



**HAL**  
open science

# Recurrent Attention Network

Yannis Bendi-Ouis, Xavier Hinaut

► **To cite this version:**

Yannis Bendi-Ouis, Xavier Hinaut. Recurrent Attention Network. Bernstein Conference 2024, Sep 2024, Frankfurt Am Main, Germany. 2024. hal-04721042

**HAL Id: hal-04721042**

**<https://hal.science/hal-04721042v1>**

Submitted on 4 Oct 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Recurrent Attention Network

Yannis Bendi-Ouis<sup>1,2,3</sup>, Xavier Hinaut<sup>1,2,3</sup>

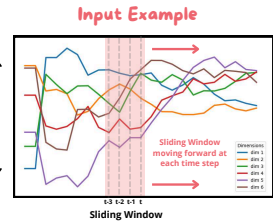
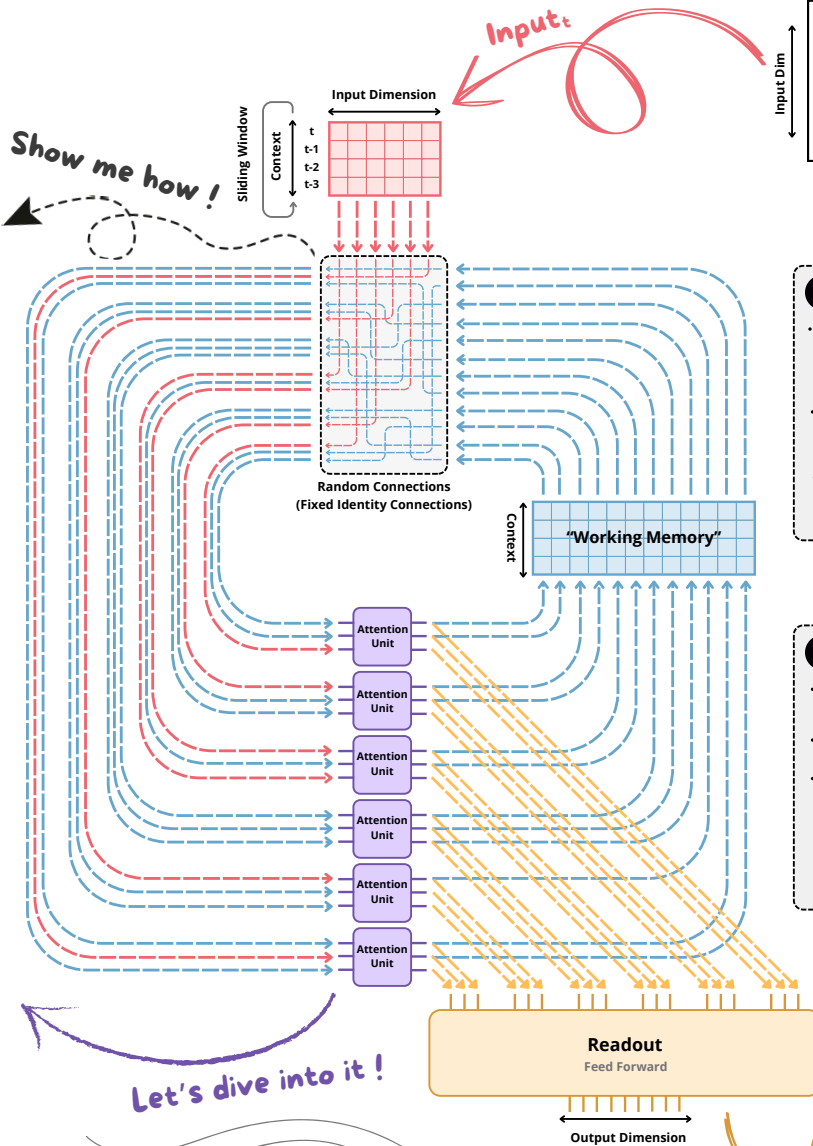
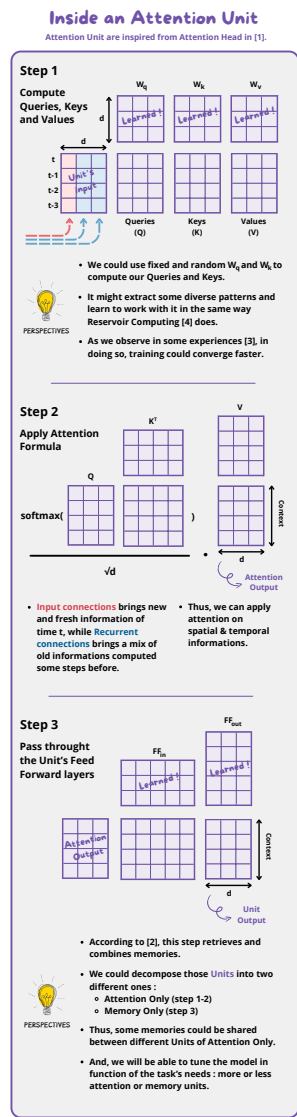
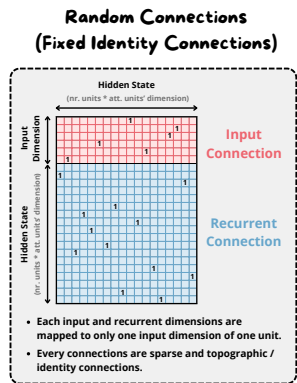
<sup>1</sup>Inria Center of Bordeaux University, Bordeaux, France

<sup>2</sup>LaBRI, Bordeaux Univ., Bordeaux INP, CNRS UMR 5800, France

<sup>3</sup>Bordeaux Univ., CNRS, IMN, UMR 5293, Bordeaux, France

This work is supported by Inria AEx BrainGPT project.  
Thanks to Experimental Inria Cluster Plafrim.

Inspired by Transformers, we're trying to  
make Reservoir Computing scalable, by using  
more complex units.



**Working Memory and Context Length**

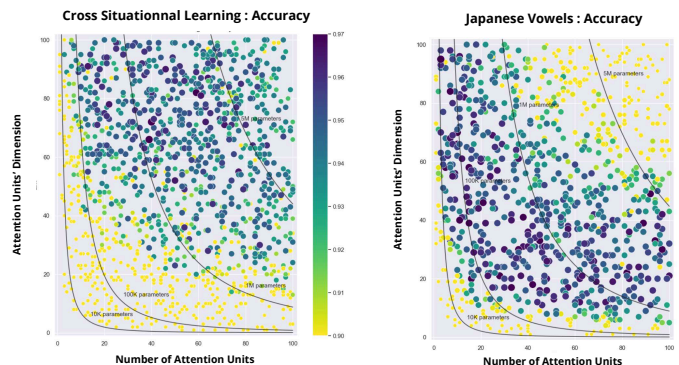
- With a "Working Memory" and a fixed context length, we can compute each step in  $O(1)$  (constant time) instead of  $O(n^2)$  (quadratic time) and still access information beyond context length thanks to recurrent connections.
- The shorter the context length is, the faster the computations. It's better to have a lot of units with a very small context length, than a few units with a very large context length.
- For example, one unit with context length of 100 takes more time than 20 units with context length of 5:  
 $100^2 > 20 \times 5^2$

**Trainable Parameters and Learning**

- Only Attention Weight ( $W_q, W_k, W_v$ ), Unit's Feed-Forward ( $FF_{in}, FF_{out}$ ) and the Readout are learned.
- We use Truncated Back-Propagation Through Time with the length of the context.
- The number of parameters trained in the network follows that formula:  
 $P = 11D^2U + UDO + 8DU + O$

where  
 $P$  = Number of Trained Parameters  
 $D$  = Attention Units' Dimension  
 $U$  = Number of Attention Units  
 $O$  = Output Dimensions

## What about the results ?



**How to understand the results ?**

- This model is a middle ground between Transformers and Reservoir Computing.
- We get closer to Transformers when the number of units is low compared to units' dimension.
- We get closer to Reservoir Computing when number of units is high compared to units' dimension.

## References

[1] Vaswani, "Attention is all you need." Advances in Neural Information Processing Systems (2017).  
 [2] Geva, "Transformer feed-forward layers are key-value memories." arXiv preprint arXiv:2012.14913 (2020).  
 [3] L ger, "Evolving Reservoirs for Meta Reinforcement Learning." International Conference on the Applications of Evolutionary Computation (2024).  
 [4] Jaeger, "The "echo state" approach to analysing and training recurrent neural networks-with an erratum note." GMD Technical Report 148.34 (2001).