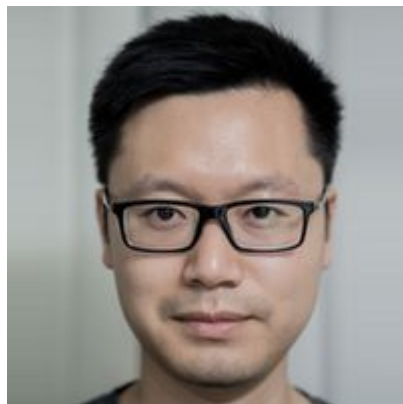


Online and Linear-Time Attention by Enforcing Monotonic Alignments

Colin Raffel



Thang Luong



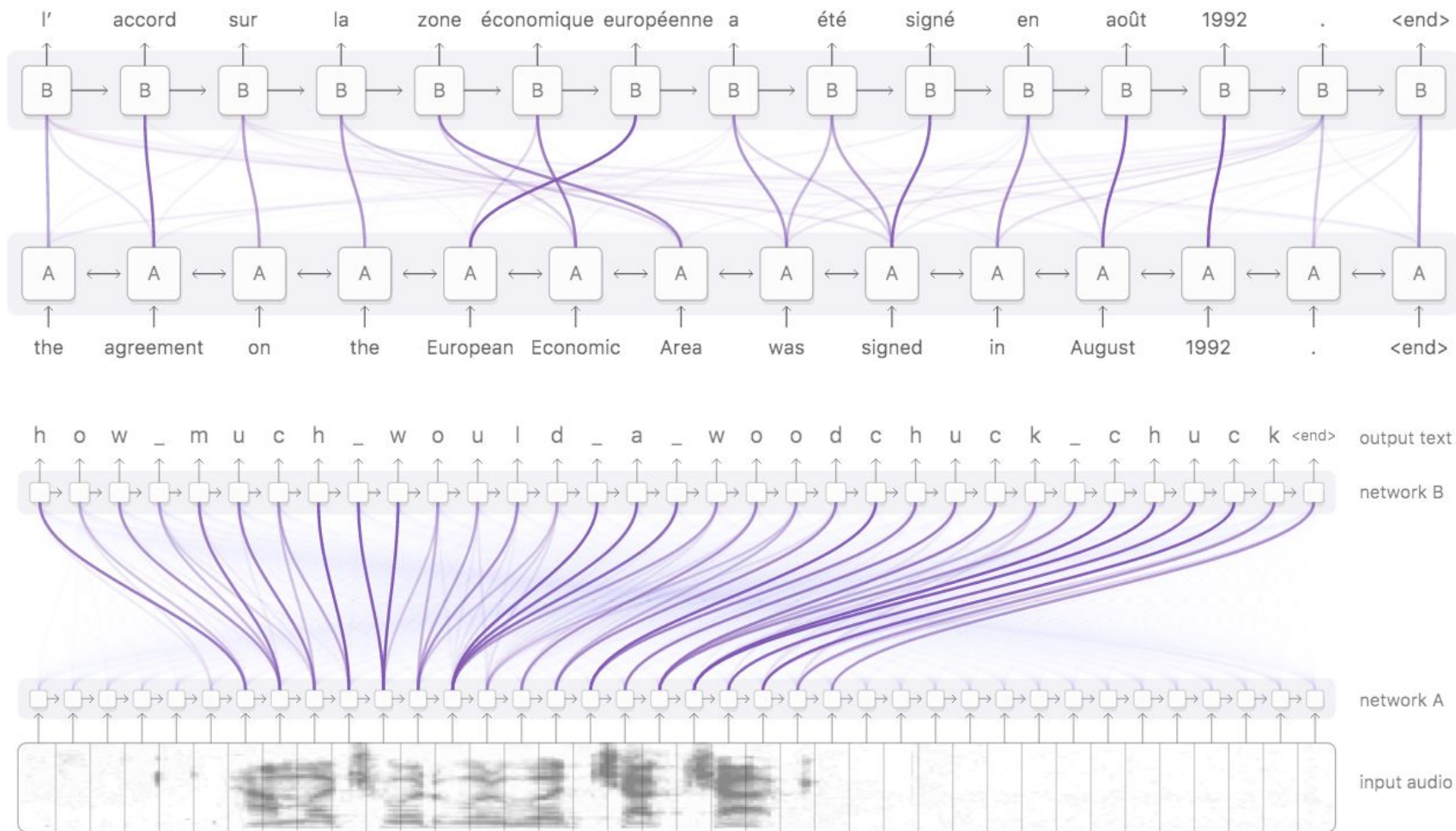
Peter J. Liu



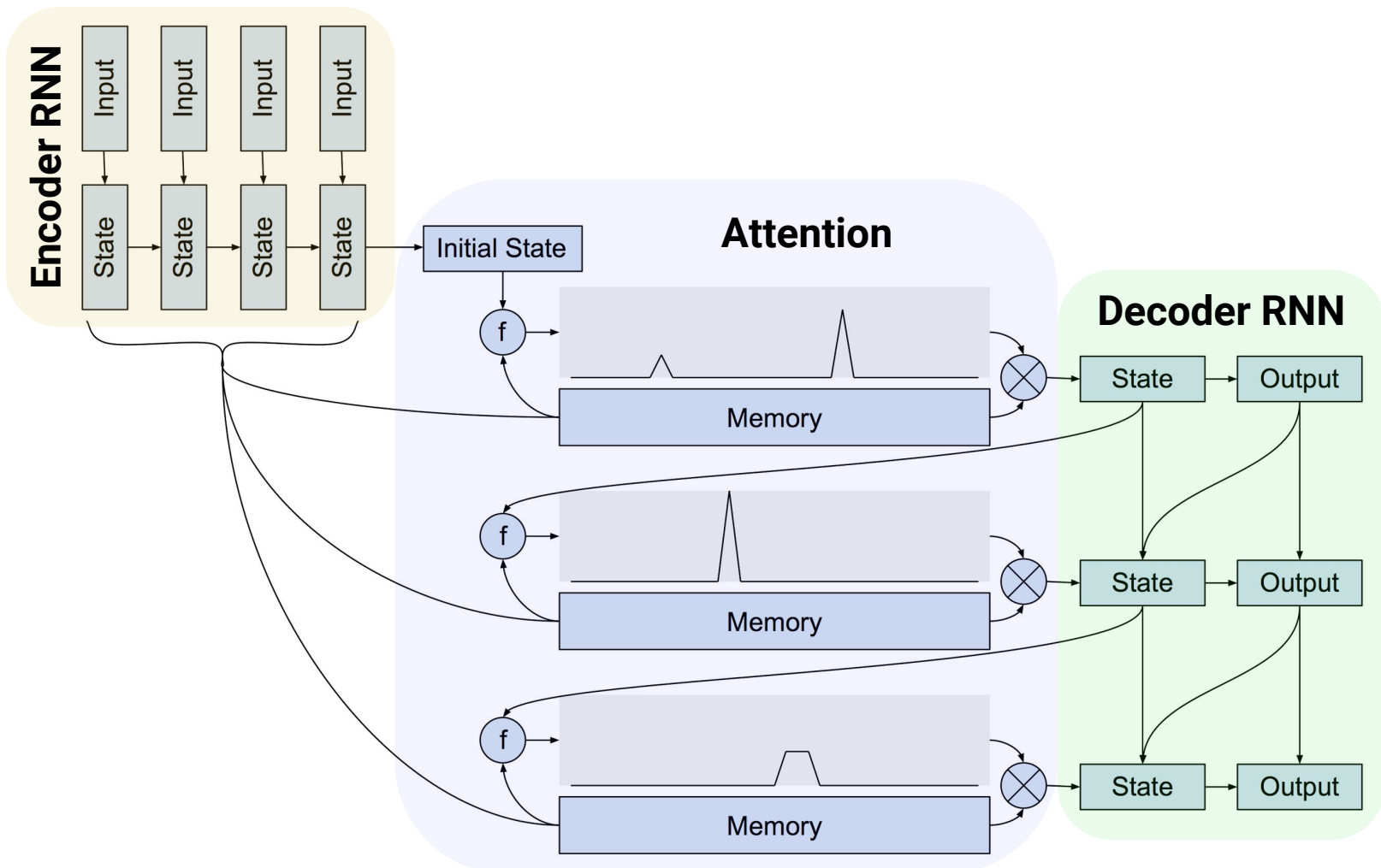
Ron J. Weiss



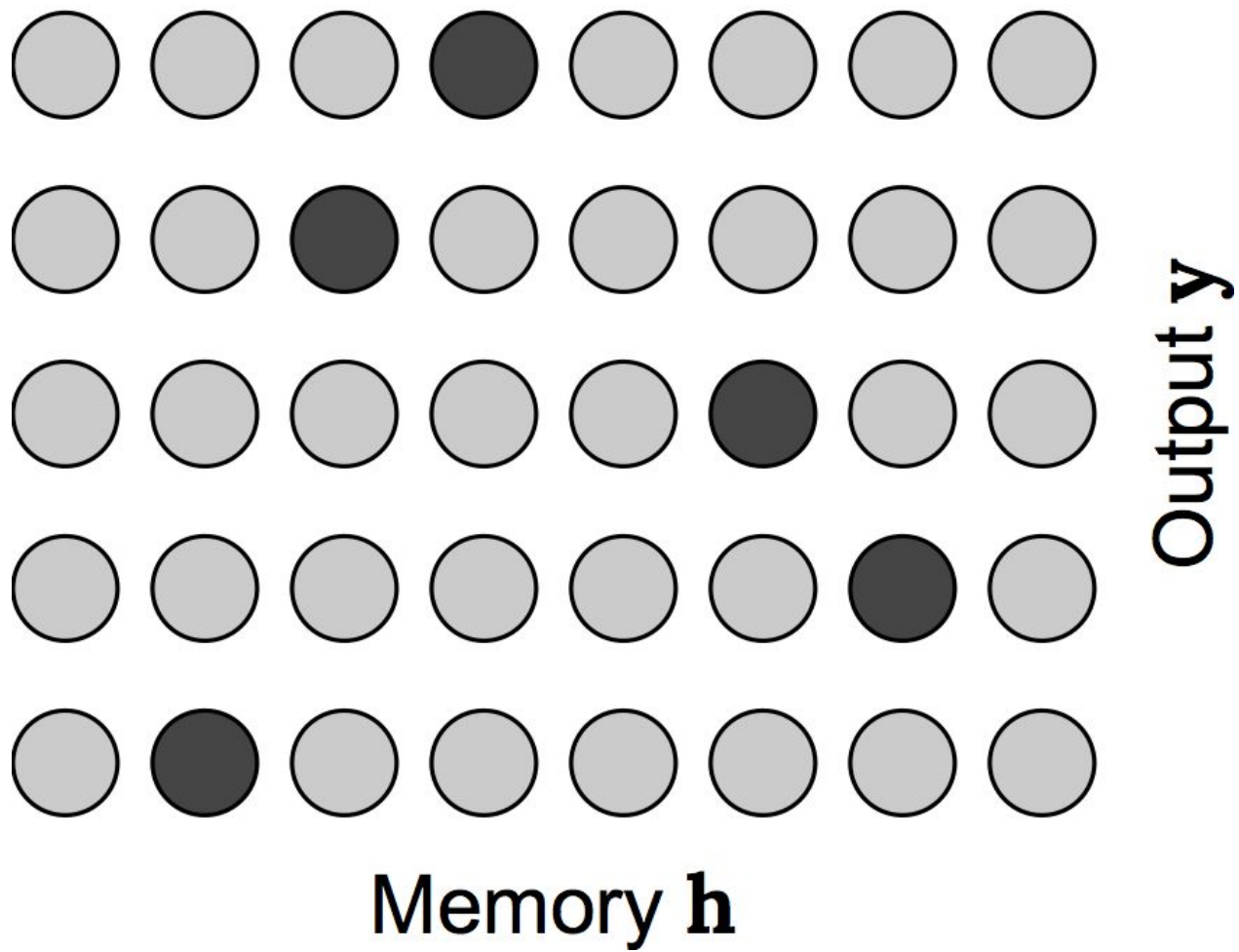
Douglas Eck

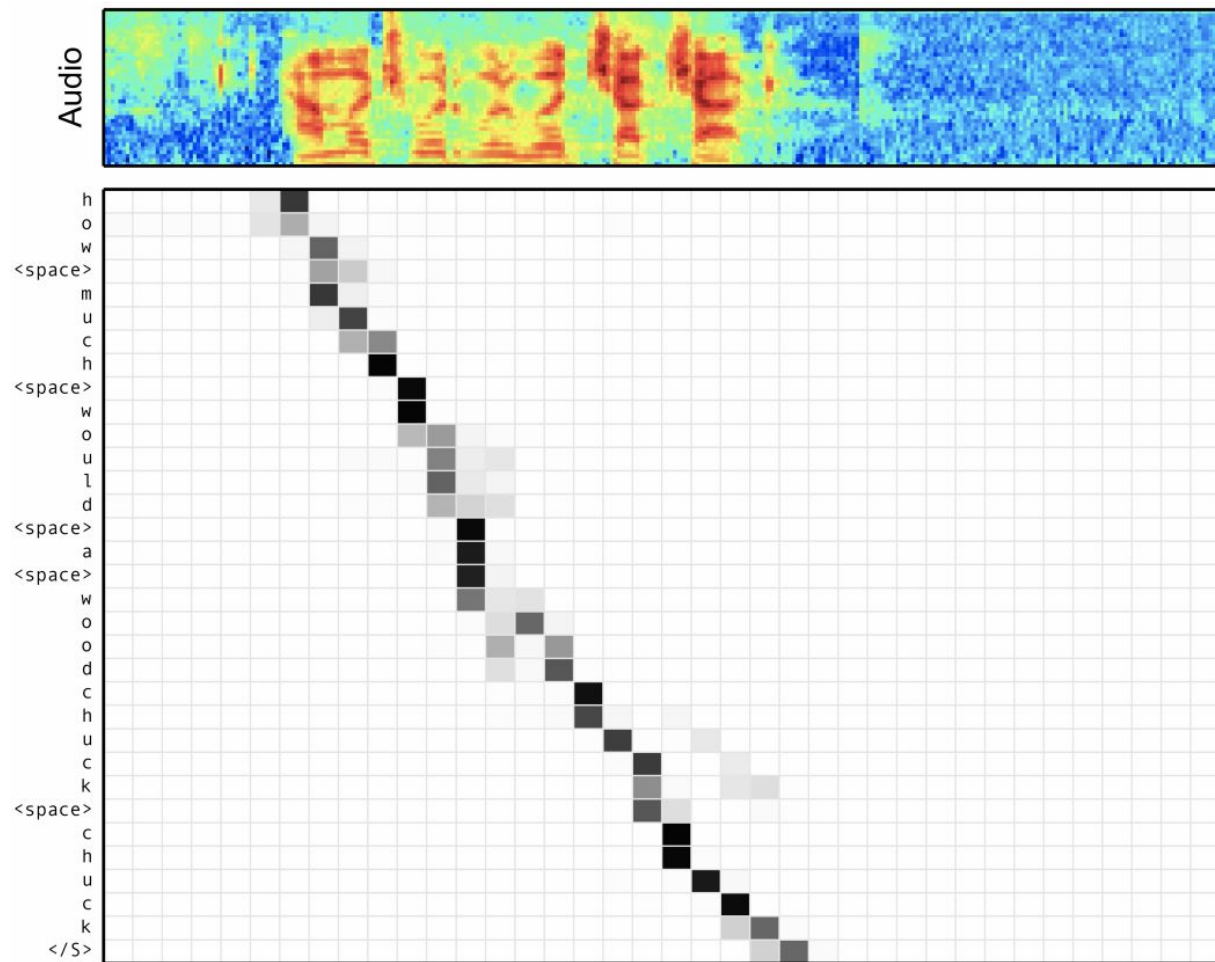


Figures from Olah & Carter, "Attention and Augmented Recurrent Neural Networks"

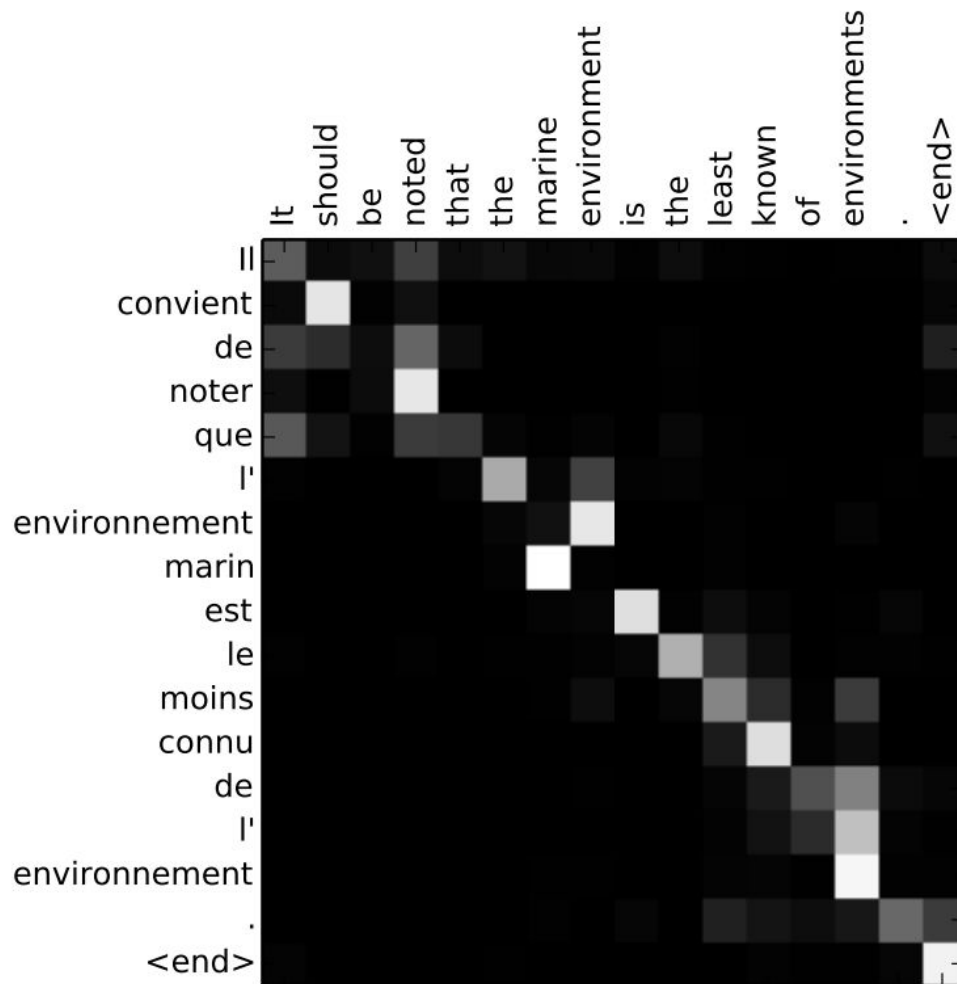


Bahdanau, Cho & Bengio, "Neural Machine Translation by Jointly Learning to Align and Translate"





Chan, Jaitly, Le & Vinyals, *"Listen, Attend and Spell"*





A(0.98)



woman(0.54)



is(0.37)



throwing(0.33)



a(0.28)



frisbee(0.37)



in(0.21)



a(0.18)

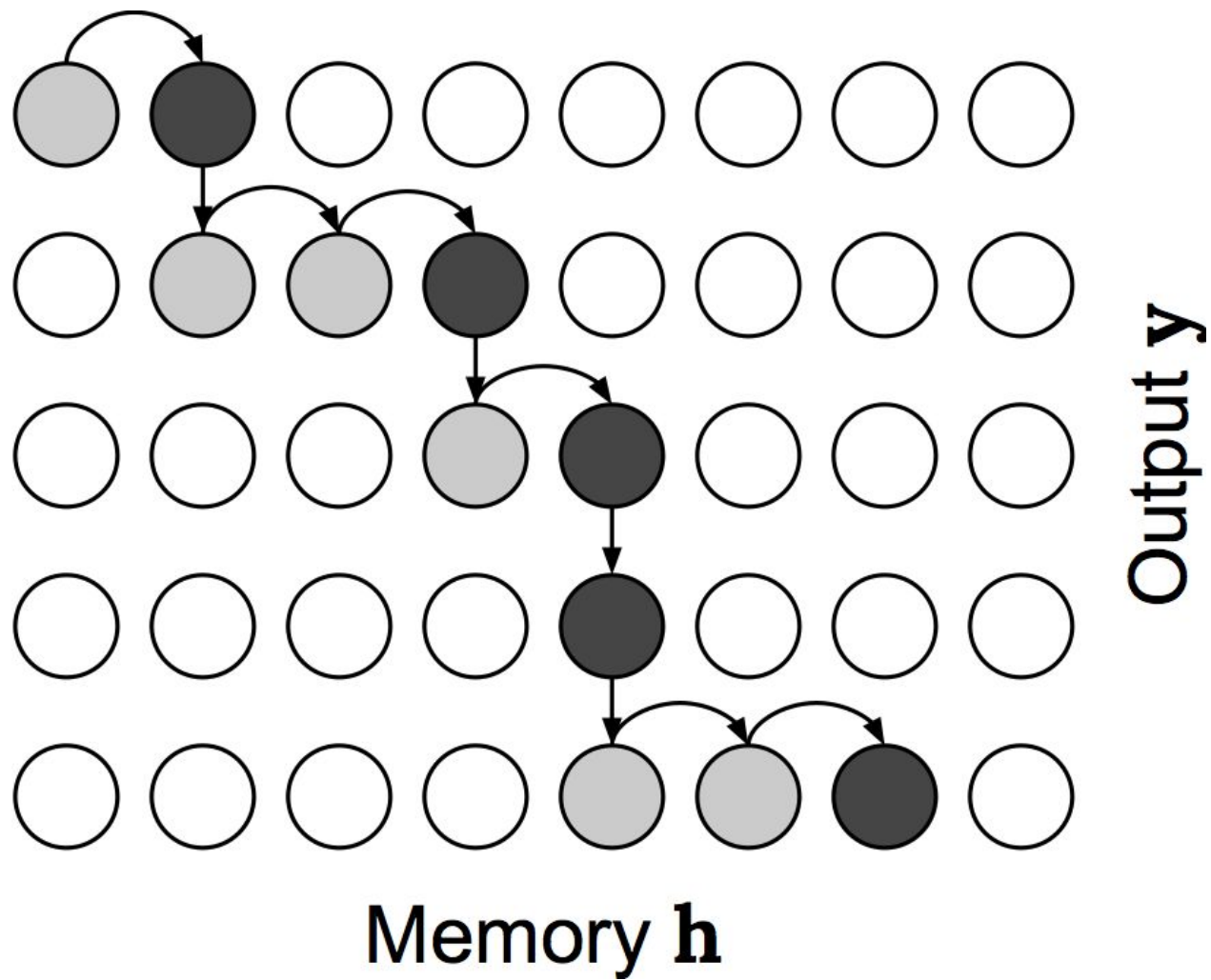


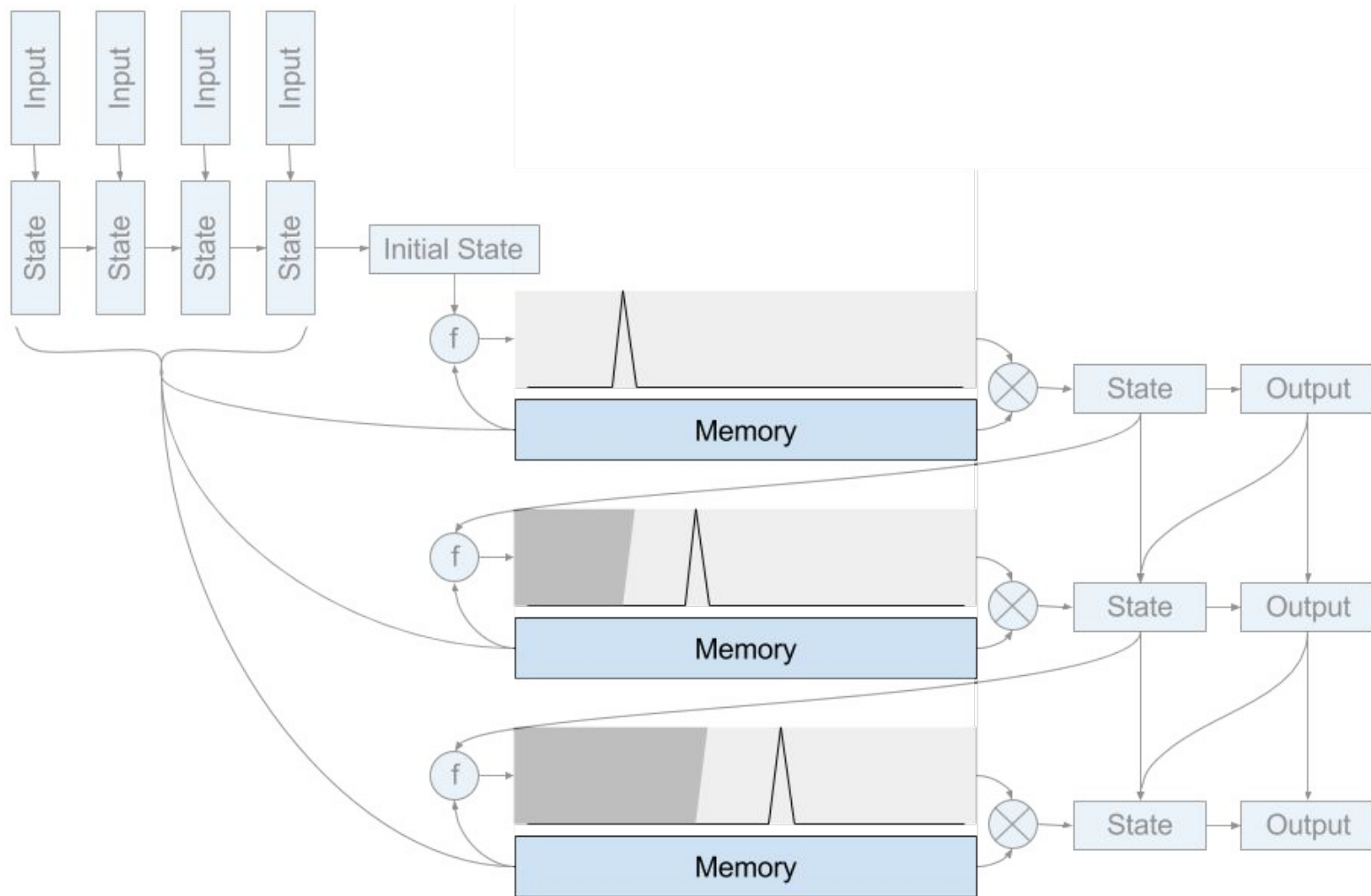
park(0.35)

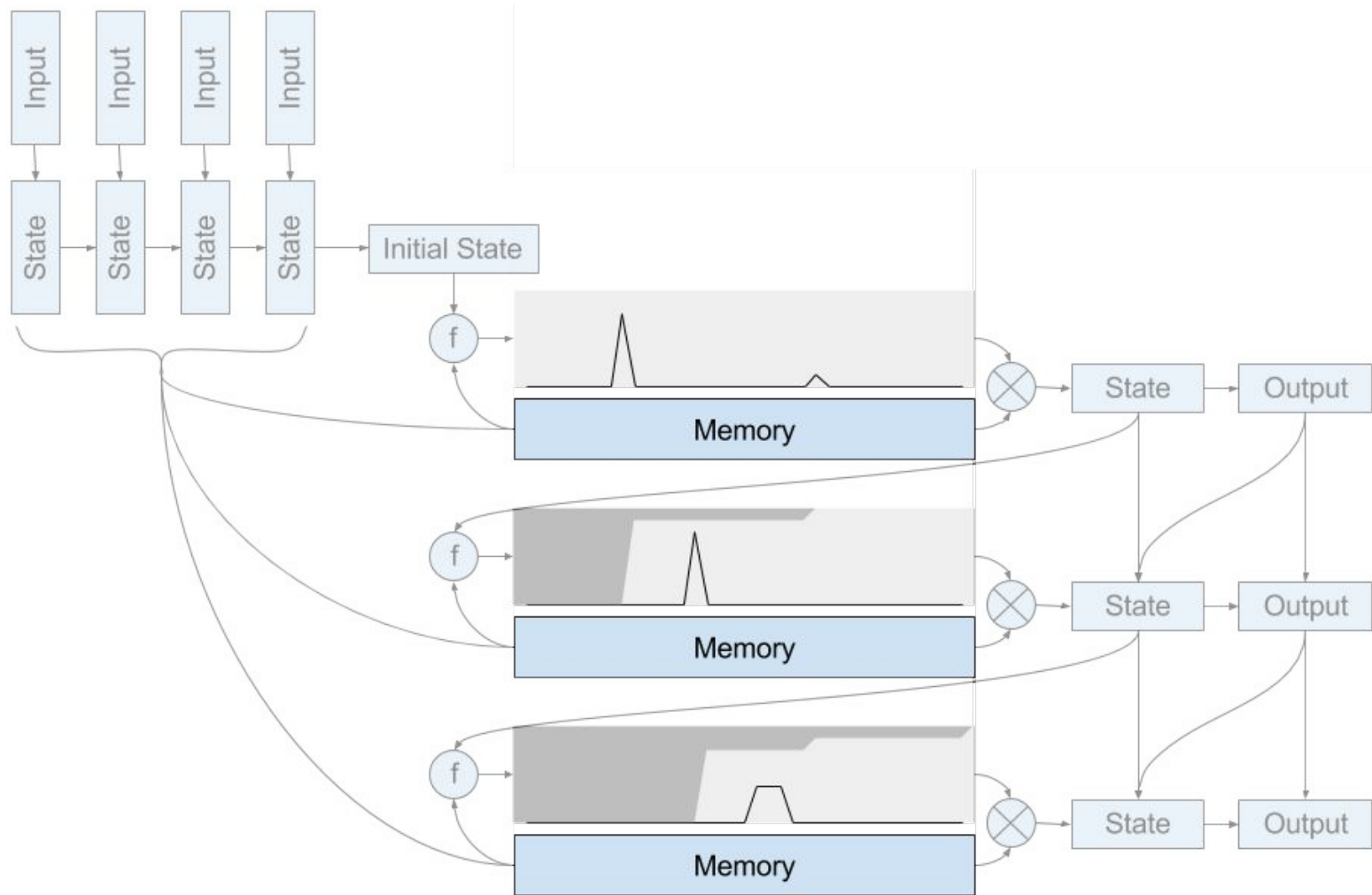


.(0.33)









$$\text{energy}_{i,j} = \text{EnergyFunction}(\text{state}_{i-1}, \text{memory}_j)$$

$$\text{attention}_{i,j} = \exp(\text{energy}_{i,j}) / \sum_{k=1}^T \exp(\text{energy}_{i,k})$$

$$\text{energy}_{i,j} = \text{EnergyFunction}(\text{state}_{i-1}, \text{memory}_j)$$

$$\text{select}_{i,j} = \sigma(\text{energy}_{i,j})$$

$$\text{attention}_{i,j} = \text{select}_{i,j} \sum_{k=1}^j \left(\text{attention}_{i-1,k} \prod_{l=k}^{j-1} (1 - \text{select}_{i,l}) \right)$$

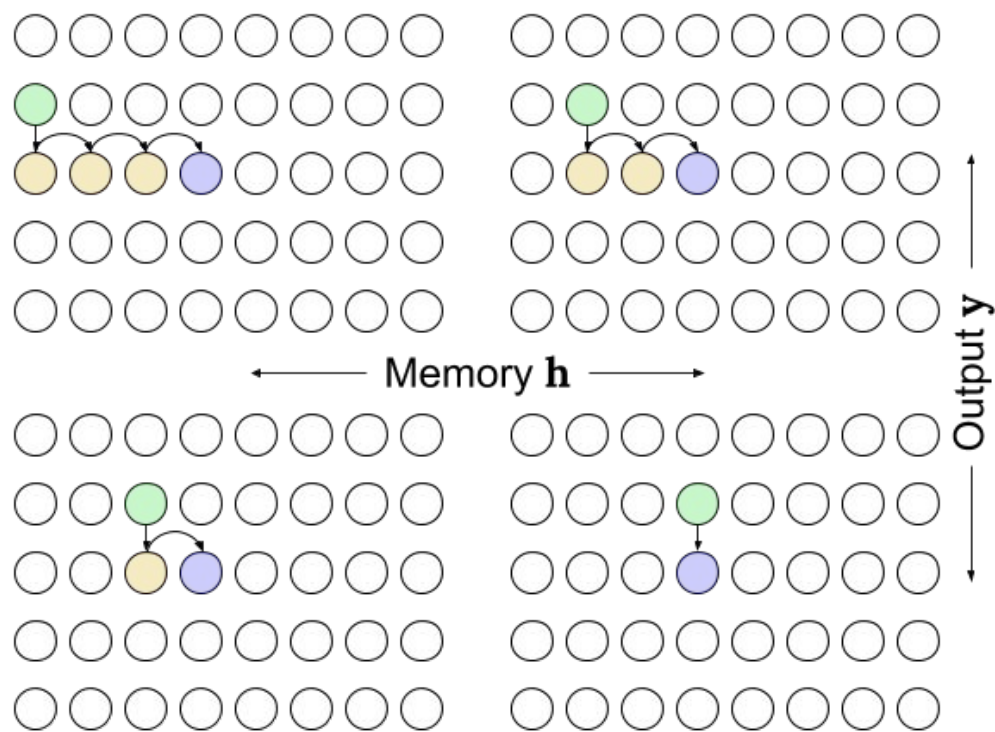
$$\text{energy}_{i,j} = \text{EnergyFunction}(\text{state}_{i-1}, \text{memory}_j)$$

$$\text{attention}_{i,j} = \exp(\text{energy}_{i,j}) / \sum_{k=1}^T \exp(\text{energy}_{i,k})$$

$$\text{energy}_{i,j} = \text{EnergyFunction}(\text{state}_{i-1}, \text{memory}_j)$$

$$\text{select}_{i,j} = \sigma(\text{energy}_{i,j})$$

$$\text{attention}_{i,j} = \text{select}_{i,j} \sum_{k=1}^j \left(\text{attention}_{i-1,k} \prod_{l=k}^{j-1} (1 - \text{select}_{i,l}) \right)$$



$$\text{attention}_{i,j} = \text{select}_{i,j} \sum_{k=1}^j \left(\text{attention}_{i-1,k} \prod_{l=k}^{j-1} (1 - \text{select}_{i,l}) \right)$$

$$\text{attention}_{i,j} = \text{select}_{i,j} \sum_{k=1}^j \left(\text{attention}_{i-1,k} \prod_{l=k}^{j-1} (1 - \text{select}_{i,l}) \right)$$

$$\text{attention}_{i,j} = \text{select}_{i,j} \left((1 - \text{select}_{i,j-1}) \frac{\text{attention}_{i,j-1}}{\text{select}_{i,j-1}} + \text{attention}_{i-1,j} \right)$$

$$\text{attention}_i = \text{select}_i \text{cumprod}(1 - \text{select}_i) \text{cumsum} \left(\frac{\text{attention}_{i-1}}{\text{cumprod}(1 - \text{select}_i)} \right)$$

$$\text{attention}_{i,j} = \text{select}_{i,j} \sum_{k=1}^j \left(\text{attention}_{i-1,k} \prod_{l=k}^{j-1} (1 - \text{select}_{i,l}) \right)$$

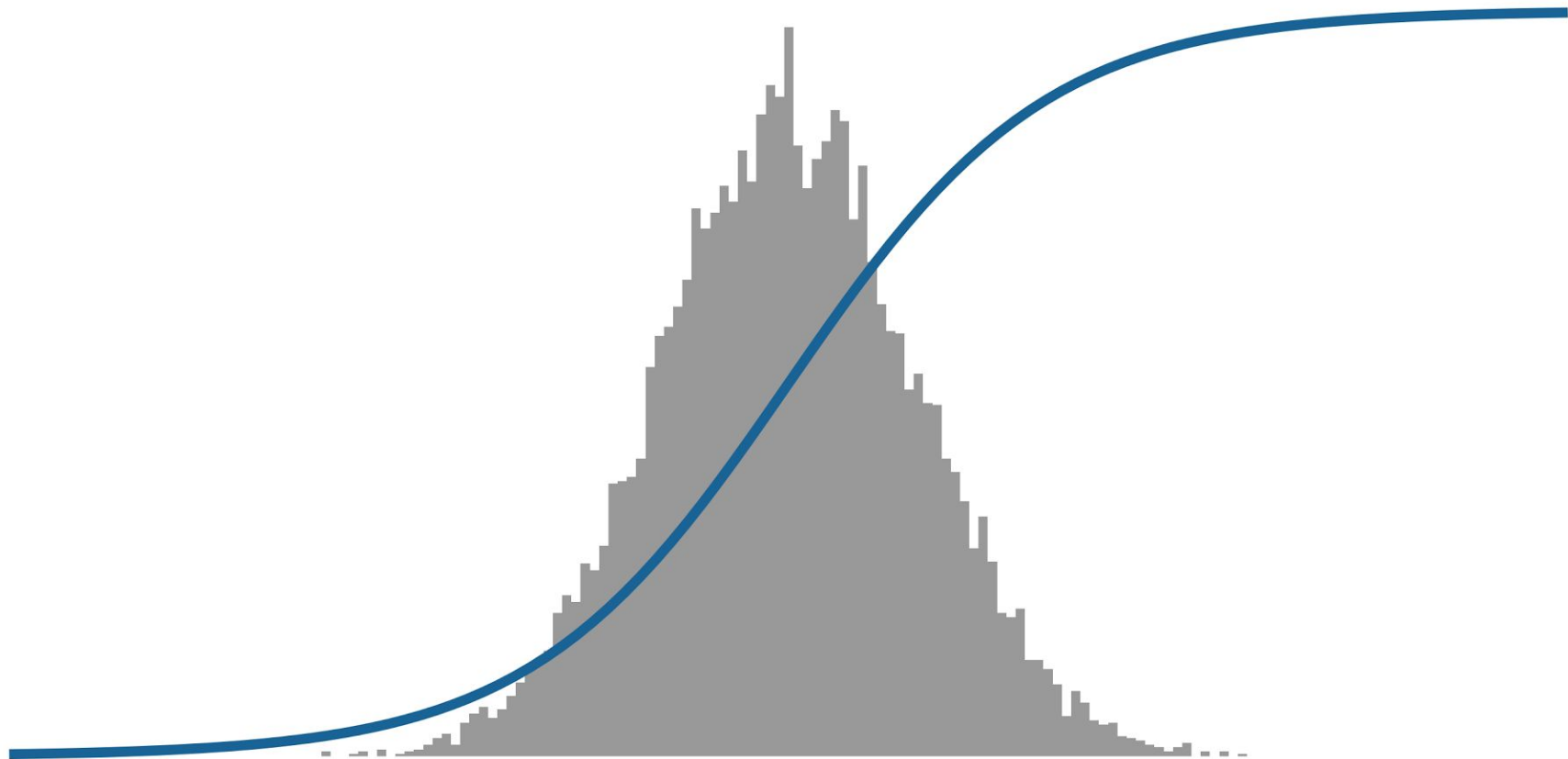
$$\text{attention}_{i,j} = \text{select}_{i,j} \left((1 - \text{select}_{i,j-1}) \frac{\text{attention}_{i,j-1}}{\text{select}_{i,j-1}} + \text{attention}_{i-1,j} \right)$$

$$\text{attention}_i = \text{select}_i \text{cumprod}(1 - \text{select}_i) \text{cumsum} \left(\frac{\text{attention}_{i-1}}{\text{cumprod}(1 - \text{select}_i)} \right)$$

$$\text{attention}_{i,j} = \text{select}_{i,j} \sum_{k=1}^j \left(\text{attention}_{i-1,k} \prod_{l=k}^{j-1} (1 - \text{select}_{i,l}) \right)$$

$$\text{attention}_{i,j} = \text{select}_{i,j} \left((1 - \text{select}_{i,j-1}) \frac{\text{attention}_{i,j-1}}{\text{select}_{i,j-1}} + \text{attention}_{i-1,j} \right)$$

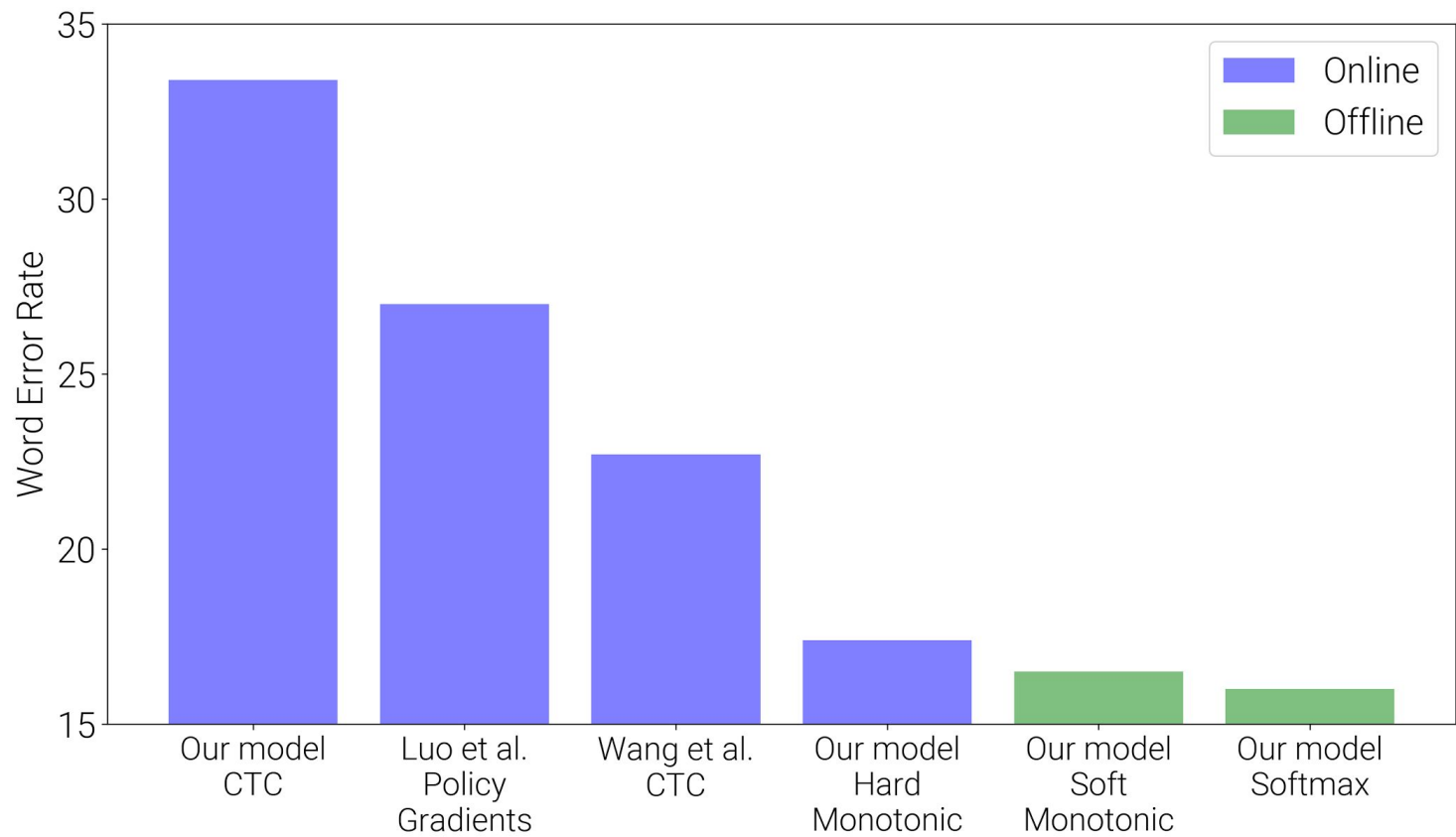
$$\text{attention}_i = \text{select}_i \text{ cumprod}(1 - \text{select}_i) \text{ cumsum} \left(\frac{\text{attention}_{i-1}}{\text{cumprod}(1 - \text{select}_i)} \right)$$



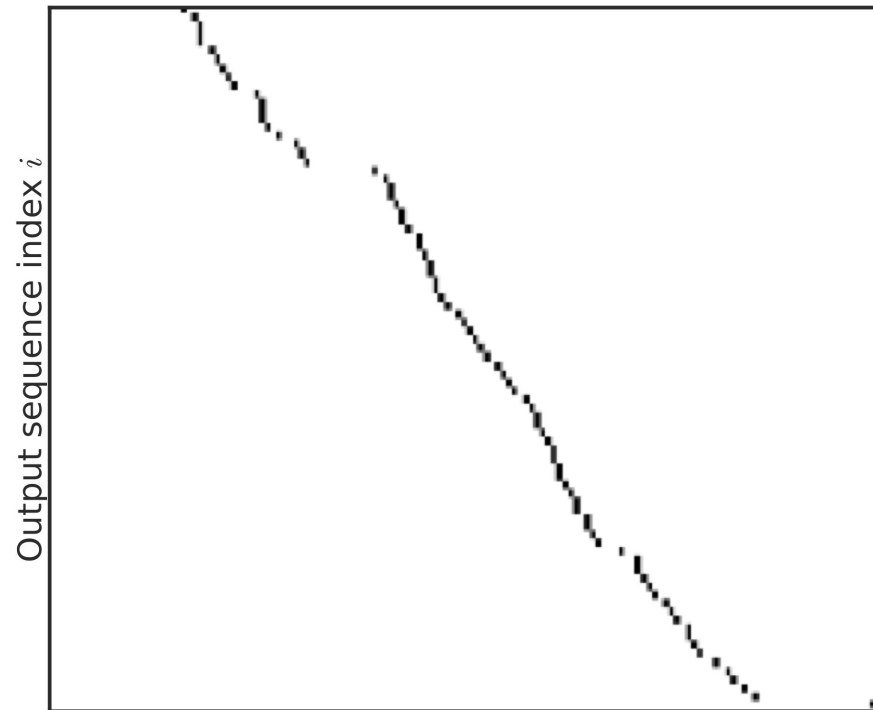
Frey, *"Continuous Sigmoidal Belief Networks Trained Using Slice Sampling"*

Salakhutdinov & Hinton, *"Semantic Hashing"*

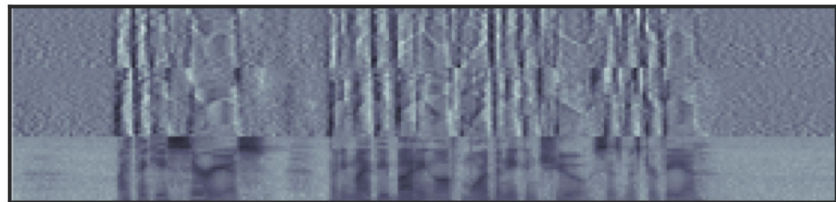
Foerster, Assael, de Freitas & Whiteson, *"Learning to Communicate with Deep Multi-Agent Reinforcement Learning"*



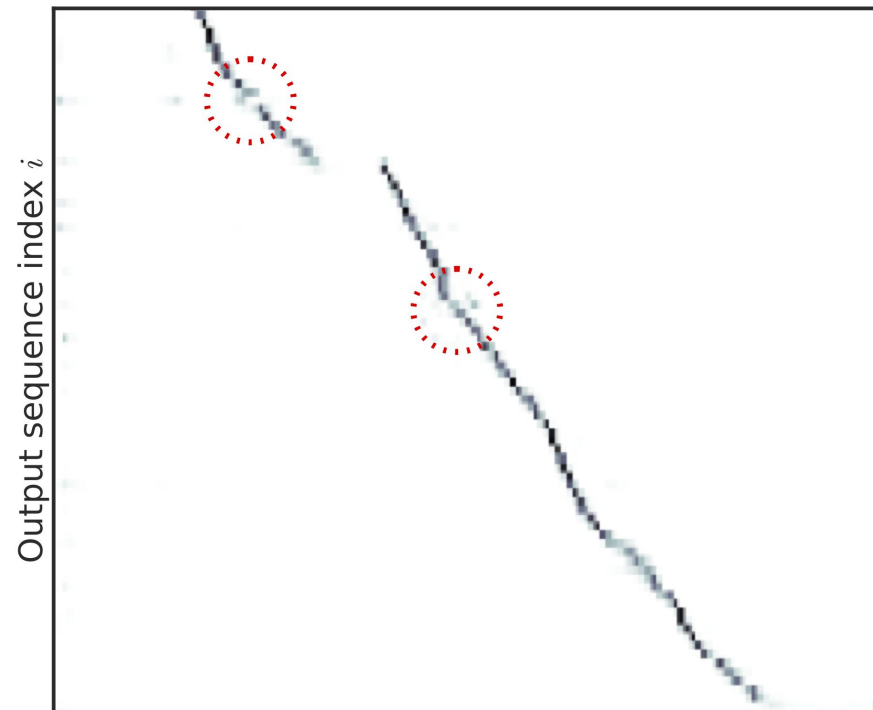
Hard Monotonic Attention



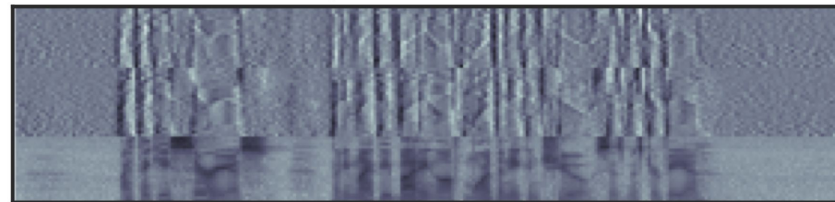
Input Features



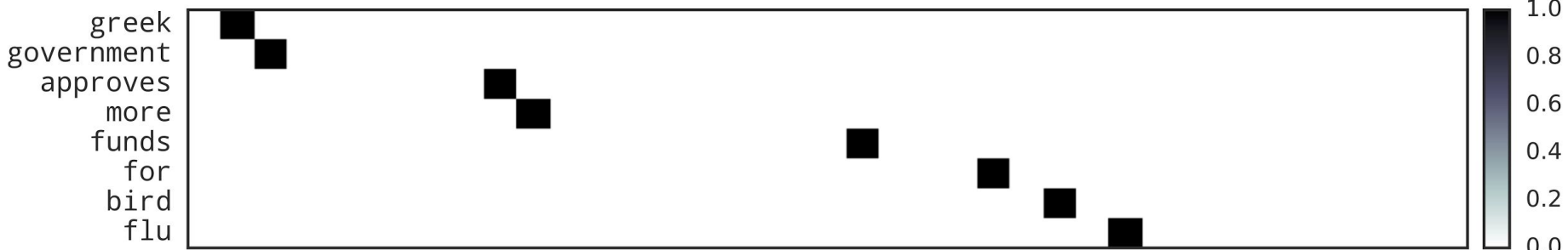
Softmax Attention



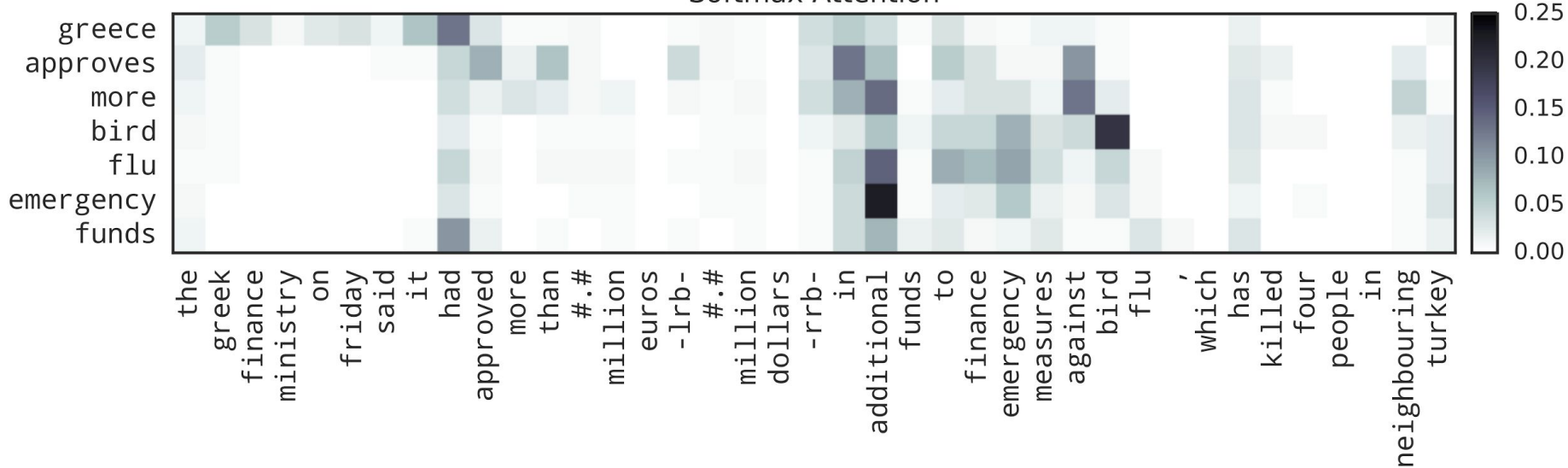
Input Features



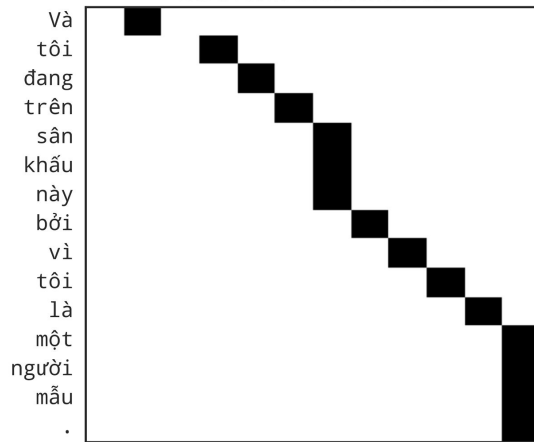
Hard Monotonic Attention



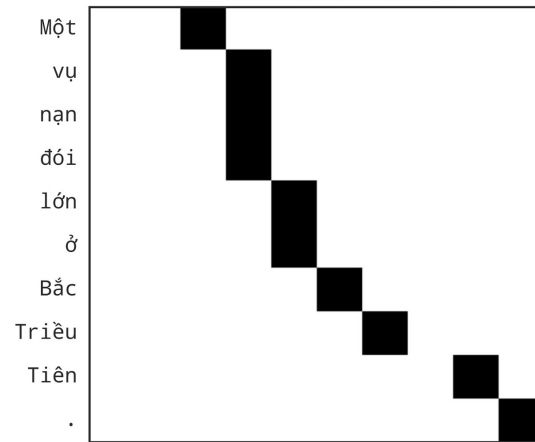
Softmax Attention



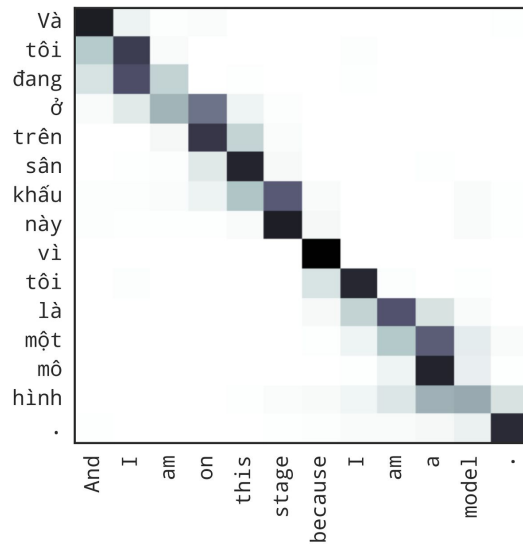
Hard Monotonic Attention



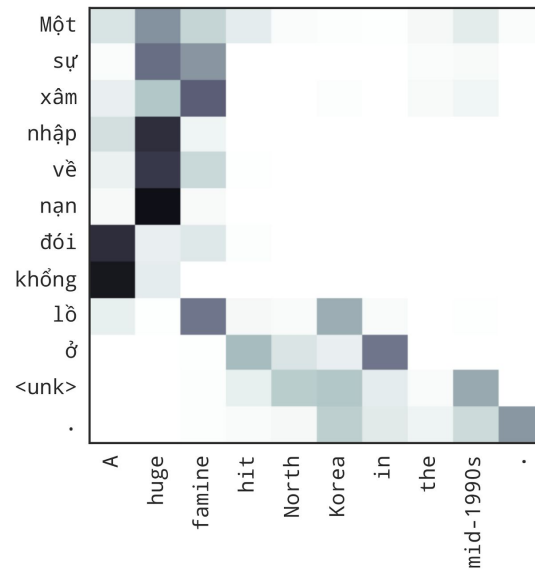
Hard Monotonic Attention



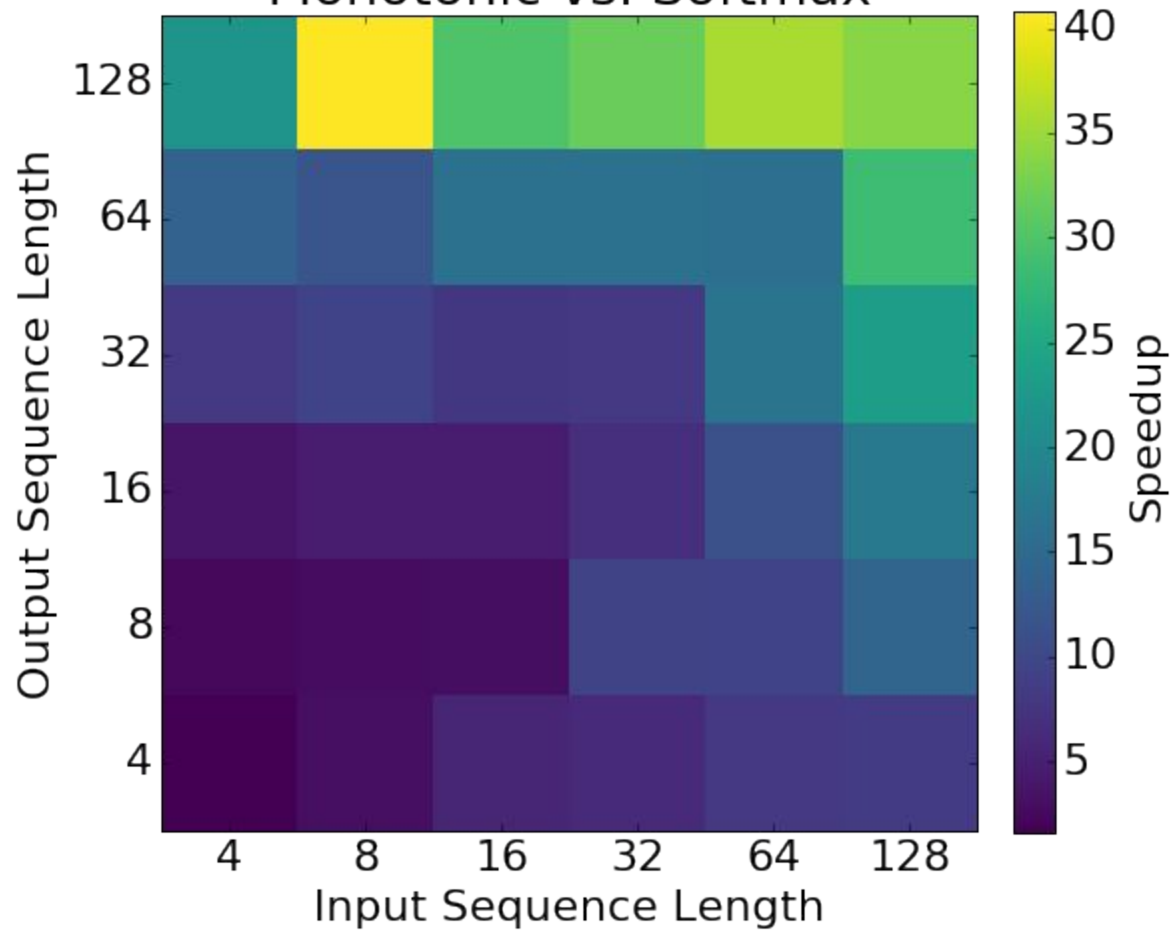
Softmax Attention



Softmax Attention



Monotonic vs. Softmax



Pointers

Implemented in **tf.contrib.seq2seq**

Additional code at <http://github.com/craffel/mad>

“Practitioner’s Guide” in Appendix G

Blog post: <http://colinraffel.com/blog>

These slides: <http://colinraffel.com/talks/icml2017online.pdf>

Poster #70 tonight

My email: craffel@gmail.com

Thanks!