How Smooth Is Attention?

Valérie Castin¹, Pierre Ablin² and Gabriel Peyré¹

¹Ecole Normale Supérieure PSL, Paris ²Apple

ENS de Lyon, 21 January 2025

• Important feature: data represented as tuples $(x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$ of tokens

Important feature: data represented as tuples
 (x₁,...,x_n) ∈ (ℝ^d)ⁿ of tokens

This is how GPT-3 tokenizes this sentence. Tokenization of text



Tokenization of images

• Important feature: data represented as tuples $(x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$ of tokens

This is how GPT-3 tokenizes this sentence. Tokenization of text



Tokenization of images

• Positional encoding encodes the order of tokens



• Important feature: data represented as tuples $(x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$ of tokens

This is how GPT-3 tokenizes this sentence. Tokenization of text



Tokenization of images

• Positional encoding encodes the order of tokens

Allows to learn local dependencies!





Architecture of Transformers:

- Attention layers: ∪_n(ℝ^d)ⁿ → ∪_n(ℝ^d)ⁿ learn dependencies between tokens
- Multilayer perceptrons: $\mathbb{R}^d \to \mathbb{R}^d$ (applied token-wise)
- Layer normalization: project each token on an ellipsis



Architecture of Transformers:

- Attention layers: ∪_n(ℝ^d)ⁿ → ∪_n(ℝ^d)ⁿ learn dependencies between tokens
- Multilayer perceptrons: $\mathbb{R}^d \to \mathbb{R}^d$ (applied token-wise)
- Layer normalization: project each token on an ellipsis

How much can the output change when slightly perturbing the input?



Architecture of Transformers:

- Attention layers: ∪_n(ℝ^d)ⁿ → ∪_n(ℝ^d)ⁿ learn dependencies between tokens
- Multilayer perceptrons: $\mathbb{R}^d \to \mathbb{R}^d$ (applied token-wise)
- Layer normalization: project each token on an ellipsis

How much can the output change when slightly perturbing the input?

controls robustness and expressive power



Architecture of Transformers:

- Attention layers: ∪_n(ℝ^d)ⁿ → ∪_n(ℝ^d)ⁿ learn dependencies between tokens
- Multilayer perceptrons: $\mathbb{R}^d \to \mathbb{R}^d$ (applied token-wise)
- Layer normalization: project each token on an ellipsis

How much can the output change when slightly perturbing the input?

- controls robustness and expressive power
- parameters are fixed



Architecture of Transformers:

- Attention layers: ∪_n(ℝ^d)ⁿ → ∪_n(ℝ^d)ⁿ learn dependencies between tokens
- Multilayer perceptrons: $\mathbb{R}^d \to \mathbb{R}^d$ (applied token-wise)
- Layer normalization: project each token on an ellipsis

How much can the output change when slightly perturbing the input?

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

• Parameters: $Q, K \in \mathbb{R}^{k \times d}$ and $V \in \mathbb{R}^{d \times d}$

- Parameters: $Q, K \in \mathbb{R}^{k \times d}$ and $V \in \mathbb{R}^{d \times d}$
- Attention of $X = (x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$ w.r.t. $z \in \mathbb{R}^d$:

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

- Parameters: $Q, K \in \mathbb{R}^{k \times d}$ and $V \in \mathbb{R}^{d \times d}$
- Attention of $X = (x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$ w.r.t. $z \in \mathbb{R}^d$:

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

• Self-attention of $X = (x_1, \ldots, x_n)$:

$$f: X \in (\mathbb{R}^d)^n \mapsto (\Gamma_X(x_1), \dots, \Gamma_X(x_n)) \in (\mathbb{R}^d)^n \qquad x_1 \overset{\bullet}{\longrightarrow} x_3$$

- Parameters: $Q, K \in \mathbb{R}^{k \times d}$ and $V \in \mathbb{R}^{d \times d}$
- Attention of $X = (x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$ w.r.t. $z \in \mathbb{R}^d$:

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

• Self-attention of $X = (x_1, \ldots, x_n)$:

$$f: X \in (\mathbb{R}^d)^n \mapsto (\Gamma_X(x_1), \dots, \Gamma_X(x_n)) \in (\mathbb{R}^d)^n$$



• Depends only on $A := K^\top Q$

- Parameters: $Q, K \in \mathbb{R}^{k \times d}$ and $V \in \mathbb{R}^{d \times d}$
- Attention of $X = (x_1, \ldots, x_n) \in (\mathbb{R}^d)^n$ w.r.t. $z \in \mathbb{R}^d$:

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

• Self-attention of $X = (x_1, \ldots, x_n)$:

$$f\colon X\in (\mathbb{R}^d)^n\mapsto ({\sf \Gamma}_X(x_1),\ldots,{\sf \Gamma}_X(x_n))\in (\mathbb{R}^d)^n$$



- Depends only on $A := K^\top Q$
- Multi-head attention: linear combination $\sum_{h=1}^{H} W^{(h)} f^{(h)}$



 $f(X)_i \coloneqq \sum_{j=1}^n P_{ij} \bigvee x_j \quad \text{with} \quad P_{ij} \coloneqq \exp(\langle Q x_i, K x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Q x_i, K x_\ell \rangle)$

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

Local Lipschitz constant

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

 $f\colon (\mathbb{R}^d)^n o (\mathbb{R}^d)^n$ 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:

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

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



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

Local Lipschitz constant

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

Local Lipschitz constant of f at X:

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

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



• Gives global guarantees:

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

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

Local Lipschitz constant

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

Local Lipschitz constant of f at X:

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

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

• Gives global guarantees:

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

• Direction of maximal change: $U = (u_1, \dots, u_n)$ first singular vector of $D_X f$

$$\frac{\|f(X+\eta U)-f(X)\|}{\|U\|} \to_{\eta\to 0} \operatorname{Lip}_X(f)$$

Presentation of the problem

• Self-attention is not globally Lipschitz continuous [Kim et al., 2021]

 $\operatorname{Lip}(f_{|B_R^n}) \geq c(A, V)R^2$

 \rightarrow set one particle at zero + spread the others

Presentation of the problem

• Self-attention is not globally Lipschitz continuous [Kim et al., 2021]

 $\operatorname{Lip}(f_{|B_R^n}) \geq c(A, V)R^2$

 \rightarrow set one particle at zero + spread the others

• Bound independent of *n* [Geshkovski et al., 2024]

 $\operatorname{Lip}(f_{|B_{R}^{n}}) \leq \|V\|_{2} (1 + 3 \|A\|_{2} R^{2}) e^{2\|A\|_{2} R^{2}}$

Presentation of the problem

• Self-attention is not globally Lipschitz continuous [Kim et al., 2021]

 $\operatorname{Lip}(f_{|B_R^n}) \geq c(A, V)R^2$

- ightarrow set one particle at zero + spread the others
- Bound independent of *n* [Geshkovski et al., 2024]

$$\operatorname{Lip}(f_{|B_{R}^{n}}) \leq \|V\|_{2} (1 + 3 \|A\|_{2} R^{2}) e^{2\|A\|_{2} R^{2}}$$

Questions

- Big discrepancy! Which bound is tighter? Dependency on *n*?
- Characterize adversarial configurations?
- Local Lipschitz constant of real data?

I - Dependency on n of the Lipschitz constant of self-attention

Theorem 1 [Castin et al., 2024]

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

Theorem 1 [Castin et al., 2024]

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

and if $V = I_d$,

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

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

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 $V = I_d$,

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

- *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 $V = I_d$,

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

- *R* fixed by layer normalization
- *n* not too large: $\operatorname{Lip}(f_{|B_P^n})$ grows like $C\sqrt{n}$
- *n* exponentially large: analyzed separately (mean-field regime)

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 $V = I_d$,

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

- Configuration for lower bound: (Ru, -Ru, ..., -Ru) or (Ru, Ru/2, ..., Ru/2) with u eigenvector of A
- Similar bound for multi-head



Plot dependency of $Lip_X(f)$ in number of tokens *n*:

• for X obtained from real text (Alice in Wonderland, AG_NEWS)



Plot dependency of $Lip_X(f)$ in number of tokens *n*:

- for X obtained from real text (Alice in Wonderland, AG_NEWS)
- for adversarial data of the lower bound



Plot dependency of $Lip_X(f)$ in number of tokens *n*:

- for X obtained from real text (Alice in Wonderland, AG_NEWS)
- for adversarial data of the lower bound
- $||D_X f||_2$ computed with power method (d = 768)



Plot dependency of $Lip_X(f)$ in number of tokens *n*:

- for X obtained from real text (Alice in Wonderland, AG_NEWS)
- for adversarial data of the lower bound
- $||D_X f||_2$ computed with power method (d = 768)
- Same theory for masked self-attention



• Growth in *Cn*^{1/4} for real data



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



- 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

The special case of learnt positional encoding



GPT-2 pretrained: magnitude of tokens R grows with number of tokens n $\operatorname{Lip}(f_{|B_R^n}) \leq R^2 \sqrt{n}$

The special case of learnt positional encoding



GPT-2 pretrained: magnitude of tokens R grows with number of tokens nLip $(f_{|B_R^n}) \leq R^2 \sqrt{n}$



Gaining intuition with the large radius regime

Large radius regime

 $R \to +\infty$ while *n* fixed

Denote $m_i = \operatorname{argmax}_{1 \le j \le n} \langle Ax_i, x_j \rangle$ (a.s. unique)

$$(D_{RX}f)(\varepsilon)
ightarrow_{R
ightarrow +\infty} (V \varepsilon_{m_1}, \dots, V \varepsilon_{m_n}) \qquad \quad \varepsilon \in (\mathbb{R}^d)^n$$

Proposition [Castin et al., 2024]

In the large radius regime:

$$\operatorname{Lip}(f_{|B_R^n}) \lesssim_{R \to +\infty} \| \mathbf{V} \|_2 \sqrt{n}$$

reached if $m_1 = \cdots = m_n$

Adversarial configurations in the large radius regime

One token x_j far away from the others s.t.

$$\forall i = 1, \dots, n, \quad \langle Ax_i, x_j \rangle \approx \max_k \langle Ax_i, x_k \rangle$$

 \rightarrow local Lipschitz constant proportional to \sqrt{n}



II - The mean-field regime (*n* exponentially large)

Self-attention $f(X) \coloneqq (\Gamma_X(x_1), \ldots, \Gamma_X(x_n))$ with

$$\Gamma_X : x \in \mathbb{R}^d \mapsto \sum_{j=1}^n p_j V_{x_j} \quad \text{with} \quad p_j := \exp(\langle Ax, x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Ax, x_\ell \rangle)$$

Self-attention $f(X) \coloneqq (\Gamma_X(x_1), \ldots, \Gamma_X(x_n))$ with

$$\Gamma_X \colon x \in \mathbb{R}^d \mapsto \sum_{j=1}^n p_j V_{x_j} \quad \text{with} \quad p_j \coloneqq \exp(\langle Ax, x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Ax, x_\ell \rangle)$$

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

Self-attention $f(X) \coloneqq (\Gamma_X(x_1), \dots, \Gamma_X(x_n))$ with

$$\Gamma_X \colon x \in \mathbb{R}^d \mapsto \sum_{j=1}^n p_j \bigvee x_j \quad \text{with} \quad p_j \coloneqq \exp(\langle Ax, x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Ax, x_\ell \rangle)$$

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

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

 ${\sf F}\colon \mu\in {\mathcal P}({\mathbb R}^d)\mapsto ({\sf \Gamma}_\mu)_\sharp\mu$ with

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

Self-attention $f(X) \coloneqq (\Gamma_X(x_1), \dots, \Gamma_X(x_n))$ with

$$\Gamma_X \colon x \in \mathbb{R}^d \mapsto \sum_{j=1}^n p_j \bigvee x_j \quad \text{with} \quad p_j \coloneqq \exp(\langle Ax, x_j \rangle) / \sum_{\ell=1}^n \exp(\langle Ax, x_\ell \rangle)$$

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

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

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

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

Covers discrete case representing X as $\frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}$

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

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

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

Lipschitz constant measured with Wasserstein distance

$$W_2(\mu,
u) \coloneqq \left(\inf_{\pi} \int |x-y|^2 \mathrm{d}\pi(x,y)\right)^{1/2}$$



where $\int \mathrm{d}\pi(x,\cdot) = \mathrm{d}\mu(x)$ and $\int \mathrm{d}\pi(\cdot,y) = \mathrm{d}\nu(y)$

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

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

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

Lipschitz constant measured with Wasserstein distance

$$W_2(\mu,
u) \coloneqq \left(\inf_{\pi} \int |x-y|^2 \mathrm{d}\pi(x,y)\right)^{1/2}$$



not optimal

optimal

where $\int d\pi(x,\cdot) = d\mu(x)$ and $\int d\pi(\cdot,y) = d\nu(y)$

$$\operatorname{Lip}(F_{|\mathcal{P}(B_R)}) \coloneqq \sup_{\mu \neq \nu \in \mathcal{P}(B_R)} \frac{W_2(F(\mu), F(\nu))}{W_2(\mu, \nu)}$$

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

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

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

Lipschitz constant measured with Wasserstein distance

$$W_2(\mu,
u)\coloneqq \left(\inf_\pi\int |x-y|^2\,\mathrm{d}\pi(x,y)
ight)^{1/2}$$

 \rightarrow upp. bound $R^2 e^{R^2}$ [Geshkovski et al., 2024]

where $\int \mathrm{d}\pi(x,\cdot) = \mathrm{d}\mu(x)$ and $\int \mathrm{d}\pi(\cdot,y) = \mathrm{d}\nu(y)$

$$\operatorname{Lip}(F_{|\mathcal{P}(B_R)}) \coloneqq \sup_{\mu \neq \nu \in \mathcal{P}(B_R)} \frac{W_2(F(\mu), F(\nu))}{W_2(\mu, \nu)}$$

Upper bound [Geshkovski et al., 2024]

$$\operatorname{Lip}(f_{|B_{R}^{n}}) \leq \| \mathbf{V} \|_{2} (1 + 3 \|\mathbf{A}\|_{2} R^{2}) e^{2 \|\mathbf{A}\|_{2} R^{2}}$$

Upper bound [Geshkovski et al., 2024]

$$\operatorname{Lip}(f_{|B_R^n}) \le \|V\|_2 (1 + 3 \|A\|_2 R^2) e^{2\|A\|_2 R^2}$$

Proposition 2 [Castin et al., 2024]

If $V = I_d$ and $n \sim e^{2\gamma R^2}$: Lip $(f_{|B_R^n}) \gtrsim \frac{\gamma}{2} R^2 e^{\gamma R^2}$

Upper bound [Geshkovski et al., 2024]

$$\operatorname{Lip}(f_{|B_R^n}) \le \| \mathbf{V} \|_2 (1 + 3 \| \mathbf{A} \|_2 R^2) e^{2 \| \mathbf{A} \|_2 R^2}$$

Proposition 2 [Castin et al., 2024]

If $V = I_d$ and $n \sim e^{2\gamma R^2}$:

$$\mathrm{Lip}(f_{|B_R^n})\gtrsim rac{\gamma}{2}R^2e^{\gamma R^2}$$

$$\begin{array}{c} 1 - p_R \\ p_R \\ Ru/2 \quad Ru \end{array}$$

Upper bound [Geshkovski et al., 2024]

$$\operatorname{Lip}(f_{|B_{R}^{n}}) \leq \| \mathbf{V} \|_{2} (1 + 3 \| \mathbf{A} \|_{2} R^{2}) e^{2 \| \mathbf{A} \|_{2} R^{2}}$$

 R^2

Proposition 2 [Castin et al., 2024]

If
$$V=I_d$$
 and $n\sim e^{2\gamma R^2}$: ${
m Lip}(f_{|B_R^n})\gtrsim rac{\gamma}{2}R^2e^{\gamma}$

Probability measures for lower bound:

$$p_R \delta_{Ru} + (1 - p_R) \delta_{Ru/2}$$
 or $p_R \delta_{Ru} + (1 - p_R) \delta_{-Ru}$

with

$$p_R = e^{-2\gamma R^2} \rightarrow_{R \rightarrow +\infty} 0$$



The mean-field regime is not practically relevant

Proposition 2 [Castin et al., 2024]

If $V = I_d$ and $n \sim e^{2\gamma R^2}$:

$$\mathrm{Lip}(f_{|B^n_R})\gtrsim rac{\gamma}{2}R^2e^{\gamma R^2}$$

In practice $2\gamma R^2 pprox 10^3
ightarrow$ not realistic



III - Masked self-attention

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

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

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

Theorem 2 [Castin et al., 2024]

The upper and lower bounds of Theorem 1 on $\operatorname{Lip}(f_{|B^n_R})$ also hold for $\operatorname{Lip}(f^m_{|B^n_R})$

Mean-field regime? \rightarrow not permutation equivariant!

Generalizing masked self-attention to measures

Mean-field self-attention $F : \mu \in \mathcal{P}_c(\mathbb{R}^d) \mapsto (\Gamma_\mu)_{\sharp} \mu$ where $\Gamma_\mu : x \in \mathbb{R}^d \mapsto \int_{\mathbb{R}^d} k(x, y) \bigvee y d\mu(y)$ with $k(x, y) \coloneqq e^{\langle Ax, y \rangle} / \int e^{\langle Ax, z \rangle} d\mu(z)$

Generalizing masked self-attention to measures

Mean-field self-attention $F : \mu \in \mathcal{P}_{c}(\mathbb{R}^{d}) \mapsto (\Gamma_{\mu})_{\sharp} \mu$ where $\Gamma_{\mu} : x \in \mathbb{R}^{d} \mapsto \int_{\mathbb{R}^{d}} k(x, y) \bigvee y \mathrm{d}\mu(y)$ with $k(x, y) \coloneqq e^{\langle Ax, y \rangle} / \int e^{\langle Ax, z \rangle} \mathrm{d}\mu(z)$

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

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

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

with

$$k_s(x,y) := e^{\langle A_{x,y}
angle} \mathbf{1}_{ au \leq s} / \int_{[0,1] imes \mathbb{R}^d} e^{\langle A_{x,y}
angle} \mathbf{1}_{ au \leq s} \mathrm{d}ar{\mu}(au,y)$$

Generalizing masked self-attention to measures

Mean-field self-attention $F : \mu \in \mathcal{P}_{c}(\mathbb{R}^{d}) \mapsto (\Gamma_{\mu})_{\sharp} \mu$ where $\Gamma_{\mu} : x \in \mathbb{R}^{d} \mapsto \int_{\mathbb{R}^{d}} k(x, y) \bigvee y \mathrm{d}\mu(y)$ with $k(x, y) \coloneqq e^{\langle Ax, y \rangle} / \int e^{\langle Ax, z \rangle} \mathrm{d}\mu(z)$

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

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

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

with

$$k_{s}(x,y) := e^{\langle \mathbf{A}x,y \rangle} \mathbf{1}_{\tau \leq s} / \int_{[0,1] \times \mathbb{R}^{d}} e^{\langle \mathbf{A}x,y \rangle} \mathbf{1}_{\tau \leq s} \mathrm{d}\bar{\mu}(\tau,y)$$

Same upper bound as unmasked mean-field self-attention!

IV - Application of the mean-field framework: modeling Transformers as PDEs

Modeling an infinitely deep Transformer as a PDE

• Simplified Transformer with only attention layers:

$$f = f^L \circ \cdots \circ f^1$$
 with $f^\ell(X) \coloneqq X + \frac{1}{L}(\Gamma^\ell_X(x_1), \dots, \Gamma^\ell_X(x_n))$

Modeling an infinitely deep Transformer as a PDE

• Simplified Transformer with only attention layers:

$$f = f^{L} \circ \cdots \circ f^{1}$$
 with $f^{\ell}(X) \coloneqq X + \frac{1}{L}(\Gamma_{X}^{\ell}(x_{1}), \dots, \Gamma_{X}^{\ell}(x_{n}))$

• Discretizes

$$\dot{x}_i(t) = \Gamma_{\mu(t)}(x_i(t)), \qquad 1 \le i \le n$$

with $\mu(t) \coloneqq \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$ and $\Gamma_{\mu(t)} \colon \mathbb{R}^d \to \mathbb{R}^d$ velocity field

Modeling an infinitely deep Transformer as a PDE

• Simplified Transformer with only attention layers:

$$f = f^{L} \circ \cdots \circ f^{1}$$
 with $f^{\ell}(X) \coloneqq X + \frac{1}{L}(\Gamma^{\ell}_{X}(x_{1}), \dots, \Gamma^{\ell}_{X}(x_{n}))$

Discretizes

$$\dot{x}_i(t) = \Gamma_{\mu(t)}(x_i(t)), \qquad 1 \le i \le n$$

with $\mu(t) \coloneqq \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i(t)}$ and $\Gamma_{\mu(t)} \colon \mathbb{R}^d \to \mathbb{R}^d$ velocity field

Corresponding PDE:

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

also on continuous measures

Tokens cluster after renormalization

For discrete tokens x_1, \ldots, x_n following

$$\dot{x}_i(t) = \Gamma_{X(t)}(x_i(t))$$

and renormalizing

$$y_i(t) := e^{-t \mathbf{V}} x_i(t)$$

 \rightarrow clusters emerge [Geshkovski et al., 2024]



Tokens cluster after renormalization

For discrete tokens x_1, \ldots, x_n following

$$\dot{x}_i(t) = \Gamma_{X(t)}(x_i(t))$$

and renormalizing

$$y_i(t) := e^{-t \mathbf{V}} x_i(t)$$

 \rightarrow clusters emerge [Geshkovski et al., 2024]



Similar phenomenon for Gaussian inputs! (Castin, Carrillo, Peyré, Ablin, in preparation)

References

```
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.
Kim, H., Papamakarios, G., and Mnih, A. (2021).
The lipschitz constant of self-attention.
In International Conference on Machine Learning, pages 5562–5571, PMLR.
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.
```