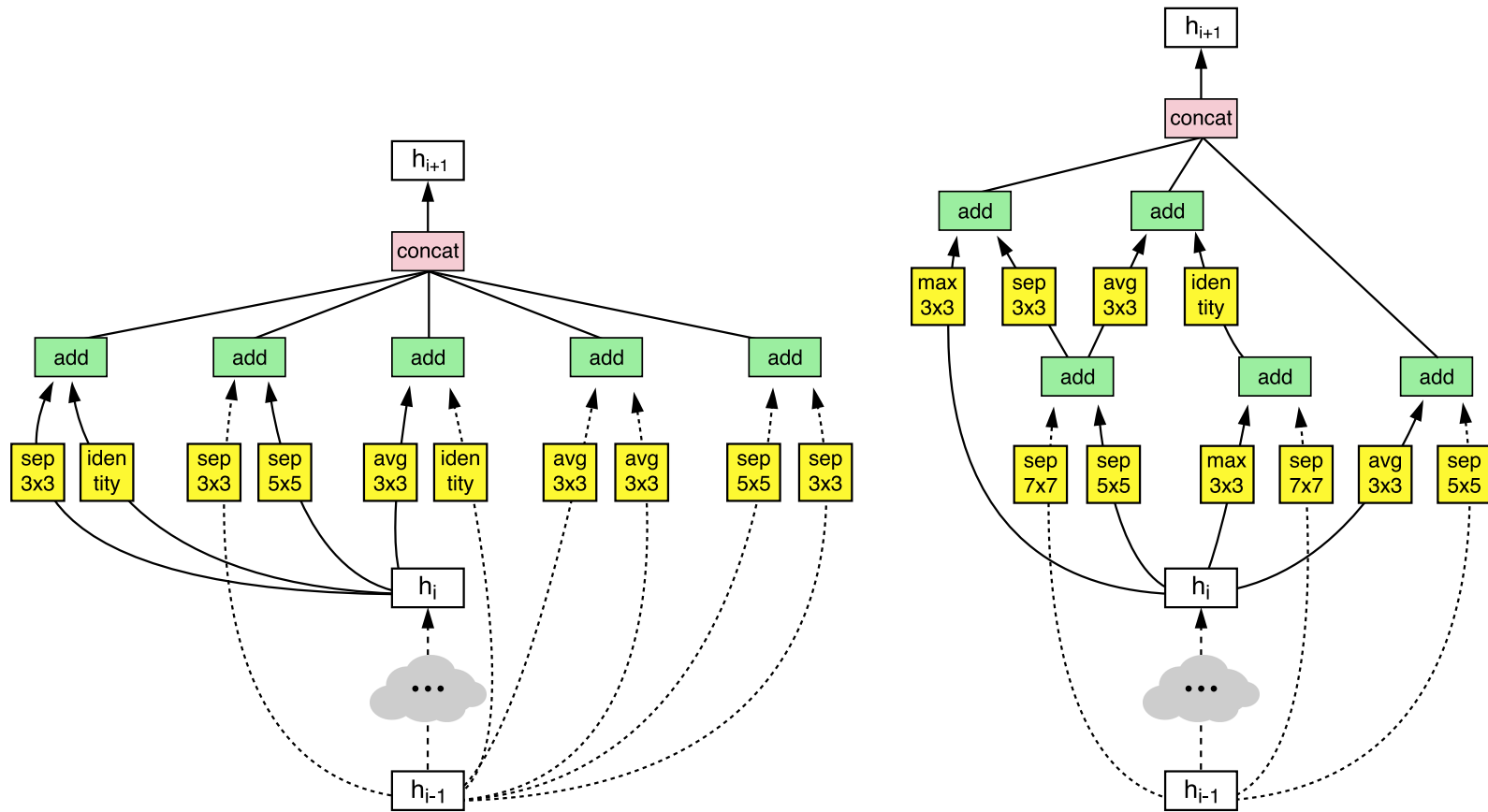
- Step 1.** Select a hidden state from  $h_i, h_{i-1}$  or from the set of hidden states created in previous blocks.
- Step 2.** Select a second hidden state from the same options as in Step 1.
- Step 3.** Select an operation to apply to the hidden state selected in Step 1.
- Step 4.** Select an operation to apply to the hidden state selected in Step 2.
- Step 5.** Select a method to combine the outputs of Step 3 and 4 to create a new hidden state.

Page 3 of "Learning Transferable Architectures for Scalable Image Recognition", <https://arxiv.org/abs/1707.07012>

- identity
- 1x7 then 7x1 convolution
- 3x3 average pooling
- 5x5 max pooling
- 1x1 convolution
- 3x3 depthwise-separable conv
- 7x7 depthwise-separable conv
- 1x3 then 3x1 convolution
- 3x3 dilated convolution
- 3x3 max pooling
- 7x7 max pooling
- 3x3 convolution
- 5x5 depthwise-separable conv

Figure 2 of "Learning Transferable Architectures for Scalable Image Recognition", <https://arxiv.org/abs/1707.07012>

The final Normal Cell and Reduction Cell chosen from 20k architectures (500GPUs, 4days).



Normal Cell

Reduction Cell

Page 3 of "Learning Transferable Architectures for Scalable Image Recognition", <https://arxiv.org/abs/1707.07012>

EfficientNet changes the search in three ways.

- Computational requirements are part of the return. Notably, the goal is to find an architecture  $m$  maximizing

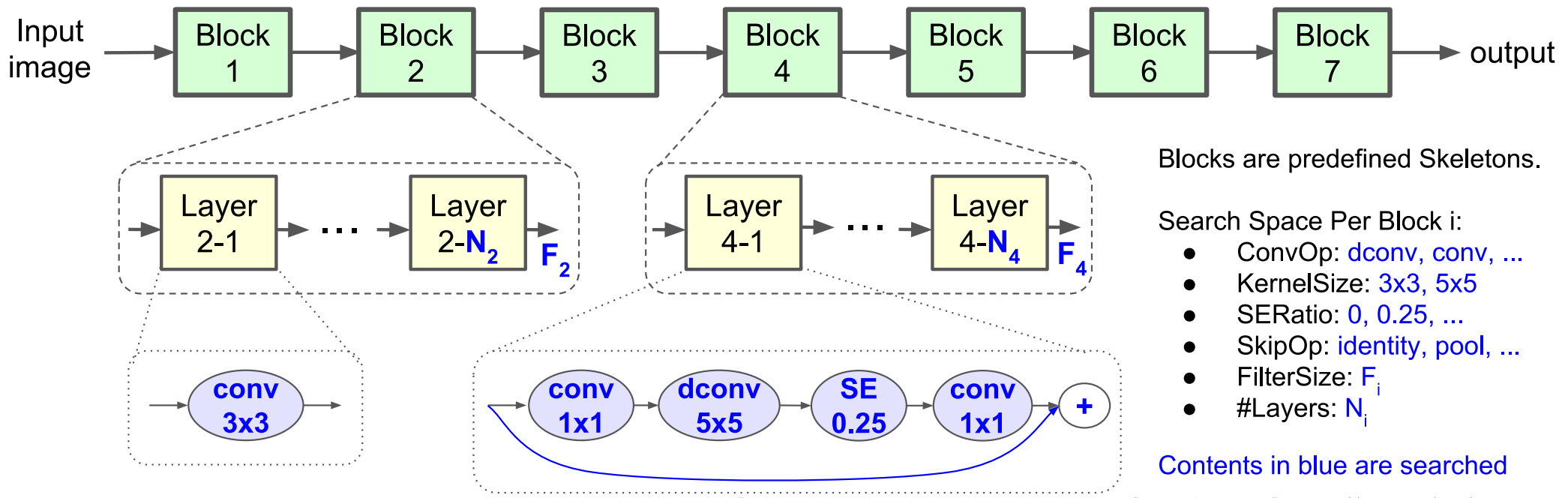
$$\text{DevelopmentAccuracy}(m) \cdot \left( \frac{\text{TargetFLOPS}=400\text{M}}{\text{FLOPS}(m)} \right)^{0.07},$$

where the constant 0.07 balances the accuracy and FLOPS (*the constant comes from an empirical observation that doubling the FLOPS brings about 5% relative accuracy gain, and  $1.05 = 2^\beta$  gives  $\beta \approx 0.0704$* ).

- Using a different search space, which allows to control kernel sizes and channels in different parts of the architecture (compared to using the same cell everywhere as in NASNet).
- Training directly on ImageNet, but only for 5 epochs.

In total, 8k model architectures are sampled, and PPO algorithm is used instead of the REINFORCE with baseline.

# EfficientNet Search



Blocks are predefined Skeletons.

Search Space Per Block  $i$ :

- ConvOp: dconv, conv, ...
- KernelSize: 3x3, 5x5
- SERatio: 0, 0.25, ...
- SkipOp: identity, pool, ...
- FilterSize:  $F_i$
- #Layers:  $N_i$

Contents in blue are searched

Figure 4 of "MnasNet: Platform-Aware Neural Architecture Search for Mobile", <https://arxiv.org/abs/1807.11626>

The overall architecture consists of 7 blocks, each described by 6 parameters – 42 parameters in total, compared to 50 parameters of the NASNet search space.

- Convolutional ops *ConvOp*: regular conv (conv), depthwise conv (dconv), and mobile inverted bottleneck conv [29].
- Convolutional kernel size *KernelSize*: 3x3, 5x5.
- Squeeze-and-excitation [13] ratio *SERatio*: 0, 0.25.
- Skip ops *SkipOp*: pooling, identity residual, or no skip.
- Output filter size  $F_i$ .
- Number of layers per block  $N_i$ .

Page 4 of "MnasNet: Platform-Aware Neural Architecture Search for Mobile", <https://arxiv.org/abs/1807.11626>

Stage $i$	Operator $\hat{\mathcal{F}}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels $\hat{C}_i$	#Layers $\hat{L}_i$
1	Conv3x3	$224 \times 224$	32	1
2	MBCConv1, k3x3	$112 \times 112$	16	1
3	MBCConv6, k3x3	$112 \times 112$	24	2
4	MBCConv6, k5x5	$56 \times 56$	40	2
5	MBCConv6, k3x3	$28 \times 28$	80	3
6	MBCConv6, k5x5	$14 \times 14$	112	3
7	MBCConv6, k5x5	$14 \times 14$	192	4
8	MBCConv6, k3x3	$7 \times 7$	320	1
9	Conv1x1 & Pooling & FC	$7 \times 7$	1280	1

Table 1 of "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks", <https://arxiv.org/abs/1905.11946>



Our goal is to model speech, using an auto-regressive model

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

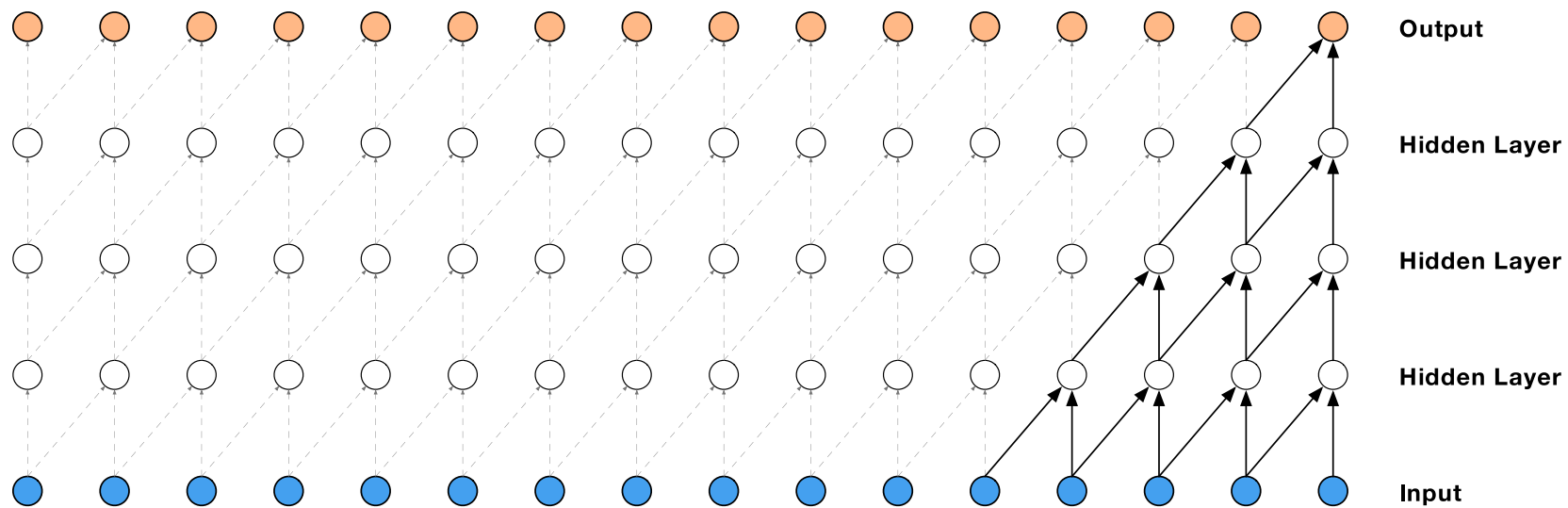


Figure 2: Visualization of a stack of causal convolutional layers.

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

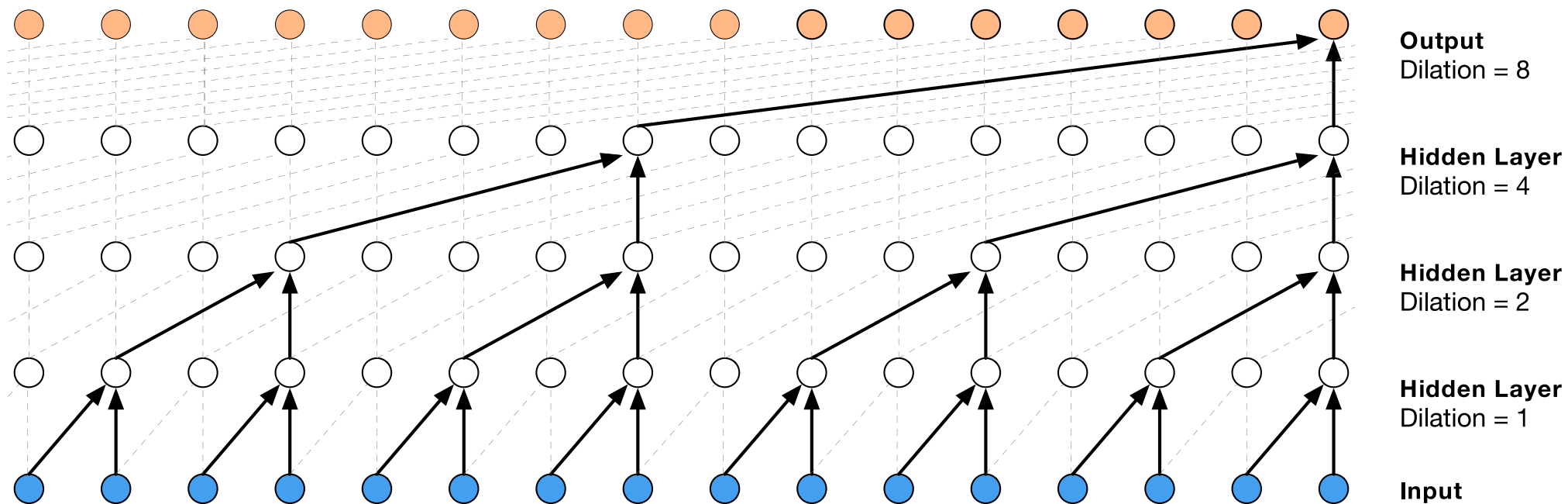


Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

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

## Output Distribution

WaveNet generates audio with 16kHz frequency and 16-bit samples. However, classification into 65 536 classes would not be tractable; instead, WaveNet adopts  $\mu$ -law transformation and quantize the samples into 256 values using

$$\text{sign}(x) \frac{\ln(1 + 255|x|)}{\ln(1 + 255)}.$$

## Gated Activation

To allow greater flexibility, the outputs of the dilated convolutions are passed through the gated activation units:

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

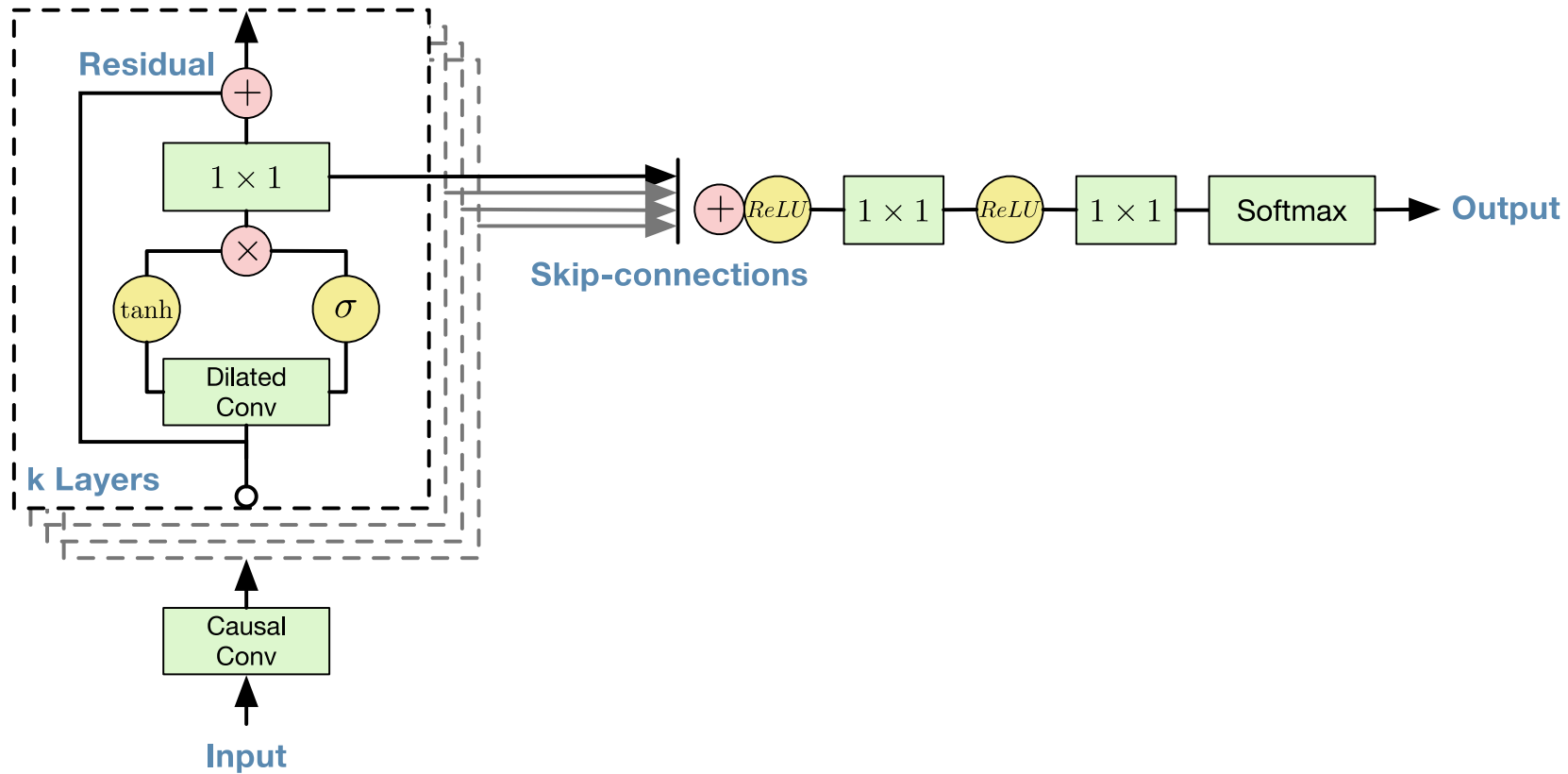


Figure 4: Overview of the residual block and the entire architecture.

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

## Global Conditioning

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

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

## 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

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

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
- filter size of a dilated convolution is 3
- residual connection has dimension 512
- gating layer uses 256+256 hidden units
- the  $1 \times 1$  output convolution produces 256 filters
- trained for 1 000 000 steps using Adam with a fixed learning rate of 0.0002

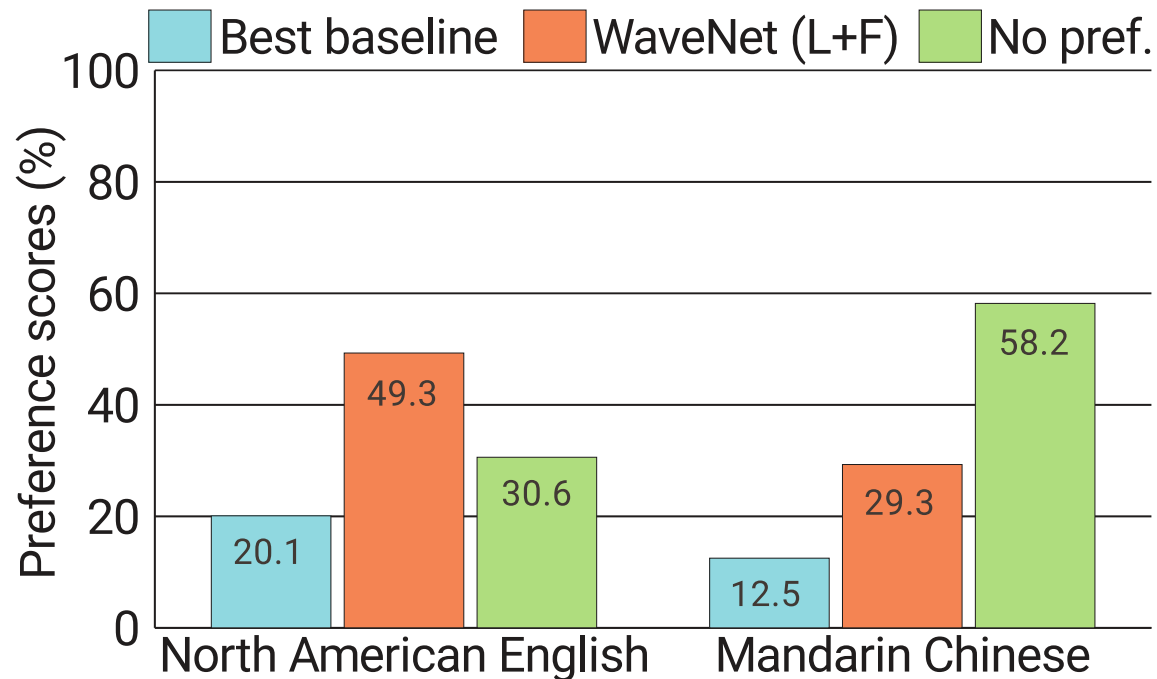


Figure 5: Subjective preference scores (%) of speech samples between (top) two baselines, (middle) two WaveNets, and (bottom) the best baseline and WaveNet. Note that LSTM and Concat correspond to LSTM-RNN-based statistical parametric and HMM-driven unit selection concatenative baseline synthesizers, 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>
ReLU	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>

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):

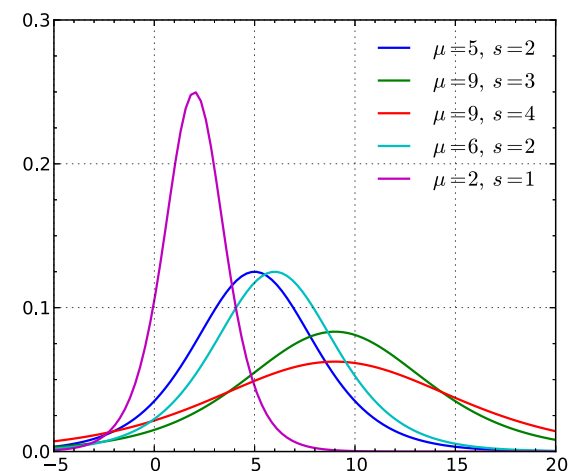
$$x \sim \sum_i \pi_i \text{Logistic}(\mu_i, s_i).$$

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, 10 mixture components are used.



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

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

Instead, we will 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  $\text{Logistic}(\mathbf{0}, \mathbf{1})$ . Then

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

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

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

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

$$\boldsymbol{\mu}_{\text{tot}} = \sum_i^N \mu^i(\mathbf{x}_{<t}^{i-1}) \cdot \left( \prod_{j>i}^N s^j(\mathbf{x}_{<t}^{j-1}) \right) \text{ and } \mathbf{s}_{\text{tot}} = \prod_i^N s^i(\mathbf{x}_{<t}^{i-1}).$$

# Probability Density Distillation

The network is trained using a **probability density distillation** using a teacher WaveNet, using KL-divergence as loss.

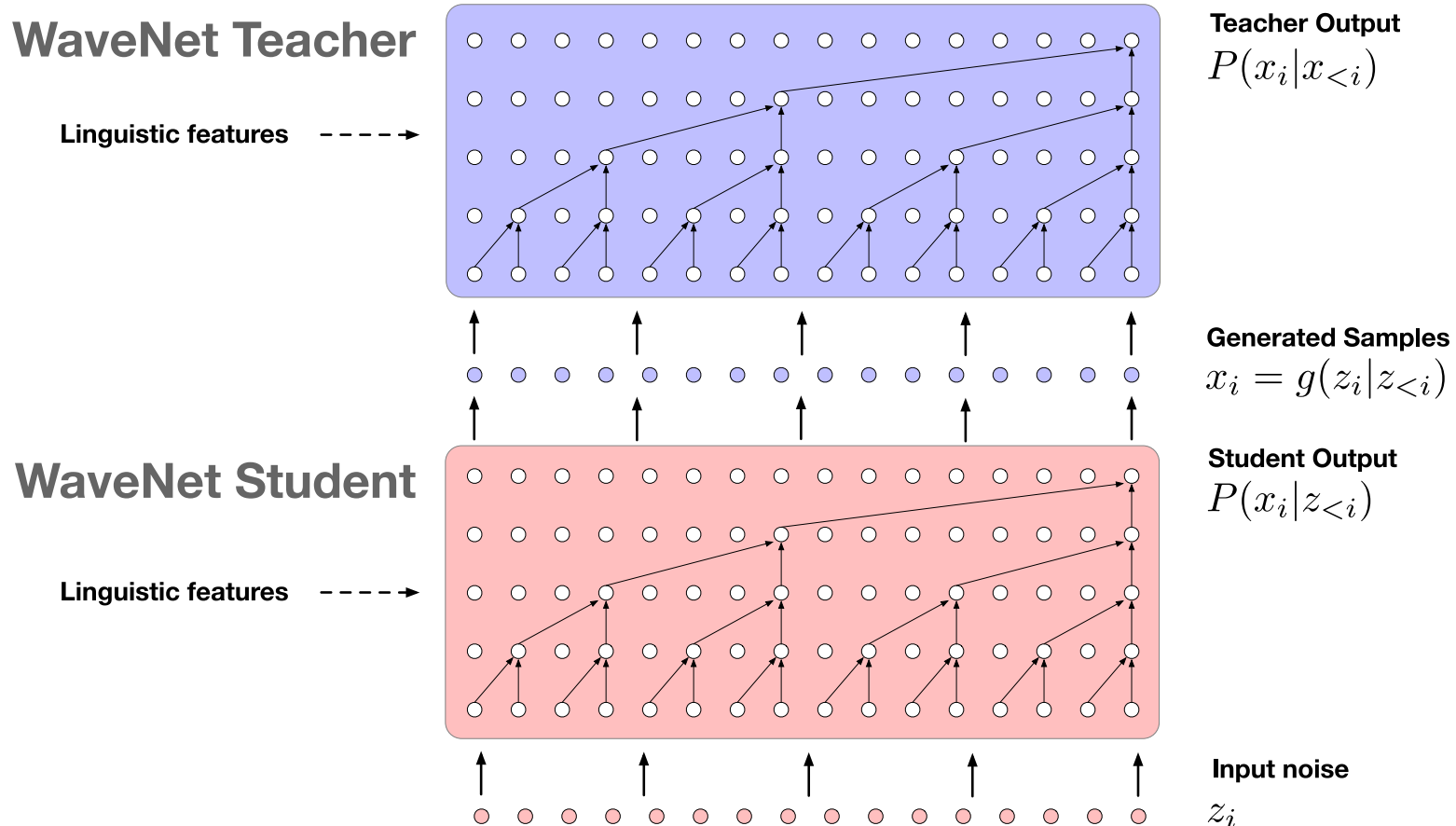


Figure 2 of "Parallel WaveNet: Fast High-Fidelity Speech Synthesis", <https://arxiv.org/abs/1711.10433>

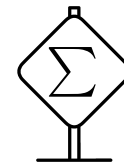
Denoting the teacher distribution as  $P_T$  and the student distribution as  $P_S$ , the loss is specifically

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

Therefore, we do not only minimize cross-entropy, but we also try to keep the entropy of the student as high as possible. That is 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.

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

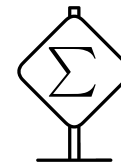


$$\begin{aligned} H(P_S) &= \mathbb{E}_{z \sim \text{Logistic}(0,1)} \left[ \sum_{t=1}^T -\ln p_S(\mathbf{x}_t | \mathbf{z}_{<t}) \right] \\ &= \mathbb{E}_{z \sim \text{Logistic}(0,1)} \left[ \sum_{t=1}^T \ln s(\mathbf{z}_{<t}, \boldsymbol{\theta}) \right] + 2T. \end{aligned}$$

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

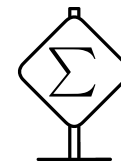
# Probability Density Distillation Details

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



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

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



We can therefore estimate  $H(P_S, P_T)$  by drawing a single sample  $\mathbf{x}$  from the student  $P_S$ , compute all  $P_T(x_t|\mathbf{x}_{<t})$  from the teacher in parallel, 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 the sample under the teacher using the original formulation.

Finally, analogously to the normal distribution, the logistic distribution offers the **reparametric trick**, which means we can differentiate  $\ln P_T(x_t|\mathbf{x}_{<t})$  with respect to both  $x_t$  and  $\mathbf{x}_{<t}$  (while the categorical distribution would be differentiable only with respect to  $\mathbf{x}_{<t}$ ).



With the 4 iterations, 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>

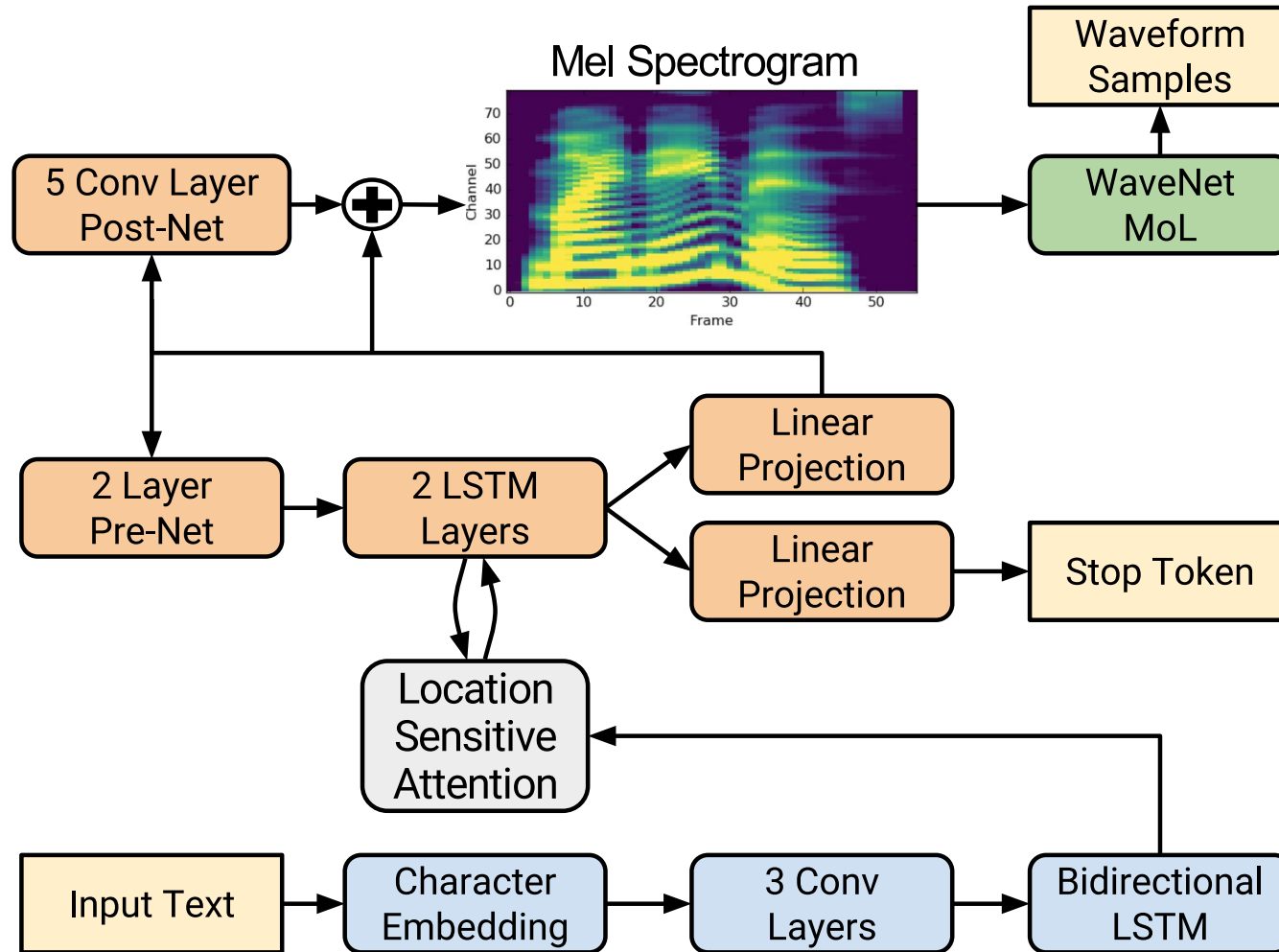


Figure 1 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

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>

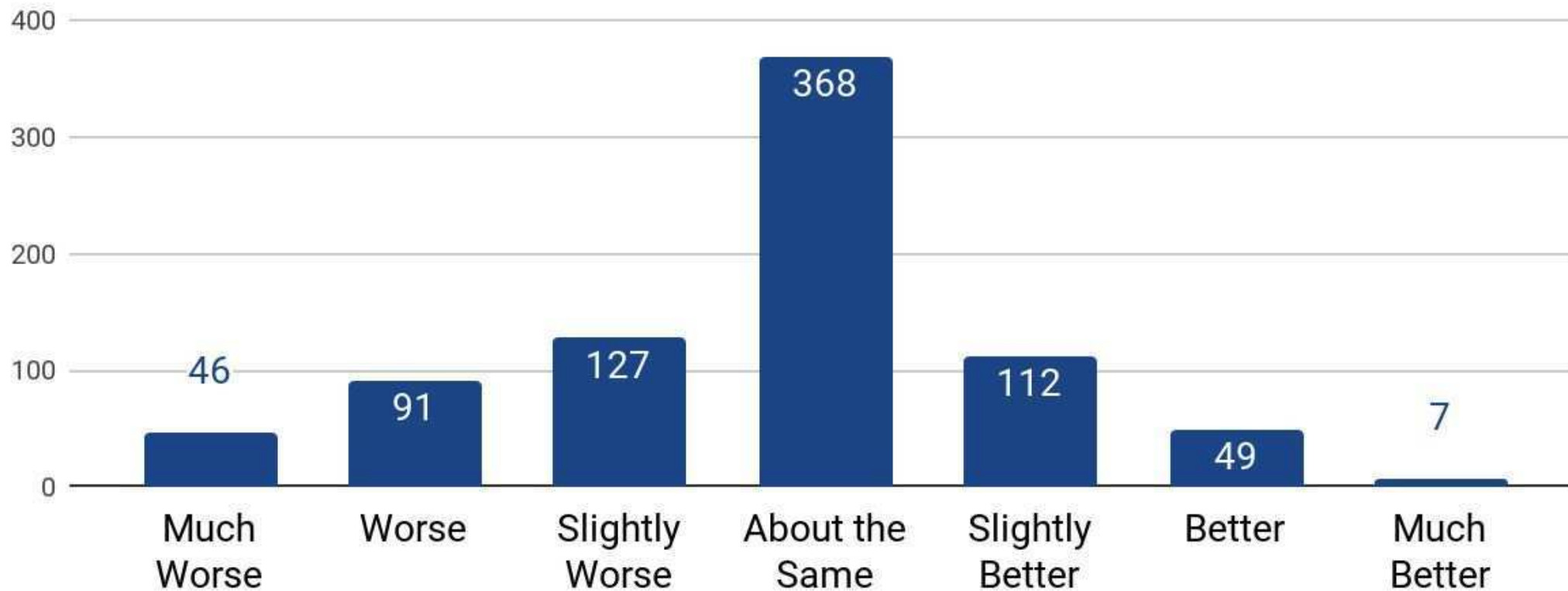


Figure 2 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

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

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.

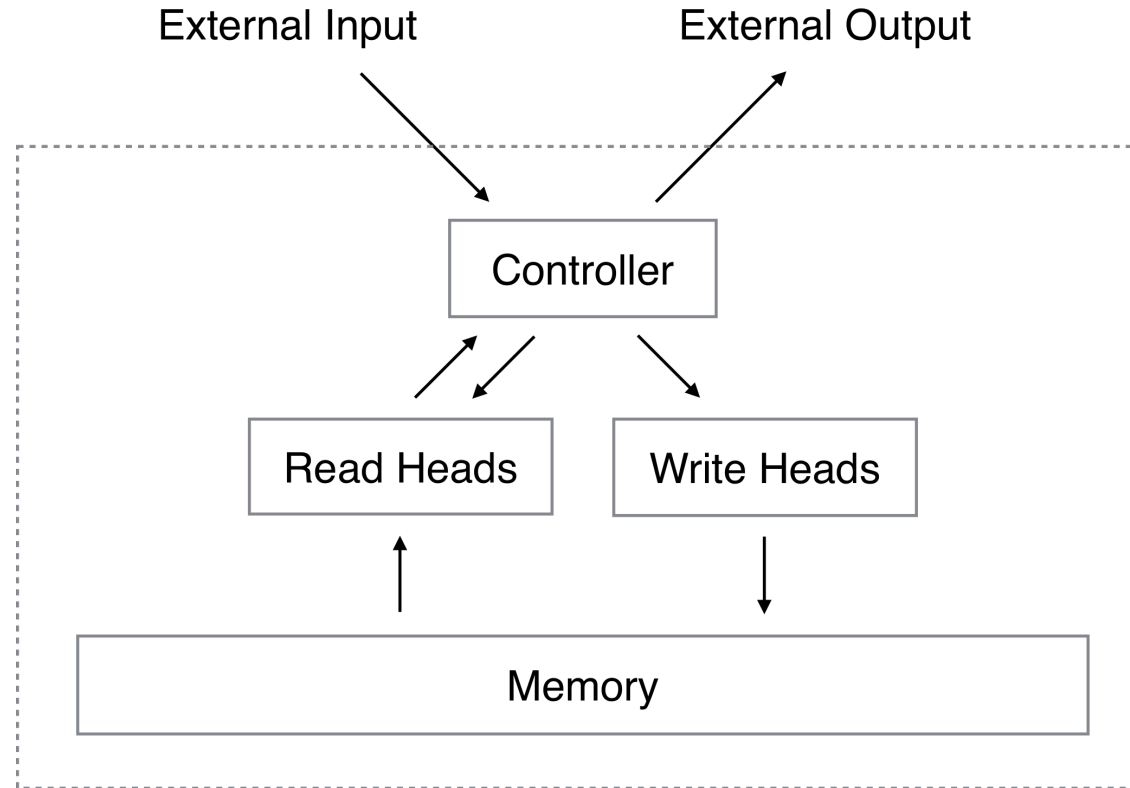


Figure 1 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>

## 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.

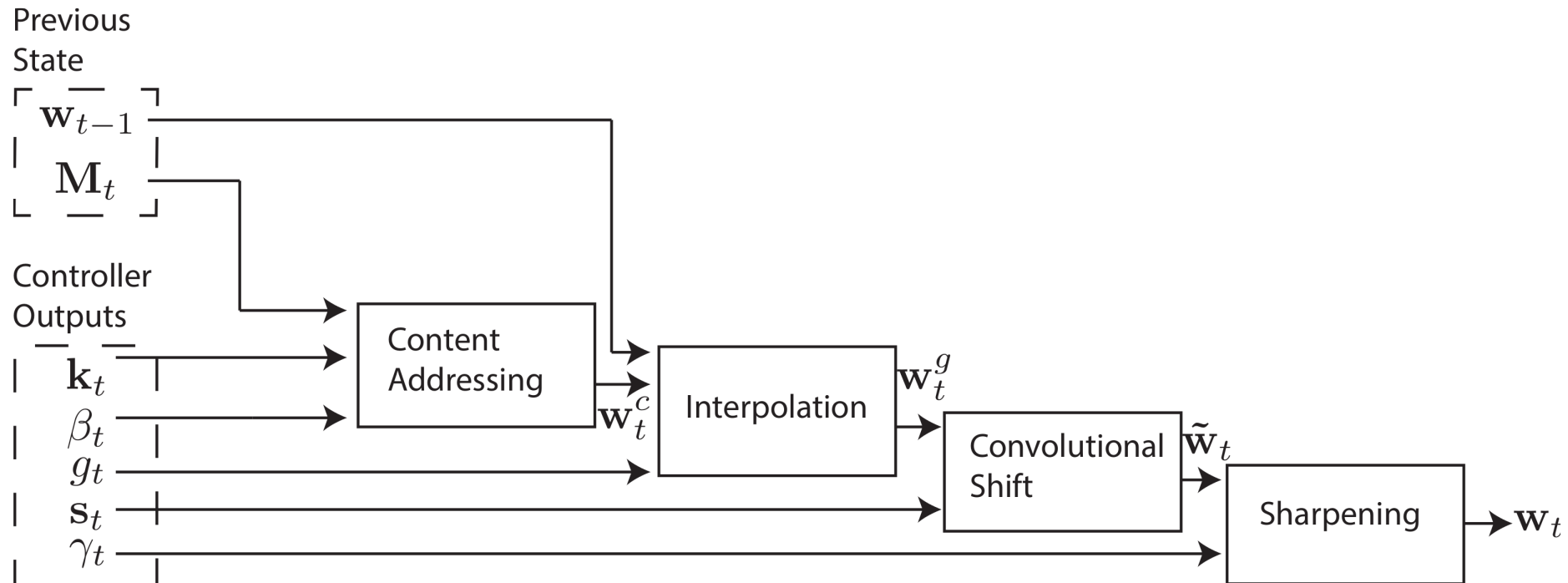


Figure 2 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>



## 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, the controller emits an *interpolation gate*  $g_t$  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*  $\gamma_t$ , 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.

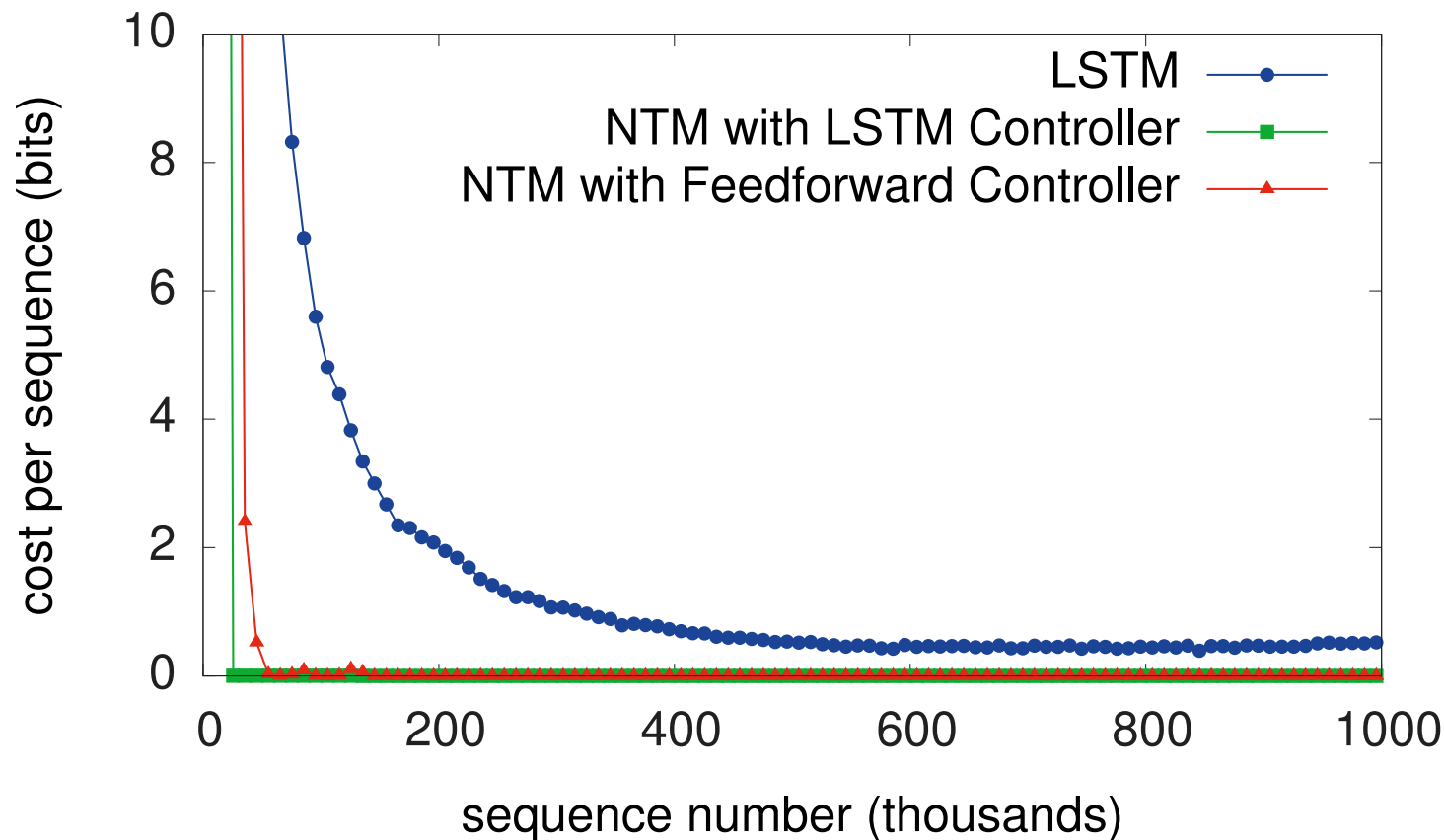
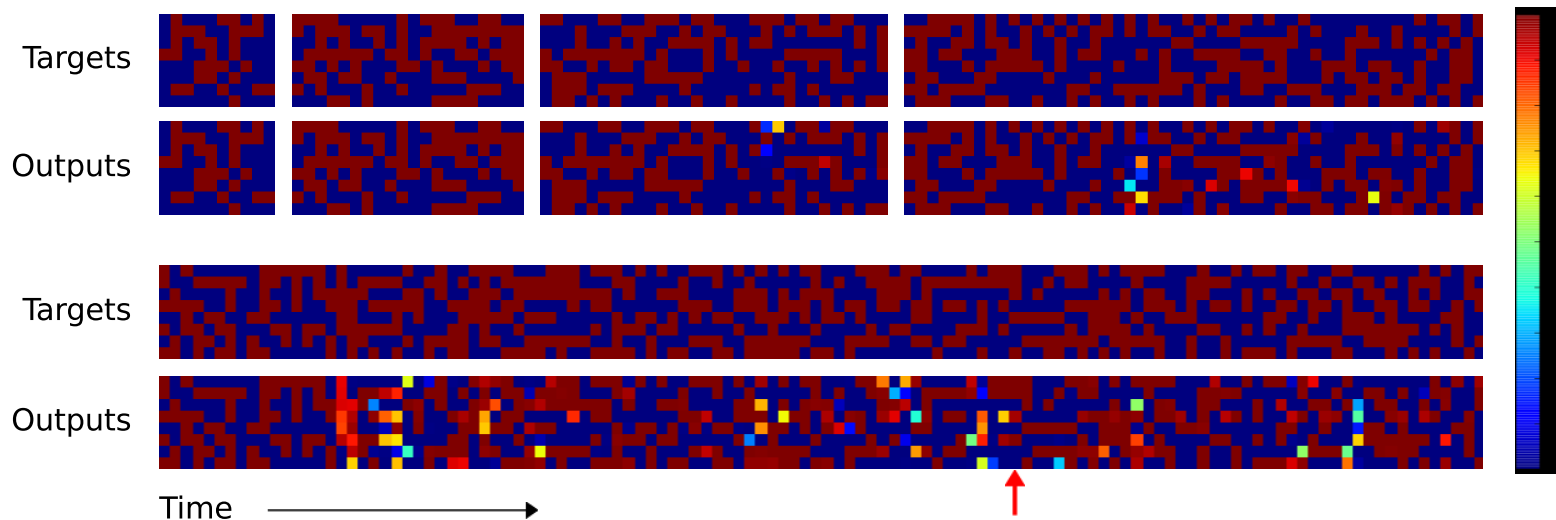
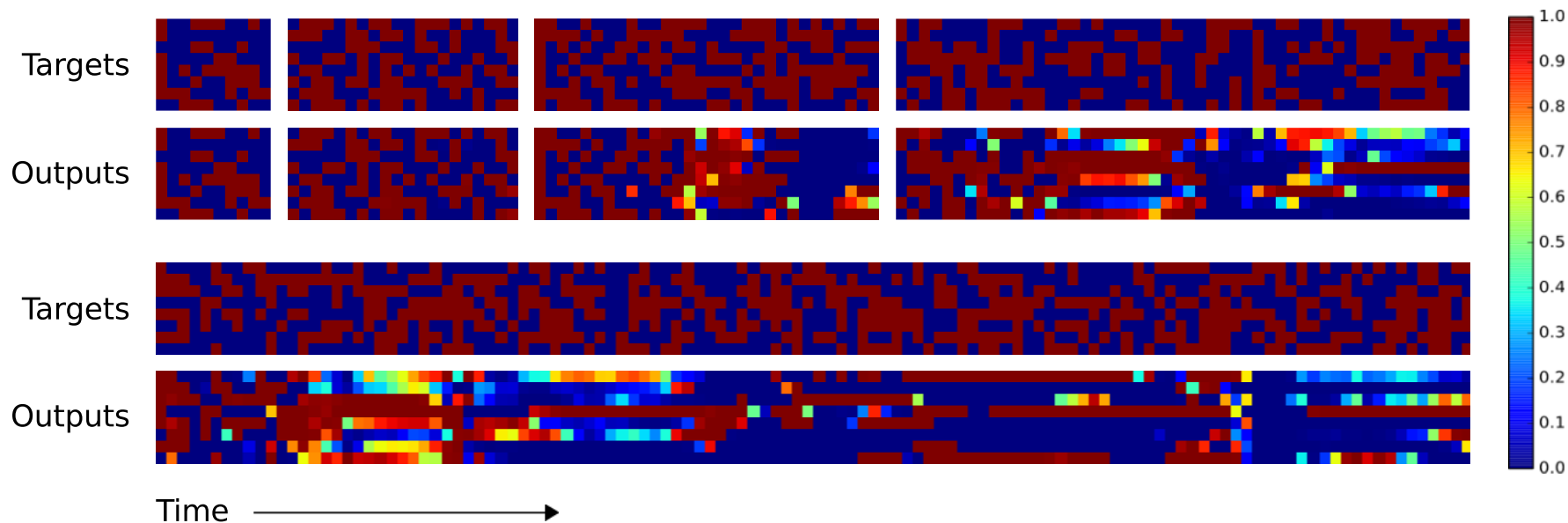


Figure 3 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>



**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>

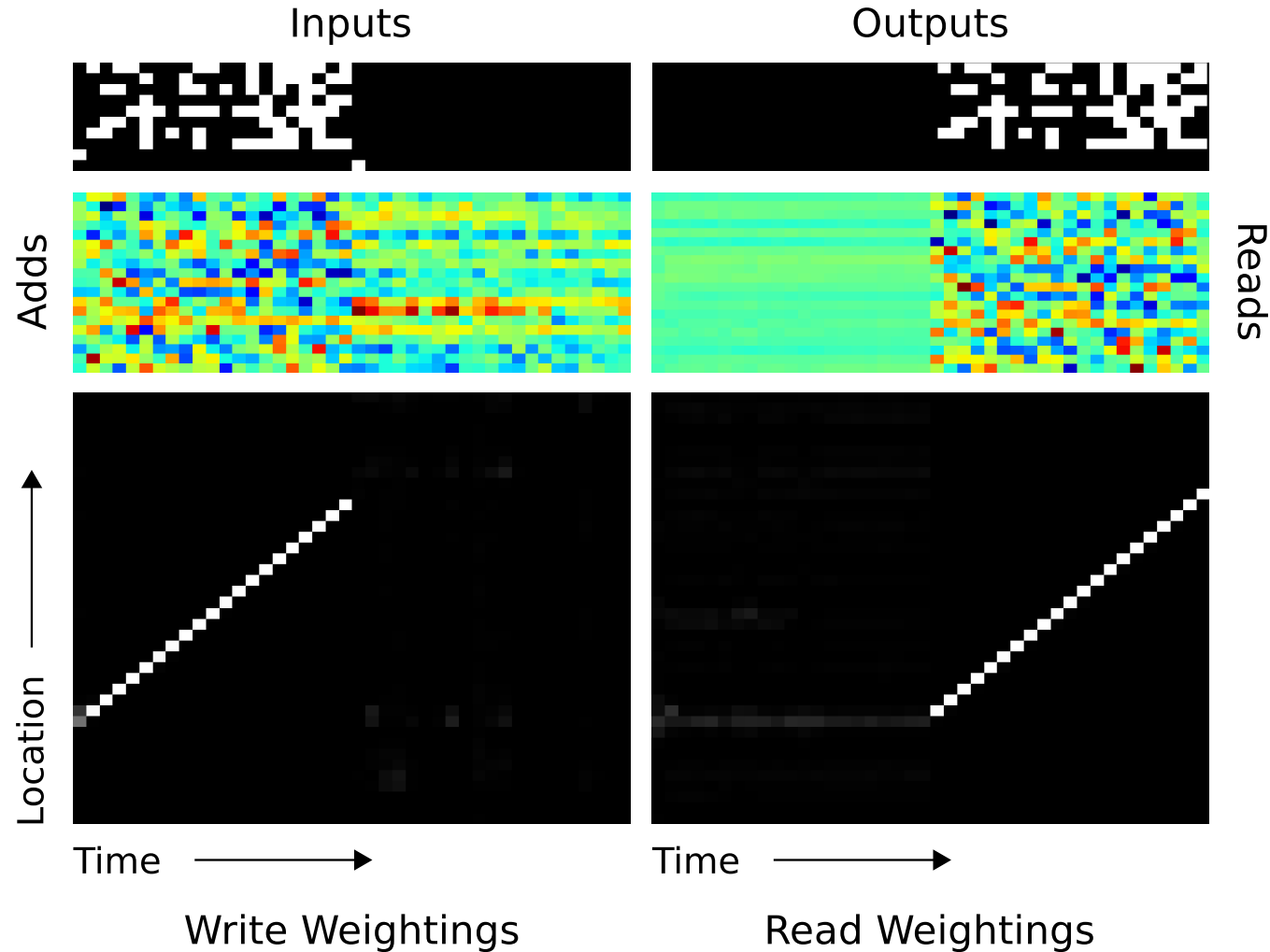


Figure 6 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>

## 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.



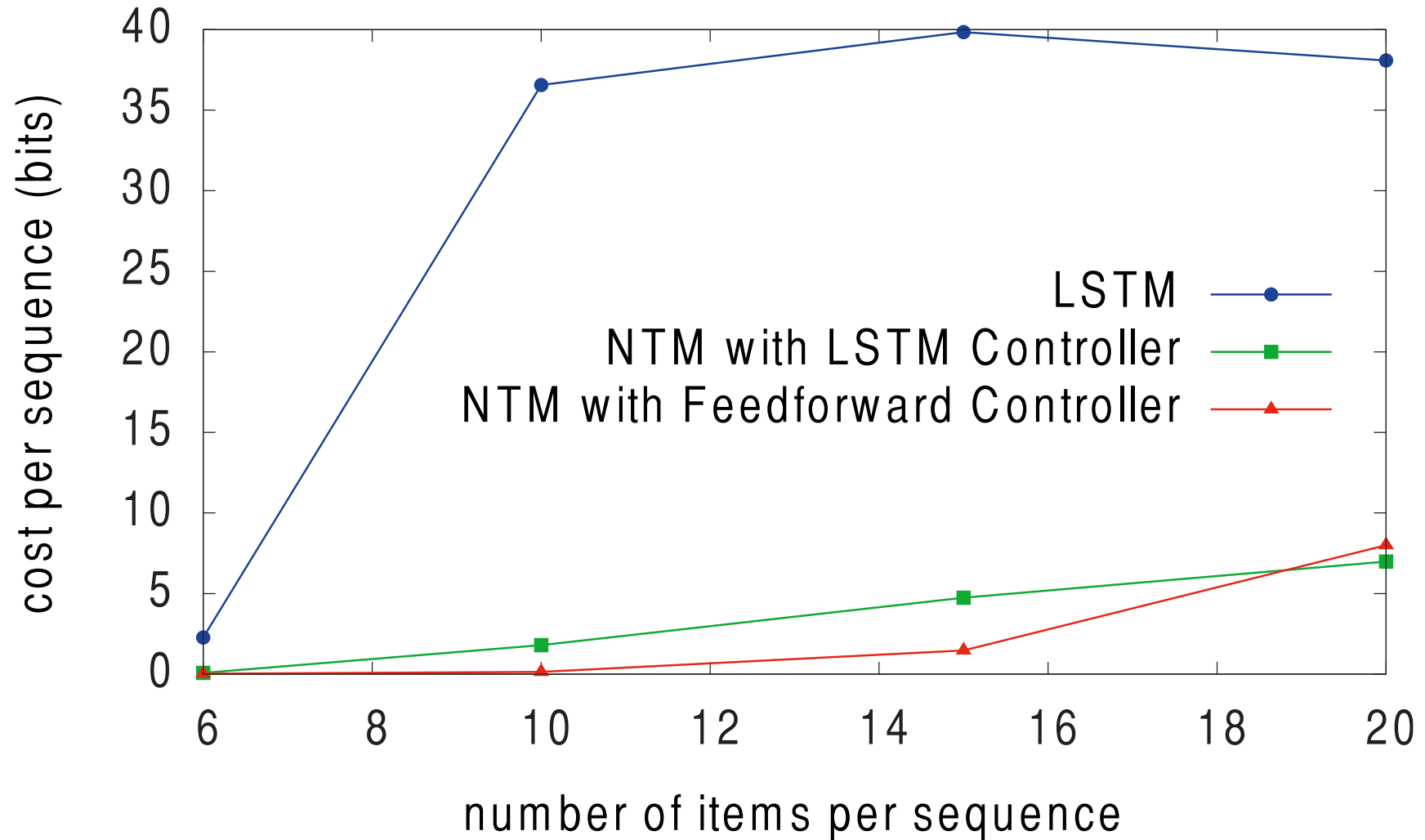


Figure 11 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>

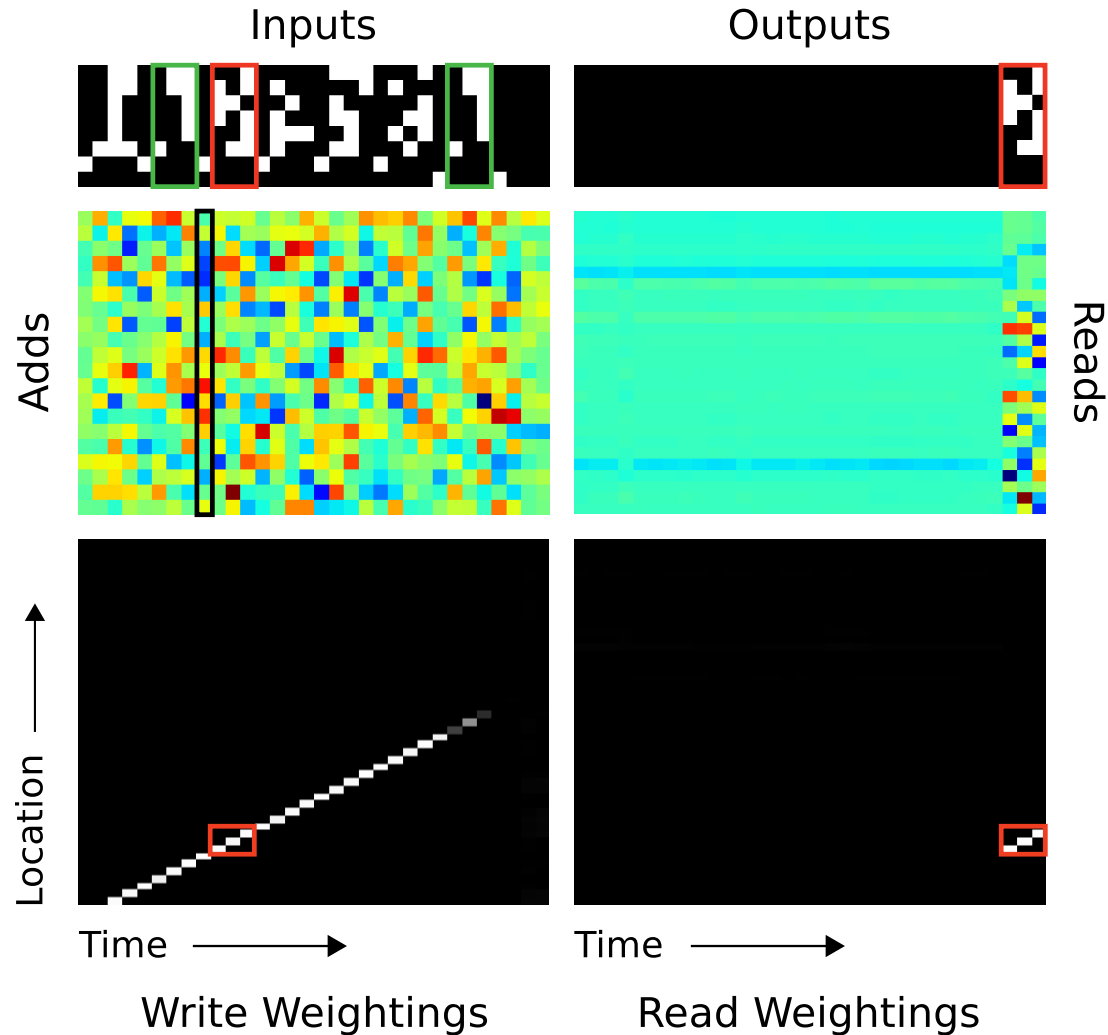


Figure 12 of "Neural Turing Machines", <https://arxiv.org/abs/1410.5401>

# Differentiable Neural Computer

NTM was later extended to a Differentiable Neural Computer.

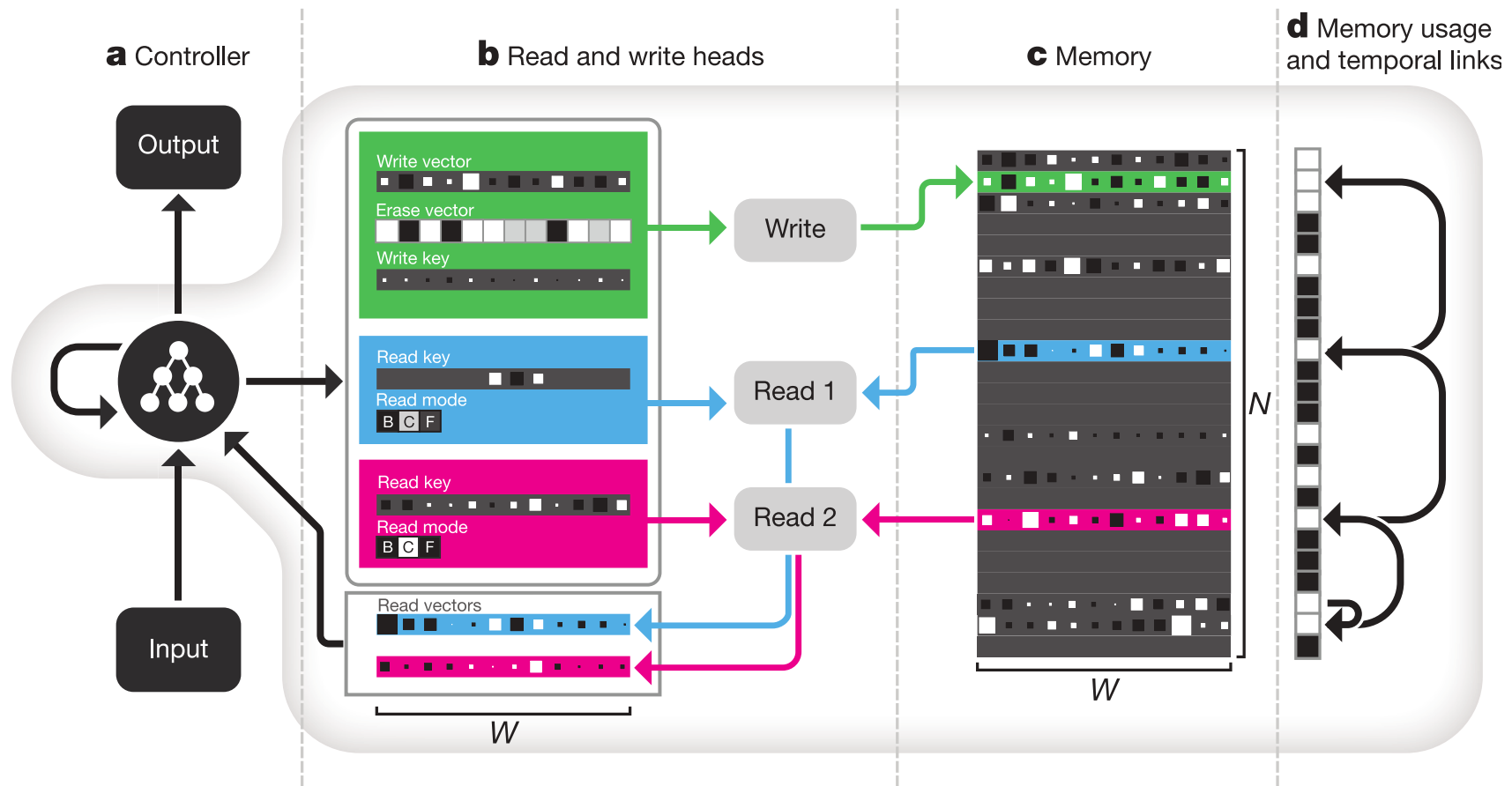


Figure 1 of "Hybrid computing using a neural network with dynamic external memory", <https://www.nature.com/articles/nature20101>

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, ...).

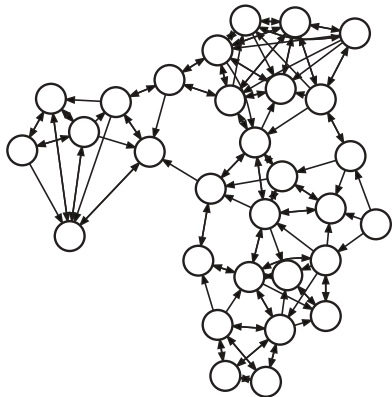
In DNC, the usage of every memory location is tracked, which enables performing dynamic allocation – at each time step, a cell with least usage can be allocated.

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.

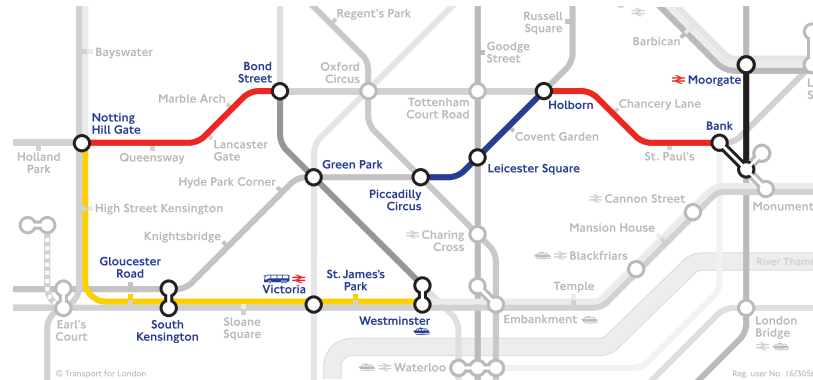
The write weighting is defined as a weighted combination of the allocation weighting and write content weighting, and read weighting is computed as a weighted combination of read content weighting, previous write weighting, and subsequent write weighting.

# Differentiable Neural Computer

**a** Random graph



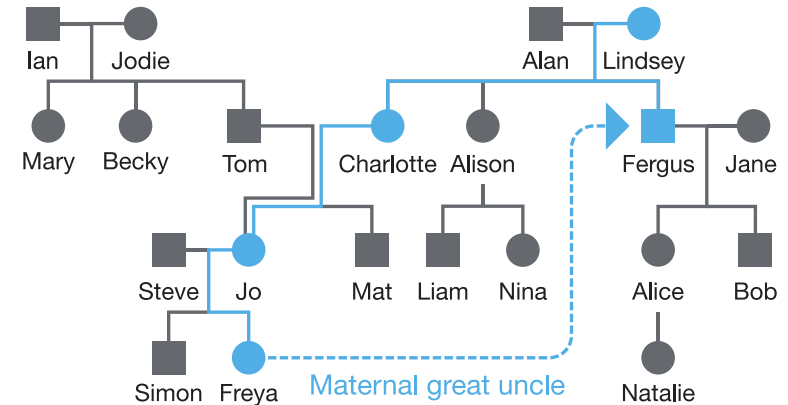
**b** London Underground



Traversal

Shortest-path

**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

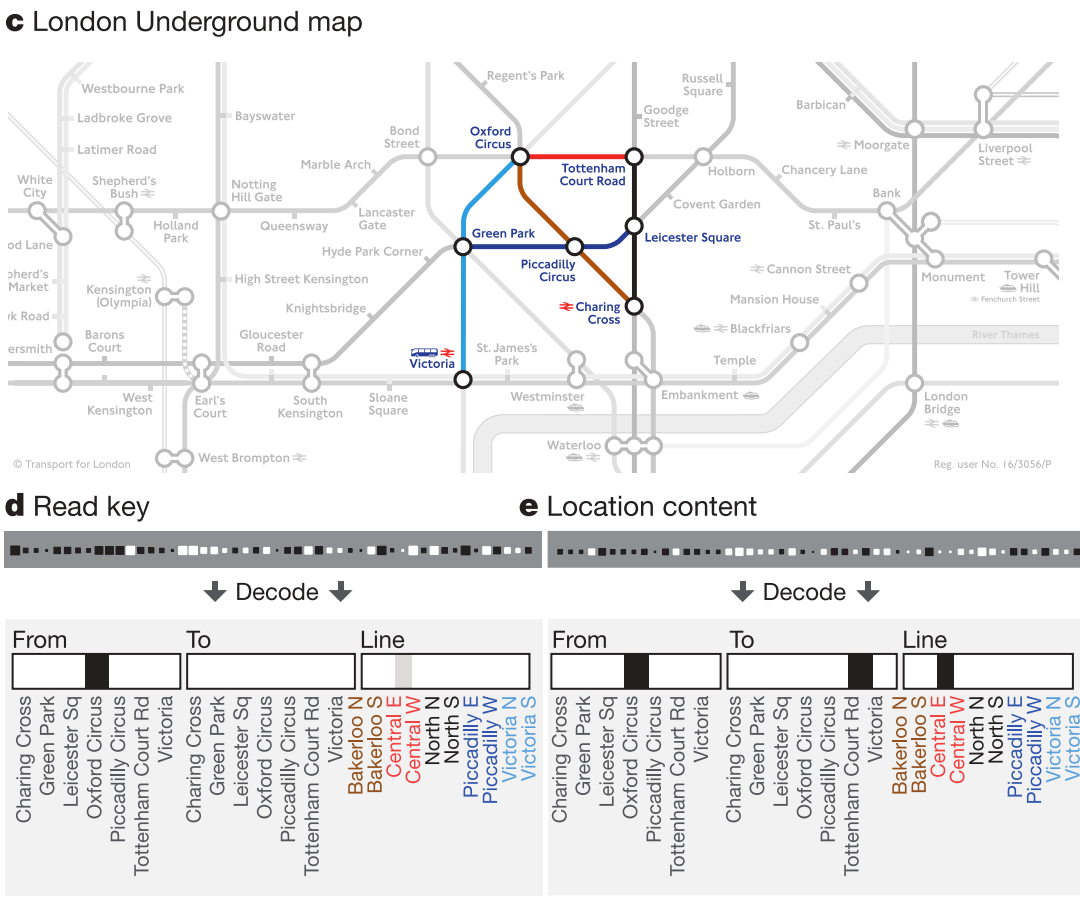
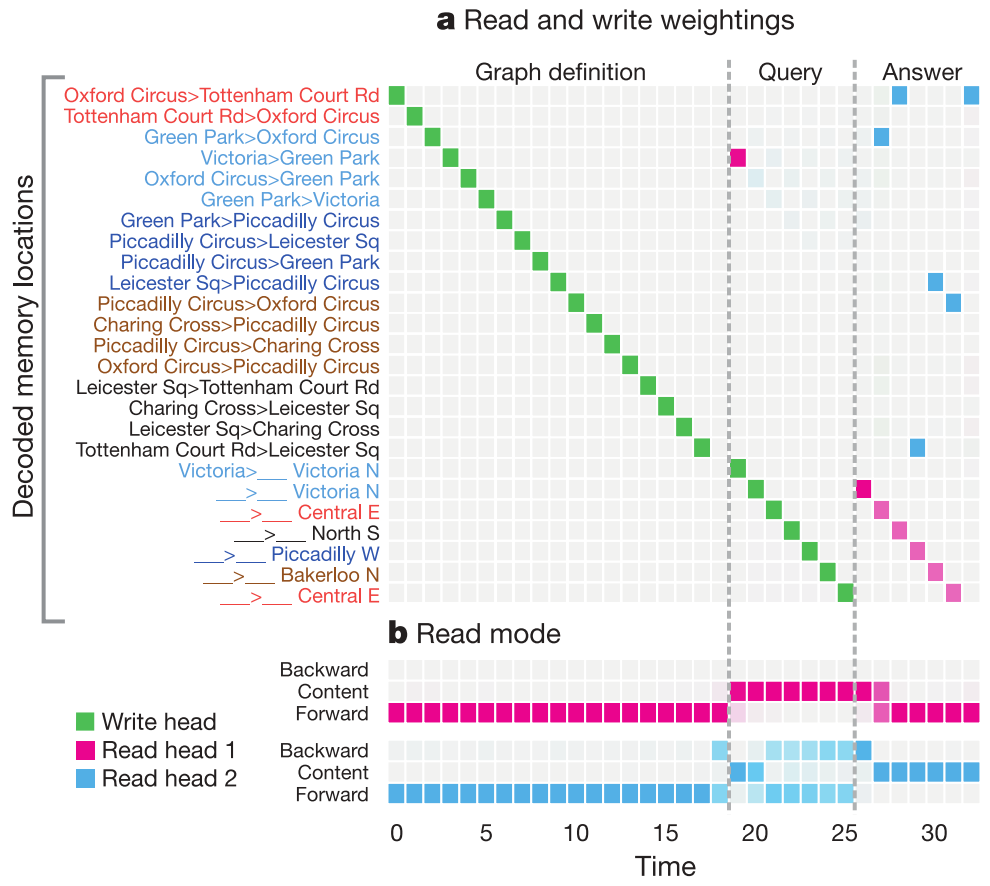


Figure 3 of "Hybrid computing using a neural network with dynamic external memory", <https://www.nature.com/articles/nature20101>

# Memory-augmented Neural Networks

External memory can be also utilized for **learning to learn**. 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.

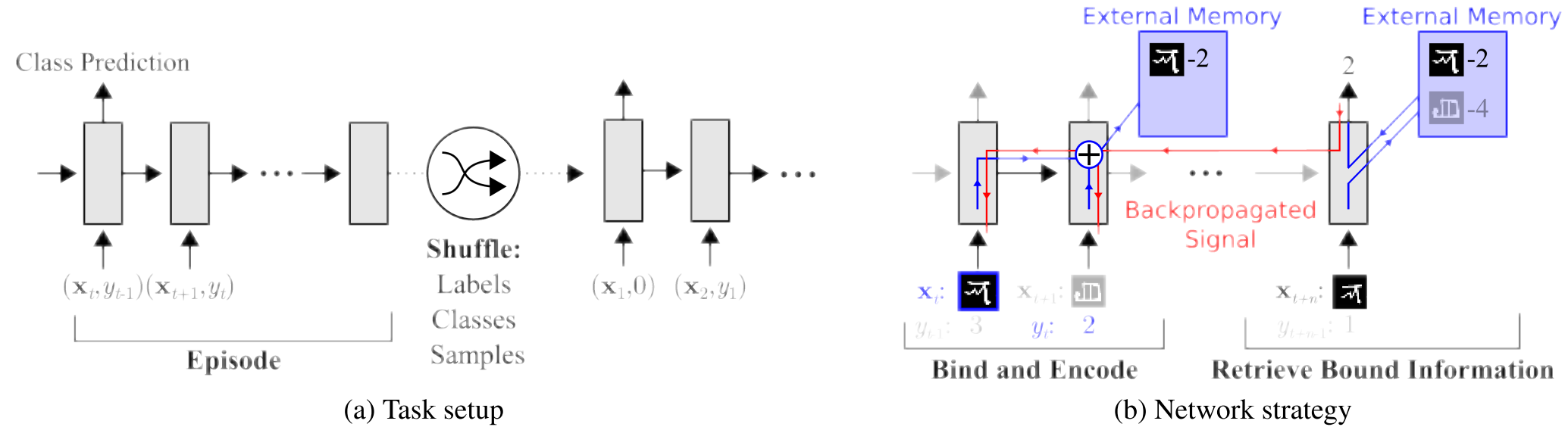


Figure 1 of "One-shot learning with Memory-Augmented Neural Networks", <https://arxiv.org/abs/1605.06065>

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

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:

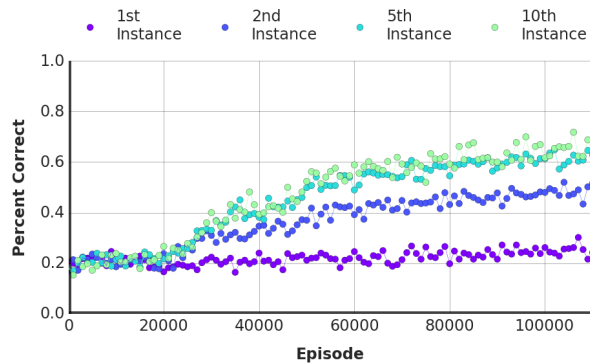
$$\mathbf{w}_t^w \leftarrow \sigma(\alpha) \mathbf{w}_{t-1}^r + (1 - \sigma(\alpha)) \mathbf{w}_{t-1}^{lu}. \quad (7)$$

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:

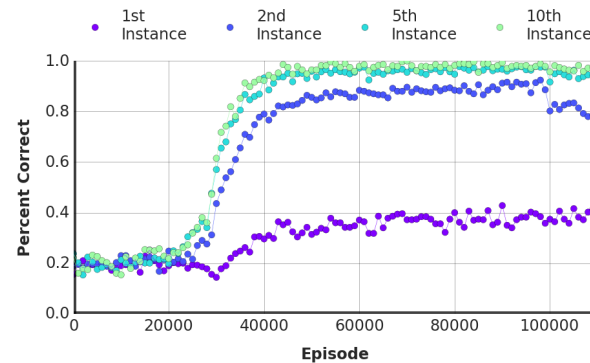
$$\mathbf{M}_t(i) \leftarrow \mathbf{M}_{t-1}(i) + w_t^w(i) \mathbf{k}_t, \forall i \quad (8)$$

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

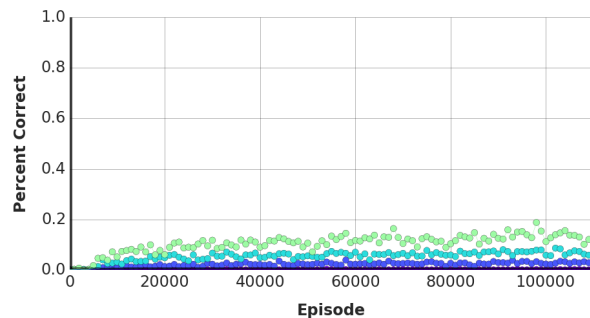




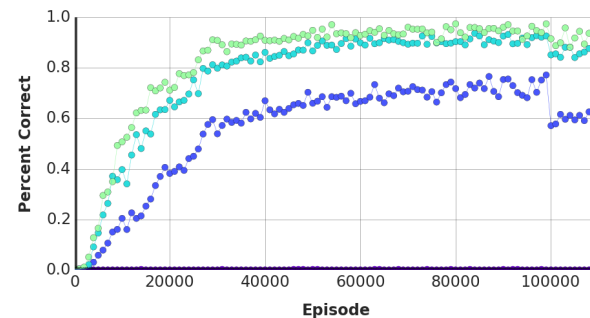
(a) LSTM, five random classes/episode, one-hot vector labels



(b) MANN, five random classes/episode, one-hot vector labels



(c) LSTM, fifteen classes/episode, five-character string labels



(d) MANN, fifteen classes/episode, five-character string labels

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>