Increasing the Sequence Length in LLMs: Smoothness and Dynamics of Self-Attention

Valérie Castin Ecole Normale Supérieure PSL, Paris

University of Amsterdam, 27 May 2025



A Unified Perspective on the Dynamics of Deep Transformers, Castin, Ablin, Carrillo, Peyré, *preprint 2025*

How Smooth Is Attention?, Castin, Ablin and Peyré, in ICML 2024

I - The dynamics of tokens processed by a Transformer



• One input sentence \rightarrow *n* tokens $X \coloneqq (x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$



- One input sentence ightarrow n tokens $X \coloneqq (x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$
- Self-attention with params Q, K, V:

$$f(x_1,\ldots,x_n) \coloneqq (\Gamma_X(x_1),\ldots,\Gamma_X(x_n))$$
 where

$$\Gamma_X(z) \coloneqq \sum_{j=1}^n p_j(z) V_{x_j} \quad ext{with} \quad p_j(z) \coloneqq e^{\langle Q z, K x_j
angle} / \sum_{\ell=1}^n e^{\langle Q z, K x_\ell
angle}$$



- One input sentence ightarrow *n* tokens $X \coloneqq (x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$
- Self-attention with params Q, K, V:

$$f(x_1,\ldots,x_n) := (\Gamma_X(x_1),\ldots,\Gamma_X(x_n))$$
 where

$$\Gamma_X(z) := \sum_{j=1}^n p_j(z) V_{x_j} \quad ext{with} \quad p_j(z) := e^{\langle Qz, K_{x_j}
angle} / \sum_{\ell=1}^n e^{\langle Qz, K_{x_\ell}
angle}$$

• Layer normalization LN: $x \mapsto \beta \odot \frac{x}{|x|} \sqrt{d}$ (applied token-wise)



- One input sentence \rightarrow *n* tokens $X \coloneqq (x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$
- Self-attention with params Q, K, V:

$$f(x_1, \dots, x_n) := (\Gamma_X(x_1), \dots, \Gamma_X(x_n)) \text{ where}$$
$$T_X(z) := \sum_{j=1}^n p_j(z) \bigvee x_j \text{ with } p_j(z) := e^{\langle Q_z, K_{X_j} \rangle} / \sum_{\ell=1}^n e^{\langle Q_z, K_{X_\ell} \rangle}$$

- Layer normalization LN: $x \mapsto \beta \odot \frac{x}{|x|} \sqrt{d}$ (applied token-wise)
- Multilayer perceptron $g : \mathbb{R}^d \to \mathbb{R}^d, x \mapsto W\sigma(Ux + b)$ (applied token-wise)



- One input sentence ightarrow n tokens $X:=(x_1,\ldots,x_n)\in (\mathbb{R}^d)^n$
- Self-attention with params Q, K, V:

$$f(x_1, \dots, x_n) := (\Gamma_X(x_1), \dots, \Gamma_X(x_n)) \text{ where}$$
$$\Gamma_X(z) := \sum_{j=1}^n p_j(z) \bigvee x_j \text{ with } p_j(z) := e^{\langle Q_z, K_{X_j} \rangle} / \sum_{\ell=1}^n e^{\langle Q_z, K_{X_\ell} \rangle}$$

• Layer normalization LN: $x \mapsto \beta \odot \frac{x}{|x|} \sqrt{d}$ (applied token-wise)

Without MLP and LN:

$$x_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \qquad 1 \le i \le n$$



$$x_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell))$$
 $1 \le i \le n$

>



$$\kappa_i(\ell+1) = \kappa_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \qquad 1 \le i \le n$$

>



$$\kappa_i(\ell+1) = \kappa_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \qquad 1 \le i \le n$$

What are the dynamics of tokens going through the Transformer? The geometry of learned representations?

• understand clustering effect

>



$$\kappa_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \qquad 1 \le i \le n$$

- understand clustering effect
- influence of parameters $Q_{\ell}, K_{\ell}, V_{\ell}$



$$x_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell))$$
 $1 \le i \le n$

- understand clustering effect
- influence of parameters $Q_{\ell}, K_{\ell}, V_{\ell}$
- effect of layer normalization



$$\kappa_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \qquad 1 \le i \le n$$

- understand clustering effect
- influence of parameters $Q_{\ell}, K_{\ell}, V_{\ell}$
- effect of layer normalization
- influence of sequence length



$$x_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell))$$
 $1 \le i \le n$

- understand clustering effect
- influence of parameters $Q_{\ell}, K_{\ell}, V_{\ell}$
- effect of layer normalization
- influence of sequence length
- compare attention variants

$$x_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \quad \longleftrightarrow \quad \dot{x}_i(t) = \Gamma_{X(t)}(x_i(t)) \quad \text{[Sander et al., 2021]}$$

$$x_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \quad \longleftrightarrow \quad \dot{x}_i(t) = \Gamma_{X(t)}(x_i(t)) \quad [\text{Sander et al., 2021}]$$

Clustering result [Geshkovski et al., 2024]

- Assume $A \succ 0$ constant and $V = I_d$
- Rescale $z_i(t) \coloneqq e^{-tV} x_i(t)$

$$x_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \quad \longleftrightarrow \quad \dot{x}_i(t) = \Gamma_{X(t)}(x_i(t)) \quad [\text{Sander et al., 2021}]$$

Clustering result [Geshkovski et al., 2024]

- Assume $A \succ 0$ constant and $V = I_d$
- Rescale $z_i(t) \coloneqq e^{-tV} x_i(t)$

Then $\forall (z_1(0), \ldots, z_n(0))$ there exists a convex polytope $K \subset \mathbb{R}^d$ such that

all $z_i(t)$ converge in $\{0\} \cup \partial K$.

$$x_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \quad \longleftrightarrow \quad \dot{x}_i(t) = \Gamma_{X(t)}(x_i(t)) \quad [\text{Sander et al., 2021}]$$

Clustering result [Geshkovski et al., 2024]

- Assume $A \succ 0$ constant and $V = I_d$
- Rescale $z_i(t) \coloneqq e^{-tV} x_i(t)$

Then $\forall (z_1(0), \ldots, z_n(0))$ there exists a convex polytope $K \subset \mathbb{R}^d$ such that

all $z_i(t)$ converge in $\{0\} \cup \partial K$.

• rescaling mimics layer normalization

$$x_i(\ell+1) = x_i(\ell) + \Gamma_{X_\ell}(x_i(\ell)) \quad \longleftrightarrow \quad \dot{x}_i(t) = \Gamma_{X(t)}(x_i(t)) \quad [\text{Sander et al., 2021}]$$

Clustering result [Geshkovski et al., 2024]

• Assume $A \succ 0$ constant and $V = I_d$

• Rescale
$$z_i(t) \coloneqq e^{-tV} x_i(t)$$

Then $\forall (z_1(0), \ldots, z_n(0))$ there exists a convex polytope $K \subset \mathbb{R}^d$ such that

all $z_i(t)$ converge in $\{0\} \cup \partial K$.

- rescaling mimics layer normalization
- typically ♯vertices(K) ≪ n



Self-attention $f(X) = (\Gamma_X(x_1), \ldots, \Gamma_X(x_n))$ where

$$\Gamma_X(z) = \sum_{j=1}^n p_j(z) V_{x_j}$$
 with $p_j(z) = \exp(\langle Az, x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Az, x_\ell \rangle)$

Self-attention $f(X) = (\Gamma_X(x_1), \ldots, \Gamma_X(x_n))$ where

$$\Gamma_X(z) = \sum_{j=1}^n p_j(z) V_{x_j} \quad \text{with} \quad p_j(z) = \exp(\langle Az, x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Az, x_\ell \rangle)$$

Permutation equivariant: $x_i \leftrightarrow x_j \Rightarrow f(X)_i \leftrightarrow f(X)_j$

Self-attention $f(X) = (\Gamma_X(x_1), \ldots, \Gamma_X(x_n))$ where

$$\Gamma_X(z) = \sum_{j=1}^n p_j(z) \bigvee x_j \quad \text{with} \quad p_j(z) = \exp(\langle Az, x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Az, x_\ell \rangle)$$

Permutation equivariant: $x_i \leftrightarrow x_j \Rightarrow f(X)_i \leftrightarrow f(X)_j$

Mean-field self-attention [Sander et al., 2022, Vuckovic et al., 2020]

 $F \colon \mu \in \mathcal{P}(\mathbb{R}^d) \mapsto (\Gamma_\mu)_{\sharp} \mu$ where

$$\Gamma_{\mu} \colon z \in \mathbb{R}^{d} \mapsto \int k(z,y) \bigvee y \mathrm{d}\mu(y) \quad ext{with} \quad k(z,y) \coloneqq e^{\langle \mathcal{A} z, y
angle} / \int e^{\langle \mathcal{A} z, y'
angle} \mathrm{d}\mu(y')$$

Self-attention $f(X) = (\Gamma_X(x_1), \ldots, \Gamma_X(x_n))$ where

$$\Gamma_X(z) = \sum_{j=1}^n p_j(z) V_{x_j}$$
 with $p_j(z) = \exp(\langle Az, x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Az, x_\ell \rangle)$

Permutation equivariant: $x_i \leftrightarrow x_j \Rightarrow f(X)_i \leftrightarrow f(X)_j$

Mean-field self-attention [Sander et al., 2022, Vuckovic et al., 2020]

$$F \colon \mu \in \mathcal{P}(\mathbb{R}^d) \mapsto (\Gamma_{\mu})_{\sharp} \mu$$
 where

$$\Gamma_{\mu} \colon z \in \mathbb{R}^{d} \mapsto \int k(z,y) V y \mathrm{d} \mu(y) \quad ext{with} \quad k(z,y) \coloneqq e^{\langle Az,y
angle} / \int e^{\langle Az,y
angle} \mathrm{d} \mu(y')$$

- Well-defined for μ compactly supported, Gaussian...
- Covers discrete case: view X as $\frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}$

Self-attention $f(X) = (\Gamma_X(x_1), \ldots, \Gamma_X(x_n))$ where

$$\Gamma_X(z) = \sum_{j=1}^n p_j(z) V_{x_j} \quad \text{with} \quad p_j(z) = \exp(\langle Az, x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Az, x_\ell \rangle)$$

Permutation equivariant: $x_i \leftrightarrow x_j \Rightarrow f(X)_i \leftrightarrow f(X)_j$

Mean-field self-attention [Sander et al., 2022, Vuckovic et al., 2020] $F: \mu \in \mathcal{P}(\mathbb{R}^d) \mapsto (\Gamma_{\mu})_{\sharp} \mu$ where $\Gamma_{\mu}: z \in \mathbb{R}^d \mapsto \int k(z, y) V y d\mu(y)$ with $k(z, y) := e^{\langle Az, y \rangle} / \int e^{\langle Az, y' \rangle} d\mu(y')$

$$\dot{x}_i(t) = \Gamma_{X(t)}(x_i(t)) \longrightarrow \partial_t \mu + \operatorname{div}(\mu \Gamma_\mu) = 0$$
 (1)

I - 1) Studying the Transformer PDE with compactly supported initial data

$$\partial_t \mu + \operatorname{div}(\mu \Gamma_\mu) = 0$$

$$\mu_0$$
 compactly supported, $\partial_t \mu + \operatorname{div}(\mu \Gamma_\mu) = 0$ (1)

 μ_0 compactly supported, $\partial_t \mu + \operatorname{div}(\mu \Gamma_\mu) = 0$ (1)

Well-posedness [Geshkovski et al., 2024, Castin et al., 2025]

- Assume A(t), V(t) continuous
- Assume supp $\mu_0 \subset B(0, R_0)$

Then (1) has a unique global weak solution μ , such that

$$\operatorname{supp} \mu(t) \subset B(0, e^{\int_0^t \| \mathbf{V}(s) \|_2 \mathrm{d}s} R_0).$$

$$\mu_0$$
 compactly supported, $\partial_t \mu + \operatorname{div}(\mu \Gamma_\mu) = 0$ (1)

Well-posedness [Geshkovski et al., 2024, Castin et al., 2025]

- Assume A(t), V(t) continuous
- Assume supp $\mu_0 \subset B(0, R_0)$

Then (1) has a unique global weak solution μ , such that

$$\mathrm{supp}\,\mu(t)\subset B(0,e^{\int_0^t\|oldsymbol{V}(s)\|_2\mathrm{d} s}R_0).$$

If supp $\nu_0 \subset B(0, R_0)$ then

$$W_{p}(\mu(t),
u(t))\leq C(t,R_{0})W_{p}(\mu_{0},
u_{0}) \quad orall p\geq 1$$

with $C(t, R_0) \propto e^{tR(t)^2}$

$$\Gamma_{\mu}(z) = \int k(z,y) V y \mathrm{d}\mu(y) \quad ext{with} \quad k(z,y) \coloneqq e^{\langle Az,y
angle} / \int e^{\langle Az,y
angle} \mathrm{d}\mu(y')$$

Central estimates for proof

- 1. $\sup_{x\in\mathbb{R}^d}|\Gamma_{\mu}(x)|\leq \|V\|_2 R$,
- 2. $\sup_{x \in \mathbb{R}^d} \|D_x \Gamma_\mu(x)\|_2 \le \|V\|_2 \|A\|_2 R^2$,
- 3. $|\Gamma_{\mu}(x) \Gamma_{\nu}(x)| \leq c(x, R) W_{\rho}(\mu, \nu)$

$$\Gamma_{\mu}(z) = \int k(z,y) V y \mathrm{d}\mu(y) \quad ext{with} \quad k(z,y) \coloneqq e^{\langle Az,y
angle} / \int e^{\langle Az,y
angle} \mathrm{d}\mu(y')$$

Central estimates for proof

- 1. $\sup_{x\in\mathbb{R}^d}|\Gamma_{\mu}(x)|\leq \|V\|_2 R$,
- 2. $\sup_{x \in \mathbb{R}^d} \|D_x \Gamma_\mu(x)\|_2 \le \|V\|_2 \|A\|_2 R^2$,
- 3. $|\Gamma_{\mu}(x) \Gamma_{\nu}(x)| \leq c(x, R) W_{\rho}(\mu, \nu)$
- Eq. 1 controls radius growth: $\operatorname{supp} \mu(t) \subset B(0, e^{\int_0^t \|V(s)\|_2 \mathrm{d}s} R_0)$

$$\Gamma_{\mu}(z) = \int k(z,y) V y \mathrm{d}\mu(y) \quad ext{with} \quad k(z,y) \coloneqq e^{\langle Az,y
angle} / \int e^{\langle Az,y
angle} \mathrm{d}\mu(y')$$

Central estimates for proof

- 1. $\sup_{x\in\mathbb{R}^d}|\Gamma_{\mu}(x)|\leq \|V\|_2 R$,
- 2. $\sup_{x \in \mathbb{R}^d} \|D_x \Gamma_\mu(x)\|_2 \le \|V\|_2 \|A\|_2 R^2$,
- 3. $|\Gamma_{\mu}(x) \Gamma_{\nu}(x)| \leq c(x, R)W_{\rho}(\mu, \nu)$
- Eq. 1 controls radius growth: $\operatorname{supp} \mu(t) \subset B(0, e^{\int_0^t || V(s)||_2 \mathrm{d} s} R_0)$
- Modular framework \rightarrow extends to attention variants!

Attention map: $\Gamma_{\mu}(z) = \int k(z, y) V y d\mu(y)$

✓ Softmax attention: $k(z, y) = e^{\langle Az, y \rangle} / \int e^{\langle Az, y' \rangle} d\mu(y')$

Attention map: $\Gamma_{\mu}(z) = \int k(z, y) V y d\mu(y)$

✓ Softmax attention: $k(z, y) = e^{\langle Az, y \rangle} / \int e^{\langle Az, y' \rangle} d\mu(y')$

✓ L2 attention: $k(z, y) = e^{-|Qz - Ky|^2} / \int e^{-|Qz - Ky'|^2} d\mu(y')$

Attention map: $\Gamma_{\mu}(z) = \int k(z, y) V y d\mu(y)$

- ✓ Softmax attention: $k(z, y) = e^{\langle Az, y \rangle} / \int e^{\langle Az, y' \rangle} d\mu(y')$
- ✓ L2 attention: $k(z, y) = e^{-|Qz Ky|^2} / \int e^{-|Qz Ky'|^2} d\mu(y')$
- ✓ Sinkhorn attention: k(z, y) is the limit $j \to +\infty$ of

$$\kappa^{0}(z,y) = e^{\langle \mathbf{A}z,y \rangle}, \quad \kappa^{j+1}(z,y) = \begin{cases} \frac{\kappa^{j}(z,y)}{\int \kappa^{j}(z,y') d\mu(y')} & \text{if } j \text{ is even}, \\ \frac{\kappa^{j}(z,y)}{\int \kappa^{j}(z',y) d\mu(z')} & \text{if } j \text{ is odd} \end{cases}$$

Attention map: $\Gamma_{\mu}(z) = \int k(z, y) V y d\mu(y)$

- ✓ Softmax attention: $k(z, y) = e^{\langle Az, y \rangle} / \int e^{\langle Az, y' \rangle} d\mu(y')$
- ✓ L2 attention: $k(z, y) = e^{-|Qz Ky|^2} / \int e^{-|Qz Ky'|^2} d\mu(y')$
- ✓ Sinkhorn attention: k(z, y) is the limit $j \to +\infty$ of

$$\kappa^{0}(z,y) = e^{\langle \mathbf{A}z,y \rangle}, \quad \kappa^{j+1}(z,y) = \begin{cases} \frac{\kappa^{j}(z,y)}{\int \kappa^{j}(z,y') d\mu(y')} & \text{if } j \text{ is even}, \\ \frac{\kappa^{j}(z,y)}{\int \kappa^{j}(z',y) d\mu(z')} & \text{if } j \text{ is odd} \end{cases}$$

✓ Masked attention
Attention map: $\Gamma_{\mu}(z) = \int k(z, y) V y d\mu(y)$

- ✓ Softmax attention: $k(z, y) = e^{\langle Az, y \rangle} / \int e^{\langle Az, y' \rangle} d\mu(y')$
- ✓ L2 attention: $k(z, y) = e^{-|Qz Ky|^2} / \int e^{-|Qz Ky'|^2} d\mu(y')$
- ✓ Sinkhorn attention: k(z, y) is the limit $j \to +\infty$ of

$$\kappa^{0}(z,y) = e^{\langle \mathbf{A}z,y \rangle}, \quad \kappa^{j+1}(z,y) = \begin{cases} \frac{\kappa^{j}(z,y)}{\int \kappa^{j}(z,y') d\mu(y')} & \text{if } j \text{ is even}, \\ \frac{\kappa^{j}(z,y)}{\int \kappa^{j}(z',y) d\mu(z')} & \text{if } j \text{ is odd} \end{cases}$$

- ✓ Masked attention
- ✓ Multihead attention: $\Gamma_{\mu} = \sum_{h=1}^{H} \Gamma_{\mu}^{(h)}$

Attention map: $\Gamma_{\mu}(z) = \int k(z, y) V y d\mu(y)$

- **X** Unnormalized attention: $k(z, y) = e^{\langle Az, y \rangle}$
- **X** Linear attention: $k(z, y) = \langle Az, y \rangle$
- **X** ReLU attention: $k(z, y) = \text{ReLU}(\langle Az, y \rangle)$

X Sigmoid attention:
$$k(z, y) = \sigma(\langle Az, y \rangle)$$

Attention map: $\Gamma_{\mu}(z) = \int k(z, y) V y d\mu(y)$

- **X** Unnormalized attention: $k(z, y) = e^{\langle Az, y \rangle}$
- **X** Linear attention: $k(z, y) = \langle Az, y \rangle$
- **X** ReLU attention: $k(z, y) = \text{ReLU}(\langle Az, y \rangle)$

X Sigmoid attention: $k(z, y) = \sigma(\langle Az, y \rangle)$

• Estimate 1

$$\sup_{x\in\mathbb{R}^d}|\Gamma_{\mu}(x)|\leq \|\boldsymbol{V}\|_2\,R$$

is not satisfied \rightarrow no global solution

Attention map: $\Gamma_{\mu}(z) = \int k(z, y) V y d\mu(y)$

- **X** Unnormalized attention: $k(z, y) = e^{\langle Az, y \rangle}$
- **X** Linear attention: $k(z, y) = \langle Az, y \rangle$
- **X** ReLU attention: $k(z, y) = \operatorname{ReLU}(\langle Az, y \rangle)$

X Sigmoid attention: $k(z, y) = \sigma(\langle Az, y \rangle)$

• Estimate 1

$$\sup_{x\in\mathbb{R}^d}|\Gamma_{\mu}(x)|\leq \|\boldsymbol{V}\|_2\,R$$

is not satisfied \rightarrow no global solution

• But LayerNorm solves the problem

I - 2) Beyond compactly supported data: the Gaussian case [Castin et al., 2025]

 $\partial_t \mu + \operatorname{div}(\mu \Gamma_\mu) = 0$

The Transformer PDE in the Gaussian case

Lemma: attention map on Gaussians

If $\mu = \mathcal{N}(\alpha, \Sigma)$ then

$$\Gamma_{\mu}(x) = V(\alpha + \Sigma A x)$$

Similar for L2 and Sinkhorn attention!

The Transformer PDE in the Gaussian case

Lemma: attention map on Gaussians

If $\mu = \mathcal{N}(\alpha, \Sigma)$ then

$$\Gamma_{\mu}(x) = V(\alpha + \Sigma A x)$$

Similar for L2 and Sinkhorn attention!

Proposition: evolution of Gaussian initial data [Castin et al., 2025]

Consider

$$\partial_t \mu + \operatorname{div}(\mu \Gamma_\mu) = 0 \tag{1}$$

Assume A, V continuous and $\mu_0 = \mathcal{N}(\alpha_0, \Sigma_0)$. Then (1) has a unique maximal solution on $[0, t_{\max})$, Gaussian for all t.

The Transformer PDE in the Gaussian case

Lemma: attention map on Gaussians

If $\mu = \mathcal{N}(\alpha, \Sigma)$ then

$$\Gamma_{\mu}(x) = V(\alpha + \Sigma A x)$$

Similar for L2 and Sinkhorn attention!

Proposition: evolution of Gaussian initial data [Castin et al., 2025]

Consider

$$\partial_t \mu + \operatorname{div}(\mu \Gamma_\mu) = 0 \tag{1}$$

Assume A, V continuous and $\mu_0 = \mathcal{N}(\alpha_0, \Sigma_0)$. Then (1) has a unique maximal solution on $[0, t_{\text{max}})$, Gaussian for all t. Denoting $\mu(t) = \mathcal{N}(\alpha(t), \Sigma(t))$:

$$\begin{cases} \dot{\alpha} = \mathbf{V} (\mathbf{I}_d + \mathbf{\Sigma} \mathbf{A}) \alpha \\ \dot{\mathbf{\Sigma}} = \mathbf{V} \mathbf{\Sigma} \mathbf{A} \mathbf{\Sigma} + \mathbf{\Sigma} \mathbf{A}^\top \mathbf{\Sigma} \mathbf{V}^\top \end{cases}$$

Proposition: closed-form analysis

Consider

$$\dot{\boldsymbol{\Sigma}} = \boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{A}\boldsymbol{\Sigma} + \boldsymbol{\Sigma}\boldsymbol{A}^{\top}\boldsymbol{\Sigma}\boldsymbol{V}^{\top}.$$

Assume

- A, V constant
- V commutes with $VA + A^{\top}V^{\top}$ and Σ_0

Then

$$\Sigma(t) = (\Sigma_0^{-1} - t(\boldsymbol{V}\boldsymbol{A} + \boldsymbol{A}^\top\boldsymbol{V}^\top))^{-1}$$

Proposition: closed-form analysis

Consider

$$\dot{\boldsymbol{\Sigma}} = \boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{A}\boldsymbol{\Sigma} + \boldsymbol{\Sigma}\boldsymbol{A}^{\top}\boldsymbol{\Sigma}\boldsymbol{V}^{\top}.$$

Assume

- A, V constant
- V commutes with $VA + A^{\top}V^{\top}$ and Σ_0

Then

$$\Sigma(t) = (\Sigma_0^{-1} - t(\boldsymbol{V}\boldsymbol{A} + \boldsymbol{A}^\top \boldsymbol{V}^\top))^{-1}$$

• If $VA + A^{\top}V^{\top} \leq 0$ the solution is global and converges to Σ^* such that $\lambda_i(VA + A^{\top}V^{\top}) < 0 \Rightarrow \lambda_i(\Sigma^*) = 0$ "clustering"

Proposition: closed-form analysis

Consider

$$\dot{\boldsymbol{\Sigma}} = \boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{A}\boldsymbol{\Sigma} + \boldsymbol{\Sigma}\boldsymbol{A}^{\top}\boldsymbol{\Sigma}\boldsymbol{V}^{\top}.$$

Assume

- A, V constant
- V commutes with $VA + A^{\top}V^{\top}$ and Σ_0

Then

$$\Sigma(t) = (\Sigma_0^{-1} - t(\boldsymbol{V}\boldsymbol{A} + \boldsymbol{A}^\top \boldsymbol{V}^\top))^{-1}$$

• If $VA + A^{\top}V^{\top} \leq 0$ the solution is global and converges to Σ^* such that $\lambda_i(VA + A^{\top}V^{\top}) < 0 \Rightarrow \lambda_i(\Sigma^*) = 0$ "clustering"

• Otherwise $\lambda_1(\Sigma(t)) \to +\infty$ in finite time

Proposition: limiting covariances have low-rank

Let $\boldsymbol{\Sigma}$ such that

 $\boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{A}\boldsymbol{\Sigma}+\boldsymbol{\Sigma}\boldsymbol{A}^{\top}\boldsymbol{\Sigma}\,\boldsymbol{V}^{\top}=\boldsymbol{0}$

Assume

- A constant and symmetric
- $V = I_d$

Then

 $\operatorname{rk} \Sigma \leq \dim \ker A + \min(\sharp \text{pos. eigvals of } A, \sharp \text{neg. eigvals of } A)$

Proposition: limiting covariances have low-rank

Let $\boldsymbol{\Sigma}$ such that

 $\boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{A}\boldsymbol{\Sigma}+\boldsymbol{\Sigma}\boldsymbol{A}^{\top}\boldsymbol{\Sigma}\,\boldsymbol{V}^{\top}=\boldsymbol{0}$

Assume

• A constant and symmetric

• $V = I_d$

Then

 $\operatorname{rk} \Sigma \leq \dim \ker A + \min(\sharp \text{pos. eigvals of } A, \sharp \text{neg. eigvals of } A)$

• If A invertible, rk limiting covariance $\leq \lceil d/2 \rceil$

Proposition: limiting covariances have low-rank

Let $\boldsymbol{\Sigma}$ such that

 $\boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{A}\boldsymbol{\Sigma}+\boldsymbol{\Sigma}\boldsymbol{A}^{\top}\boldsymbol{\Sigma}\,\boldsymbol{V}^{\top}=\boldsymbol{0}$

Assume

• A constant and symmetric

• $V = I_d$

Then

 $\operatorname{rk} \Sigma \leq \dim \ker A + \min(\sharp \text{pos. eigvals of } A, \sharp \text{neg. eigvals of } A)$

- If A invertible, rk limiting covariance $\leq \lceil d/2 \rceil$
- Numerically: holds for V random

Plot the covariance evolution for d = 2:



Comparing attention variants

- Softmax, L2, Sinkhorn attention \rightarrow similar behavior when converge
- L2 does not blow-up in finite time \rightarrow more regular

Comparing attention variants

- Softmax, L2, Sinkhorn attention \rightarrow similar behavior when converge
- L2 does not blow-up in finite time \rightarrow more regular



Comparing attention variants

- Softmax, L2, Sinkhorn attention \rightarrow similar behavior when converge
- L2 does not blow-up in finite time \rightarrow more regular



Open question: beyond Gaussians?

- Mean-field attention and the Transformer PDE generalize dynamics to infinitely many tokens
- Attention variants should have a normalized kernel to induce well-posed dynamics
- The Gaussian case reduces to a system of matrix equations \rightarrow theoretical and numerical analysis, variety of behaviors, clustering
- LayerNorm changes a lot the dynamics
- What about next-token prediction? (Masked attention dynamics)
- From discrete to continuous time, what changes?

II - The Lipschitz constant of self-attention [Castin et al., 2024]



How much can the output of a Transformer change when slightly perturbing the input?



How much can the output of a Transformer change when slightly perturbing the input?

controls robustness and expressive power



How much can the output of a Transformer change when slightly perturbing the input?

- controls robustness and expressive power
- parameters are fixed



How much can the output of a Transformer change when slightly perturbing the input?

- controls robustness and expressive power
- parameters are fixed
- we analyze only one attention layer



How much can the output of a Transformer change when slightly perturbing the input?

- controls robustness and expressive power
- parameters are fixed
- we analyze only one attention layer

Does the robustness of an input depend on the sequence length?

Measuring regularity with the local Lipschitz constant

 $f\colon (\mathbb{R}^d)^n o (\mathbb{R}^d)^n$ self-attention

Local Lipschitz constant

Norm on $(\mathbb{R}^d)^n$: $||X||^2 \coloneqq \sum_{i=1}^n |x_i|^2$

Measuring regularity with the local Lipschitz constant

 $f\colon (\mathbb{R}^d)^n o (\mathbb{R}^d)^n$ self-attention

Local Lipschitz constant

Norm on $(\mathbb{R}^d)^n$: $||X||^2 \coloneqq \sum_{i=1}^n |x_i|^2$

Local Lipschitz constant of f at X:

$$\operatorname{Lip}_{X}(f) := \|D_{X}f\|_{2} = \sup_{\|\varepsilon\|=1} \|D_{X}f(\varepsilon)\|$$

where $D_X f \colon (\mathbb{R}^d)^n o (\mathbb{R}^d)^n$ Jacobian of f

Measuring regularity with the local Lipschitz constant

 $f\colon (\mathbb{R}^d)^n o (\mathbb{R}^d)^n$ self-attention

Local Lipschitz constant

Norm on $(\mathbb{R}^d)^n$: $||X||^2 \coloneqq \sum_{i=1}^n |x_i|^2$

Local Lipschitz constant of f at X:

$$\operatorname{Lip}_X(f) := \|D_X f\|_2 = \sup_{\|\varepsilon\|=1} \|D_X f(\varepsilon)\|$$

where $D_X f \colon (\mathbb{R}^d)^n o (\mathbb{R}^d)^n$ Jacobian of f

Gives global guarantees:

$$\sup_{X\neq Y\in B_R^n}\frac{\|f(X)-f(Y)\|}{\|X-Y\|}=\sup_{X\in B_R^n}\operatorname{Lip}_X(f)$$

Theorem 1 [Castin et al., 2024]

$$\operatorname{Lip}(f_{|B_R^n}) \leq \sqrt{3} \| \mathbf{V} \|_2 \left(\| \mathbf{A} \|_2^2 R^4 (4n+1) + n \right)^{1/2} \approx R^2 \sqrt{n}$$

Theorem 1 [Castin et al., 2024]

$$\operatorname{Lip}(f_{|B_R^n}) \leq \sqrt{3} \| \mathbf{V} \|_2 \left(\| \mathbf{A} \|_2^2 R^4 (4n+1) + n \right)^{1/2} \approx R^2 \sqrt{n}$$

and if $V = I_d$,

$$\operatorname{Lip}(f_{|B_R^n}) \geq rac{1}{1+(n-1)e^{-2R^2 oldsymbol{\gamma}}} \sqrt{n-1}$$

where $R^2 \gamma \approx 10^{2-3}$ in practical Transformers.

Theorem 1 [Castin et al., 2024]

$$\operatorname{Lip}(f_{|B_R^n}) \leq \sqrt{3} \| \mathbf{V} \|_2 \left(\| \mathbf{A} \|_2^2 R^4 (4n+1) + n \right)^{1/2} \approx R^2 \sqrt{n}$$

and if $V = I_d$,

$$\operatorname{Lip}(f_{|B_R^n}) \geq \frac{1}{1+(n-1)e^{-2R^2\gamma}}\sqrt{n-1}$$

where $R^2 \gamma \approx 10^{2-3}$ in practical Transformers.

• *R* fixed by layer normalization

Theorem 1 [Castin et al., 2024]

$$\operatorname{Lip}(f_{|B_R^n}) \leq \sqrt{3} \| \mathbf{V} \|_2 \left(\| \mathbf{A} \|_2^2 R^4 (4n+1) + n \right)^{1/2} \approx R^2 \sqrt{n}$$

and if $V = I_d$,

$$\operatorname{Lip}(f_{|B_R^n}) \geq rac{1}{1+(n-1)e^{-2R^2 \gamma}} \sqrt{n-1}$$

where $R^2 \gamma \approx 10^{2-3}$ in practical Transformers.

- *R* fixed by layer normalization
- *n* not too large: $\operatorname{Lip}(f_{|B_R^n})$ grows like $C\sqrt{n}$

Theorem 1 [Castin et al., 2024]

$$\operatorname{Lip}(f_{|B_R^n}) \leq \sqrt{3} \| \mathbf{V} \|_2 \left(\| \mathbf{A} \|_2^2 R^4 (4n+1) + n \right)^{1/2} \approx R^2 \sqrt{n}$$

and if $V = I_d$,

$$\operatorname{Lip}(f_{|B_R^n}) \geq rac{1}{1+(n-1)e^{-2R^2 \gamma}} \sqrt{n-1}$$

where $R^2 \gamma \approx 10^{2-3}$ in practical Transformers.

- *R* fixed by layer normalization
- *n* not too large: $\operatorname{Lip}(f_{|B_R^n})$ grows like $C\sqrt{n}$
- mean-field bound: $\operatorname{Lip}(f_{|B_R^n}) \leq c(A, V)R^2 e^{1\|A\|_2 R^2}$

Experiments: typical case and worst case



• Growth in *Cn*^{1/4} for real data

Experiments: typical case and worst case



- Growth in *Cn*^{1/4} for real data
- Growth in $C\sqrt{n}$ for adv. data \rightarrow matches lower bound

Experiments: typical case and worst case



- Growth in *Cn*^{1/4} for real data
- Growth in $C\sqrt{n}$ for adv. data \rightarrow matches lower bound

Obstacle to Lipschitz attention (GPT-40 context window: 128k)

Local Lipschitz constant of real vs. adversarial data
Thank you!

References

- Castin, V., Ablin, P., Carrillo, J. A., and Peyré, G. (2025).
 A unified perspective on the dynamics of deep transformers.
 In arXiv preprint arXiv:2501.18322.
- Castin, V., Ablin, P., and Peyré, G. (2024).
 How smooth is attention?
 In *ICML 2024*.
- Geshkovski, B., Letrouit, C., Polyanskiy, Y., and Rigollet, P. (2024). The emergence of clusters in self-attention dynamics. Advances in Neural Information Processing Systems, 36.
- Sander, M. E., Ablin, P., Blondel, M., and Peyré, G. (2022).
 Sinkformers: Transformers with doubly stochastic attention.
 In International Conference on Artificial Intelligence and Statistics, pages 3515–3530. PMLR.
 - Vuckovic, J., Baratin, A., and Combes, R. T. d. (2020). A mathematical theory of attention. *arXiv preprint arXiv:2007.02876.*

III - Masked self-attention

Masked self-attention $f^m \colon (\mathbb{R}^d)^n \to (\mathbb{R}^d)^n$ such that

 $f^m(X)_i := f(x_1,\ldots,x_i)_i$

with params $A, V \in \mathbb{R}^{d \times d}$

Mean-field regime? \rightarrow not permutation equivariant!

Generalizing masked self-attention to measures

Mean-field self-attention $F \colon \mu \in \mathcal{P}_c(\mathbb{R}^d) \mapsto (\Gamma_{\mu})_{\sharp} \mu$ where

$$\Gamma_{\mu} \colon x \in \mathbb{R}^{d} \mapsto \int_{\mathbb{R}^{d}} k(x,y) \bigvee y \mathrm{d}\mu(y) \quad ext{with} \quad k(x,y) \coloneqq e^{\langle Ax,y
angle} / \int e^{\langle Ax,z
angle} \mathrm{d}\mu(z)$$

Generalizing masked self-attention to measures

Mean-field self-attention $F \colon \mu \in \mathcal{P}_{c}(\mathbb{R}^{d}) \mapsto (\Gamma_{\mu})_{\sharp} \mu$ where

$$ar{f}_{\mu} \colon x \in \mathbb{R}^d \mapsto \int_{\mathbb{R}^d} k(x,y) V y \mathrm{d} \mu(y) \quad ext{with} \quad k(x,y) \coloneqq e^{\langle Ax,y
angle} / \int e^{\langle Ax,z
angle} \mathrm{d} \mu(z)$$

Mean-field masked self-attention [Castin et al., 2024]

Replace $\mu \in \mathcal{P}_c(\mathbb{R}^d)$ by $\bar{\mu} \in \mathcal{P}_c([0,1] \times \mathbb{R}^d)$:

$$F^m \colon \bar{\mu} \mapsto \left(\mathsf{\Gamma}_{\bar{\mu}}
ight)_{\sharp} \bar{\mu} \quad \text{where} \quad \mathsf{\Gamma}_{\bar{\mu}}(s, x) \coloneqq \left(s, \int_{[0,1] imes \mathbb{R}^d} \bigvee y k_s(x, y) \mathrm{d}\bar{\mu}(\tau, y)
ight)$$

with

$$k_{s}(x,y) \coloneqq e^{\langle Ax,y
angle} \mathbb{1}_{ au \leq s} / \int_{[0,1] imes \mathbb{R}^{d}} e^{\langle Ax,y
angle} \mathbb{1}_{ au \leq s} \mathrm{d}ar{\mu}(au,y)$$