

Attention!

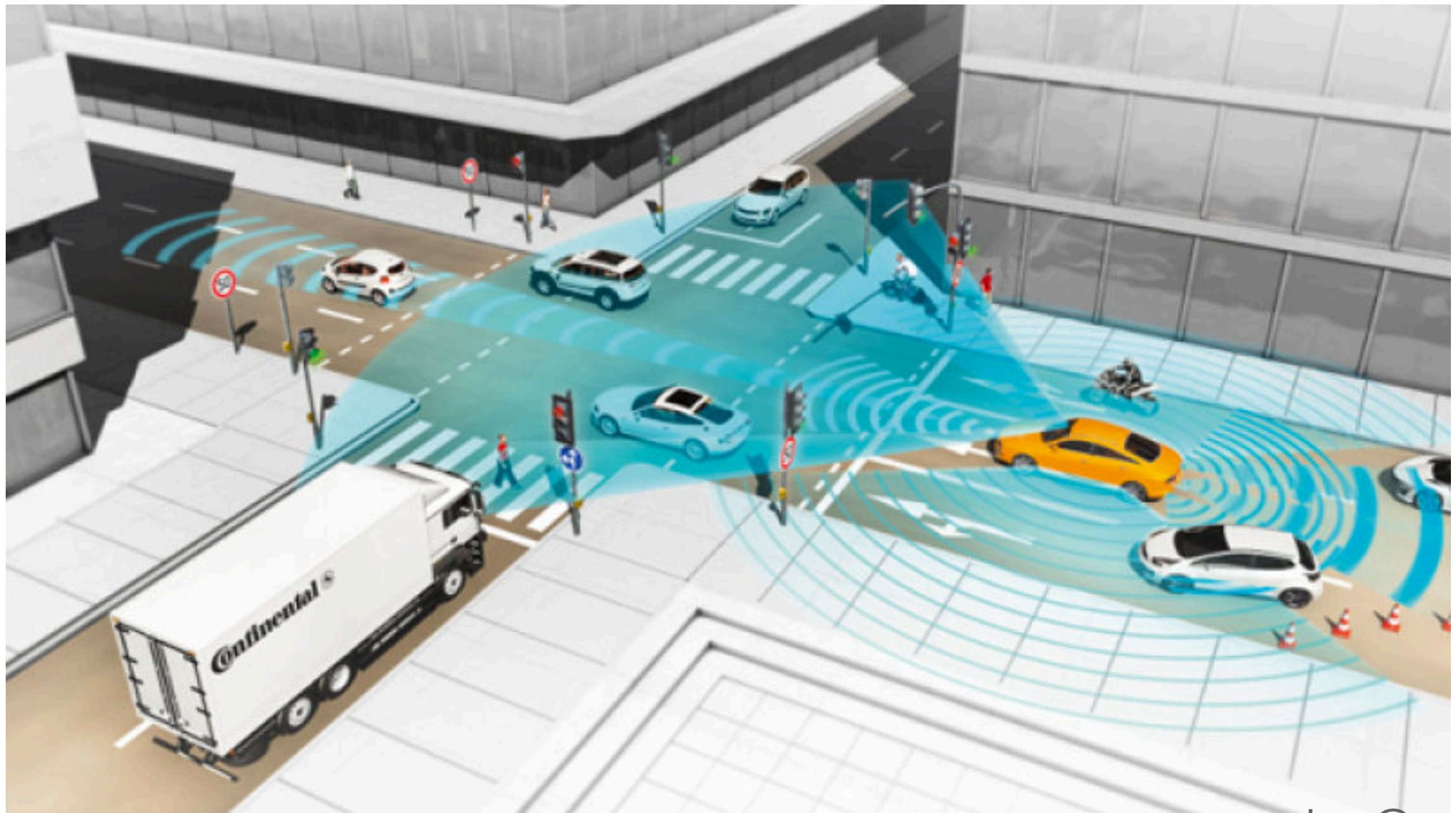
...and its use in Deep Learning

Anders Huss

Stockholm AI - Reading Group
2017-08-23

What about it?

Self driving cars



What about it?

Machine Translation

“narrow down the focus to the most relevant parts”

“recentrez l'attention sur les éléments les plus pertinents”

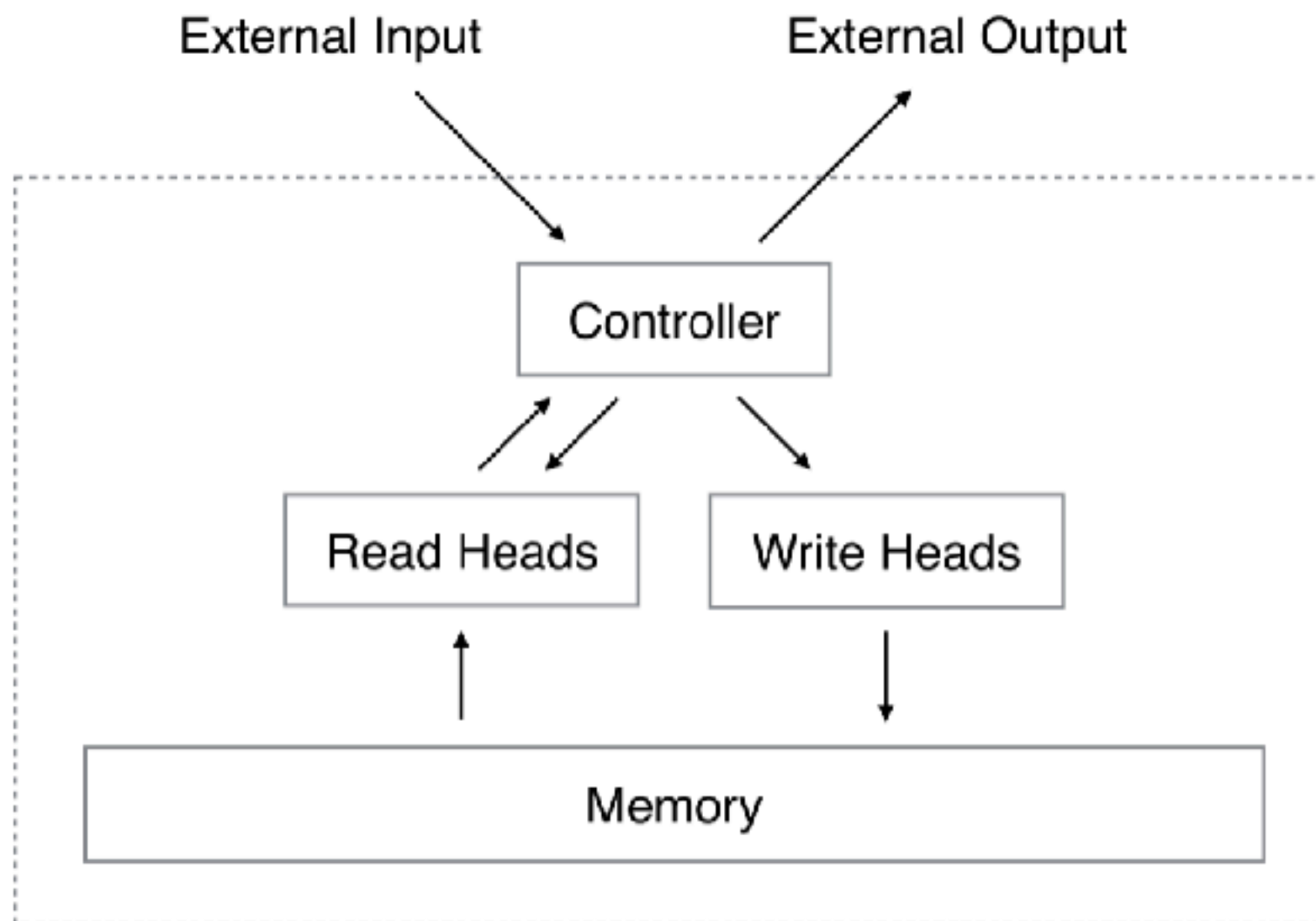
What about it?

Solving Problems (Logical Reasoning)

*“many times complexity
is a collection of
connected simplicities”*

What about it?

Solving Problems (Logical Reasoning)



Agenda

- Problem Formulation
- Spatial Attention
 - Example/Benchmark
 - Applications
- Beyond Spatial Attention
 - Differentiable Neural Computers
- Discussion and Questions

Iterative Application of Attention (RNN)

Regular RNN/LSTM/GRU...

$$y^t, s^t = f(x^t, s^{t-1}) \qquad y^{1:T}, s^{1:T} = RNN_f(x^{1:T}, s^0)$$

With Attention (**Z**)

$$y^t, s^t = g(x^t, s^{t-1}, \mathbf{Z}) \qquad y^{1:T}, s^{1:T} = RNN_g(x^{1:T}, s^0, \mathbf{Z})$$

E.g. by *altering input* to the regular RNN/GRU/LSTM...

$$y^t, s^t = f(x'^t, s^{t-1}), \quad x'^t = a(x^t, s^{t-1}, \mathbf{Z})$$

Attention Mechanisms

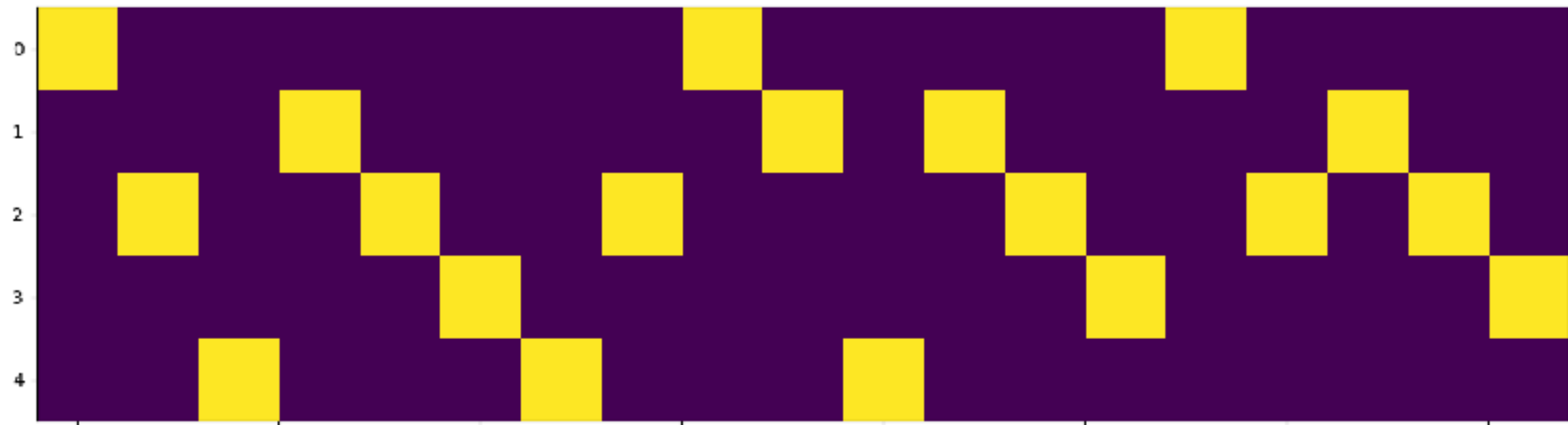
At each time step t we predict a **weighting** over input \mathbf{Z} and map it to an attention encoding z^t

$$y^t, s^t = f(x'^t, s^{t-1})$$

$$x'^t = a(x^t, s^{t-1}, \mathbf{Z}) = [x^t, z^t(x^t, s^{t-1}, \mathbf{Z})]$$

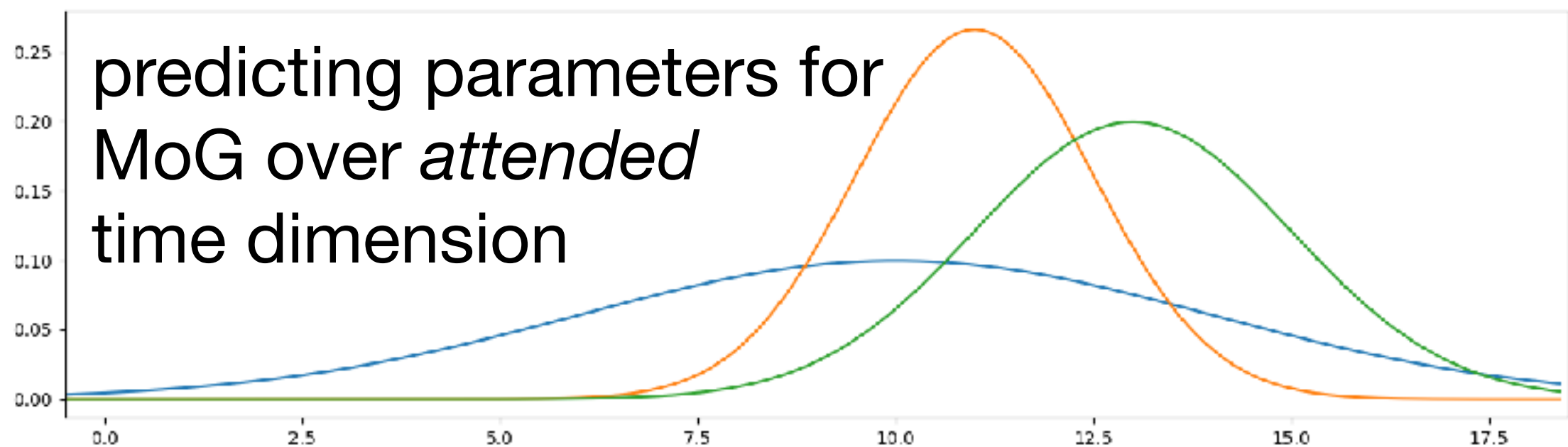
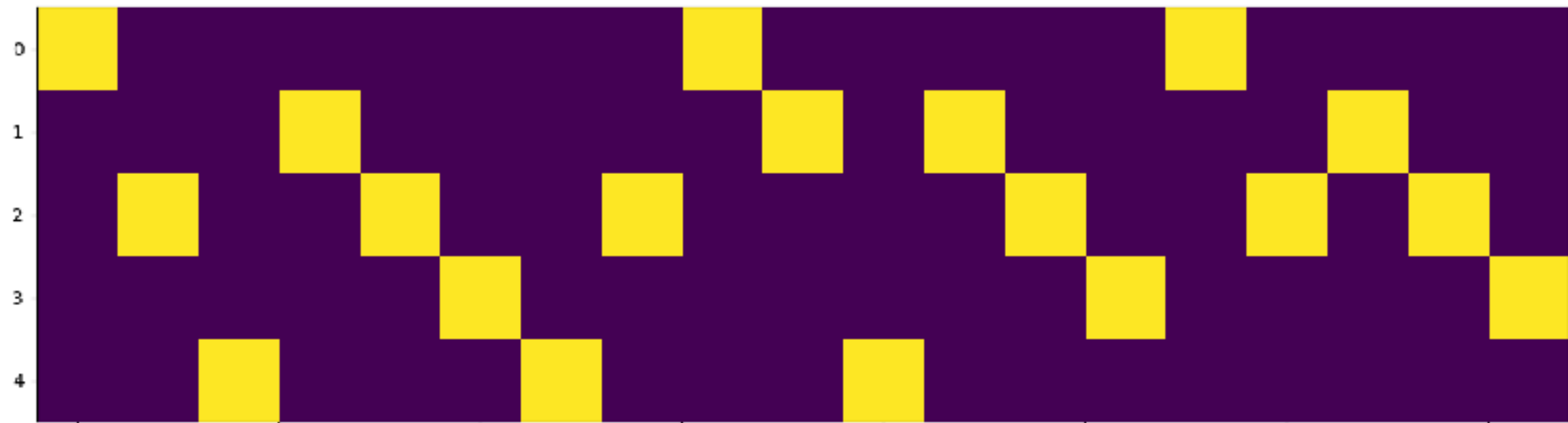
There are many options for how to do this...

Spatial Attention: 1D

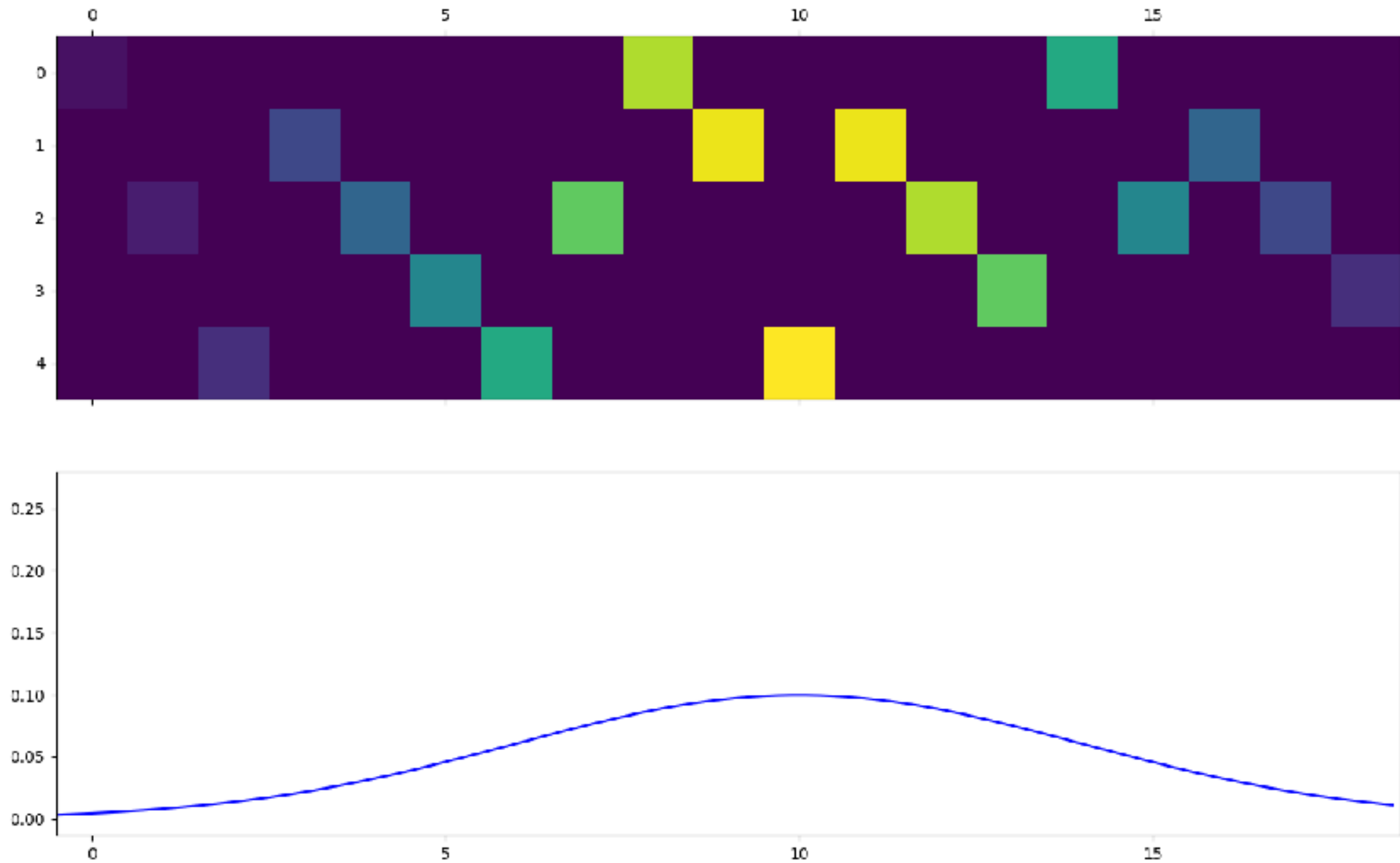


At each time step t we predict a **weighting** over input sequence \mathbf{Z} and map it to an attention encoding z^t

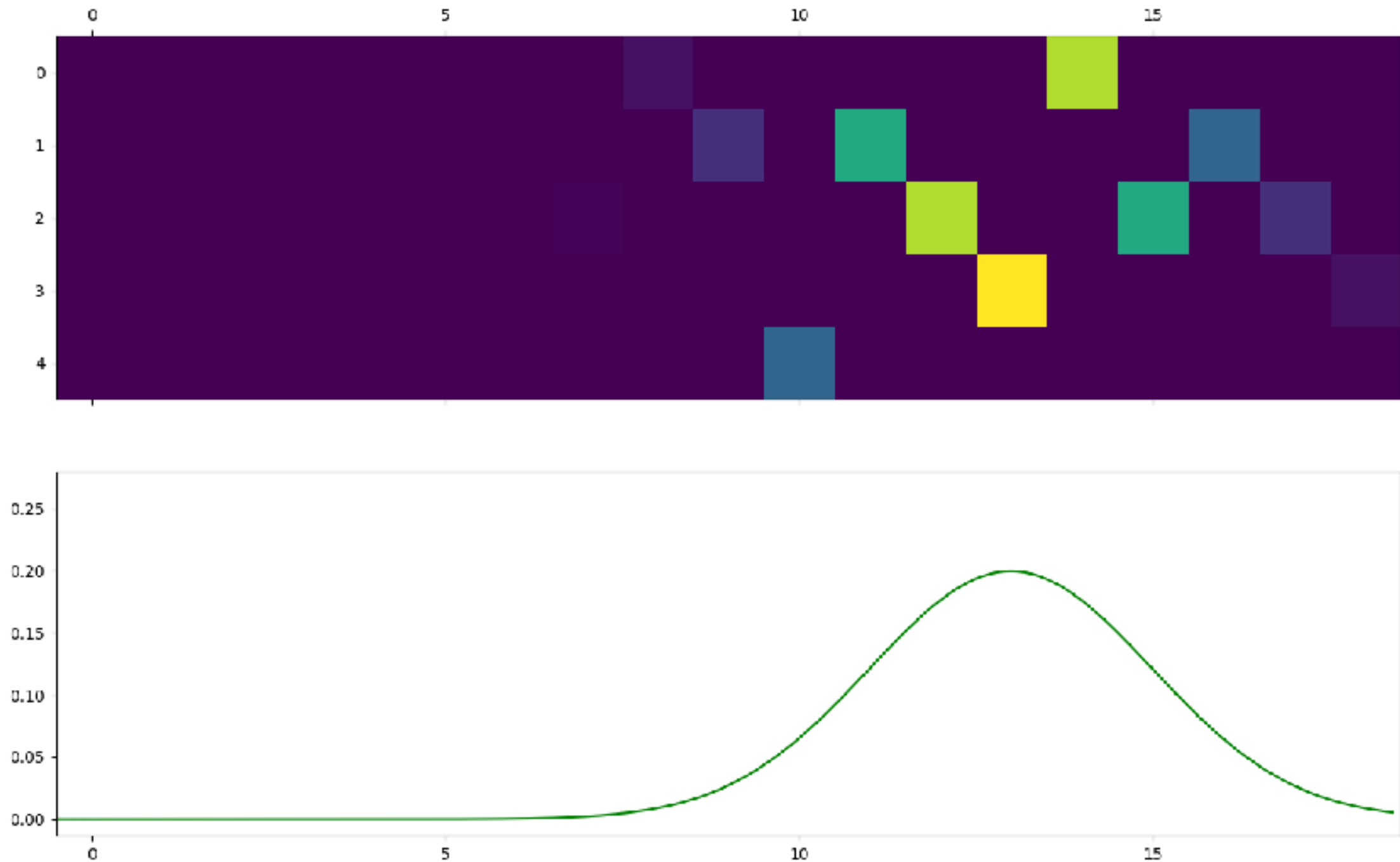
Spatial Attention: 1D



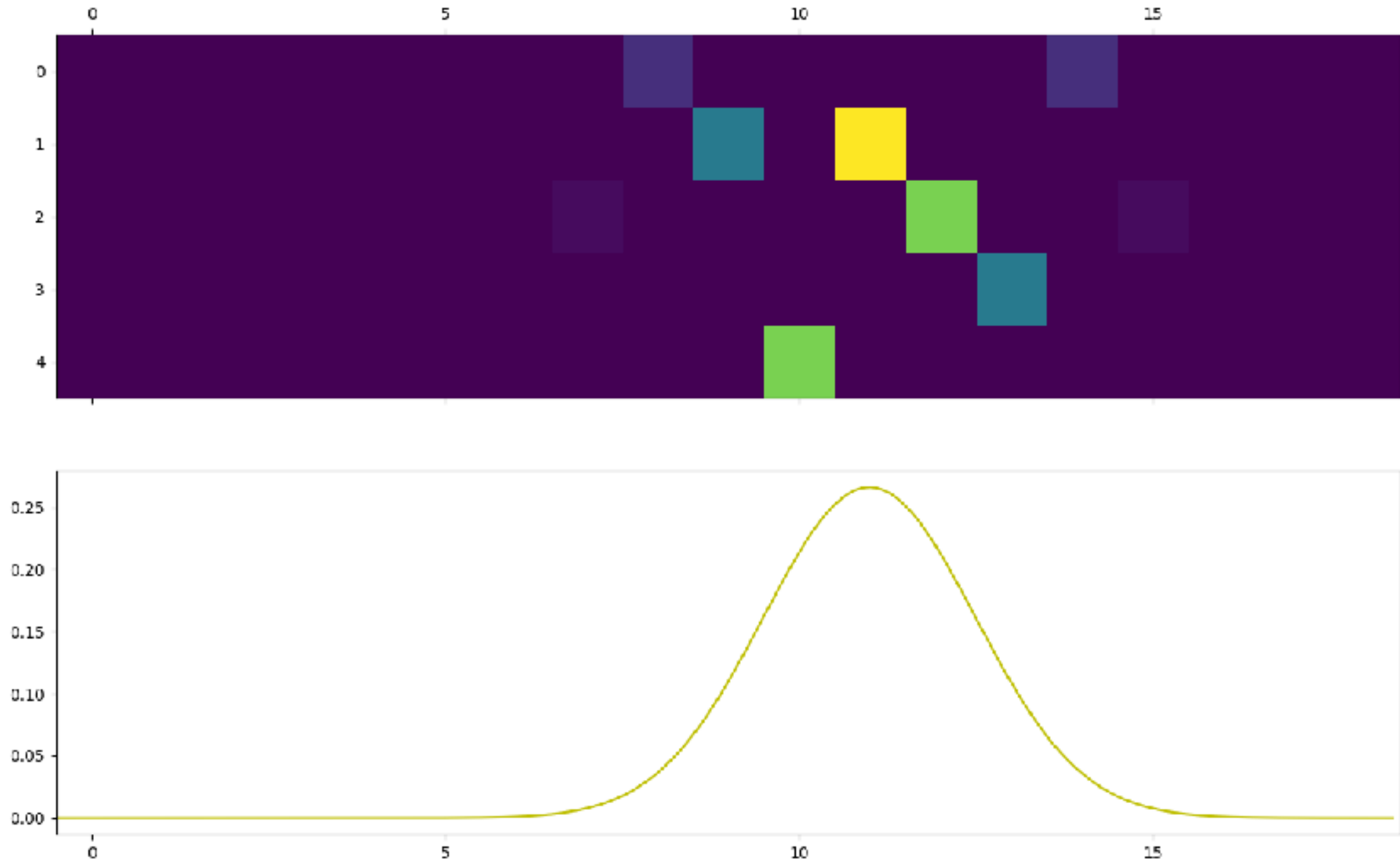
Spatial Attention: 1D



Spatial Attention: 1D

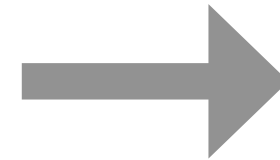
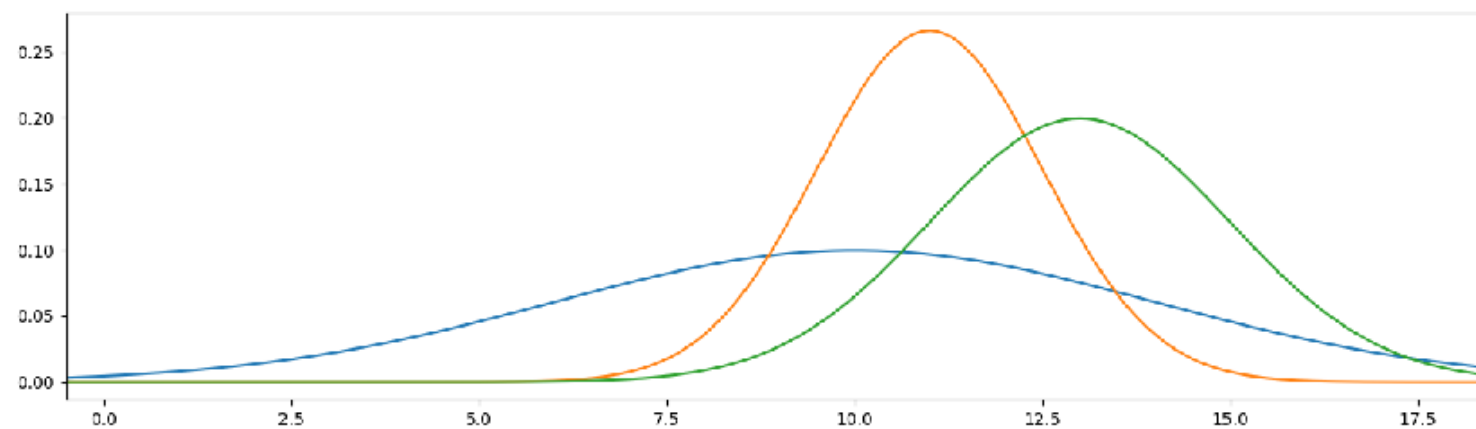
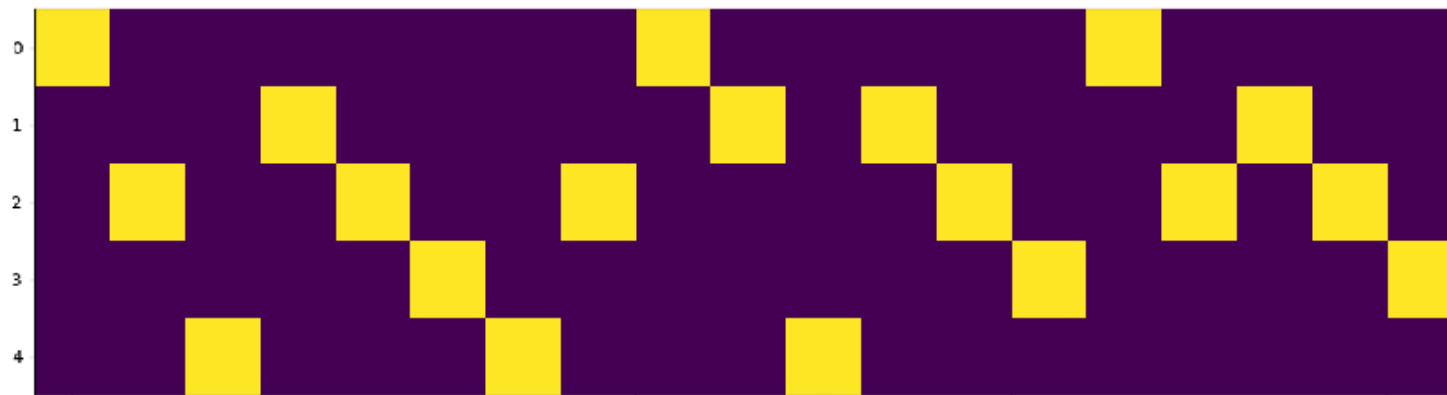


Spatial Attention: 1D



Spatial Attention: 1D

Z

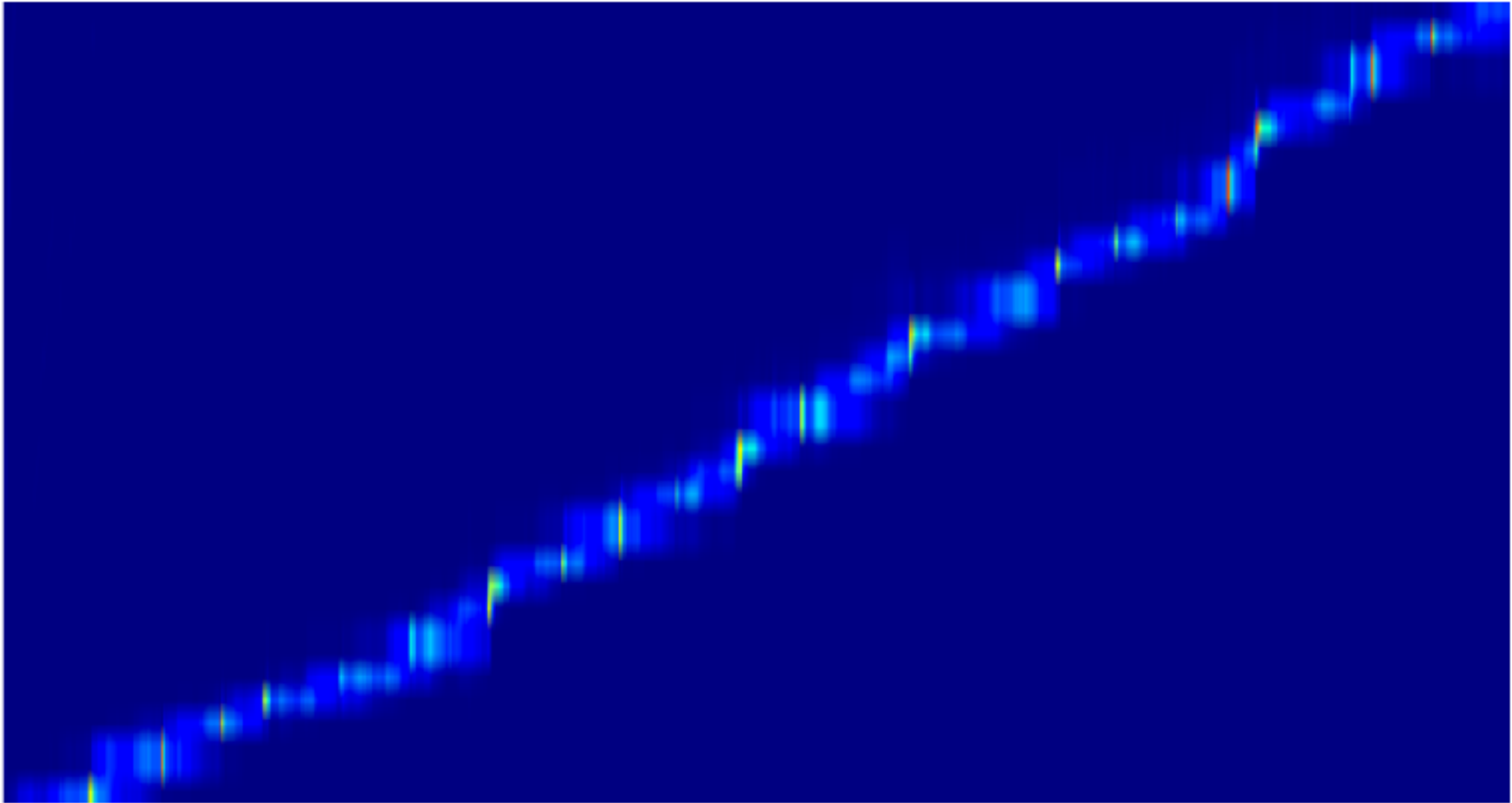


z^t



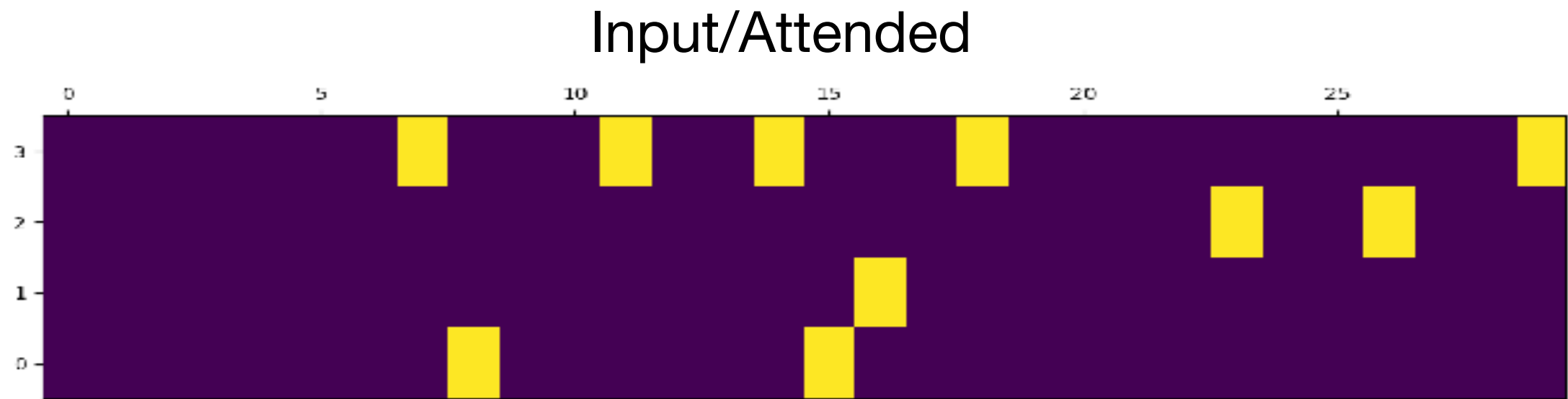
Application: Handwriting synthesis

Thought that the muster from

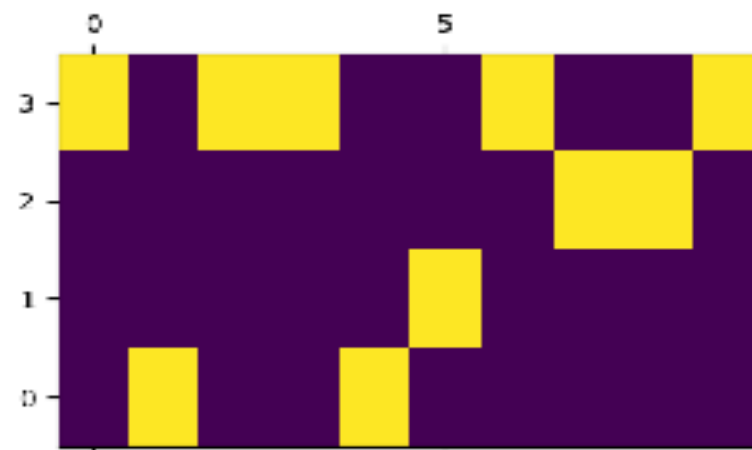


Thought that the muster from

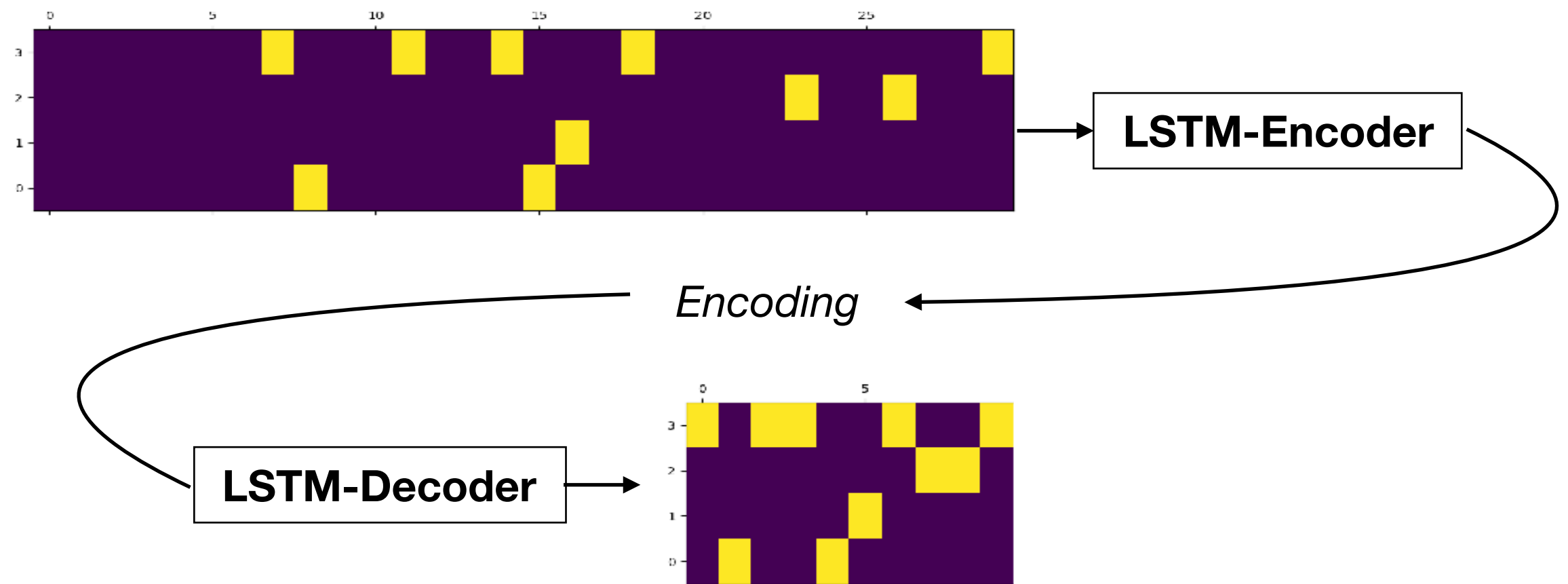
Canonical Example of Alignment/Selection Problem



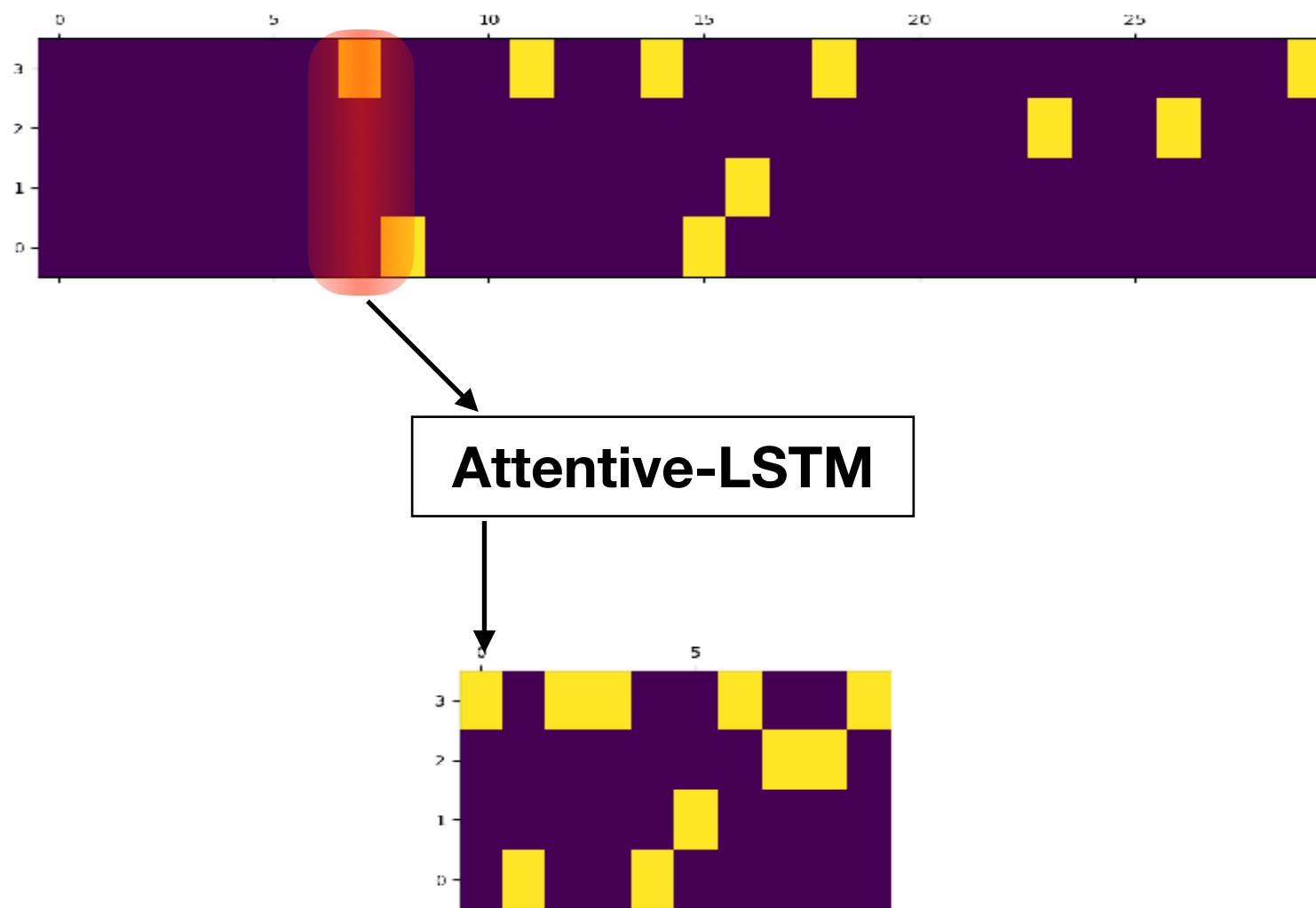
Target



Sequence-to-Sequence

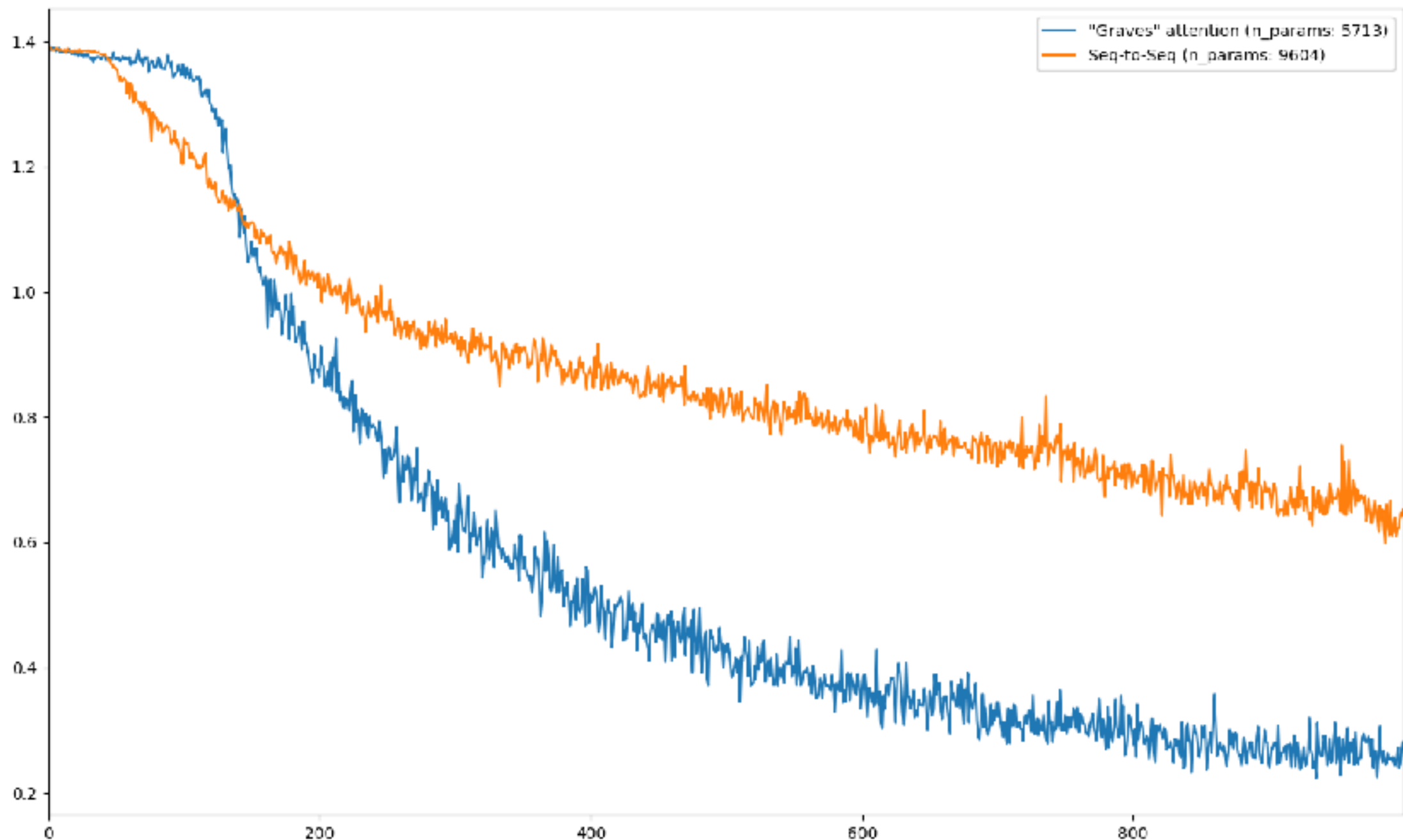


Attention



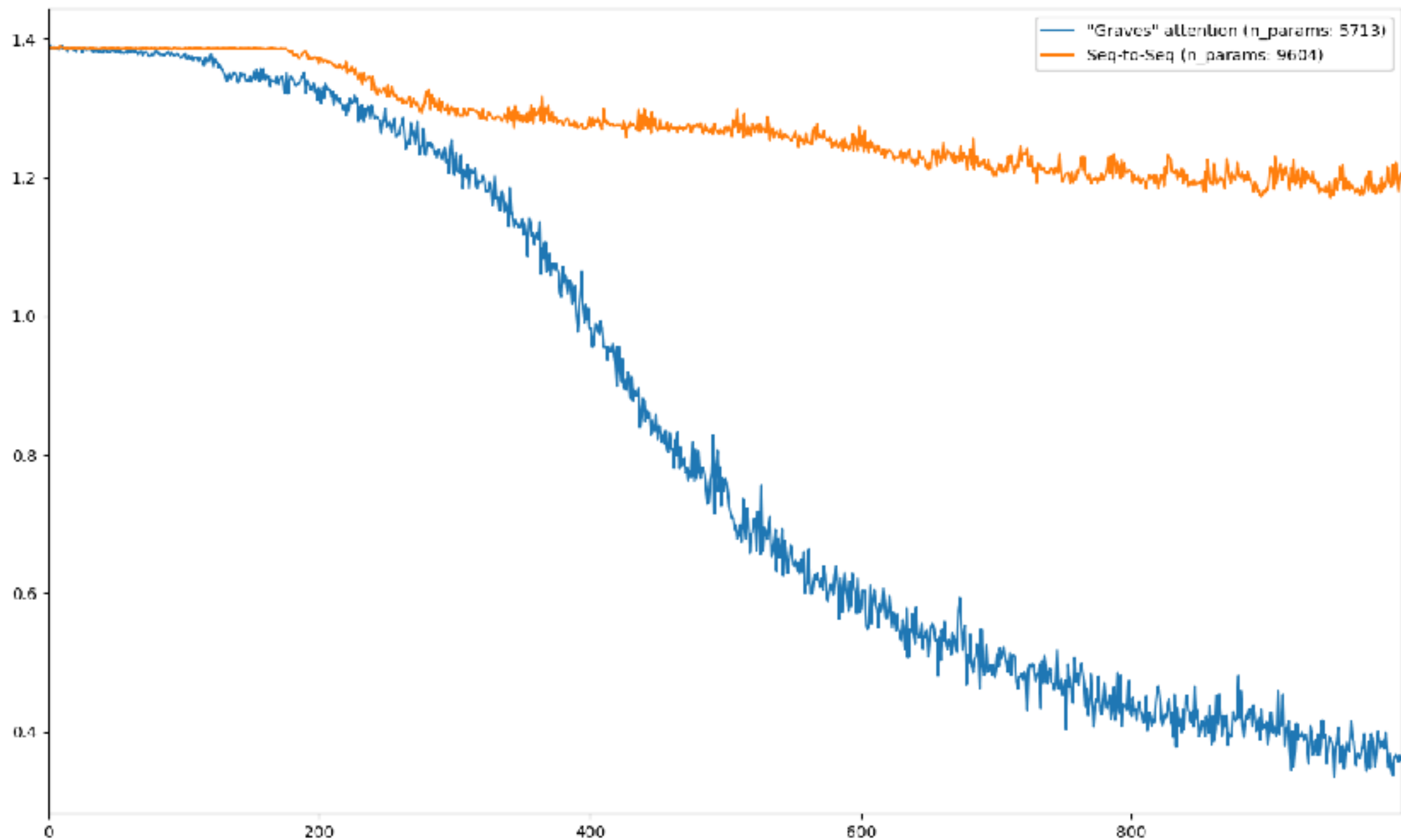
Seq-to-Seq vs Attention

attended sequence length: 30, target sequence length: 10



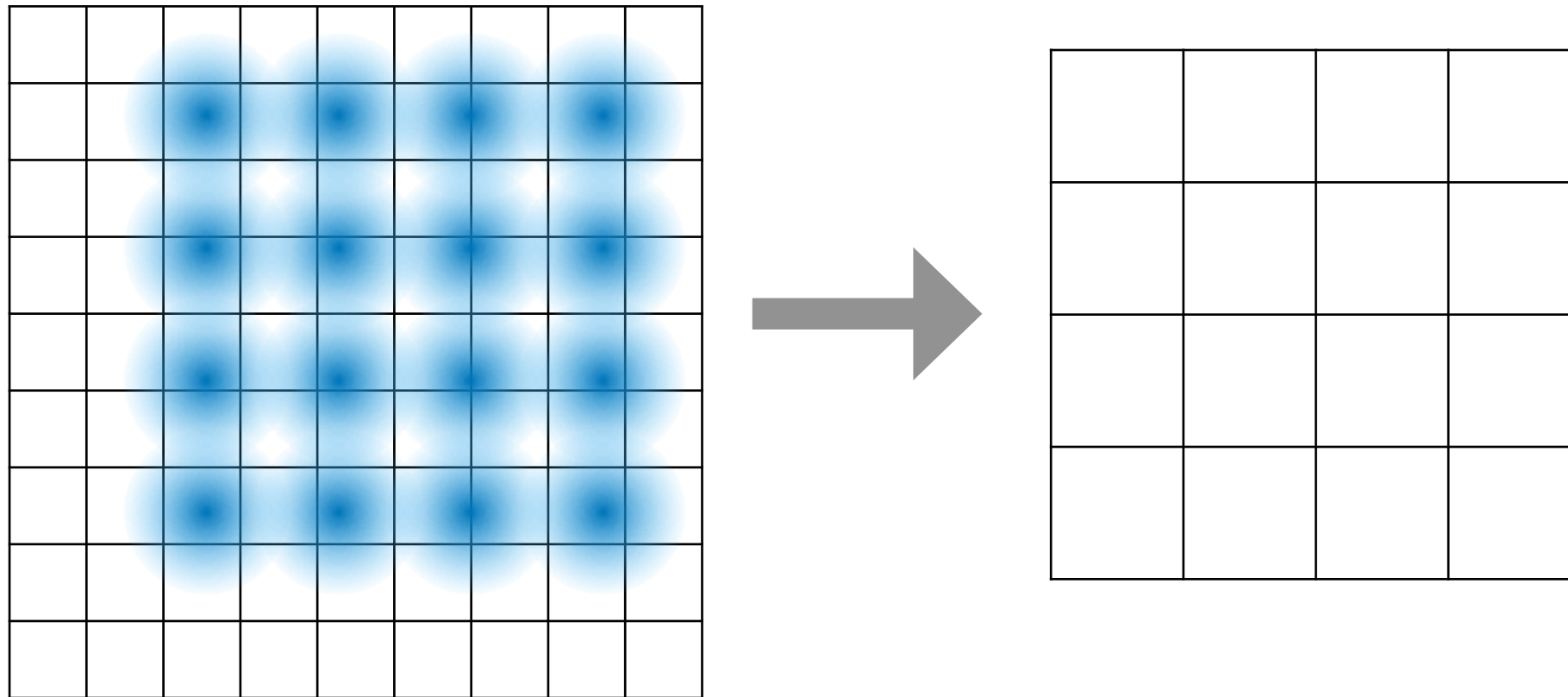
Seq-to-Seq vs Attention

attended sequence length: 60, target sequence length: 20



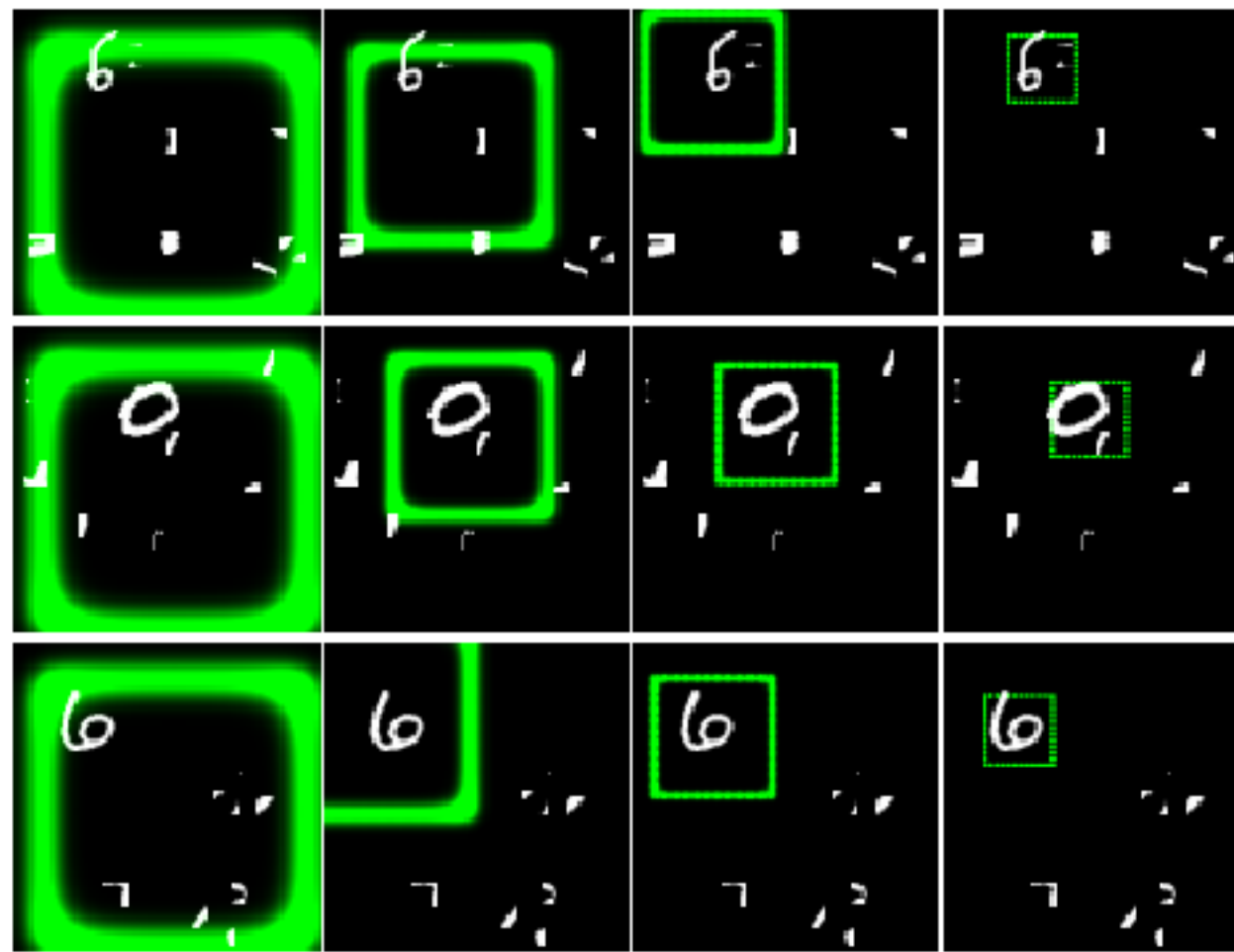
Spatial Attention: 2D

“zoom and shift” focus



Spatial Attention: 2D

“zoom and shift” focus



Spatial Attention: 2D

“zoom and shift” focus



A woman is throwing a **frisbee** in a park.

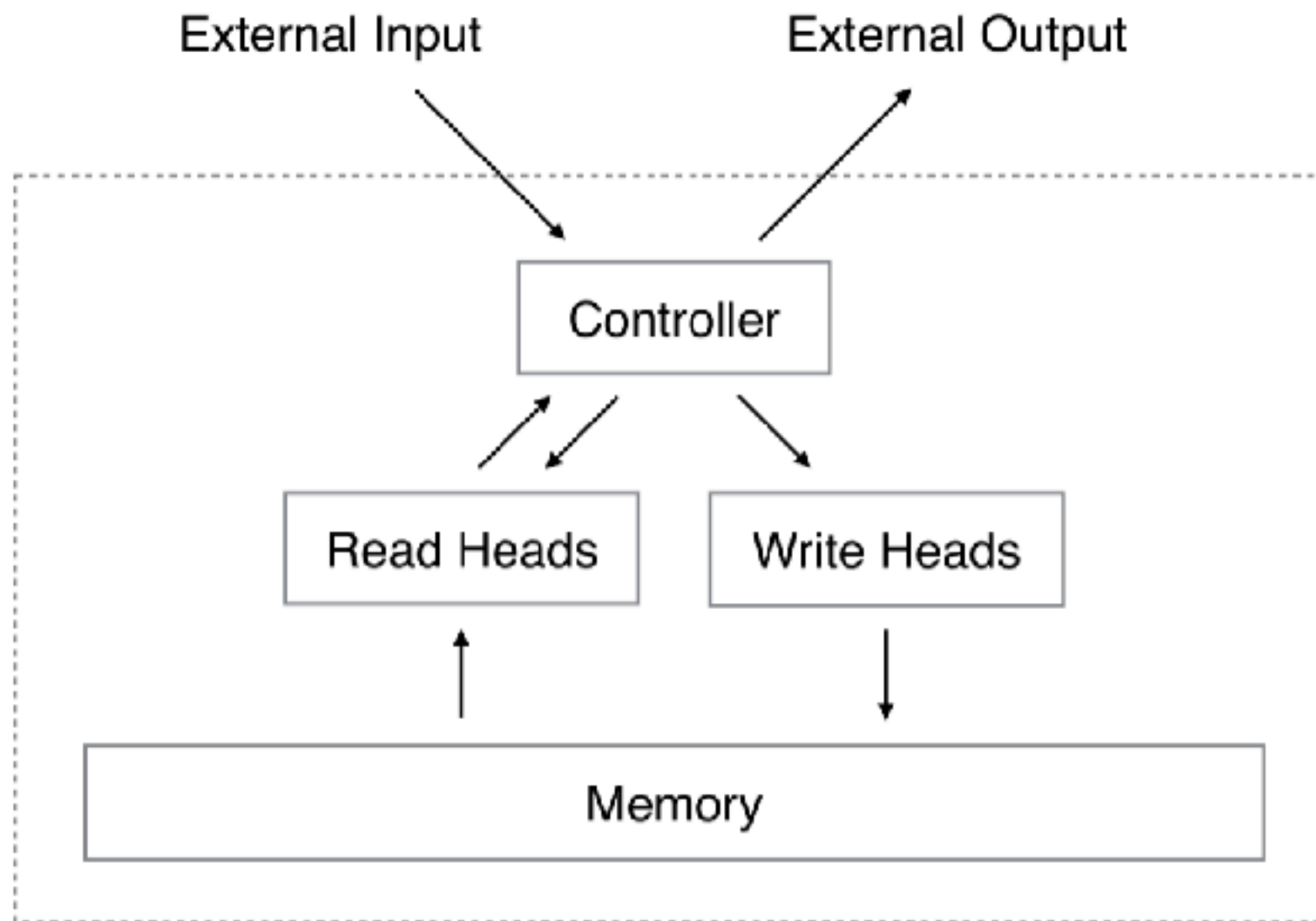
Differentiability

- Very convenient - we can train with SGD and Backprop as usual
- Obvious limitations regarding efficiency and size of the attended

Beyond Spatial Attention

- With attention, we can selectively make use of specific pieces of information when we need it.
- “Information retrieval problems”
- Add support for *storing* information at certain locations -> “Neural Touring Machine”

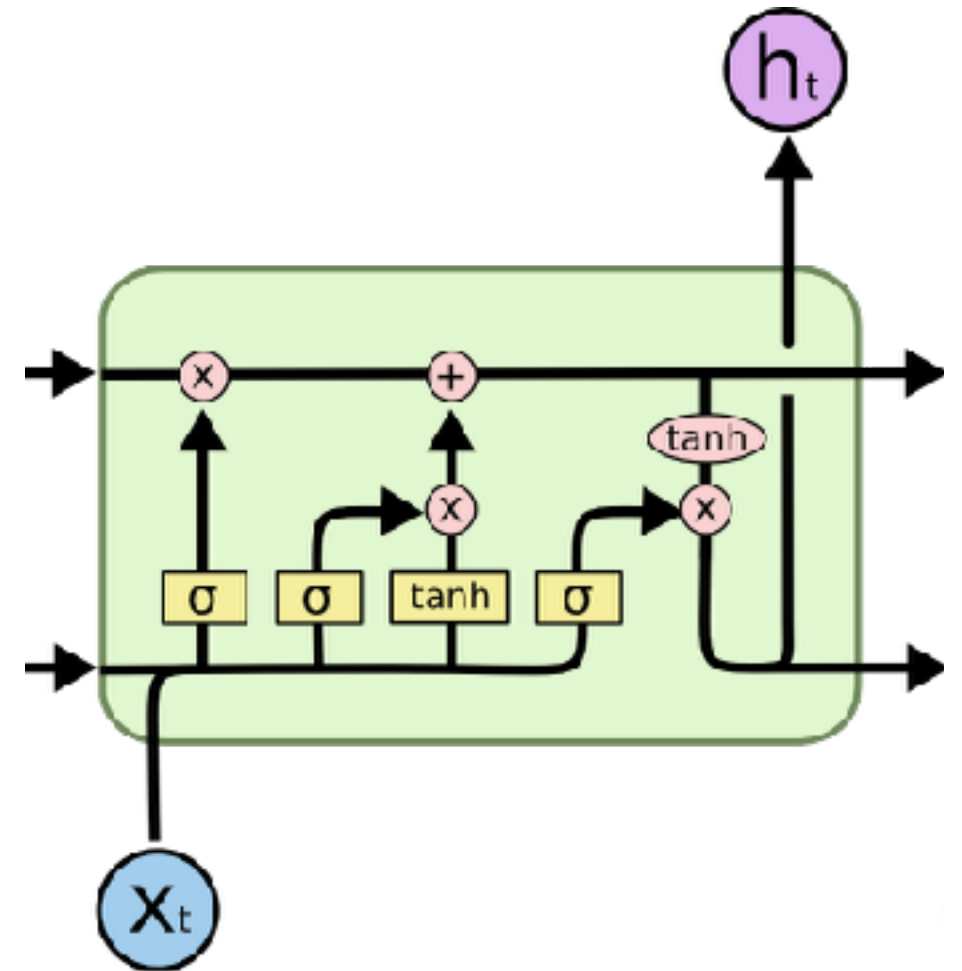
Neural Turing Machine



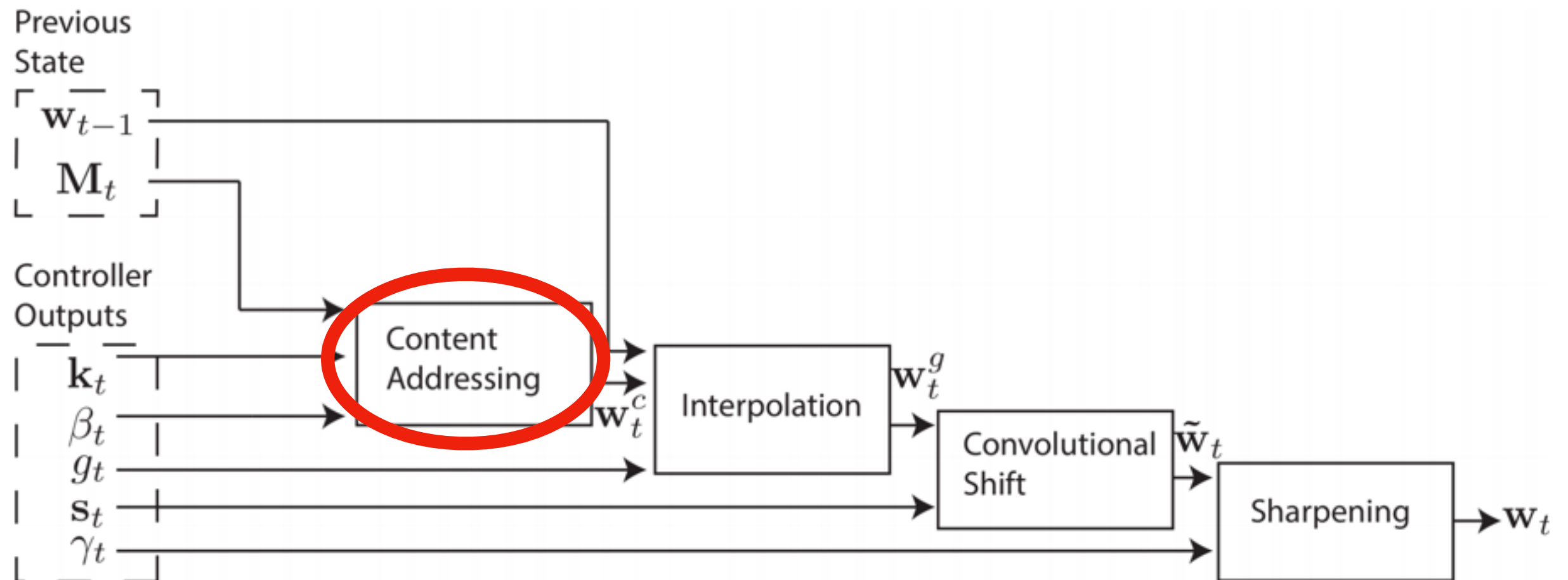
Neural Touring Machine

Comparison LSTM

- Similar concept
- NTM considerable more sophisticated (tailored) read/write mechanisms
- size of NTM memory can be increase without increasing number of parameters to learn (grows quadratically in LSTM)



Neural Touring Machine



Questions/Discussion