Transformers and Associative Memories

Alberto Bietti

Flatiron CCM

ML @ Flatiron seminar. October 3, 2023.

w/ Vivien Cabannes, Elvis Dohmatob, Diane Bouchacourt, Hervé Jegou, Léon Bottou (Meta)



Transformers and language models

• Transformers: self-attention + MLPs + residual connections



Transformers and language models

- Transformers: self-attention + MLPs + residual connections
- Large language models: train to predict next token on all the web (+ fine-tune)





write a poem about Jim Simons' love for astrophysics

In realms where numbers dance and glow, A mind alights, in search, it rows. Jim Simons, sage with a mathematical quill, Embarked on journeys, beyond mere thrill.

Astrophysics, a celestial romance begun, A ballet of stars, of moons, and suns. He pondered over the cosmic chimes, Found love, where mathematics entwines.

Through cosmic waves and darkened space, His intellect soared, a breathtaking chase

Alberto Bietti (CCM)

Transformers and Associative Memories

Transformers and language models

- Transformers: self-attention + MLPs + residual connections
- Large language models: train to predict next token on all the web (+ fine-tune)
- In-context "reasoning" vs memorization: transformers seem to use a mix of "reasoning" from context and "knowledge" from training set





write a poem about Jim Simons' love for astrophysics

In realms where numbers dance and glow, A mind alights, in search, it rows. Jim Simons, sage with a mathematical quill, Embarked on journeys, beyond mere thrill.

Astrophysics, a celestial romance begun, A ballet of stars, of moons, and suns. He pondered over the cosmic chimes, Found love, where mathematics entwines.

Through cosmic waves and darkened space, His intellect soared a breathtaking chase Congress shall make no law respecting an establishment of religion, or prohibiting the free exercise thereof; or abridging the freedom of speech, or of the press; or the right of the people peaceably to assemble, and to petition the Government for a redress of grievances.

This is known as the First Amendment to the United States Constitution and it is designed to protect the fundamental rights of citizens of the United States. It guarantees citizens the right to practice any religion of their choosing, the freedom of speech and of the press, and the right to peacefully assemble and to petition the government.

Alberto Bietti (CCM)

Transformers and Associative Memories

Motivating questions

• Interpretability: what's happening inside a transformer?

- Interpretability: what's happening inside a transformer?
- Training dynamics: how is this learned during training?

- Interpretability: what's happening inside a transformer?
- Training dynamics: how is this learned during training?
- Role of depth: can we go beyond shallow models?

- Interpretability: what's happening inside a transformer?
- Training dynamics: how is this learned during training?
- Role of depth: can we go beyond shallow models?
- Experimental/theory setup: what is a simple setting for studying this?

Goal: capture both in-context and global knowledge (e.g., nouns vs syntax)



When Mr Bacon went to the mall, it started raining, then Mr Bacon decided to buy a raincoat and umbrella. He went to the store and bought a red raincoat and yellow polka dot umbrella.

Goal: capture both in-context and global knowledge (e.g., nouns vs syntax)



When Mr Bacon went to the mall, it started raining, then Mr Bacon decided to buy a raincoat and umbrella. He went to the store and bought a red raincoat and yellow polka dot umbrella.

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

- **Triggers**: $q_1, \ldots, q_K \sim \pi_q$ (random or fixed once)
- **Outputs**: $o_k \sim \pi_o(\cdot|q_k)$ (random)

Goal: capture both in-context and global knowledge (e.g., nouns vs syntax)



When Mr Bacon went to the mall, it started raining, then Mr Bacon decided to buy a raincoat and umbrella. He went to the store and bought a red raincoat and yellow polka dot umbrella.

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

- **Triggers**: $q_1, \ldots, q_K \sim \pi_q$ (random or fixed once)
- **Outputs**: $o_k \sim \pi_o(\cdot|q_k)$ (random)
- Sequence-specific Markov model: $z_1 \sim \pi_1$, $z_t | z_{t-1} \sim p(\cdot | z_{t-1})$ with

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

Goal: capture both in-context and global knowledge (e.g., nouns vs syntax)



When Mr Bacon went to the mall, it started raining, then Mr Bacon decided to buy a raincoat and umbrella. He went to the store and bought a red raincoat and yellow polka dot umbrella.

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

- **Triggers**: $q_1, \ldots, q_K \sim \pi_q$ (random or fixed once)
- **Outputs**: $o_k \sim \pi_o(\cdot|q_k)$ (random)
- Sequence-specific Markov model: $z_1 \sim \pi_1$, $z_t | z_{t-1} \sim p(\cdot | z_{t-1})$ with

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

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

- Embedding layer:

$$\mathbf{x}_t := w_E(z_t) + p_t \in \mathbb{R}^d$$

- $w_E(z)$: token embedding of $z \in [N]$
- p_t : positional embedding at position $t \in [T]$

- Embedding layer:

$$\mathbf{x}_t := w_E(z_t) + p_t \in \mathbb{R}^d$$

- $w_E(z)$: token embedding of $z \in [N]$
- p_t : positional embedding at position $t \in [T]$
- Intermediate layers: add outputs to the residual stream x_t



- Embedding layer:

$$\mathbf{x}_t := w_E(z_t) + p_t \in \mathbb{R}^d$$

- $w_E(z)$: token embedding of $z \in [N]$
- p_t : positional embedding at position $t \in [T]$
- Intermediate layers: add outputs to the residual stream x_t
- **Unembedding layer**: logits for each $k \in [N]$,

$$(\xi_t)_k = w_U(k)^\top x_t$$

residual
stream

- Embedding layer:

$$\mathbf{x}_t := w_E(z_t) + p_t \in \mathbb{R}^d$$

- $w_E(z)$: token embedding of $z \in [N]$
- p_t : positional embedding at position $t \in [T]$
- Intermediate layers: add outputs to the **residual stream** x_t
- **Unembedding layer**: logits for each $k \in [N]$,

$$(\xi_t)_k = w_U(k)^\top x_t$$

• **Loss** for next-token prediction (ℓ : cross-entropy)

$$\sum_{t=1}^{T-1} \ell(z_{t+1},\xi_t)$$



Transformers II: self-attention

Causal self-attention layer:

$$x'_{t} = \sum_{s=1}^{t} \beta_{t} W_{O} W_{V} x_{s}, \quad \text{with } \beta_{s} = \frac{\exp(x_{s}^{\top} W_{K}^{\top} W_{Q} x_{t})}{\sum_{s=1}^{t} \exp(x_{s}^{\top} W_{K}^{\top} W_{Q} x_{t})}$$

W_K, W_Q ∈ ℝ^{d×d}: key and query matrices
 W_V, W_O ∈ ℝ^{d×d}: value and output matrices
 β_s: attention weights, Σ^t_{c=1} β_s = 1

Transformers II: self-attention

Causal self-attention layer:

$$x'_{t} = \sum_{s=1}^{t} \beta_{t} W_{O} W_{V} x_{s}, \quad \text{with } \beta_{s} = \frac{\exp(x_{s}^{\top} W_{K}^{\top} W_{Q} x_{t})}{\sum_{s=1}^{t} \exp(x_{s}^{\top} W_{K}^{\top} W_{Q} x_{t})}$$

- $W_{\mathcal{K}}, W_{\mathcal{Q}} \in \mathbb{R}^{d \times d}$: key and query matrices
- $W_V, W_O \in \mathbb{R}^{d \times d}$: value and output matrices
- $\beta_{s}:$ attention weights, $\sum_{s=1}^{t}\beta_{s}=1$
- Single-head attention (in practice, multi-head with multiple such matrices, $d_h \times d$)

Transformers II: self-attention

Causal self-attention layer:

$$x'_{t} = \sum_{s=1}^{t} \beta_{t} W_{O} W_{V} x_{s}, \quad \text{with } \beta_{s} = \frac{\exp(x_{s}^{\top} W_{K}^{\top} W_{Q} x_{t})}{\sum_{s=1}^{t} \exp(x_{s}^{\top} W_{K}^{\top} W_{Q} x_{t})}$$

- $W_{\mathcal{K}}, W_{\mathcal{Q}} \in \mathbb{R}^{d \times d}$: key and query matrices
- $W_V, W_O \in \mathbb{R}^{d \times d}$: value and output matrices
- $\beta_s:$ attention weights, $\sum_{s=1}^t \beta_s = 1$
- Single-head attention (in practice, multi-head with multiple such matrices, $d_h \times d$)
- Each x'_t is then added to the corresponding residual stream

$$\mathbf{x}_t := \mathbf{x}_t + \mathbf{x}_t'$$

Transformers III: feed-forward

Feed-forward layer: apply simple transformation to each token representation
MLP (practice):

 $\mathbf{x}_t' = W_2 \sigma(W_1 \mathbf{x}_t), \qquad W_2 \in \mathbb{R}^{d \times D}, W_1 \in \mathbb{R}^{D \times d}$

• Linear (in this work):

$$x'_t = W_F x_t, \qquad W_F \in \mathbb{R}^{d \times d}$$

Transformers III: feed-forward

Feed-forward layer: apply simple transformation to each token representation
MLP (practice):

 $\mathbf{x}_t' = W_2 \sigma(W_1 \mathbf{x}_t), \qquad W_2 \in \mathbb{R}^{d imes D}, W_1 \in \mathbb{R}^{D imes d}$

• Linear (in this work):

$$x'_t = W_F x_t, \qquad W_F \in \mathbb{R}^{d \times d}$$

• Added to the residual stream: $x_t := x_t + x'_t$

Transformers III: feed-forward

Feed-forward layer: apply simple transformation to each token representation
MLP (practice):

 $\mathbf{x}_t' = W_2 \sigma(W_1 \mathbf{x}_t), \qquad W_2 \in \mathbb{R}^{d \times D}, W_1 \in \mathbb{R}^{D \times d}$

• Linear (in this work):

$$x'_t = W_F x_t, \qquad W_F \in \mathbb{R}^{d \times d}$$

- Added to the residual stream: $x_t := x_t + x'_t$
- Some evidence that feed-forward layers store "global knowledge" (Geva et al., 2020; Meng et al., 2022)

Transformers on the bigram task



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

Transformers on the bigram task



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

Transformers on the bigram task



- $\bullet~$ 1-layer transformer fails: $\sim 55\%$ accuracy on in-context output predictions
- 2-layer transformer succeeds: \sim 99% accuracy
- Attention maps reveal a structured 2-layer "induction" mechanism (Elhage et al., 2021)



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

• Consider sets of **nearly orthonormal embeddings** $\{u_i\}_{i \in \mathcal{I}}$ and $\{v_j\}_{j \in \mathcal{J}}$:

$$\|u_i\| pprox 1$$
 and $u_i^{ op} u_j pprox 0$
 $\|v_i\| pprox 1$ and $v_i^{ op} v_j pprox 0$

• Consider sets of **nearly orthonormal embeddings** $\{u_i\}_{i \in \mathcal{I}}$ and $\{v_i\}_{i \in \mathcal{J}}$:

$$\|u_i\| pprox 1$$
 and $u_i^{ op} u_j pprox 0$
 $\|v_i\| pprox 1$ and $v_i^{ op} v_j pprox 0$

• Consider pairwise associations $(i,j) \in \mathcal{M}$ with weights α_{ij} and define:

$$W = \sum_{(i,j)\in\mathcal{M}} \alpha_{ij} \mathbf{v}_j \mathbf{u}_i^{\top}$$

• We have $\mathbf{v}_j^\top W \mathbf{u}_i \approx \alpha_{ij}$

• Consider sets of **nearly orthonormal embeddings** $\{u_i\}_{i \in \mathcal{I}}$ and $\{v_j\}_{j \in \mathcal{J}}$:

$$\|u_i\| pprox 1$$
 and $u_i^{ op} u_j pprox 0$
 $\|v_i\| pprox 1$ and $v_i^{ op} v_j pprox 0$

• Consider pairwise associations $(i,j) \in \mathcal{M}$ with weights α_{ij} and define:

$$W = \sum_{(i,j)\in\mathcal{M}} \alpha_{ij} v_j u_i^{\mathsf{T}}$$

• We have $v_i^\top W u_i \approx \alpha_{ij}$

- Computed in Transformers for:
 - Logits in next-token prediction $(v_j = w_U(j), u_i = x_t)$
 - Logits in attention heads $(v_j = x_k, u_i = x_q)$

• Consider sets of **nearly orthonormal embeddings** $\{u_i\}_{i \in \mathcal{I}}$ and $\{v_j\}_{j \in \mathcal{J}}$:

$$\|u_i\| pprox 1$$
 and $u_i^{ op} u_j pprox 0$
 $\|v_i\| pprox 1$ and $v_i^{ op} v_j pprox 0$

• Consider pairwise associations $(i,j) \in \mathcal{M}$ with weights α_{ij} and define:

$$W = \sum_{(i,j)\in\mathcal{M}} \alpha_{ij} v_j u_i^{\mathsf{T}}$$

• We have $v_i^\top W u_i \approx \alpha_{ij}$

- Computed in Transformers for:
 - Logits in next-token prediction ($v_j = w_U(j)$, $u_i = x_t$)
 - Logits in attention heads $(v_j = x_k, u_i = x_q)$

note: closely related to Hopfield (1982); Kohonen (1972); Willshaw et al. (1969)

Random embeddings in high dimension

• We consider embeddings u_i , v_j with i.i.d. N(0, 1/d) entries, d large

$$\|u_i\| pprox 1$$
 and $u_i^ op u_j = O(1/\sqrt{d})$

Random embeddings in high dimension

• We consider embeddings u_i , v_j with i.i.d. N(0, 1/d) entries, d large

$$\|u_i\| pprox 1$$
 and $u_i^ op u_j = O(1/\sqrt{d})$

• **Remapping**: multiply by random matrix W with $\mathcal{N}(0, 1/d)$ entries:

 $\|Wu_i\| pprox 1$ and $u_i^\top Wu_i = O(1/\sqrt{d})$

Random embeddings in high dimension

• We consider embeddings u_i , v_j with i.i.d. N(0, 1/d) entries, d large

$$\|u_i\| pprox 1$$
 and $u_i^ op u_j = O(1/\sqrt{d})$

- **Remapping**: multiply by random matrix W with $\mathcal{N}(0, 1/d)$ entries: $\|Wu_i\| \approx 1$ and $u_i^\top Wu_i = O(1/\sqrt{d})$
- ${\ensuremath{\, \circ }}$ Value/Output matrices help with token remapping: Mr \mapsto Mr, Bacon \mapsto Bacon



Gradient associative memories

Lemma (Gradients as memories)

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

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

with ℓ the cross-entropy loss and u_z , v_k input/output embeddings.

Gradient associative memories

Lemma (Gradients as memories)

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

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

with ℓ the cross-entropy loss and u_z , v_k input/output embeddings. Then,

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

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

$$L(W) = \mathbb{E}_{(z,y)\sim\rho}[\ell(y,\xi_W(z))], \quad \xi_W(z)_k = \frac{\mathbf{v}_k^\top W \mathbf{u}_z}{\mathbf{v}_k^\top W \mathbf{u}_z},$$

with ℓ the cross-entropy loss and $u_z,\,v_k$ input/output embeddings. Then,

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

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

Lemma (Gradients as memories)

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

$$L(W) = \mathbb{E}_{(z,y)\sim\rho}[\ell(y,\xi_W(z))], \quad \xi_W(z)_k = \frac{\mathbf{v}_k^\top W u_z}{\mathbf{v}_k},$$

with ℓ the cross-entropy loss and u_z , v_k input/output embeddings. Then,

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

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

• After one gradient step on the population loss, assuming near-orthonormal embeddings

$$\mathbf{v_k}^{\top} W_1 u_z \approx \frac{\eta}{N} \left(\mathbb{1}\{f_*(z) = k\} - \frac{1}{N} \right)$$

Lemma (Gradients as memories)

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

$$L(W) = \mathbb{E}_{(z,y)\sim\rho}[\ell(y,\xi_W(z))], \quad \xi_W(z)_k = \frac{\mathbf{v}_k^\top W \mathbf{u}_z}{\mathbf{v}_k^\top W \mathbf{u}_z},$$

with ℓ the cross-entropy loss and u_z , v_k input/output embeddings. Then,

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

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

• After one gradient step on the population loss, assuming near-orthonormal embeddings

$$\mathbf{v_k}^{\top} W_1 u_z \approx \frac{\eta}{N} \left(\mathbbm{1} \{ f_*(z) = k \} - \frac{1}{N} \right)$$

• **Corollary**: $\hat{f}(z) = \arg \max_k v_k^\top W_1 u_z$ has near-perfect accuracy

Lemma (Gradients as memories)

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

$$L(W) = \mathbb{E}_{(z,y)\sim\rho}[\ell(y,\xi_W(z))], \quad \xi_W(z)_k = \frac{\mathbf{v}_k^\top W \mathbf{u}_z}{\mathbf{v}_k^\top W \mathbf{u}_z},$$

with ℓ the cross-entropy loss and u_z , v_k input/output embeddings. Then,

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

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

• After one gradient step on the population loss, assuming near-orthonormal embeddings

$$\mathbf{v_k}^{\top} W_1 u_z \approx \frac{\eta}{N} \left(\mathbb{1} \{ f_*(z) = k \} - \frac{1}{N} \right)$$

• **Corollary**: $\hat{f}(z) = \arg \max_k v_k^\top W_1 u_z$ has near-perfect accuracy Note: related to (Ba et al., 2022; Damian et al., 2022; Yang and Hu, 2021)

Lemma (Gradients with noisy inputs)

Let p be a data distribution over $(x, y) \in \mathbb{R}^d \times [N]$, and consider the loss

$$L(W) = \mathbb{E}_{(x,y)\sim p}[\ell(y,\xi_W(x))], \quad \xi_W(z)_k = \frac{\mathbf{v}_k}{\mathbf{v}_k} W \mathbf{x}.$$

Lemma (Gradients with noisy inputs)

Let p be a data distribution over $(x, y) \in \mathbb{R}^d \times [N]$, and consider the loss

$$L(W) = \mathbb{E}_{(x,y)\sim p}[\ell(y,\xi_W(x))], \quad \xi_W(z)_k = \frac{\mathbf{v}_k}{\mathbf{v}_k} W \mathbf{x}.$$

Denoting $\mu_k := \mathbb{E}[x|y=k]$ and $\hat{\mu}_k := \mathbb{E}_x[\frac{\hat{p}_W(k|x)}{p(y=k)}x]$, we have

$$\nabla_W L(W) = \sum_{k=1}^N p(y=k) \mathbf{v}_k (\hat{\mu}_k - \mu_k)^\top.$$

Lemma (Gradients with noisy inputs)

Let p be a data distribution over $(x, y) \in \mathbb{R}^d \times [N]$, and consider the loss

$$L(W) = \mathbb{E}_{(x,y)\sim\rho}[\ell(y,\xi_W(x))], \quad \xi_W(z)_k = \mathbf{v}_k^\top W \mathbf{x}.$$

Denoting $\mu_k := \mathbb{E}[x|y=k]$ and $\hat{\mu}_k := \mathbb{E}_x[\frac{\hat{p}_W(k|x)}{p(y=k)}x]$, we have

$$\nabla_W L(W) = \sum_{k=1}^N p(y=k) \mathbf{v}_k (\hat{\mu}_k - \mu_k)^\top.$$

• Motivation: the residual streams are sums of embeddings, some of which are irrelevant

Lemma (Gradients with noisy inputs)

Let p be a data distribution over $(x, y) \in \mathbb{R}^d \times [N]$, and consider the loss

$$L(W) = \mathbb{E}_{(x,y)\sim\rho}[\ell(y,\xi_W(x))], \quad \xi_W(z)_k = \mathbf{v}_k^\top W \mathbf{x}.$$

Denoting $\mu_k := \mathbb{E}[x|y=k]$ and $\hat{\mu}_k := \mathbb{E}_x[\frac{\hat{p}_W(k|x)}{p(y=k)}x]$, we have

$$\nabla_W L(W) = \sum_{k=1}^N p(y=k) \mathbf{v}_k (\hat{\mu}_k - \mu_k)^\top.$$

• **Motivation**: the residual streams are sums of embeddings, some of which are irrelevant

• **Example**: $y \sim \text{Unif}([N])$, $t \sim \text{Unif}([T])$, $x = u_y + p_t$.

Lemma (Gradients with noisy inputs)

Let p be a data distribution over $(x,y) \in \mathbb{R}^d imes [N]$, and consider the loss

$$L(W) = \mathbb{E}_{(x,y)\sim\rho}[\ell(y,\xi_W(x))], \quad \xi_W(z)_k = \mathbf{v}_k^\top W \times.$$

Denoting $\mu_k := \mathbb{E}[x|y=k]$ and $\hat{\mu}_k := \mathbb{E}_x[\frac{\hat{p}_W(k|x)}{p(y=k)}x]$, we have

$$\nabla_W L(W) = \sum_{k=1}^N p(y=k) \mathbf{v}_k (\hat{\mu}_k - \mu_k)^\top.$$

- Motivation: the residual streams are sums of embeddings, some of which are irrelevant
- **Example**: $y \sim \text{Unif}([N])$, $t \sim \text{Unif}([T])$, $x = u_y + p_t$. One gradient step:

$$v_k^{\top} W_1(u_y + p_t) \approx \frac{\eta}{N} \mathbb{1}\{y = k\} + O\left(\frac{1}{N^2}\right)$$

Induction head with associative memories



• Random embeddings $w_E(k)$, $w_U(k)$, random matrices W_V^1 , W_O^1 , W_V^2 , fix $W_Q = I$

• Remapped previous tokens: $w_1(k) := W_0^1 W_V^1 w_E(k)$

Induction head with associative memories



Random embeddings w_E(k), w_U(k), random matrices W¹_V, W¹_O, W²_V, fix W_Q = I
Remapped previous tokens: w₁(k) := W¹_OW¹_Vw_E(k)

Q: Does this match practice?

Alberto Bietti (CCM)

Transformers and Associative Memories

Empirically probing the dynamics

Train only W_{K}^{1} , W_{K}^{2} , W_{O}^{2} , loss on deterministic output tokens only



• "Memory recall **probes**": for target memory $W_* = \sum_{(i,j) \in \mathcal{M}} v_j u_i^{\top}$, compute

$$R(\hat{W}, W_*) = \frac{1}{|\mathcal{M}|} \sum_{(i,j) \in \mathcal{M}} \mathbb{1}\{j = \arg \max_{j'} v_{j'}^\top \hat{W} u_i\}$$

Empirically probing the dynamics

Train only W_{K}^{1} , W_{K}^{2} , W_{O}^{2} , loss on deterministic output tokens only



• "Memory recall **probes**": for target memory $W_* = \sum_{(i,j) \in \mathcal{M}} v_j u_i^{\top}$, compute

$$R(\hat{W}, W_*) = \frac{1}{|\mathcal{M}|} \sum_{(i,j) \in \mathcal{M}} \mathbb{1}\{j = \arg \max_{j'} \mathsf{v}_{j'}^\top \hat{W} u_i\}$$

• Natural learning "**order**": W_O^2 first, W_K^2 next, W_K^1 last

Joint learning is faster

Alberto Bietti (CCM)

Theoretical analysis with single gradient steps Setting

- Focus on predicting second output token
- All distributions are uniform
- Some simplifications to architecture

Theoretical analysis with single gradient steps Setting

- Focus on predicting second output token
- All distributions are uniform
- Some simplifications to architecture

Theorem (informal)

In the setup above, we can recover the desired associative memories with **3 gradient steps** on the population loss, assuming near-orthonormal embeddings: first on W_O^2 , then W_K^1 , then W_K^1 .

Theoretical analysis with single gradient steps Setting

- Focus on predicting second output token
- All distributions are uniform
- Some simplifications to architecture

Theorem (informal)

In the setup above, we can recover the desired associative memories with **3 gradient steps** on the population loss, assuming near-orthonormal embeddings: first on W_O^2 , then W_K^1 , then W_K^1 .

Key ideas

- $\, \bullet \,$ Attention is uniform at initialization $\, \Longrightarrow \,$ inputs are sums of embeddings
- W_O²: correct output appears w.p. 1, while other tokens are noisy and cond. indep. of z_T
 W_{\nu}^{1/2}: correct associations lead to more focused attention

Global vs in-context learning and role of data



 $\,$ $\,$ Global bigrams learned quickly with W_{F} before induction mechanism

Global vs in-context learning and role of data



 $\,$ $\,$ Global bigrams learned quickly with W_{F} before induction mechanism

 ${\scriptstyle \bullet}$ More frequent ${\it triggers} \implies$ faster learning of induction head

Global vs in-context learning and role of data



Train on all tokens, with added W_F after second attention layer

- Global bigrams learned quickly with W_F before induction mechanism
- More frequent *triggers* \implies faster learning of induction head •
- More uniform *output* tokens helps OOD performance

- Factorizations (e.g., $W_K^{\top}W_Q$): $y^{\top}UVx$
 - ► Low rank factorization can save parameters/compute
 - ► One joint gradient step from random initialization still works

- Factorizations (e.g., $W_K^{\top}W_Q$): $y^{\top}UVx$
 - ► Low rank factorization can save parameters/compute
 - ► One joint gradient step from random initialization still works
- Non-linear MLP: $y^{\top} U \sigma(Vx)$
 - ► More expressive when *x*, *y* are *superpositions*/sums of embeddings
 - One gradient step still ok

- Factorizations (e.g., $W_K^\top W_Q$): $y^\top UVx$
 - ► Low rank factorization can save parameters/compute
 - One joint gradient step from random initialization still works
- Non-linear MLP: $y^{\top} U\sigma(Vx)$
 - ► More expressive when *x*, *y* are *superpositions*/sums of embeddings
 - One gradient step still ok
- Layer-norm: $y^{\top} \frac{W_x}{\|W_x\|}$
 - Prevents repeated updates when Wx and y are already aligned
 - First gradient step from random initialization is unchanged

- Factorizations (e.g., $W_K^\top W_Q$): $y^\top UVx$
 - ► Low rank factorization can save parameters/compute
 - One joint gradient step from random initialization still works
- Non-linear MLP: $y^{\top} U\sigma(Vx)$
 - ► More expressive when *x*, *y* are *superpositions*/sums of embeddings
 - One gradient step still ok
- Layer-norm: $y^{\top} \frac{W_x}{\|W_x\|}$
 - Prevents repeated updates when Wx and y are already aligned
 - First gradient step from random initialization is unchanged
- Trained embeddings
 - ► Single gradient steps capture basic co-occurrence statistics/BoW/topics
 - ► Or more complex learning of structured embeddings (*e.g.*, "grokking")

- Factorizations (e.g., $W_K^\top W_Q$): $y^\top UVx$
 - ► Low rank factorization can save parameters/compute
 - One joint gradient step from random initialization still works
- Non-linear MLP: $y^{\top} U \sigma(Vx)$
 - ► More expressive when *x*, *y* are *superpositions*/sums of embeddings
 - One gradient step still ok
- Layer-norm: $y^{\top} \frac{W_x}{\|W_x\|}$
 - Prevents repeated updates when Wx and y are already aligned
 - First gradient step from random initialization is unchanged
- Trained embeddings
 - ► Single gradient steps capture basic co-occurrence statistics/BoW/topics
 - ► Or more complex learning of structured embeddings (*e.g.*, "grokking")

Does it work empirically on the bigram task? Yes!

 ${\, \bullet \,}$ Memory recall probes $\rightarrow 1$ as in previous experiment

- Factorizations (e.g., $W_K^\top W_Q$): $y^\top UVx$
 - ► Low rank factorization can save parameters/compute
 - One joint gradient step from random initialization still works
- Non-linear MLP: $y^{\top} U\sigma(Vx)$
 - ► More expressive when *x*, *y* are *superpositions*/sums of embeddings
 - One gradient step still ok
- Layer-norm: $y^{\top} \frac{Wx}{\|Wx\|}$
 - Prevents repeated updates when Wx and y are already aligned
 - First gradient step from random initialization is unchanged
- Trained embeddings
 - ► Single gradient steps capture basic co-occurrence statistics/BoW/topics
 - ► Or more complex learning of structured embeddings (*e.g.*, "grokking")

Does it work empirically on the bigram task? Yes!

- ${\scriptstyle \bullet }$ Memory recall probes $\rightarrow 1$ as in previous experiment
- But: adding heads and layers loses identifiability

•
$$z_i \sim p(z)$$
, $y_i = f^*(z_i)$, n samples: $S_n = \{z_1, \dots, z_n\}$, $0/1$ loss:
 $L(\hat{f}_n) = \mathbb{P}(y \neq \hat{f}_n(z))$

•
$$z_i \sim p(z), y_i = f^*(z_i), n \text{ samples: } S_n = \{z_1, \dots, z_n\}, 0/1 \text{ loss:}$$

 $L(\hat{f}_n) = \mathbb{P}(y \neq \hat{f}_n(z))$

• Zipf law: $p(z) \propto z^{-lpha}$ (up to permutation)

•
$$z_i \sim p(z), y_i = f^*(z_i), n \text{ samples: } S_n = \{z_1, \dots, z_n\}, 0/1 \text{ loss:}$$

 $L(\hat{f}_n) = \mathbb{P}(y \neq \hat{f}_n(z))$

• Zipf law: $p(z) \propto z^{-lpha}$ (up to permutation)

• Random embeddings $u_z, v_y \in \mathbb{R}^d$ with $\mathcal{N}(0, 1/d)$ entries

•
$$z_i \sim p(z), y_i = f^*(z_i), n \text{ samples: } S_n = \{z_1, \dots, z_n\}, 0/1 \text{ loss:}$$

 $L(\hat{f}_n) = \mathbb{P}(y \neq \hat{f}_n(z))$

• Zipf law: $p(z) \propto z^{-lpha}$ (up to permutation)

- Random embeddings $u_z, v_y \in \mathbb{R}^d$ with $\mathcal{N}(0, 1/d)$ entries
- Hutter (2021): with infinite memory, $L(\hat{f}_n) \lesssim n^{-\frac{\alpha-1}{\alpha}}$

•
$$z_i \sim p(z), y_i = f^*(z_i), n \text{ samples: } S_n = \{z_1, \dots, z_n\}, 0/1 \text{ loss:}$$

 $L(\hat{f}_n) = \mathbb{P}(y \neq \hat{f}_n(z))$

• Zipf law: $p(z) \propto z^{-lpha}$ (up to permutation)

- Random embeddings $u_z, v_y \in \mathbb{R}^d$ with $\mathcal{N}(0, 1/d)$ entries
- Hutter (2021): with infinite memory, $L(\hat{f}_n) \lesssim n^{-\frac{\alpha-1}{\alpha}}$

Theorem (Cabannes, Dohmatob, B., 2023, informal) Consider the estimator $\hat{f}_{n,d}(x) = \arg \max_{v} v_{v}^{\top} W_{n,d} u_{z}$, with $W_{n,d} = \sum_{z=1}^{N} q(z) v_{f^{*}(z)} u_{z}^{\top}$.

•
$$z_i \sim p(z), y_i = f^*(z_i), n \text{ samples: } S_n = \{z_1, \dots, z_n\}, 0/1 \text{ loss:}$$

 $L(\hat{f}_n) = \mathbb{P}(y \neq \hat{f}_n(z))$

• Zipf law: $p(z) \propto z^{-lpha}$ (up to permutation)

- Random embeddings $u_z, v_y \in \mathbb{R}^d$ with $\mathcal{N}(0, 1/d)$ entries
- Hutter (2021): with infinite memory, $L(\hat{f}_n) \lesssim n^{-\frac{\alpha-1}{\alpha}}$

Theorem (Cabannes, Dohmatob, B., 2023, informal) Consider the estimator $\hat{f}_{n,d}(x) = \arg \max_y v_y^\top W_{n,d} u_z$, with $W_{n,d} = \sum_{z=1}^N q(z) v_{f^*(z)} u_z^\top$. **a** For $q(z) = \sum_i \mathbb{1}\{z = z_i\}$: $L(\hat{f}_{n,d}) \lesssim n^{-\frac{\alpha-1}{\alpha}} + d^{-\frac{\alpha-1}{2\alpha}}$

•
$$z_i \sim p(z), y_i = f^*(z_i), n \text{ samples: } S_n = \{z_1, \dots, z_n\}, 0/1 \text{ loss:}$$

 $L(\hat{f}_n) = \mathbb{P}(y \neq \hat{f}_n(z))$

• Zipf law: $p(z) \propto z^{-lpha}$ (up to permutation)

- Random embeddings $u_z, v_y \in \mathbb{R}^d$ with $\mathcal{N}(0, 1/d)$ entries
- Hutter (2021): with infinite memory, $L(\hat{f}_n) \lesssim n^{-\frac{\alpha-1}{\alpha}}$

Theorem (Cabannes, Dohmatob, B., 2023, informal) Consider the estimator $\hat{f}_{n,d}(x) = \arg \max_y v_y^\top W_{n,d}u_z$, with $W_{n,d} = \sum_{z=1}^N q(z)v_{f^*(z)}u_z^\top$. ① For $q(z) = \sum_i \mathbb{1}\{z = z_i\}$: $L(\hat{f}_{n,d}) \lesssim n^{-\frac{\alpha-1}{\alpha}} + d^{-\frac{\alpha-1}{2\alpha}}$ ② For $q(z) = \mathbb{1}\{z \in S_n\}$, and $d \gg N$: $L(\hat{f}_{n,d}) \lesssim n^{-\frac{\alpha-1}{\alpha}} + d^{-k}$ for any k

•
$$z_i \sim p(z), y_i = f^*(z_i), n \text{ samples: } S_n = \{z_1, \dots, z_n\}, 0/1 \text{ loss:}$$

 $L(\hat{f}_n) = \mathbb{P}(y \neq \hat{f}_n(z))$

• Zipf law: $p(z) \propto z^{-lpha}$ (up to permutation)

- Random embeddings $u_z, v_y \in \mathbb{R}^d$ with $\mathcal{N}(0, 1/d)$ entries
- Hutter (2021): with infinite memory, $L(\hat{f}_n) \lesssim n^{-\frac{\alpha-1}{\alpha}}$

Theorem (Cabannes, Dohmatob, B., 2023, informal) Consider the estimator $\hat{f}_{n,d}(x) = \arg \max_y v_y^\top W_{n,d} u_z$, with $W_{n,d} = \sum_{z=1}^N q(z) v_{f^*(z)} u_z^\top$. 1) For $q(z) = \sum_i \mathbb{1}\{z = z_i\}$: $L(\hat{f}_{n,d}) \lesssim n^{-\frac{\alpha-1}{\alpha}} + d^{-\frac{\alpha-1}{2\alpha}}$ 2) For $q(z) = \mathbb{1}\{z \in S_n\}$, and $d \gg N$: $L(\hat{f}_{n,d}) \lesssim n^{-\frac{\alpha-1}{\alpha}} + d^{-k}$ for any k3) For $q(z) = \mathbb{1}\{z \text{ seen at least s times in } S_n\}$: $L(\hat{f}_{n,d}) \lesssim n^{-\frac{\alpha-1}{\alpha}} + d^{-\alpha+1}$

Scaling laws with optimization algorithms

$$W_{n,d} = \sum_{z=1}^{N} q(z) v_{f^*(z)} u_z^{\top}$$

Different algorithms lead to different memory schemes q(z):

Scaling laws with optimization algorithms

$$W_{n,d} = \sum_{z=1}^{N} \frac{q(z)}{v_{f^*(z)}} u_z^{\top}$$

Different algorithms lead to different memory schemes q(z):

• One step of SGD with large batch: $q(z) \approx p(z)$

Scaling laws with optimization algorithms

$$W_{n,d} = \sum_{z=1}^{N} q(z) v_{f^*(z)} u_z^{\mathsf{T}}$$

Different algorithms lead to different memory schemes q(z):

- One step of SGD with large batch: $q(z) \approx p(z)$
- SGD with batch size one + large step-size, $d \gg N$: $q(z) \approx 1$
Scaling laws with optimization algorithms

$$W_{n,d} = \sum_{z=1}^{N} q(z) v_{f^*(z)} u_z^{\mathsf{T}}$$

Different algorithms lead to different memory schemes q(z):

- One step of SGD with large batch: $q(z) \approx p(z)$
- SGD with batch size one + large step-size, $d \gg N$: $q(z) \approx 1$
- For $d \leq N$, smaller step-sizes can help later in training

Scaling laws with optimization algorithms

$$\mathcal{N}_{n,d} = \sum_{z=1}^{N} \frac{q(z)}{v_{f^*(z)}} u_z^{\top}$$

Different algorithms lead to different memory schemes q(z):

- One step of SGD with large batch: $q(z) \approx p(z)$
- SGD with batch size one + large step-size, $d \gg N$: $q(z) \approx 1$
- For $d \leq N$, smaller step-sizes can help later in training
- Adam and layer-norm help with practical settings (large batch sizes + smaller step-size)

Scaling laws with optimization algorithms

$$W_{n,d} = \sum_{z=1}^{N} \frac{q(z)}{v_{f^*(z)}} u_z^{\top}$$

Different algorithms lead to different memory schemes q(z):

- One step of SGD with large batch: $q(z) \approx p(z)$
- SGD with batch size one + large step-size, $d \gg N$: $q(z) \approx 1$
- For $d \leq N$, smaller step-sizes can help later in training
- Adam and layer-norm help with practical settings (large batch sizes + smaller step-size)



Discussion and next steps

Summary

- Bigram model: simple but rich toy model for discrete data
- Transformer weights as associative memories
- Learning via few top-down gradient steps
- Better algorithms help for better scaling laws for heavy-tailed data

Discussion and next steps

Summary

- Bigram model: simple but rich toy model for discrete data
- Transformer weights as associative memories
- Learning via few top-down gradient steps
- Better algorithms help for better scaling laws for heavy-tailed data

Future directions

- More complex "reasoning" mechanisms, links with "emergence"
- Learning dynamics: multiple gradient steps? joint training? embeddings?
- Applications: interpretability, model editing, factual recall, efficient fine-tuning
- LLM large-width scalings (links with μP)
- Replace weights by hash tables?? (a.k.a. Leon's dream)

Discussion and next steps

Summary

- Bigram model: simple but rich toy model for discrete data
- Transformer weights as associative memories
- Learning via few top-down gradient steps
- Better algorithms help for better scaling laws for heavy-tailed data

Future directions

- More complex "reasoning" mechanisms, links with "emergence"
- Learning dynamics: multiple gradient steps? joint training? embeddings?
- Applications: interpretability, model editing, factual recall, efficient fine-tuning
- LLM large-width scalings (links with μP)
- Replace weights by hash tables?? (a.k.a. Leon's dream)

Thank you!

References I

- 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.
- 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.
- 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.
- M. Hutter. Learning curve theory. arXiv preprint arXiv:2102.04074, 2021.
- T. Kohonen. Correlation matrix memories. *IEEE Transactions on Computers*, 1972.

References II

- K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and editing factual associations in GPT. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 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.
- 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.