

# Blending Complementary Memory Systems in Hybrid Quadratic-Linear Transformers



Kazuki Irie<sup>1</sup>

Morris Yau<sup>2</sup>

Samuel J. Gershman<sup>1,3</sup>

<sup>1</sup>Department of Psychology & Center for Brain Science, Harvard University, USA   <sup>2</sup>MIT CSAIL, USA   <sup>3</sup>Kempner Institute for the Study of Natural and Artificial Intelligence, USA

## Background

We have **two types** of transformers (“memory mechanisms”) with **complementary strengths**:

- softmax/quadratic attention = **key-value memory**
- linear attention / fast weight programming = **fast weight memory**

Property	Key-value memory	Fast weight memory
Complexity	quadratic	<b>linear</b>
Context length	bounded	<b>unbounded</b>
Retrieval precision	<b>high</b>	low
Expressivity	low	<b>high</b>

*How to get the best of both worlds?*

**Ideal solution (dream!)**: devise a clean/principled single memory mechanism with all these nice properties!

**Engineer’s solution (this work)**: Combine existing solutions! = build a hybrid system.

Remark: somewhat reminiscent of a **neuroscience hypothesis** (**Complementary Learning Systems**; McClelland et al. 1995): the brain achieves two memory functions that are **normally incompatible** (remembering specifics vs. extracting generalities) through a **division of labor** between separate episodic and semantic memory systems.

There have been many recent works on hybrid sequence models combining transformer and RNN-like models: layer-wise (e.g., xLSTM), or intra-layer (e.g., BASED, Infini-attention) hybrid approaches.

→ Our focus is **intra-layer hybrids** of **sliding window attention** and **DeltaNet**. Since our hybrid combines two types of “transformers”, Query/Key/Value projections/variables can be **nicely shared between the two sub-systems**.

## Quick reminder about DeltaNet (ICML 2021)

$$\begin{aligned}[q_t, k_t, v_t, \beta_t] &= \mathbf{W}_{\text{slow}} \mathbf{x}_t \\ \mathbf{W}_t &= \mathbf{W}_{t-1} + \sigma(\beta_t)(v_t - \mathbf{W}_{t-1}\phi(k_t)) \otimes \phi(k_t) \\ \mathbf{y}_t &= \mathbf{W}_t \phi(q_t)\end{aligned}$$

- More **expressive** than quadratic attention (Grazzi+ 2024)
- while still supporting **parallel training** (Yang+ 2024)

## Summary

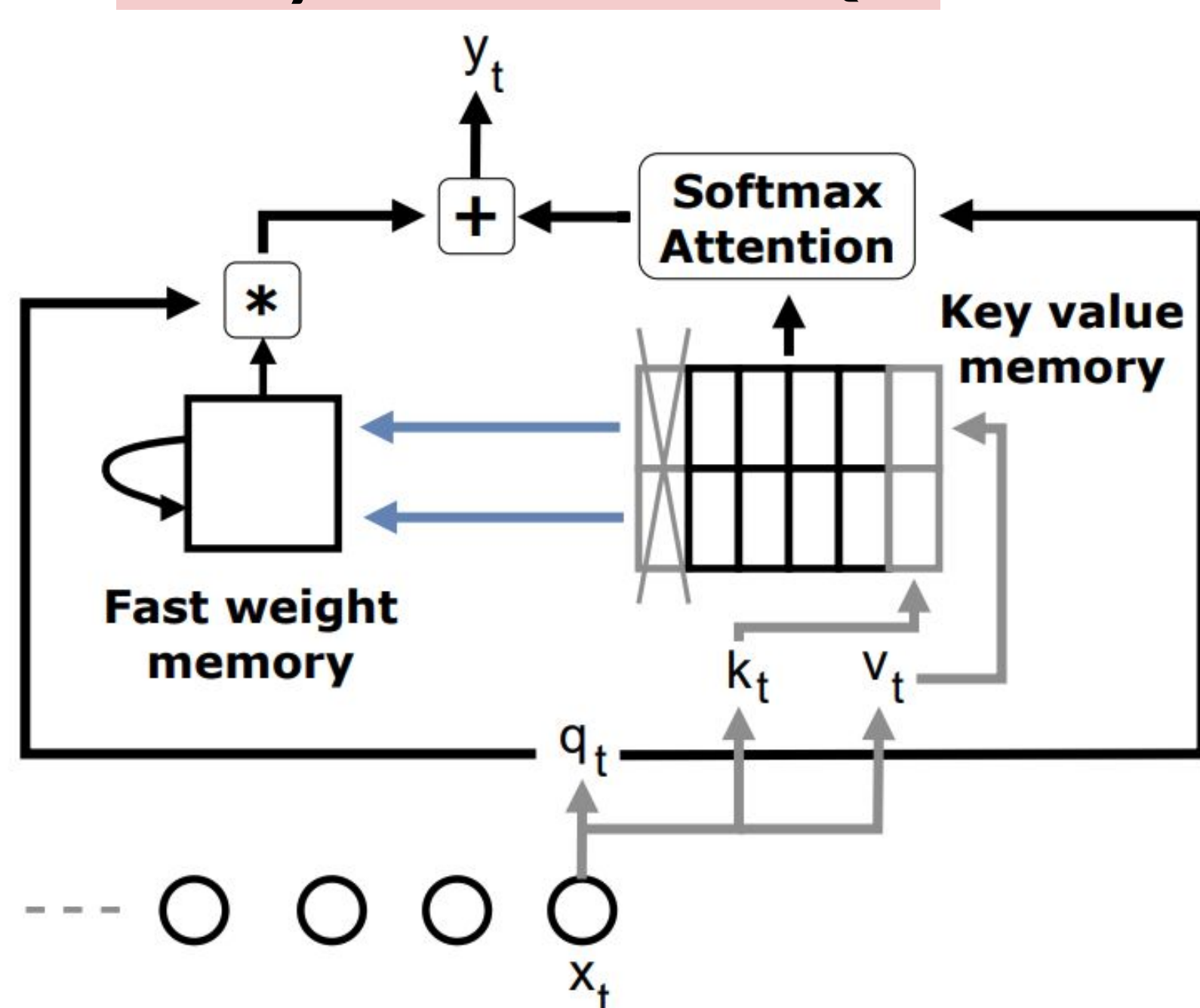
*How to build **intra-layer hybrid quadratic-linear attention models** while leveraging complementary strengths of their components?*

- We examine **design choices** for passing information between the two sub-systems, including the timing of information flow and the method used to combine their outputs (i.e., gating).
- We advocate for **expressivity** as an important, yet often overlooked, aspect of hybrid designs. Consequently, the “time-synchronous” hybrid approach performs best.

## HQLT Model Designs

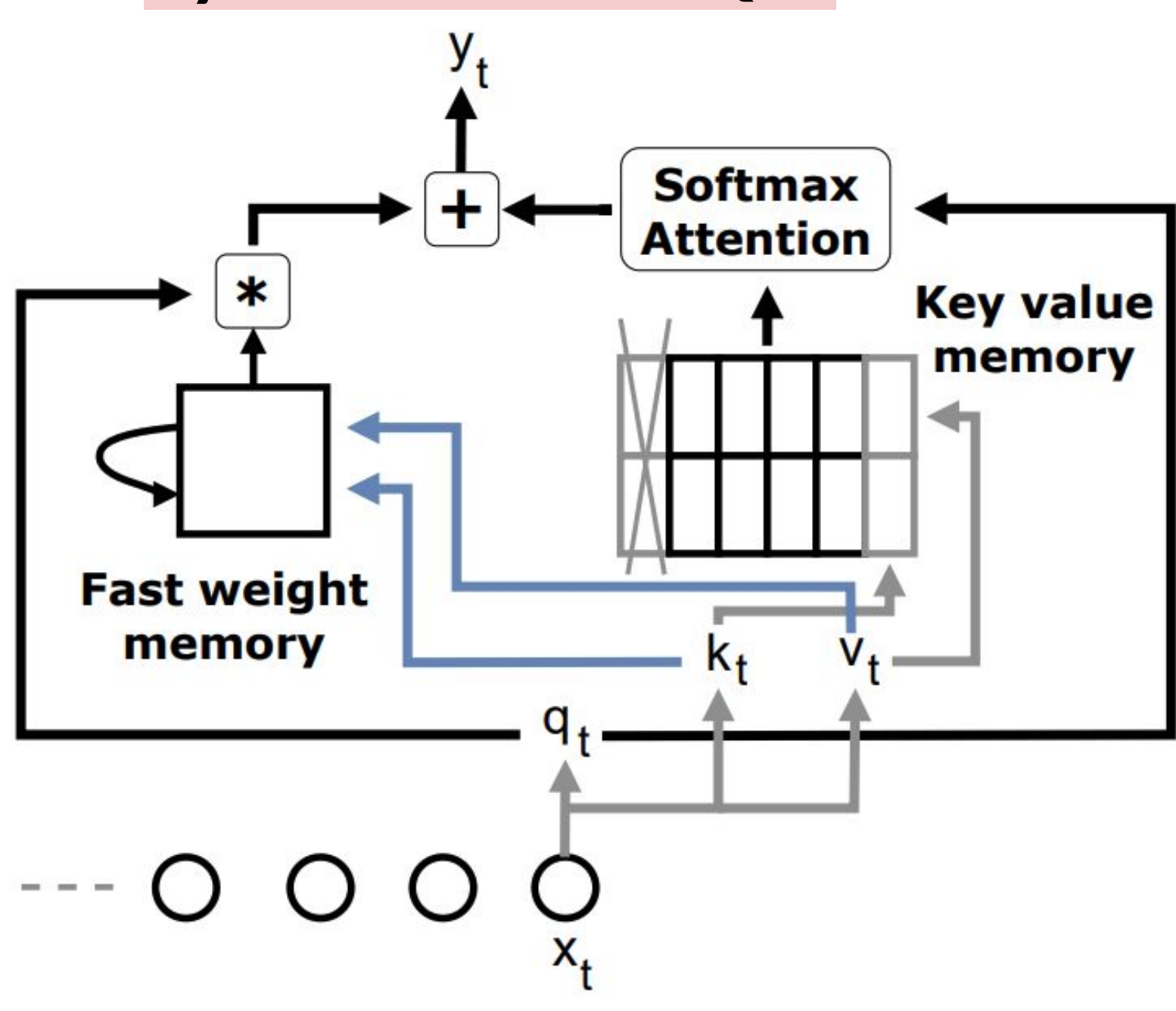
*Different division of labor between the two sub-systems*

### Delayed-Stream HQLT



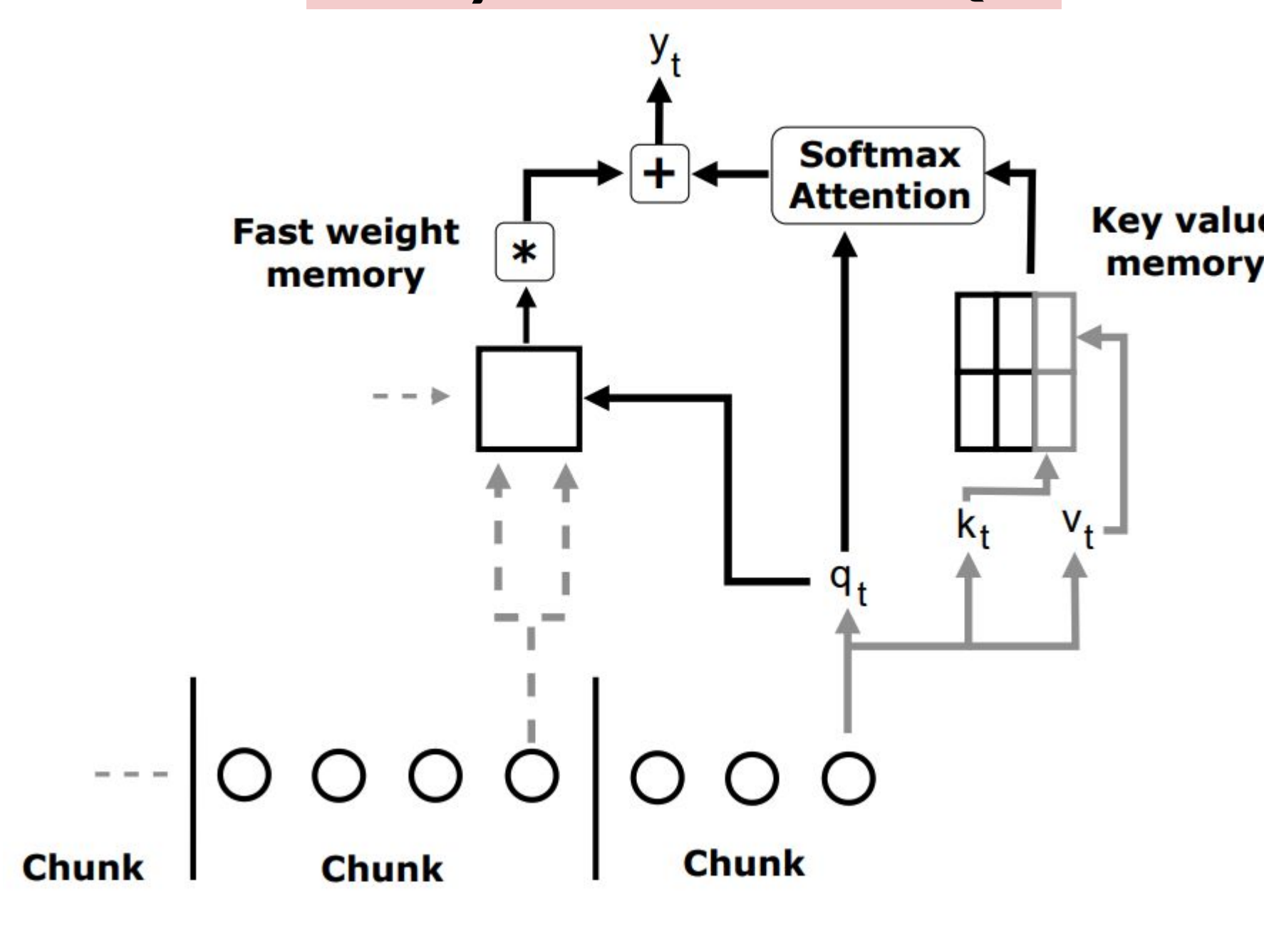
Old key/value vectors falling outside of the sliding window are fed to FW-memory (*clean division of labors*)

### Synchronous HQLT



Key/value vectors are synchronously fed to both KV-memory and FW-memory (*as we’ll see this is crucial to leverage the full advantage of FW-memory*)

### Delayed-Chunk HQLT



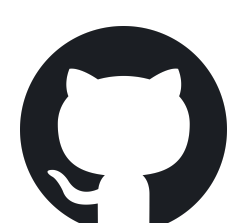
Delayed variant with a “chunk size”: KV-memory contents are flashed at the beginning of each chunk (*natural from the chunk-wise training view; just adding softmax within chunk*)

## General language modelling tasks

*Sanity checks!*

- Performance is mostly comparable between different HQLT variants (slight advantage for Delayed-Stream for the 1.3B model)
- Enlarging the window size also has a limited impact for these general LM tasks
- Use of vanilla linear attention instead of DeltaNet yields significant performance degradation

Remark: Our code is based on “**flame**” for training and “**fla**” for model implementation.



[github.com/kazuki-irie/hybrid-memory](https://github.com/kazuki-irie/hybrid-memory)

Model	Wiki. ppl ↓	LMB. ppl ↓	LMB. acc ↑	PIQA acc ↑	Hella. acc_n ↑	Wino. acc ↑	ARC-e acc ↑	ARC-c acc_n ↑	Avg.
<b>340M params</b>									
Transformer++	26.5	34.9	<b>33.9</b>	67.6	41.0	53.7	60.2	29.0	47.6
DeltaNet	27.6	35.0	32.8	67.1	40.8	52.6	58.5	28.8	46.8
<b>HQLT</b>									
Delayed-Stream	26.4	33.1	33.6	<b>67.9</b>	42.1	51.8	59.4	29.5	47.4
Delayed-Chunk	26.7	29.9	33.5	66.8	42.3	50.9	61.1	<b>30.6</b>	47.5
Synchronous	<b>26.3</b>	<b>29.4</b>	33.3	66.2	<b>42.7</b>	<b>53.8</b>	<b>61.5</b>	29.4	<b>47.8</b>
<b>1.3B params</b>									
Transformer++	19.8	17.9	42.6	71.0	50.3	55.8	65.2	33.2	53.0
DeltaNet	20.6	19.9	39.3	70.1	49.5	52.5	68.5	34.2	52.3
<b>HQLT</b>									
Delayed-Stream	20.0	16.5	<b>43.5</b>	70.7	<b>51.6</b>	56.0	<b>69.3</b>	<b>36.0</b>	<b>54.5</b>
Delayed-Chunk	20.2	16.3	41.3	71.8	50.9	55.0	67.9	35.2	53.7
Synchronous	<b>19.8</b>	<b>15.9</b>	42.8	<b>72.0</b>	51.5	<b>56.1</b>	68.1	33.1	53.9
<b>Ablations with 340M params</b>									
HQLT Synchronous	27.7	34.8	32.5	67.0	41.2	52.4	60.9	28.9	47.2
sum mixing	26.3	29.4	33.3	66.2	42.7	53.8	61.5	29.4	47.8
dynamic vector mixing (25M)	26.6	27.8	35.6	68.1	41.7	51.2	61.4	30.1	48.0
w. Linear Attn. (no DeltaNet)	26.7	27.7	35.4	67.0	42.2	52.7	59.8	29.0	47.7
window × 2 = 128	27.0	28.0	35.9	66.3	41.3	53.2	60.1	29.0	47.6
window × 4 = 256									
window × 8 = 512									

## Synthetic algorithmic tasks (expressivity evaluation)

Expectably, only the “Synchronous” variant can leverage the expressivity advantage of DeltaNet

Model	Parity acc ↑	Mod Arith acc ↑
Transformer [7]	2.2	3.1
Mamba [7]	100.0	24.1
DeltaNet [7]	100.0	97.1
<b>HQLT</b>		
Delayed-Stream	3.3	27.8
Delayed-Chunk	2.8	1.4
Synchronous	<b>100.0</b>	<b>97.0</b>
w. Linear Attn.	2.5	44.5

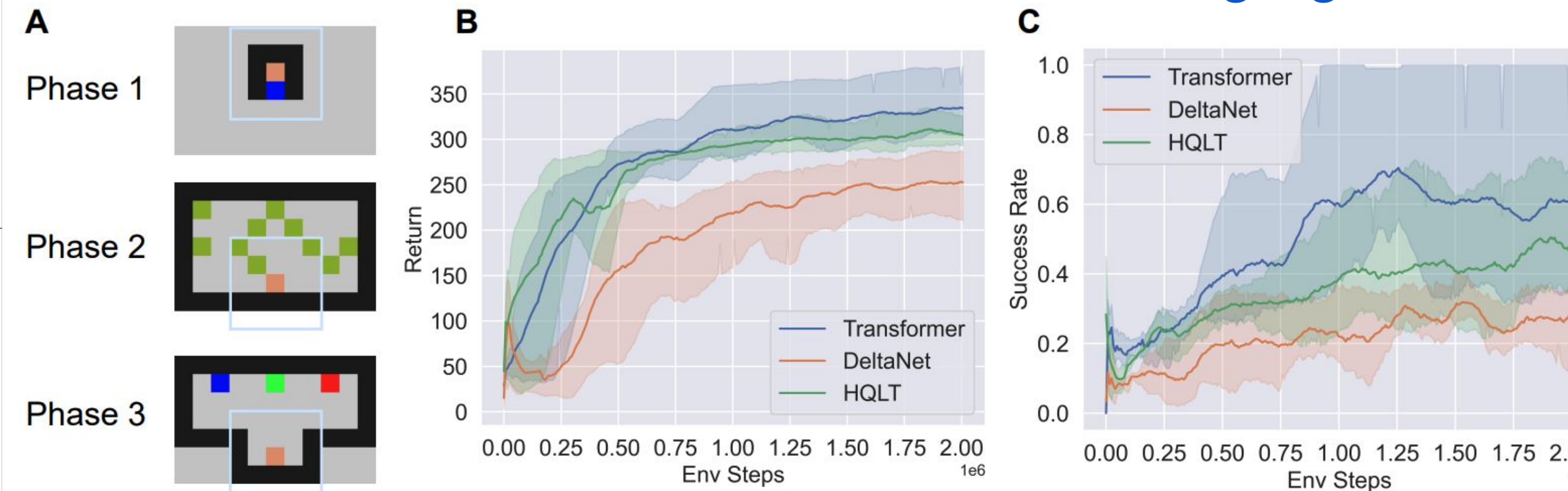
## Retrieval tasks

- Large window size is needed for these retrieval intensive tasks  
- Naive “sum mixing” is not enough for HQLT (context-dependent scalar or vector gating helps)

	Window size	SWDE acc ↑	SQuAD acc ↑	FDA acc ↑	Avg.
<b>340M params</b>					
Transformer++	2048	<b>44.9</b>	36.9	<b>52.3</b>	<b>44.7</b>
Transformer++	1024	30.4	25.5	31.2	29.0
DeltaNet	-	18.5	25.2	8.6	17.4
<b>HQLT Synchronous</b>					
sum mixer	64	13.3	26.2	12.6	17.4
dynamic scalar	64	21.1	27.7	11.4	20.1
dynamic vector	64	20.0	28.4	10.9	19.8
<b>1.3B params</b>					
window × 2	128	16.9	30.5	17.9	21.8
window × 4	256	18.6	35.6	15.1	23.1
window × 8	512	22.9	35.7	17.3	25.3
window × 16	1024	34.0	<b>37.1</b>	50.1	40.4

## Reinforcement learning in partially observable environment

*Evaluation outside of the language domain*



**Task**: there are three phases: (1) observe a color (here blue), (2) collect green apples (distraction task!), (3) reach the right color (if so, counted as success). Partial observation (light blue box) is limited around the agent (brown). Total sequence len = 750

**Result**: small window size of 64 improves HQLT over DeltaNet

## Discussions

- Overall promising results but retrieval performance is not satisfactory (still requires a large sliding window size)
- Designing a more sophisticated communication between the two sub-systems may address this limitation