Some Memory Mechanisms in Transformers and their Emergence

Alberto Bietti

Flatiron Institute, Simons Foundation

Memory and Vision workshop. ICCV 2025.





Some Memory Mechanisms in Transformers and their Emergence

Alberto Bietti

Flatiron Institute, Simons Foundation

Memory and Vision workshop. ICCV 2025.

w/ V. Cabannes, E. Dohmatob, D. Bouchacourt, H. Jégou, L. Bottou (Meta AI), E. Nichani, Z. Wang, J. Lee (Princeton), L. Chen, D. Wu, J. Bruna (NYU), D. Hsu (Columbia)

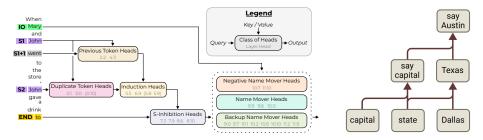




Mechanisms inside Transformer LLMs

Reasoning over context

- Circuits of attention heads (Elhage et al., 2021; Olsson et al., 2022; Wang et al., 2022)
- A "biology" of circuits found by mechanistic interpretability (e.g., Anthropic, 2025)



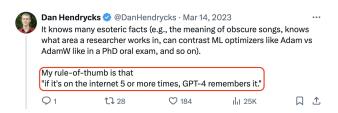
Mechanisms inside Transformer LLMs

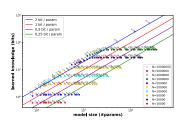
Reasoning over context

- Circuits of attention heads (Elhage et al., 2021; Olsson et al., 2022; Wang et al., 2022)
- A "biology" of circuits found by mechanistic interpretability (e.g., Anthropic, 2025)

Knowledge storage

- Memorization, factual recall, parameter scaling
 - ► (Geva et al., 2020; Meng et al., 2022; Allen-Zhu and Li, 2024)
- Learn rules that help higher-level reasoning





Mechanisms inside Transformer LLMs

Reasoning over context

- Circuits of attention heads (Elhage et al., 2021; Olsson et al., 2022; Wang et al., 2022)
- A "biology" of circuits found by mechanistic interpretability (e.g., Anthropic, 2025)

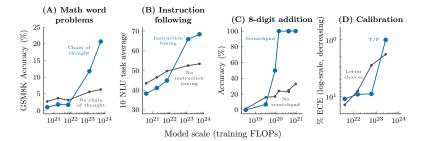
Knowledge storage

- Memorization, factual recall, parameter scaling
 - ► (Geva et al., 2020; Meng et al., 2022; Allen-Zhu and Li, 2024)
- Learn rules that help higher-level reasoning

Q: How do these arise during training?

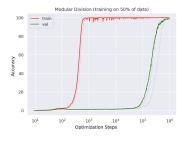
Emergence and training dynamics

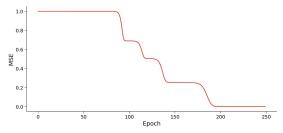
Emergent capabilities and sequential learning (Wei et al. 2022; Power et al., 2022; Saxe et al. 2013)



Emergence and training dynamics

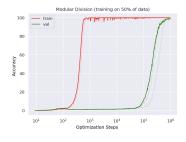
Emergent capabilities and sequential learning (Wei et al. 2022; Power et al., 2022; Saxe et al. 2013)

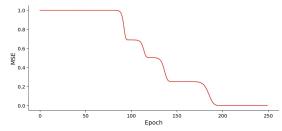




Emergence and training dynamics

Emergent capabilities and sequential learning (Wei et al. 2022; Power et al., 2022; Saxe et al. 2013)





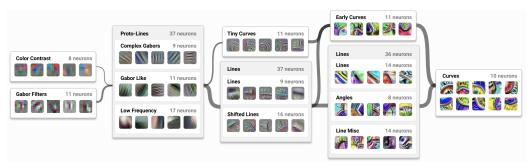
Why do we care about understanding?

- Highlight role of data, training algorithm, architecture
- ⇒ improve training methodology

What about computer vision?

A mix of continuous and discrete (?)

- Low-level processing usually involves continuous operations (e.g., Gabor filters)
- Some evidence for learned circuits
- Models also involve discrete semantic concepts and knowledge (especially VLMs)



(Cammarata et al., 2021)

What about computer vision?

A mix of continuous and discrete (?)

- Low-level processing usually involves continuous operations (e.g., Gabor filters)
- Some evidence for learned circuits
- Models also involve discrete semantic concepts and knowledge (especially VLMs)

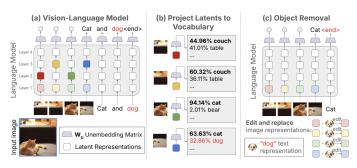


(Templeton et al., 2024)

What about computer vision?

A mix of continuous and discrete (?)

- Low-level processing usually involves continuous operations (e.g., Gabor filters)
- Some evidence for learned circuits
- Models also involve discrete semantic concepts and knowledge (especially VLMs)



(Jiang et al., 2025)

Input: sequence of discrete tokens $(z_1, \ldots, z_T) \in [N]^T$

Input: sequence of discrete tokens $(z_1, \ldots, z_T) \in [N]^T$

Embeddings

• input e_z , positional p_t , output u_y , in \mathbb{R}^d

Input: sequence of discrete tokens $(z_1, \ldots, z_T) \in [N]^T$

Embeddings

• input e_z , positional p_t , output u_y , in \mathbb{R}^d

Residual streams (Elhage et al., 2021)

- embed each token $z_t \in [N]$ as $x_t := e_{z_t} + p_t$
- (causal) self-attention $x_t := x_t + \mathsf{MHSA}(x_t, x_{1:t})$



$$\mathsf{MHSA}(\mathbf{x}_t, \mathbf{x}_{1:t}) = \sum_{h=1}^H \sum_{s=1}^t \beta_s^h W_O^{h\top} W_V^h \mathbf{x}_s, \quad \text{ with } \beta_s^h = \frac{\exp(\mathbf{x}_s^\top W_K^{h\top} W_Q^h \mathbf{x}_t)}{\sum_{s=1}^t \exp(\mathbf{x}_s^\top W_K^{h\top} W_Q^h \mathbf{x}_t)}$$

where W_K , W_Q , W_V , $W_O \in \mathbb{R}^{d_h \times d}$ (key/query/value/output matrices)

Input: sequence of discrete tokens $(z_1, \ldots, z_T) \in [N]^T$

Embeddings

• input e_z , positional p_t , output u_y , in \mathbb{R}^d

Residual streams (Elhage et al., 2021)

- embed each token $z_t \in [N]$ as $x_t := e_{z_t} + p_t$
- (causal) self-attention $x_t := x_t + MHSA(x_t, x_{1:t})$
- feed-forward $x_t := x_t + \mathsf{MLP}(x_t)$



$$\mathsf{MLP}(\mathbf{x}_t) = V^{\top} \sigma(U\mathbf{x}_t)$$

where $U, V \in \mathbb{R}^{m \times d}$, often m = 4d

Input: sequence of discrete tokens $(z_1, \ldots, z_T) \in [N]^T$

Embeddings

• input e_z , positional p_t , output u_y , in \mathbb{R}^d

Residual streams (Elhage et al., 2021)

- embed each token $z_t \in [N]$ as $x_t := e_{z_t} + p_t$
- (causal) self-attention $x_t := x_t + \mathsf{MHSA}(x_t, x_{1:t})$
- feed-forward $x_t := x_t + MLP(x_t)$
- residual stream x_t is a sum of embeddings/"features"



Input: sequence of discrete tokens $(z_1, \ldots, z_T) \in [N]^T$

Embeddings

• input e_z , positional p_t , output u_y , in \mathbb{R}^d

Residual streams (Elhage et al., 2021)

- embed each token $z_t \in [N]$ as $x_t := e_{z_t} + p_t$
- (causal) self-attention $x_t := x_t + \mathsf{MHSA}(x_t, x_{1:t})$
- feed-forward $x_t := x_t + MLP(x_t)$
- residual stream x_t is a sum of embeddings/"features"

Next-token prediction

cross-entropy loss

$$\sum_{t < T} \ell(z_{t+1}; (\underline{u_j}^\top x_t)_j)$$



Outline

1 Weights as memories

2 Context as memory

• Consider sets of nearly orthonormal embeddings $\{e_z\}_{z\in\mathcal{Z}}$ and $\{u_y\}_{y\in\mathcal{Y}}$:

$$\|e_z\| \approx 1$$
 and $e_z^{\top} e_{z'} \approx 0$
 $\|u_y\| \approx 1$ and $u_y^{\top} u_{y'} \approx 0$

• Consider sets of nearly orthonormal embeddings $\{e_z\}_{z\in\mathcal{Z}}$ and $\{u_y\}_{y\in\mathcal{Y}}$:

$$\|e_z\| \approx 1$$
 and $e_z^{\top} e_{z'} \approx 0$
 $\|u_y\| \approx 1$ and $u_y^{\top} u_{y'} \approx 0$

$$W = \sum_{\mathbf{z} \in \mathcal{M}} \mathbf{u_{f(\mathbf{z})}} \mathbf{e_{\mathbf{z}}}^{\top}$$

• Consider sets of nearly orthonormal embeddings $\{e_z\}_{z\in\mathcal{Z}}$ and $\{u_y\}_{y\in\mathcal{Y}}$:

$$\|e_z\| \approx 1$$
 and $e_z^{\top} e_{z'} \approx 0$
 $\|u_y\| \approx 1$ and $u_y^{\top} u_{y'} \approx 0$

$$W = \sum_{\mathbf{z} \in \mathcal{M}} \mathbf{u}_{f(\mathbf{z})} \mathbf{e}_{\mathbf{z}}^{\top} \implies \mathbf{u}_{y}^{\top} W \mathbf{e}_{\mathbf{z}} \approx \mathbb{1} \{ y = f(\mathbf{z}) \}$$

• Consider sets of nearly orthonormal embeddings $\{e_z\}_{z\in\mathcal{Z}}$ and $\{u_y\}_{y\in\mathcal{Y}}$:

$$\|e_z\| \approx 1$$
 and $e_z^{\top} e_{z'} \approx 0$
 $\|u_y\| \approx 1$ and $u_y^{\top} u_{y'} \approx 0$

$$W = \sum_{\mathbf{z} \in \mathcal{M}} \mathbf{u}_{f(\mathbf{z})} \mathbf{e}_{\mathbf{z}}^{\top} \implies \mathbf{u}_{y}^{\top} W \mathbf{e}_{\mathbf{z}} \approx \mathbb{1} \{ y = f(\mathbf{z}) \}$$

- Examples in Transformers:
 - ► Logits in attention heads: $x_k^\top W_{KQ} x_q$
 - ► Logits in next-token prediction: $u_y^\top \dot{U}\sigma(Vx_t)$ or $u_y^\top W_{OV}x_k$

• Consider sets of nearly orthonormal embeddings $\{e_z\}_{z\in\mathcal{Z}}$ and $\{u_y\}_{y\in\mathcal{Y}}$:

$$\|e_z\| \approx 1$$
 and $e_z^{\top} e_{z'} \approx 0$
 $\|u_y\| \approx 1$ and $u_y^{\top} u_{y'} \approx 0$

$$W = \sum_{\mathbf{z} \in \mathcal{M}} \mathbf{u}_{f(\mathbf{z})} \mathbf{e}_{\mathbf{z}}^{\top} \implies \mathbf{u}_{\mathbf{y}}^{\top} W \mathbf{e}_{\mathbf{z}} \approx \mathbb{1}\{y = f(\mathbf{z})\}$$

- Examples in Transformers:
 - ▶ Logits in attention heads: $x_k^\top W_{KQ} x_q$
 - ► Logits in next-token prediction: $u_v^\top U \sigma(V x_t)$ or $u_v^\top W_{OV} x_k$
- For random embeddings, capacity \approx number of parameters
 - ► See Cabannes et al. (2024); Nichani et al. (2025), extends to MLPs

• Consider sets of nearly orthonormal embeddings $\{e_z\}_{z\in\mathcal{Z}}$ and $\{u_y\}_{y\in\mathcal{Y}}$:

$$\|e_z\| \approx 1$$
 and $e_z^{\top} e_{z'} \approx 0$
 $\|u_y\| \approx 1$ and $u_y^{\top} u_{y'} \approx 0$

$$W = \sum_{z \in \mathcal{M}} u_{f(z)} e_z^{\top} \implies u_y^{\top} W e_z \approx \mathbb{1}\{y = f(z)\}$$

- Examples in Transformers:
 - ▶ Logits in attention heads: $x_k^\top W_{KQ} x_q$
 - ► Logits in next-token prediction: $u_v^\top U \sigma(V x_t)$ or $u_v^\top W_{OV} x_k$
- ullet For random embeddings, capacity pprox number of parameters
 - ► See Cabannes et al. (2024); Nichani et al. (2025), extends to MLPs
- Related to Hopfield (1982); Kohonen (1972); Willshaw et al. (1969); Iscen et al. (2017)

Lemma (Gradients as memories, B. et al., 2023)

Let p be a data distribution over $(z, y) \in [N]^2$, and consider the loss

$$L(W) = \mathbb{E}_{(z,y)\sim p}[\ell(y,F_W(z))], \quad F_W(z)_k = \mathbf{u_k}^\top W \mathbf{e_z},$$

with ℓ the cross-entropy loss and $e_z,~u_k$ input/output embeddings.

Lemma (Gradients as memories, B. et al., 2023)

Let p be a data distribution over $(z, y) \in [N]^2$, and consider the loss

$$L(W) = \mathbb{E}_{(z,y)\sim p}[\ell(y,F_W(z))], \quad F_W(z)_k = \frac{\mathbf{u}_k}{\mathbf{v}_k} W \mathbf{e}_z,$$

with ℓ the **cross-entropy loss** and e_z , u_k input/output embeddings. Then,

$$\nabla L(W) = \sum_{k=1}^{K} \mathbb{E}_{z} [(\hat{p}_{W}(y=k|z) - p(y=k|z)) \mathbf{u}_{k} \mathbf{e}_{z}^{\top}]$$

Lemma (Gradients as memories, B. et al., 2023)

Let p be a data distribution over $(z,y) \in [N]^2$, and consider the loss

$$L(W) = \mathbb{E}_{(z,y)\sim p}[\ell(y,F_W(z))], \quad F_W(z)_k = \frac{\mathbf{u}_k}{\mathbf{v}_k} W \mathbf{e}_z,$$

with ℓ the **cross-entropy loss** and e_z , u_k input/output embeddings. Then,

$$\nabla L(W) = \sum_{k=1}^{K} \mathbb{E}_{z} [(\hat{p}_{W}(y=k|z) - p(y=k|z)) \mathbf{u}_{k} \mathbf{e}_{z}^{\top}]$$

• Example: $z \sim \text{Unif}([N])$, $y = f_*(z)$

Lemma (Gradients as memories, B. et al., 2023)

Let p be a data distribution over $(z,y) \in [N]^2$, and consider the loss

$$L(W) = \mathbb{E}_{(z,y)\sim p}[\ell(y,F_W(z))], \quad F_W(z)_k = \frac{\mathbf{u}_k}{\mathbf{v}_k} W \mathbf{e}_z,$$

with ℓ the **cross-entropy loss** and e_z , u_k input/output embeddings. Then,

$$\nabla L(W) = \sum_{k=1}^{K} \mathbb{E}_{z} [(\hat{p}_{W}(y=k|z) - p(y=k|z)) \mathbf{u}_{k} e_{z}^{\top}]$$

- **Example**: $z \sim \text{Unif}([N])$, $y = f_*(z)$
 - ► After **one gradient step** on the population loss, assuming near-orthonormal embeddings

$$W_1 = \frac{\eta}{N} \sum_{z,k} \left(\mathbb{1}\{f_*(z) = k\} - \frac{1}{N} \right) \mathbf{u}_k \mathbf{e}_z^\top \quad \Longrightarrow \quad \mathbf{u}_k^\top W_1 \mathbf{e}_z \approx \frac{\eta}{N} \left(\mathbb{1}\{f_*(z) = k\} - \frac{1}{N} \right)$$

Lemma (Gradients as memories, B. et al., 2023)

Let p be a data distribution over $(z, y) \in [N]^2$, and consider the loss

$$L(W) = \mathbb{E}_{(z,y)\sim p}[\ell(y,F_W(z))], \quad F_W(z)_k = \frac{\mathbf{u_k}^{\top} W e_z}{\mathbf{e}_z},$$

with ℓ the **cross-entropy loss** and e_z , u_k input/output embeddings. Then,

$$\nabla L(W) = \sum_{k=1}^{K} \mathbb{E}_{\mathbf{z}}[(\hat{p}_{W}(y=k|\mathbf{z}) - p(y=k|\mathbf{z})) \mathbf{u}_{k} \mathbf{e}_{\mathbf{z}}^{\top}]$$

- **Example**: $z \sim \text{Unif}([N])$, $y = f_*(z)$
 - ► After **one gradient step** on the population loss, assuming near-orthonormal embeddings

$$W_1 = \frac{\eta}{N} \sum_{z,k} \left(\mathbb{1}\{f_*(z) = k\} - \frac{1}{N} \right) \mathbf{u}_k \mathbf{e}_z^\top \quad \Longrightarrow \quad \mathbf{u}_k^\top W_1 \mathbf{e}_z \approx \frac{\eta}{N} \left(\mathbb{1}\{f_*(z) = k\} - \frac{1}{N} \right)$$

► Corollary: $\hat{f}(z) = \arg\max_k \frac{u_k}{u_k} W_1 e_z$ has near-perfect accuracy

Lemma (Gradients as memories, B. et al., 2023)

Let p be a data distribution over $(z, y) \in [N]^2$, and consider the loss

$$L(W) = \mathbb{E}_{(z,y)\sim p}[\ell(y,F_W(z))], \quad F_W(z)_k = \frac{\mathbf{u}_k}{\mathbf{v}_k} W \mathbf{e}_z,$$

with ℓ the **cross-entropy loss** and e_z , u_k input/output embeddings. Then,

$$\nabla L(W) = \sum_{k=1}^{K} \mathbb{E}_{z} [(\hat{p}_{W}(y=k|z) - p(y=k|z)) \mathbf{u}_{k} e_{z}^{\top}]$$

- **Example**: $z \sim \text{Unif}([N])$, $y = f_*(z)$
 - ► After **one gradient step** on the population loss, assuming near-orthonormal embeddings

$$W_1 = \frac{\eta}{N} \sum_{z,k} \left(\mathbb{1}\{f_*(z) = k\} - \frac{1}{N} \right) u_k e_z^\top \quad \Longrightarrow \quad u_k^\top W_1 e_z \approx \frac{\eta}{N} \left(\mathbb{1}\{f_*(z) = k\} - \frac{1}{N} \right)$$

- ▶ Corollary: $\hat{f}(z) = \arg\max_k \frac{u_k}{u_k} W_1 e_z$ has near-perfect accuracy
- More generally, replace u_k by "backward" vector

Lemma (Gradients as memories, B. et al., 2023)

Let p be a data distribution over $(z, y) \in [N]^2$, and consider the loss

$$L(W) = \mathbb{E}_{(z,y)\sim p}[\ell(y,F_W(z))], \quad F_W(z)_k = \frac{\mathbf{u}_k}{\mathbf{v}_k} W \mathbf{e}_z,$$

with ℓ the **cross-entropy loss** and e_z , u_k input/output embeddings. Then,

$$\nabla L(W) = \sum_{k=1}^{K} \mathbb{E}_{z} [(\hat{p}_{W}(y=k|z) - p(y=k|z)) \mathbf{u}_{k} \mathbf{e}_{z}^{\top}]$$

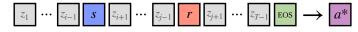
- Example: $z \sim \text{Unif}([N])$, $y = f_*(z)$
 - lacktriangledown After one gradient step on the population loss, assuming near-orthonormal embeddings

$$W_1 = \frac{\eta}{N} \sum_{z,k} \left(\mathbb{1}\{f_*(z) = k\} - \frac{1}{N} \right) \mathbf{u}_k \mathbf{e}_z^\top \quad \Longrightarrow \quad \mathbf{u}_k^\top W_1 \mathbf{e}_z \approx \frac{\eta}{N} \left(\mathbb{1}\{f_*(z) = k\} - \frac{1}{N} \right)$$

- ► Corollary: $\hat{f}(z) = \arg\max_k u_k^\top W_1 e_z$ has near-perfect accuracy
- More generally, replace u_k by "backward" vector

Note: related to (Ba et al., 2022; Damian et al., 2022; Oymak et al., 2023; Yang and Hu, 2021)

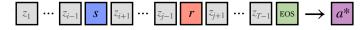
Application to factual recall: toy model



The capital of France is Paris

- $s \in S$: subject token
- $r \in \mathcal{R}$: relation token
- $a^*(s,r) \in \mathcal{A}_r$: attribute/fact to be stored
- $z_i \in \mathcal{N}$: noise tokens

Application to factual recall: toy model



The capital of France is Paris

- $s \in S$: subject token
- $r \in \mathcal{R}$: relation token
- $a^*(s,r) \in \mathcal{A}_r$: attribute/fact to be stored
- $z_i \in \mathcal{N}$: noise tokens

Q: How do Transformers solve this?

Two mechanisms

- One-layer Transformer, with or without MLP, random embeddings
- Embedding dimension d, head dimension d_h , MLP width m, H heads

Theorem (Nichani et al., 2025, informal)

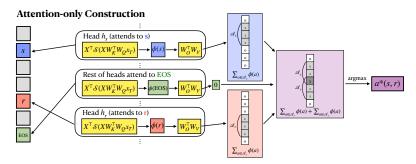
A one-layer Transformer with or without MLP can achieve perfect factual recall with $\tilde{O}(SR)$ parameters. (note: attention-only needs large enough d)

Two mechanisms

- One-layer Transformer, with or without MLP, random embeddings
- Embedding dimension d, head dimension d_h , MLP width m, H heads

Theorem (Nichani et al., 2025, informal)

A one-layer Transformer with or without MLP can achieve perfect factual recall with $\tilde{O}(SR)$ parameters. (note: attention-only needs large enough d)



Two mechanisms

- One-layer Transformer, with or without MLP, random embeddings
- Embedding dimension d, head dimension d_h , MLP width m, H heads

Theorem (Nichani et al., 2025, informal)

A one-layer Transformer with or without MLP can achieve perfect factual recall with $\tilde{O}(SR)$ parameters. (note: attention-only needs large enough d)

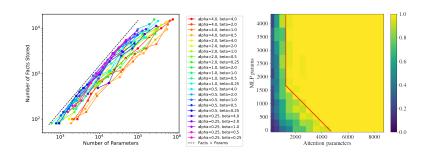
Attention + MLP Construction $\begin{array}{c} \vdots \\ \hline d/dd_h|\text{heads attends to s} \\ \hline S \\ \hline Rest of heads attend to EOS \\ \hline X^TS(XW_K^TW_{Q^XT}) \rightarrow \emptyset(S) \rightarrow W_O^TW_V \end{array}$ Rest of heads attend to EOS $\begin{array}{c} X^TS(XW_K^TW_{Q^XT}) \rightarrow \emptyset(S) \rightarrow W_O^TW_V \end{array}$ $\begin{array}{c} W_0^TW_V \rightarrow \overline{\emptyset}(S) \\ \overline{\emptyset}(S) \rightarrow W_O^TW_V \end{array}$ $\begin{array}{c} W_0^TW_V \rightarrow \overline{\emptyset}(S) \rightarrow W_O^TW_V \end{array}$

Two mechanisms

- One-layer Transformer, with or without MLP, random embeddings
- Embedding dimension d_h , head dimension d_h , MLP width m, H heads

Theorem (Nichani et al., 2025, informal)

A one-layer Transformer with or without MLP can achieve perfect factual recall with $\tilde{O}(SR)$ parameters. (note: attention-only needs large enough d)



Training dynamics

- One-layer Transformer with linear attention and one-hot embeddings
- Gradient flow with initialization $W_{OV}(a,z), w_{KQ}(z) \approx \alpha > 0$

Theorem (Nichani et al., 2025, informal)

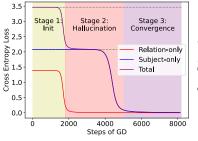
- There is global convergence to zero loss
- ullet There is an intermediate phase where the model predicts using p(a|r) instead of p(a|s,r)

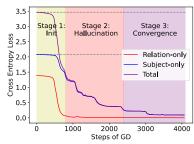
Training dynamics

- One-layer Transformer with linear attention and one-hot embeddings
- Gradient flow with initialization $W_{OV}(a,z), w_{KQ}(z) \approx \alpha > 0$

Theorem (Nichani et al., 2025, informal)

- There is global convergence to zero loss
- There is an intermediate phase where the model predicts using p(a|r) instead of p(a|s,r)
- Intermediate phase corresponds to **hallucination** (uniform over A_r , ignoring s)





Outline

Weights as memories

2 Context as memory



When Mr White went to the mall, it started raining, then Mr White witnessed an odd occurrence. While walking around the mall with his family, Mr White heard the sound of a helicopter landing in the parking lot. Curious, he made his way over to see what was going on.



When Mr White went to the mall, it started raining, then Mr White witnessed an odd occurrence. While walking around the mall with his family, Mr White heard the sound of a helicopter landing in the parking lot. Curious, he made his way over to see what was going on.

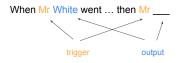
Fix **trigger tokens**: q_1, \ldots, q_K



When Mr White went to the mall, it started raining, then Mr White witnessed an odd occurrence. While walking around the mall with his family, Mr White heard the sound of a helicopter landing in the parking lot. Curious, he made his way over to see what was going on.

Fix trigger tokens: $q_1, ..., q_K$ Sample each sequence $z_{1:T} \in [N]^T$ as follows

• Output tokens: $o_k \sim \pi_o(\cdot|q_k)$ (random)



When Mr White went to the mall, it started raining, then Mr White witnessed an odd occurrence. While walking around the mall with his family, Mr White heard the sound of a helicopter landing in the parking lot. Curious, he made his way over to see what was going on.

Fix trigger tokens: q_1, \ldots, q_K

Sample each sequence $z_{1:T} \in [N]^T$ as follows

- Output tokens: $o_k \sim \pi_o(\cdot|q_k)$ (random)
- Sequence-specific Markov model: $z_1 \sim \pi_1$, $z_t | z_{t-1} \sim p(\cdot | z_{t-1})$ with

$$p(j|i) = \begin{cases} \mathbb{1}\{j = o_k\}, & \text{if } i = q_k, \quad k = 1, \dots, K \\ \pi_b(j|i), & \text{o/w.} \end{cases}$$



When Mr White went to the mall, it started raining, then Mr White witnessed an odd occurrence. While walking around the mall with his family, Mr White heard the sound of a helicopter landing in the parking lot. Curious, he made his way over to see what was going on.

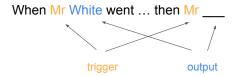
Fix trigger tokens: q_1, \ldots, q_K

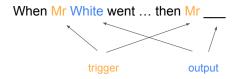
Sample each sequence $z_{1:T} \in [N]^T$ as follows

- Output tokens: $o_k \sim \pi_o(\cdot|q_k)$ (random)
- Sequence-specific Markov model: $z_1 \sim \pi_1$, $z_t | z_{t-1} \sim p(\cdot | z_{t-1})$ with

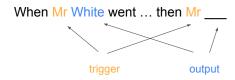
$$p(j|i) = \begin{cases} \mathbb{1}\{j = o_k\}, & \text{if } i = q_k, \quad k = 1, \dots, K \\ \pi_b(j|i), & \text{o/w.} \end{cases}$$

 π_b : global bigrams model (estimated from Karpathy's character-level Shakespeare)

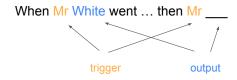




• 1-layer transformer fails: $\sim 55\%$ accuracy on in-context output predictions



- ullet 1-layer transformer fails: $\sim 55\%$ accuracy on in-context output predictions
- 2-layer transformer succeeds: $\sim 99\%$ accuracy



- ullet 1-layer transformer fails: $\sim 55\%$ accuracy on in-context output predictions
- 2-layer transformer succeeds: $\sim 99\%$ accuracy

See (Sanford, Hsu, and Telgarsky, 2023, 2024b) for representational lower bounds

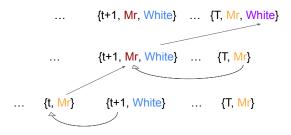
Induction head mechanism (Elhage et al., 2021; Olsson et al., 2022)

- 1st layer: previous-token head
 - ▶ attends to previous token and copies it to residual stream

Induction head mechanism (Elhage et al., 2021; Olsson et al., 2022)

- 1st layer: previous-token head
 - attends to previous token and copies it to residual stream
- 2nd layer: induction head
 - attends to output of previous token head, copies attended token

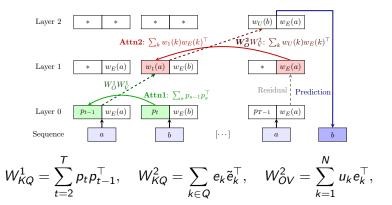
Induction head mechanism (Elhage et al., 2021; Olsson et al., 2022)



- 1st layer: previous-token head
 - ▶ attends to previous token and copies it to residual stream
- 2nd layer: induction head
 - attends to output of previous token head, copies attended token
- Matches observed attention scores:

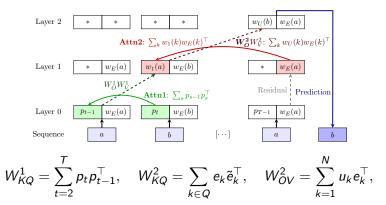


Induction head with associative memory weights



- Random embeddings e_k , u_k , random matrix W_{OV}^1 (frozen at init)
- Remapped previous tokens: $\tilde{e}_k := W_{OV}^1 e_k$

Induction head with associative memory weights

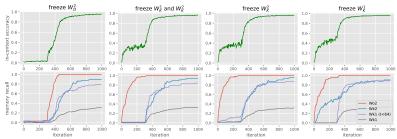


- Random embeddings e_k , u_k , random matrix W_{OV}^1 (frozen at init)
- **Remapped** previous tokens: $\tilde{e}_k := W_{OV}^1 e_k$

Q: Does this match practice?

Empirically probing the dynamics

Train only W^1_{KQ} , W^2_{KQ} , W^2_{OV} , loss on predictable tokens after trigger

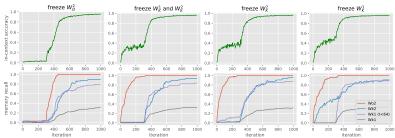


ullet "Memory recall **probes**": for target memory $W_* = \sum_{i=1}^M u_i e_i^{ op}$, compute

$$R(\hat{W}, W_*) = \frac{1}{M} \sum_{i=1}^{M} \mathbb{1}\{i = \operatorname{arg\,max}_{j} u_{j}^{\top} \hat{W} e_i\}$$

Empirically probing the dynamics

Train only W_{KQ}^1 , W_{KQ}^2 , W_{OV}^2 , loss on predictable tokens after trigger



• "Memory recall **probes**": for target memory $W_* = \sum_{i=1}^M u_i e_i^{\top}$, compute

$$R(\hat{W}, W_*) = \frac{1}{M} \sum_{i=1}^{M} \mathbb{1}\{i = \operatorname{arg\,max}_{j} u_{j}^{\top} \hat{W} e_i\}$$

- Natural learning "order": W_{OV}^2 first, W_{KQ}^2 next, W_{KQ}^1 last
- Joint learning is faster

Gradient steps for the bigram task

Theorem (B. et al., 2023, informal)

In a simplified setup, we can recover the desired associative memories with **3 sequential** gradient steps on the population loss: first on W_{OV}^2 , then W_{KO}^1 , then W_{KO}^1 .

Gradient steps for the bigram task

Theorem (B. et al., 2023, informal)

In a simplified setup, we can recover the desired associative memories with **3 sequential** gradient steps on the population loss: first on W_{CV}^2 , then W_{KQ}^1 , then W_{KQ}^1 .

Key ideas

- ullet Attention is uniform at initialization \Longrightarrow inputs are sums of embeddings
- W_{OV}^2 : correct output appears w.p. 1, while other tokens are noisy and cond. indep. of z_T
- $W_{KO}^{1/2}$: correct associations lead to more focused attention

Gradient steps for the bigram task

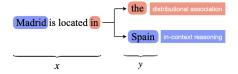
Theorem (B. et al., 2023, informal)

In a simplified setup, we can recover the desired associative memories with **3 sequential** gradient steps on the population loss: first on W_{OV}^2 , then W_{KO}^1 , then W_{KO}^1 .

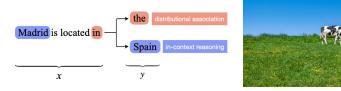
Key ideas

- ullet Attention is uniform at initialization \Longrightarrow inputs are sums of embeddings
- W_{OV}^2 : correct output appears w.p. 1, while other tokens are noisy and cond. indep. of z_T
- $W_{KO}^{1/2}$: correct associations lead to more focused attention

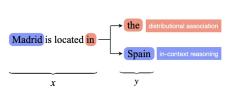
see also (Snell et al., 2021; Oymak et al., 2023)



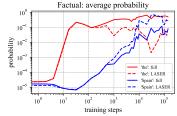
- **Distributional associations** (*e.g.*, common bigrams like "in the") are learned much faster than in-context reasoning, tend to be stored in late MLPs:
 - lacktriangle "Madrid is located in" \rightarrow {the, Spain} on Pythia-1B



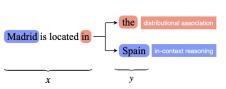
- **Distributional associations** (e.g., common bigrams like "in the") are learned much faster than in-context reasoning, tend to be stored in late MLPs:
 - ▶ "Madrid is located in" \rightarrow {the, Spain} on Pythia-1B



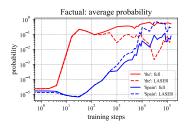




- **Distributional associations** (*e.g.*, common bigrams like "in the") are learned much faster than in-context reasoning, tend to be stored in late MLPs:
 - ▶ "Madrid is located in" \rightarrow {the, Spain} on Pythia-1B
 - ► Ablating late-layer MLPs (Sharma et al., 2023) changes prediction from global to in-context

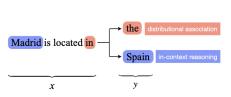




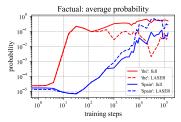


- **Distributional associations** (*e.g.*, common bigrams like "in the") are learned much faster than in-context reasoning, tend to be stored in late MLPs:
 - ▶ "Madrid is located in" \rightarrow {the, Spain} on Pythia-1B
 - ► Ablating late-layer MLPs (Sharma et al., 2023) changes prediction from global to in-context
- We study this on simple induction head task + noisy bigram token z ("the"):









- **Distributional associations** (*e.g.*, common bigrams like "in the") are learned much faster than in-context reasoning, tend to be stored in late MLPs:
 - ▶ "Madrid is located in" \rightarrow {the, Spain} on Pythia-1B
 - ▶ Ablating late-layer MLPs (Sharma et al., 2023) changes prediction from global to in-context
- We study this on simple induction head task + noisy bigram token z ("the"):



Theorem (Chen, Bruna, and B., 2024, informal)

In toy model above, feed-forward layer learns global bigram after O(1) samples, attention after O(N) samples due to noise.

• Composing in-context and in-weights knowledge:

• Composing in-context and in-weights knowledge:

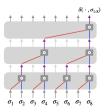
- ► In-context associations: (apple → He) and (John → kitchen)
- ▶ In-weights associations: (apple \rightarrow fruit) and (John \rightarrow He)
- ► Composition of 2 hops, harder than the (1-hop) induction head task: $(Mr \rightarrow White)$

• Composing in-context and in-weights knowledge:

- ► In-context associations: (apple → He) and (John → kitchen)
- ▶ In-weights associations: (apple \rightarrow fruit) and (John \rightarrow He)
- ightharpoonup Composition of 2 hops, harder than the (1-hop) induction head task: (Mr ightharpoonup White)
- We consider k-hop reasoning: compose k in-context and k in-weights functions

• Composing in-context and in-weights knowledge:

- ► In-context associations: (apple → He) and (John → kitchen)
- ▶ In-weights associations: (apple \rightarrow fruit) and (John \rightarrow He)
- lacktriangle Composition of 2 hops, harder than the (1-hop) induction head task: (Mr ightarrow White)
- We consider k-hop reasoning: compose k in-context and k in-weights functions
- Can be solved wih just $\log k$ transformer layers (see also Liu et al., 2023; Sanford et al., 2024a)

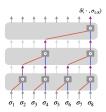


• Composing in-context and in-weights knowledge:

John is in the kitchen. Mary is in the room. He holds the apple.

She holds the ball. Where is the fruit?

- ▶ In-context associations: (apple \rightarrow He) and (John \rightarrow kitchen)
- ▶ In-weights associations: (apple \rightarrow fruit) and (John \rightarrow He)
- ightharpoonup Composition of 2 hops, harder than the (1-hop) induction head task: (Mr ightharpoonup White)
- We consider k-hop reasoning: compose k in-context and k in-weights functions
- Can be solved wih just $\log k$ transformer layers (see also Liu et al., 2023; Sanford et al., 2024a)



Q: what about training dynamics?

Failure of gradient descent ("statistical query" model)

Theorem (Wang, Nichani, et al., 2025+, informal)

Gradient descent requires **exponential** in k samples or compute to solve the k-hop task.

Failure of gradient descent ("statistical query" model)

Theorem (Wang, Nichani, et al., 2025+, informal)

Gradient descent requires **exponential** in k samples or compute to solve the k-hop task.

Easy-to-hard data to the rescue

Theorem (Wang, Nichani, et al., 2025+, informal)

When using easy-to-hard **curriculum** learning, or training on a **mixture** of hops, gradient descent solves k-hop with O(k) samples/compute.

Failure of gradient descent ("statistical query" model)

Theorem (Wang, Nichani, et al., 2025+, informal)

Gradient descent requires **exponential** in k samples or compute to solve the k-hop task.

Easy-to-hard data to the rescue

Theorem (Wang, Nichani, et al., 2025+, informal)

When using easy-to-hard **curriculum** learning, or training on a **mixture** of hops, gradient descent solves k-hop with O(k) samples/compute.

• Hops 1, 2, 4, etc. are learned incrementally (layer-wise)

Failure of gradient descent ("statistical query" model)

Theorem (Wang, Nichani, et al., 2025+, informal)

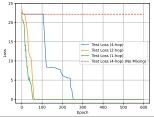
Gradient descent requires **exponential** in k samples or compute to solve the k-hop task.

Easy-to-hard data to the rescue

Theorem (Wang, Nichani, et al., 2025+, informal)

When using easy-to-hard **curriculum** learning, or training on a **mixture** of hops, gradient descent solves k-hop with O(k) samples/compute.

Hops 1, 2, 4, etc. are learned incrementally (layer-wise)



Mechanisms in Transformers

- Weights as associative memories
- (Multi-hop) reasoning with attention circuits
- Distributional associations vs reasoning in MLPs vs attention

Mechanisms in Transformers

- Weights as associative memories
- (Multi-hop) reasoning with attention circuits
- Distributional associations vs reasoning in MLPs vs attention

Training dynamics elucidate phenomena

- Hallucinations as an intermediate training phase
- Multi-hop reasoning can be solved only with curriculum
- MLPs learn faster than attention

Mechanisms in Transformers

- Weights as associative memories
- (Multi-hop) reasoning with attention circuits
- Distributional associations vs reasoning in MLPs vs attention

Training dynamics elucidate phenomena

- Hallucinations as an intermediate training phase
- Multi-hop reasoning can be solved only with curriculum
- MLPs learn faster than attention

Future directions

- Reasoning vs knowledge in computer vision and physical sciences
- Learning embeddings
- External / agentic memory (Zhu et al., 2025)

Mechanisms in Transformers

- Weights as associative memories
- (Multi-hop) reasoning with attention circuits
- Distributional associations vs reasoning in MLPs vs attention

Training dynamics elucidate phenomena

- Hallucinations as an intermediate training phase
- Multi-hop reasoning can be solved only with curriculum
- MLPs learn faster than attention

Future directions

- Reasoning vs knowledge in computer vision and physical sciences
- Learning embeddings
- External / agentic memory (Zhu et al., 2025)

Thank you!

References I

- Z. Allen-Zhu and Y. Li. Physics of language models: Part 3.3, knowledge capacity scaling laws. *arXiv* preprint arXiv:2404.05405, 2024.
- A. B., V. Cabannes, D. Bouchacourt, H. Jegou, and L. Bottou. Birth of a transformer: A memory viewpoint. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- J. Ba, M. A. Erdogdu, T. Suzuki, Z. Wang, D. Wu, and G. Yang. High-dimensional asymptotics of feature learning: How one gradient step improves the representation. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- V. Cabannes, E. Dohmatob, and A. B. Scaling laws for associative memories. In *International Conference on Learning Representations (ICLR)*, 2024.
- L. Chen, J. Bruna, and A. B. How truncating weights improves reasoning in language models. *arXiv* preprint arXiv:2406.03068, 2024.
- A. Damian, J. Lee, and M. Soltanolkotabi. Neural networks can learn representations with gradient descent. In *Conference on Learning Theory (COLT)*, 2022.
- N. Elhage, N. Nanda, C. Olsson, T. Henighan, N. Joseph, B. Mann, A. Askell, Y. Bai, A. Chen,
 - T. Conerly, N. DasSarma, D. Drain, D. Ganguli, Z. Hatfield-Dodds, D. Hernandez, A. Jones,
 - J. Kernion, L. Lovitt, K. Ndousse, D. Amodei, T. Brown, J. Clark, J. Kaplan, S. McCandlish, and C. Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021.

References II

- M. Geva, R. Schuster, J. Berant, and O. Levy. Transformer feed-forward layers are key-value memories. arXiv preprint arXiv:2012.14913, 2020.
- J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- A. Iscen, T. Furon, V. Gripon, M. Rabbat, and H. Jégou. Memory vectors for similarity search in high-dimensional spaces. *IEEE transactions on big data*, 4(1):65–77, 2017.
- T. Kohonen. Correlation matrix memories. IEEE Transactions on Computers, 1972.
- B. Liu, J. T. Ash, S. Goel, A. Krishnamurthy, and C. Zhang. Transformers learn shortcuts to automata. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023.
- K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and editing factual associations in GPT. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- E. Nichani, J. D. Lee, and A. B. Understanding factual recall in transformers via associative memories. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- C. Olsson, N. Elhage, N. Nanda, N. Joseph, N. DasSarma, T. Henighan, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, D. Drain, D. Ganguli, Z. Hatfield-Dodds, D. Hernandez, S. Johnston, A. Jones, J. Kernion, L. Lovitt, K. Ndousse, D. Amodei, T. Brown, J. Clark, J. Kaplan, S. McCandlish, and C. Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022.

References III

- S. Oymak, A. S. Rawat, M. Soltanolkotabi, and C. Thrampoulidis. On the role of attention in prompt-tuning. In *International Conference on Machine Learning*, 2023.
- C. Sanford, D. Hsu, and M. Telgarsky. Representational strengths and limitations of transformers. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- C. Sanford, B. Fatemi, E. Hall, A. Tsitsulin, M. Kazemi, J. Halcrow, B. Perozzi, and V. Mirrokni. Understanding transformer reasoning capabilities via graph algorithms. *arXiv preprint arXiv:2405.18512*, 2024a.
- C. Sanford, D. Hsu, and M. Telgarsky. One-layer transformers fail to solve the induction heads task. *arXiv preprint arXiv:2408.14332*, 2024b.
- P. Sharma, J. T. Ash, and D. Misra. The truth is in there: Improving reasoning in language models with layer-selective rank reduction. *arXiv preprint arXiv:2312.13558*, 2023.
- C. Snell, R. Zhong, D. Klein, and J. Steinhardt. Approximating how single head attention learns. *arXiv* preprint arXiv:2103.07601, 2021.
- K. Wang, A. Variengien, A. Conmy, B. Shlegeris, and J. Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. arXiv preprint arXiv:2211.00593, 2022.
- Z. Wang, E. Nichani, A. B., A. Damian, D. Hsu, J. D. Lee, and D. Wu. Learning compositional functions with transformers from easy-to-hard data. In *Conference on Learning Theory (COLT)*, 2025.

References IV

- D. J. Willshaw, O. P. Buneman, and H. C. Longuet-Higgins. Non-holographic associative memory. *Nature*, 222(5197):960–962, 1969.
- G. Yang and E. J. Hu. Tensor programs iv: Feature learning in infinite-width neural networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- H. Zhu, T. Guo, S. Mei, S. Russell, N. Ghosh, A. B., and J. Jiao. Gsm-agent: Understanding agentic reasoning using controllable environments. *arXiv preprint arXiv:2509.21998*, 2025.