

# Transformers are Universal in Context Learners

Gabriel Peyré



Takashi  
Furuya



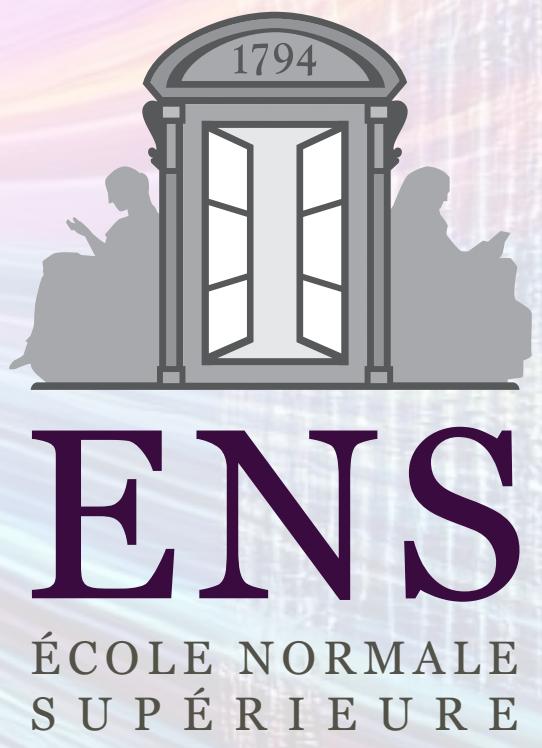
Maarten de  
Hoop



Valérie  
Castin



Pierre  
Ablin



ÉCOLE NORMALE  
SUPÉRIEURE

# Transformers and attention mechanism

Le lycée Marcelin Berthelot étant situé sur le parcours touristique de « la boucle de la Marne », est connu de tous ceux qui ont visité les environs de Paris. « Ah, c'est cet immense bâtiment moderne » dit-on.

**Tokenize**

Le lycée Marcelin Berthelot étant situé sur le parcours touristique de « la boucle de la Marne », est connu de tous ceux qui ont visité les environs de Paris. « Ah, c'est cet immense bâtiment moderne » dit-on.

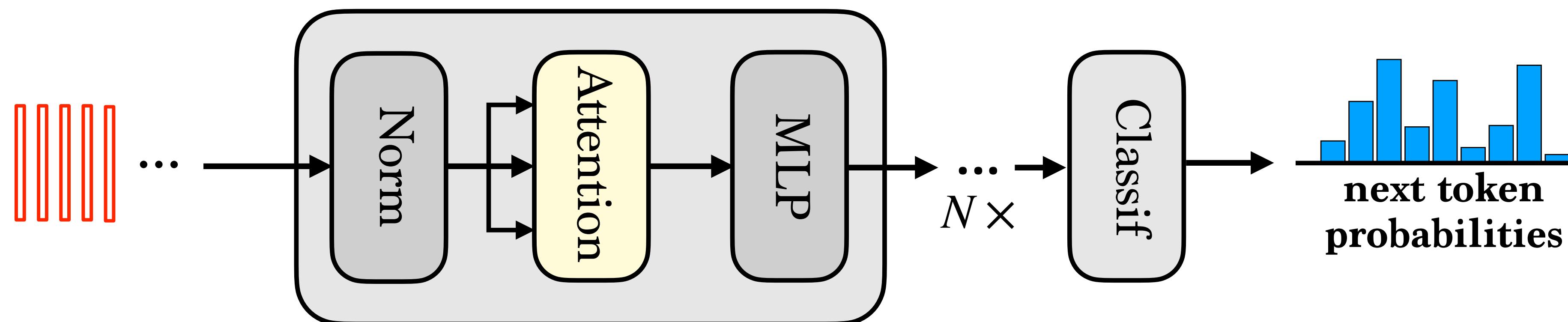
**Token  
encoding**

**Positional  
encoding**

$x_1$   
 $x_2$

...

**Points cloud**  
 $\{x_i\}_i$



**(Unmasked) Attention layer**

The diagram shows the computation of the attention weight between tokens  $x_i$  and  $x_j$ . Token  $x_i$  is projected into a query vector  $Qx_i$  and a key vector  $Kx_i$ . Token  $x_j$  is projected into a value vector  $Vx_j$ . The attention weight  $\tilde{x}_i$  is calculated as the dot product of the query vector and the transpose of the key vector, scaled by the exponential of the dot product, divided by the square root of the dimension. The formula is:

$$\tilde{x}_i := \sum_j \frac{e^{\langle Qx_i, Kx_j \rangle}}{\sum_\ell e^{\langle Qx_i, Kx_\ell \rangle}} Vx_j$$

# Transformers and attention mechanism

Le lycée Marcelin Berthelot étant situé sur le parcours touristique de « la boucle de la Marne », est connu de tous ceux qui ont visité les environs de Paris. « Ah, c'est cet immense bâtiment moderne » dit-on.

**Tokenize**

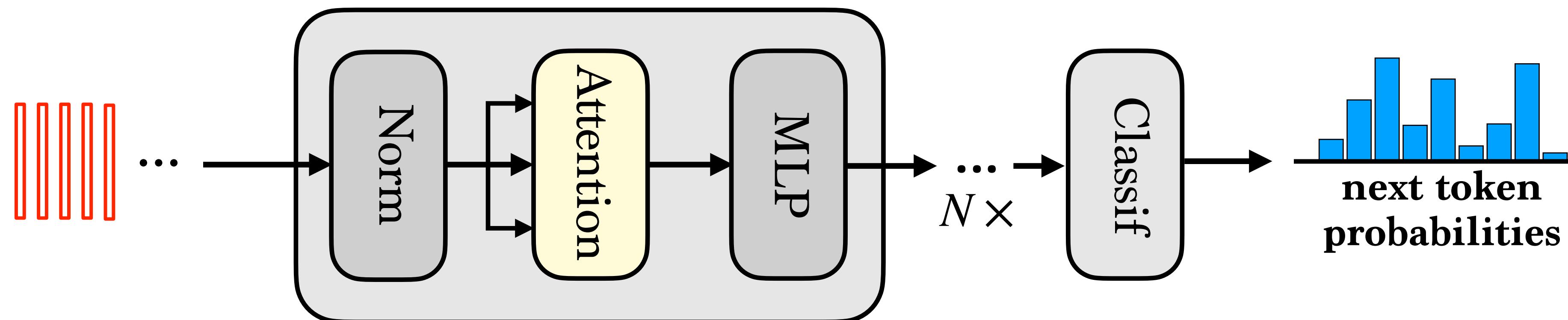
Le lycée Marcelin Berthelot étant situé sur le parcours touristique de « la boucle de la Marne », est connu de tous ceux qui ont visité les environs de Paris. « Ah, c'est cet immense bâtiment moderne » dit-on.

**Token  
encoding**

**Positional  
encoding**

$x_1$   
 $x_2$

**Points cloud**  
 $\{x_i\}_i$



**(Unmasked) Attention layer**

The equation defines the output of the Attention layer for token  $i$  as:

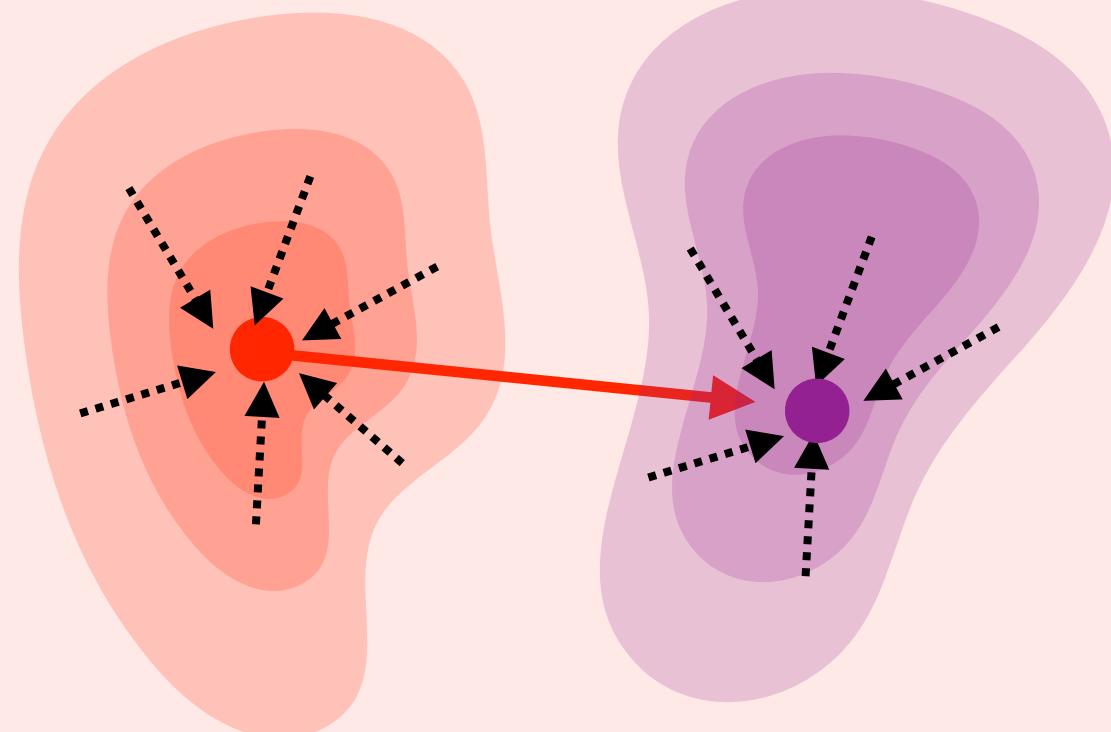
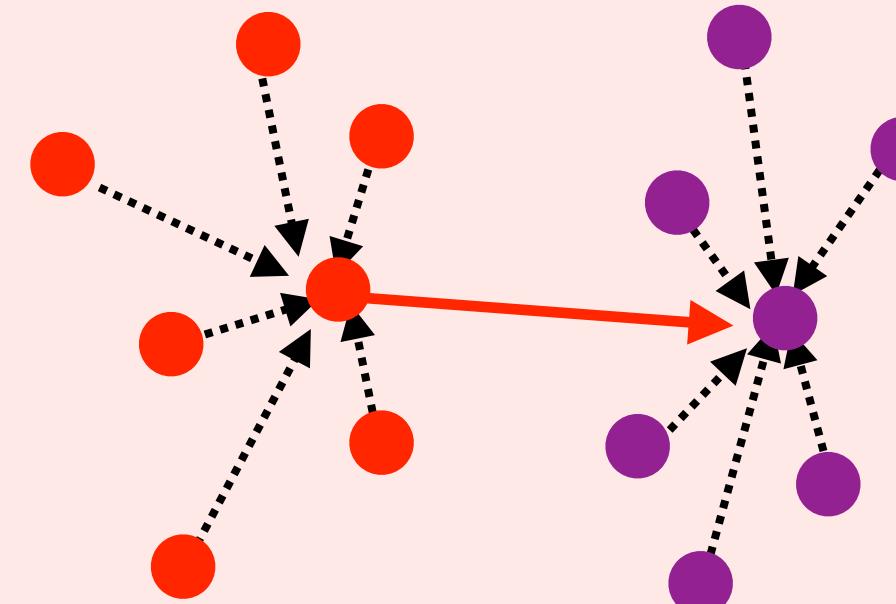
$$\tilde{x}_i := \sum_j \frac{e^{\langle Qx_i, Kx_j \rangle}}{\sum_\ell e^{\langle Qx_i, Kx_\ell \rangle}} Vx_j$$

where  $x_i$  is the query vector,  $Kx_j$  are the key vectors for other tokens, and  $Vx_j$  are the value vectors for those tokens.

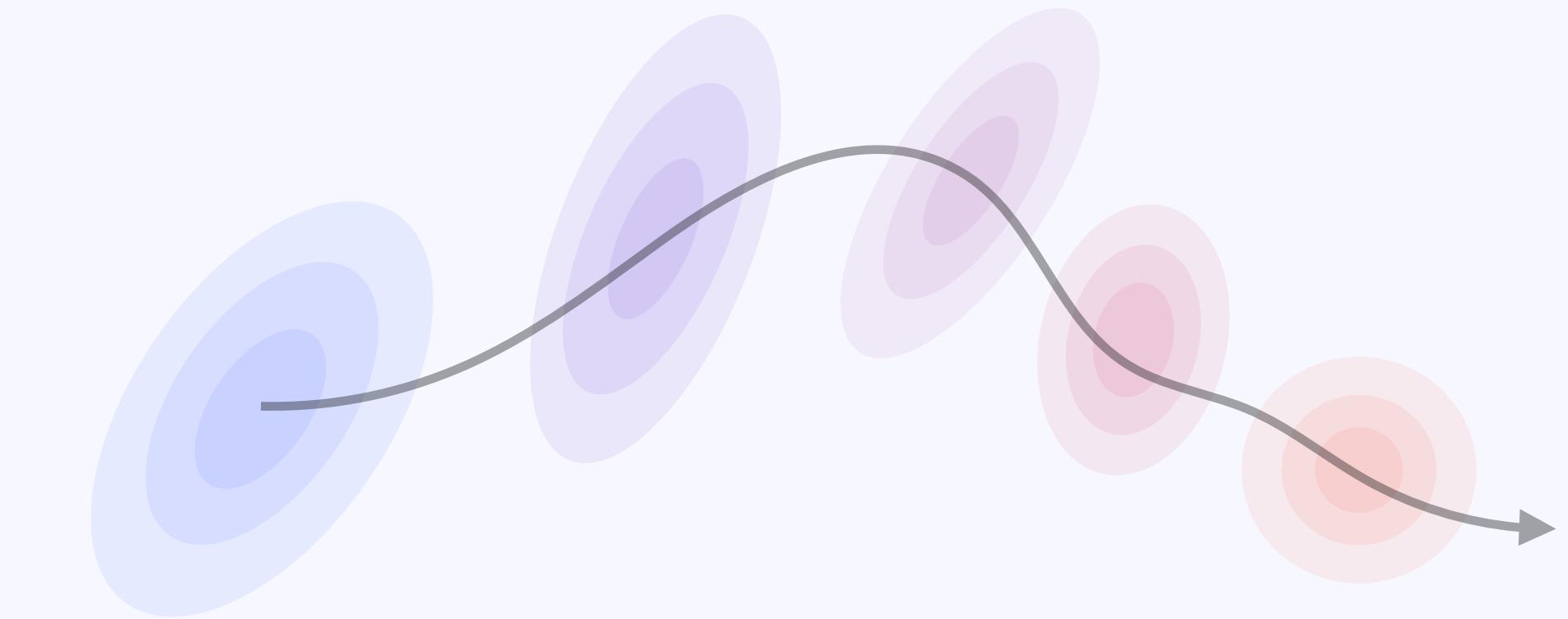
**Understanding**

Arbitrary number of tokens  
Arbitrary number of layers  
Expressivity

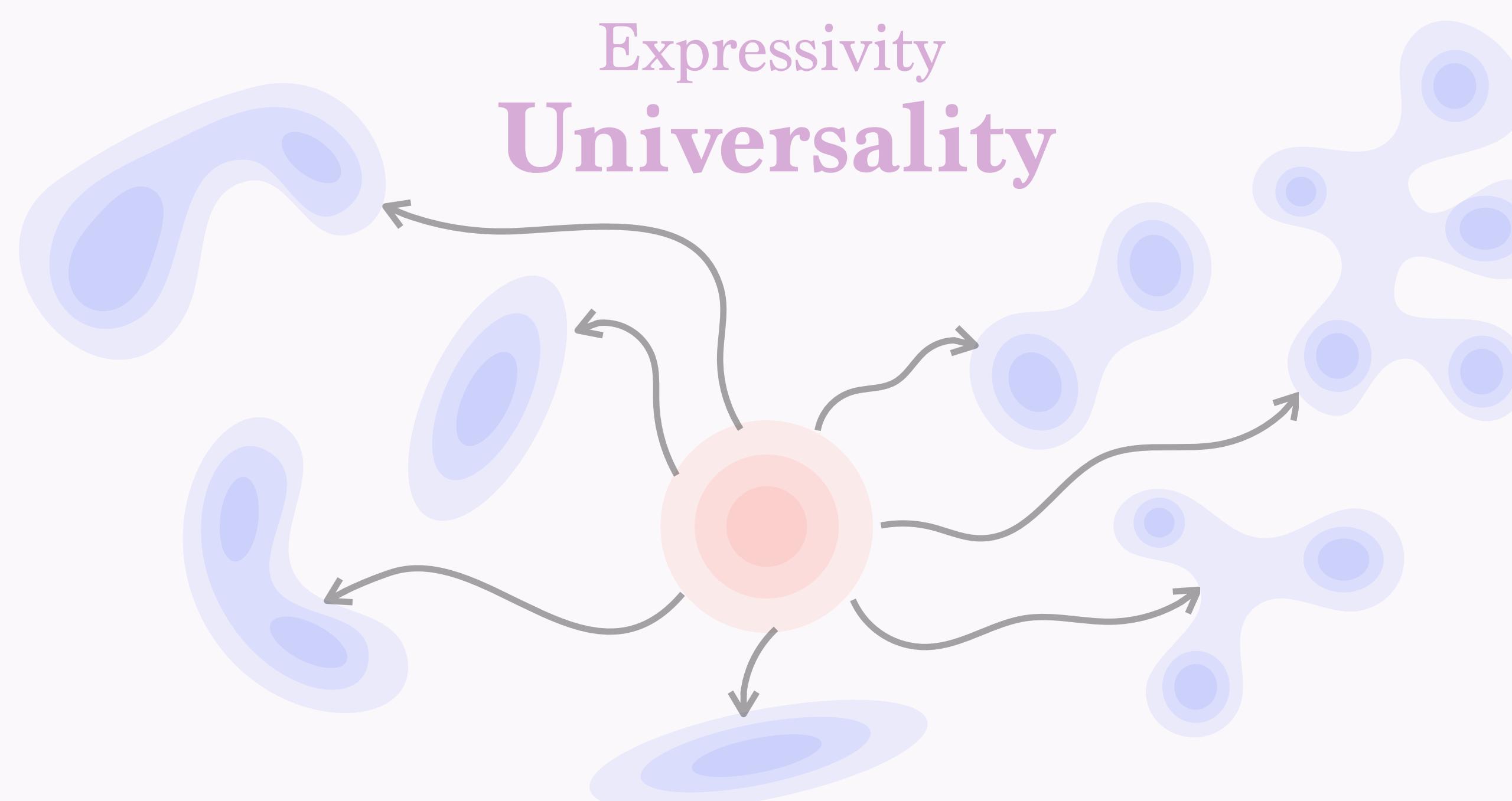
Arbitrary number of tokens  
**In Context Mappings over Measures**



Arbitrary number of layers  
**Smoothness and PDE's**



Expressivity  
**Universality**

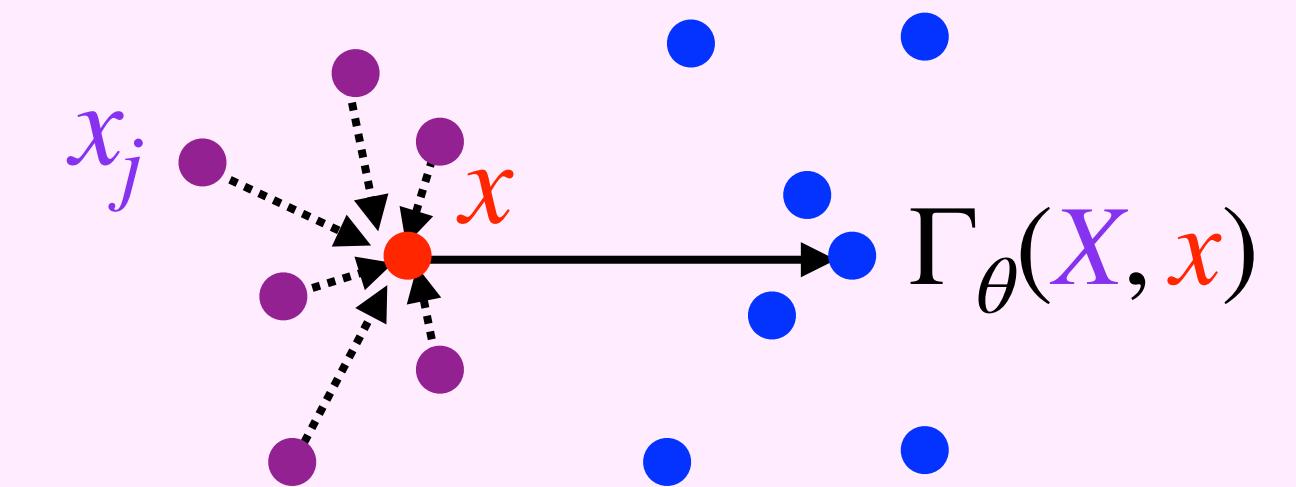


# Attention as In-context Mapping

Point clouds:  $X := \{x_i\}_{i=1}^n$

In-context mapping:  
parameters  $\theta := (Q, K, V)$

$$\Gamma_\theta[\textcolor{violet}{X}](\textcolor{red}{x}) := \sum_j \frac{e^{\langle Q\textcolor{red}{x}, Kx_j \rangle}}{\sum_\ell e^{\langle Q\textcolor{red}{x}, Kx_\ell \rangle}} V\textcolor{violet}{x}_j$$

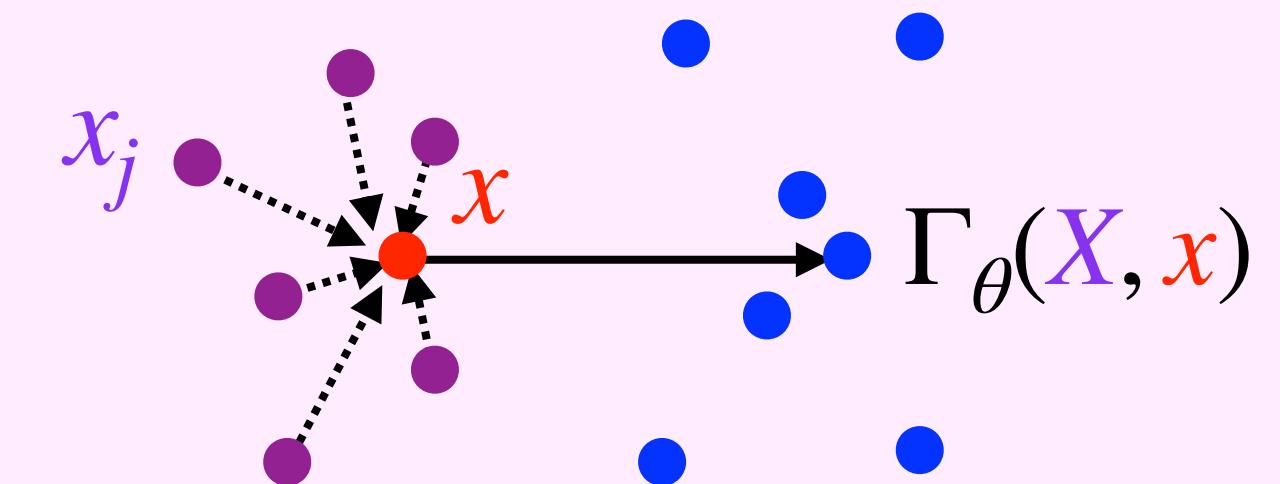


# Attention as In-context Mapping

Point clouds:  $X := \{x_i\}_{i=1}^n$

**In-context mapping:**  
parameters  $\theta := (Q, K, V)$

$$\Gamma_\theta[\mathbf{X}](\mathbf{x}) := \sum_j \frac{e^{\langle Q\mathbf{x}, K\mathbf{x}_j \rangle}}{\sum_\ell e^{\langle Q\mathbf{x}, K\mathbf{x}_\ell \rangle}} V\mathbf{x}_j$$



**Single-head attention layer:**  $X \mapsto \{\Gamma_\theta[\mathbf{X}](x_i)\}_{i=1}^n$

**Multi-head attention layer:**  $X \mapsto \left\{ \sum_{h=1}^H W_h \Gamma_{\theta_h}[\mathbf{X}](x_i) \right\}_{i=1}^n$

$$\begin{array}{|c|c|c|} \hline K_1, Q_1, V_1 \\ \hline K_2, Q_2, V_2 \\ \hline \dots \\ \hline \end{array}$$

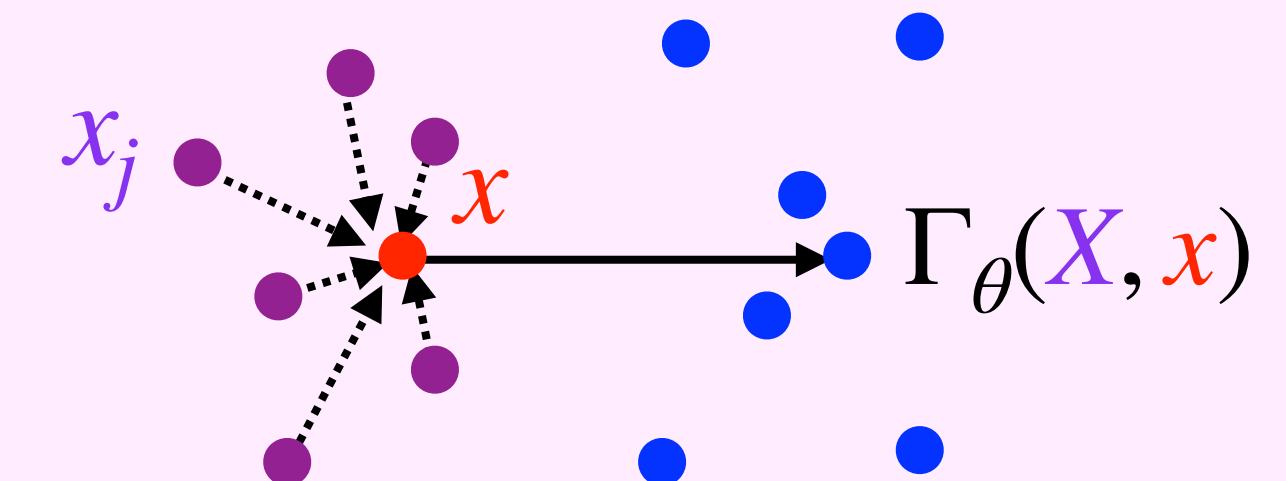
$$\begin{array}{|c|c|c|c|} \hline W_1 & W_2 & \dots & \\ \hline \end{array}$$

# Attention as In-context Mapping

Point clouds:  $X := \{x_i\}_{i=1}^n$

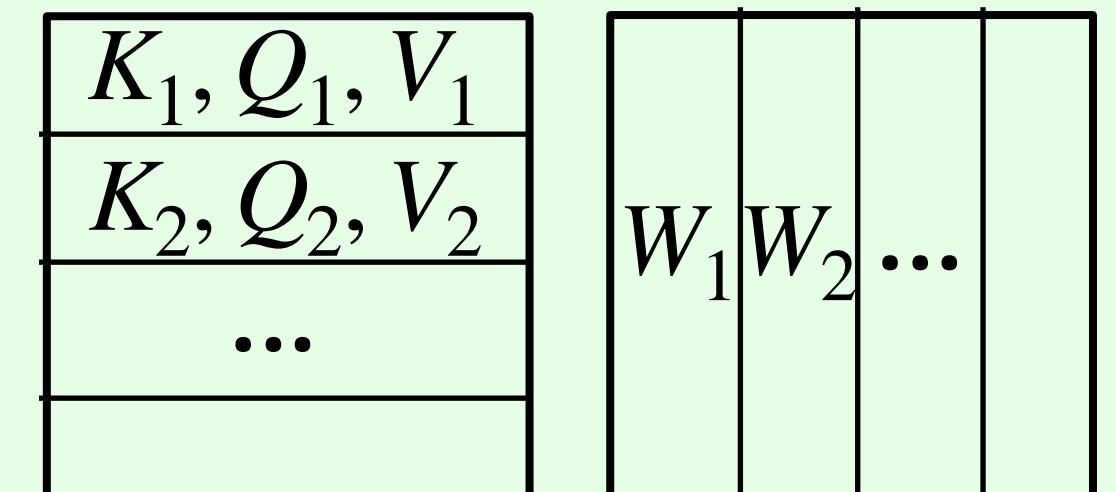
**In-context mapping:**  
parameters  $\theta := (Q, K, V)$

$$\Gamma_\theta[\mathbf{X}](\mathbf{x}) := \sum_j \frac{e^{\langle Q\mathbf{x}, K\mathbf{x}_j \rangle}}{\sum_\ell e^{\langle Q\mathbf{x}, K\mathbf{x}_\ell \rangle}} V\mathbf{x}_j$$



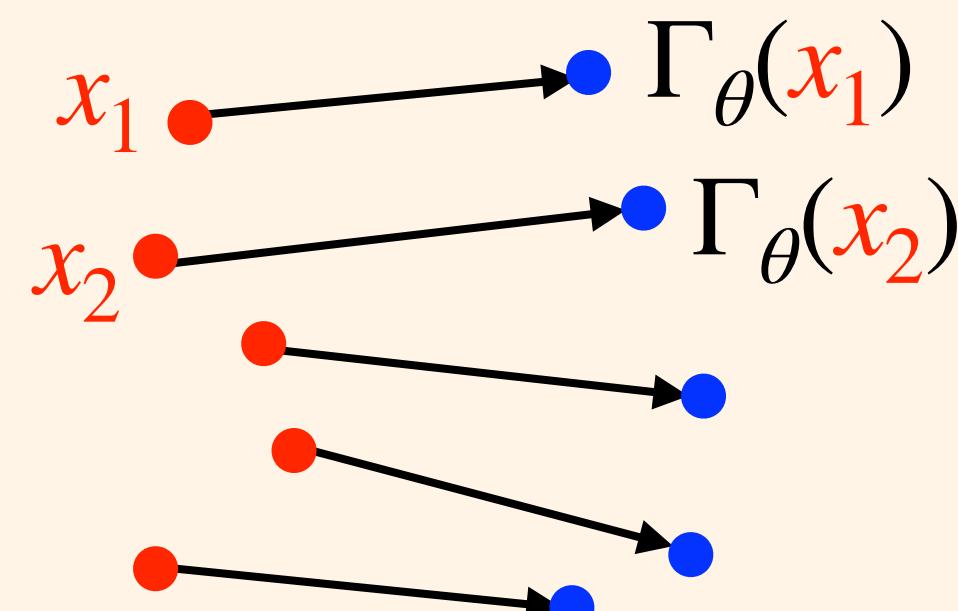
**Single-head attention layer:**  $X \mapsto \{\Gamma_\theta[\mathbf{X}](x_i)\}_{i=1}^n$

**Multi-head attention layer:**  $X \mapsto \left\{ \sum_{h=1}^H W_h \Gamma_{\theta_h}[\mathbf{X}](x_i) \right\}_{i=1}^n$



**Context-free layers:**  $X \mapsto \{\Gamma_\theta(x_i)\}_{i=1}^n$

Multi-layer perceptron:  $\Gamma_\theta(x) := x + \theta_1 \text{ReLU}(\theta_2 x)$

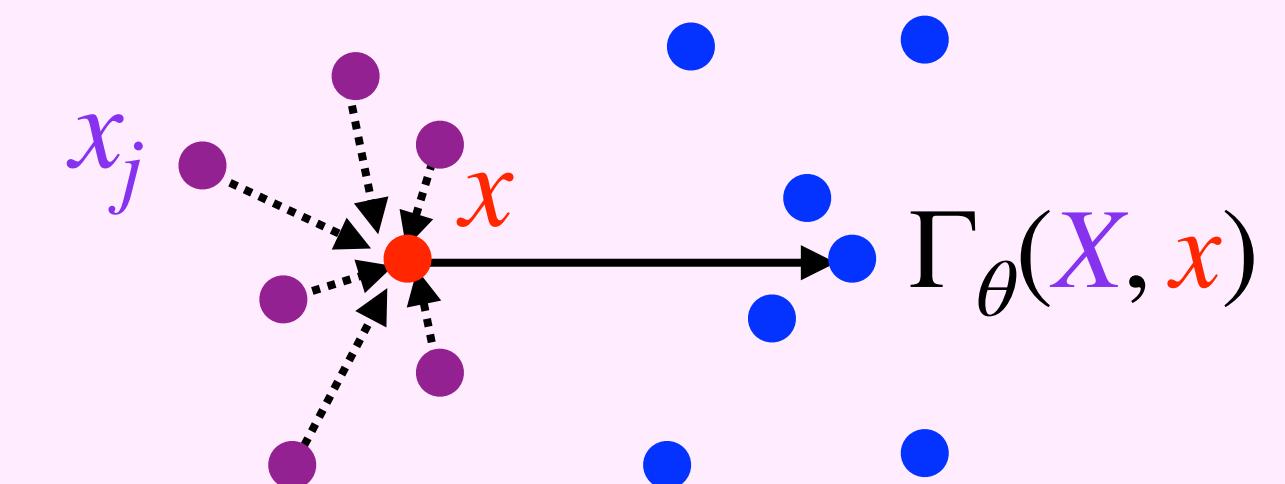


# Attention as In-context Mapping

Point clouds:  $X := \{x_i\}_{i=1}^n$

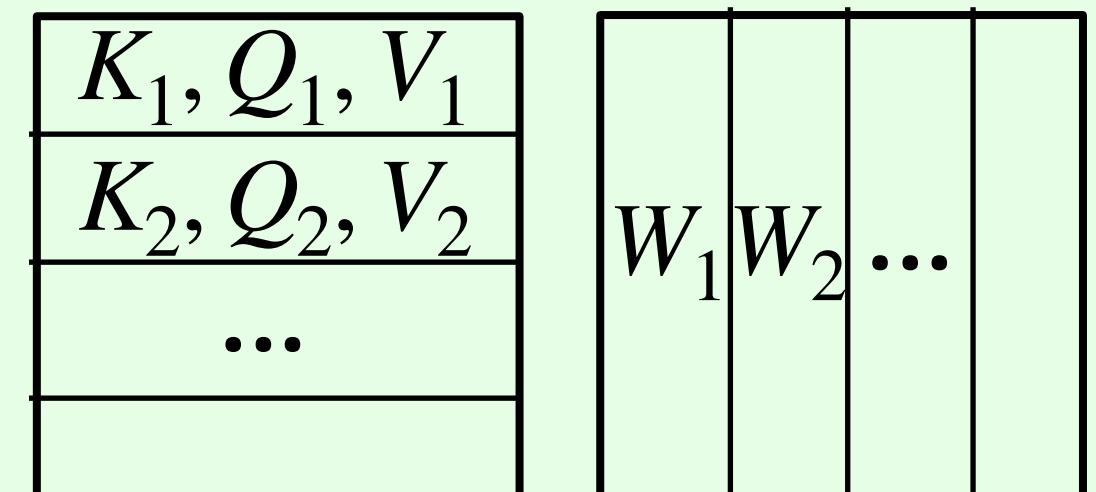
**In-context mapping:**  
parameters  $\theta := (Q, K, V)$

$$\Gamma_\theta[\mathbf{X}](\mathbf{x}) := \sum_j \frac{e^{\langle Q\mathbf{x}, K\mathbf{x}_j \rangle}}{\sum_\ell e^{\langle Q\mathbf{x}, K\mathbf{x}_\ell \rangle}} V\mathbf{x}_j$$



**Single-head attention layer:**  $X \mapsto \{\Gamma_\theta[\mathbf{X}](x_i)\}_{i=1}^n$

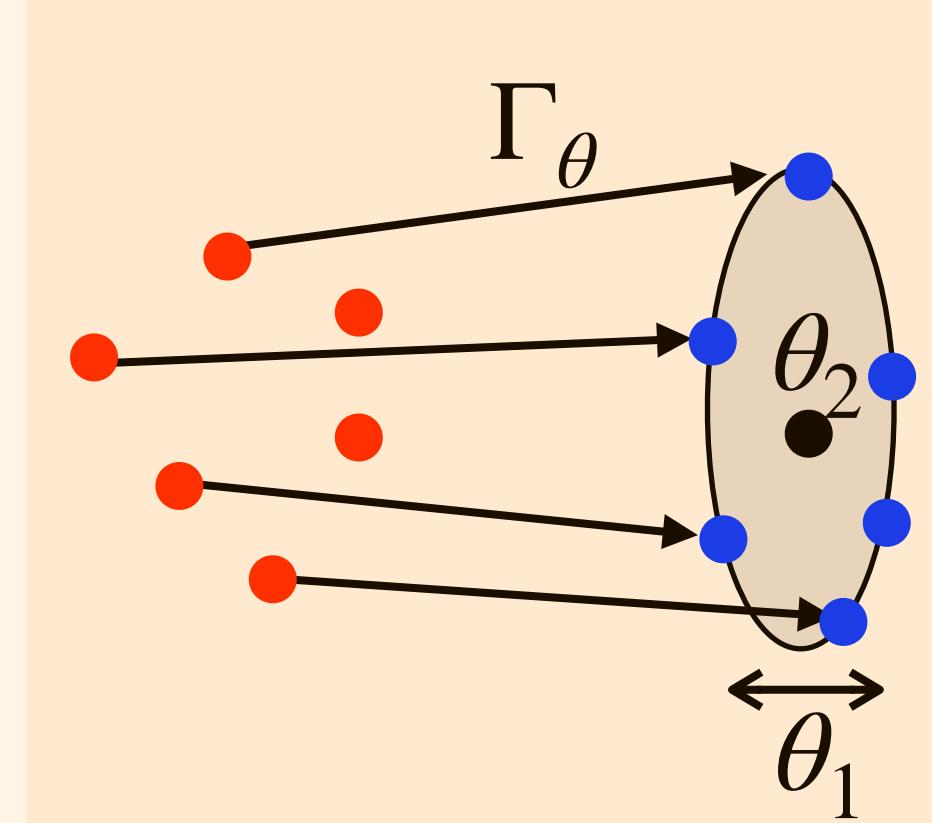
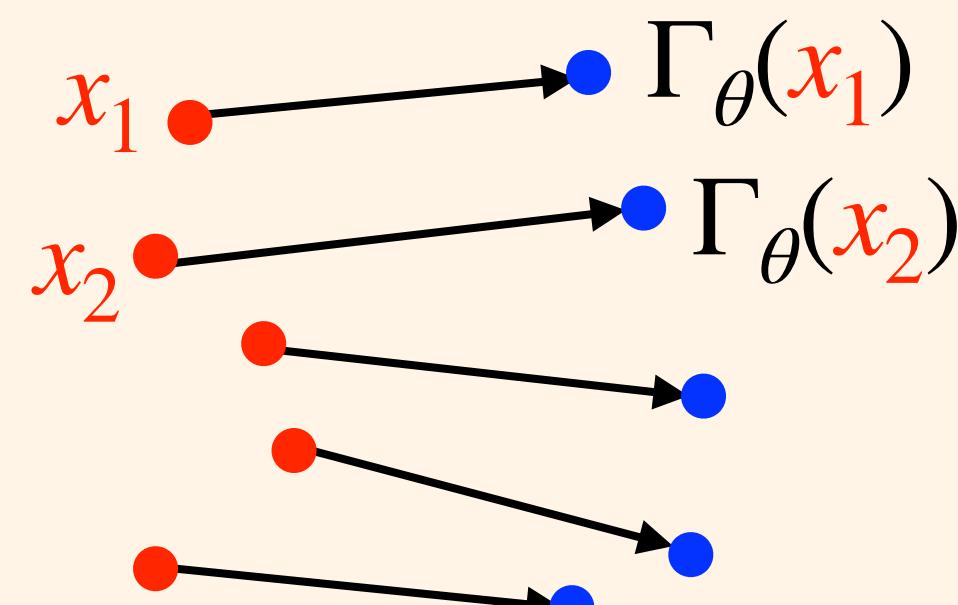
**Multi-head attention layer:**  $X \mapsto \left\{ \sum_{h=1}^H W_h \Gamma_{\theta_h}[\mathbf{X}](x_i) \right\}_{i=1}^n$



**Context-free layers:**  $X \mapsto \{\Gamma_\theta(x_i)\}_{i=1}^n$

Multi-layer perceptron:  $\Gamma_\theta(x) := x + \theta_1 \text{ReLU}(\theta_2 x)$

Layer norm:  $\Gamma_\theta(x) := \theta_1 \odot \frac{x}{\|x\|} + \theta_2$

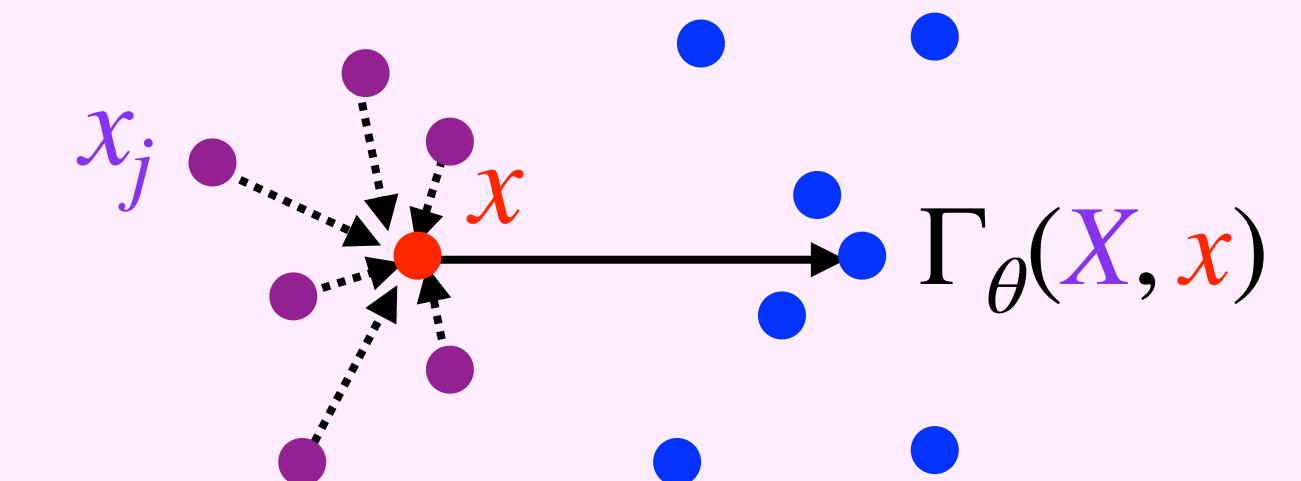


# Attention as In-context Mapping

Point clouds:  $X := \{x_i\}_{i=1}^n$

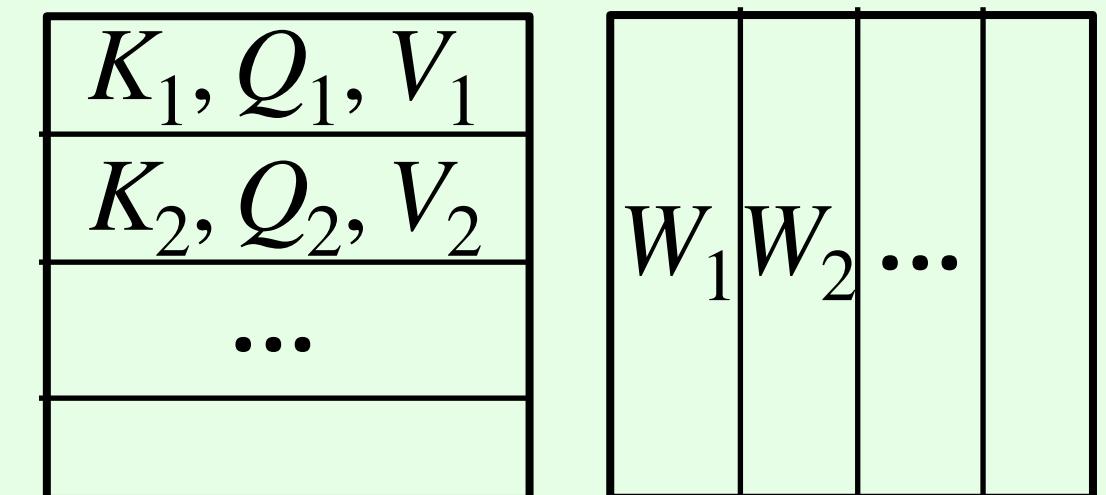
**In-context mapping:**  
parameters  $\theta := (Q, K, V)$

$$\Gamma_\theta[\mathbf{X}](\mathbf{x}) := \sum_j \frac{e^{\langle Q\mathbf{x}, K\mathbf{x}_j \rangle}}{\sum_\ell e^{\langle Q\mathbf{x}, K\mathbf{x}_\ell \rangle}} V\mathbf{x}_j$$



**Single-head attention layer:**  $X \mapsto \{\Gamma_\theta[\mathbf{X}](x_i)\}_{i=1}^n$

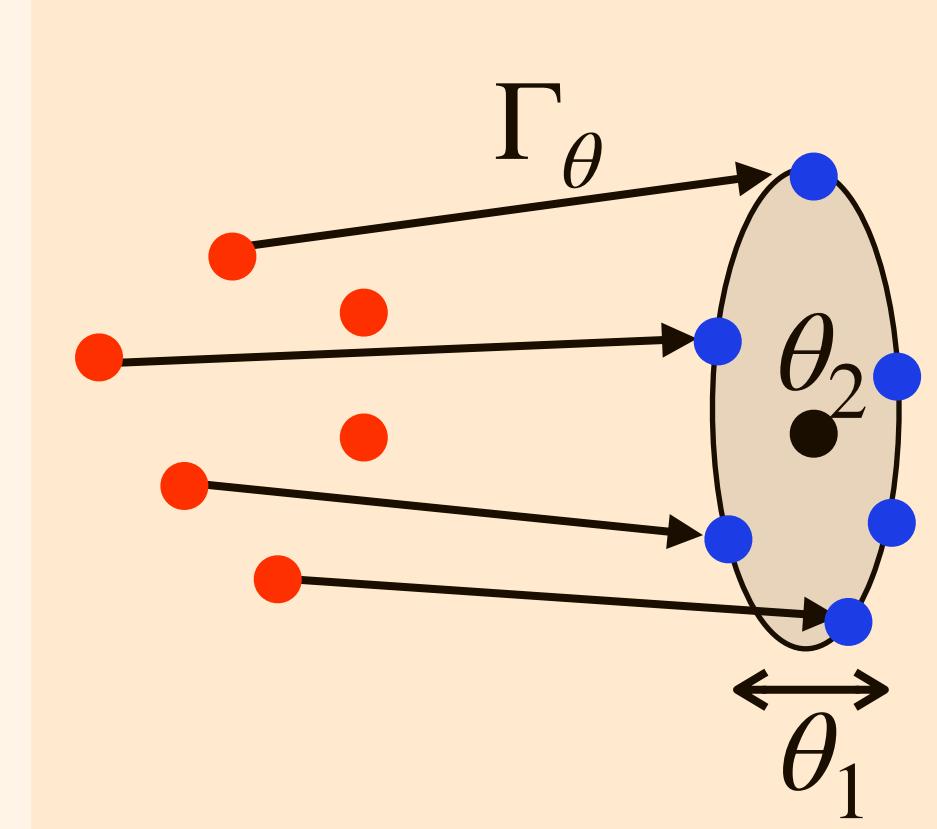
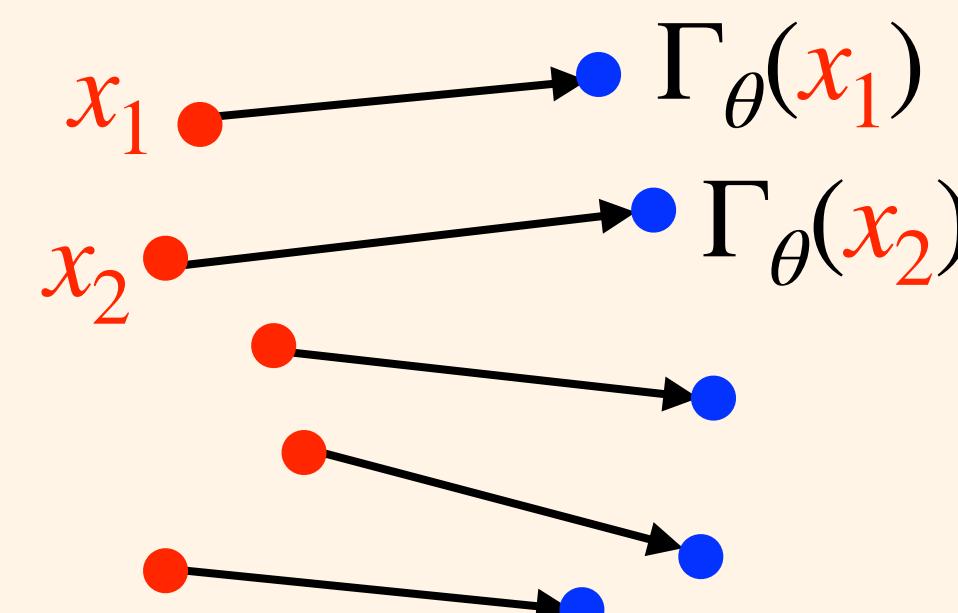
**Multi-head attention layer:**  $X \mapsto \left\{ \sum_{h=1}^H W_h \Gamma_{\theta_h}[\mathbf{X}](x_i) \right\}_{i=1}^n$



**Context-free layers:**  $X \mapsto \{\Gamma_\theta(x_i)\}_{i=1}^n$

Multi-layer perceptron:  $\Gamma_\theta(x) := x + \theta_1 \text{ReLU}(\theta_2 x)$

Layer norm:  $\Gamma_\theta(x) := \theta_1 \odot \frac{x}{\|x\|} + \theta_2$



Transformer  $\equiv$  composition of in-context and context-free layers.

# Attentions Operating over Measures

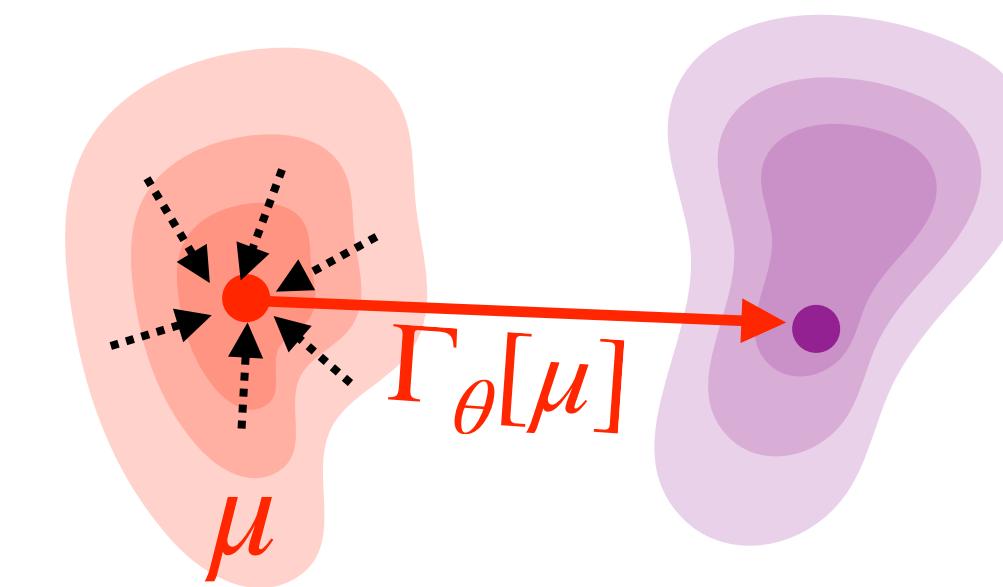
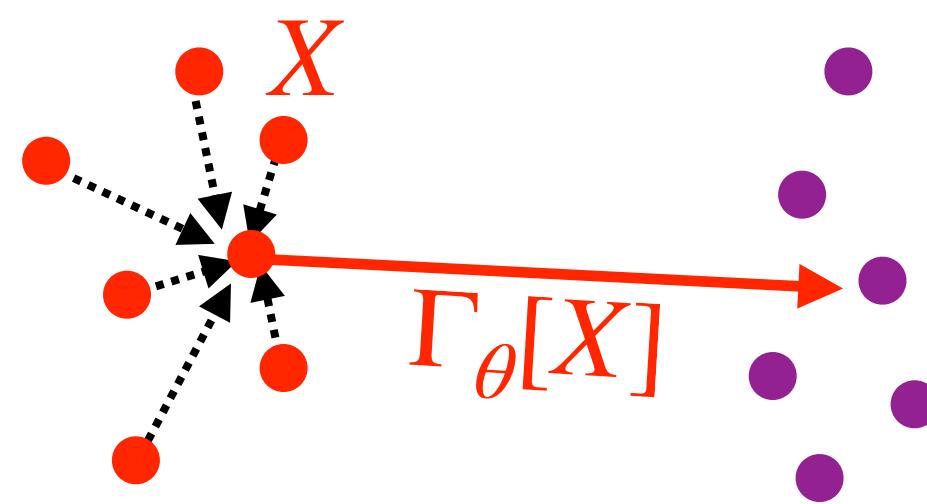
Number  $n$  of token is arbitrary.

(Unmasked) attention is permutation invariant.

$$\Gamma_\theta[X](x) := \sum_j \frac{e^{\langle Qx, Kx_j \rangle}}{\sum_\ell e^{\langle Qx, Kx_\ell \rangle}} Vx_j$$

$$\boxed{\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}}$$

$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky \rangle} d\mu(y)} Vy d\mu(y)$$

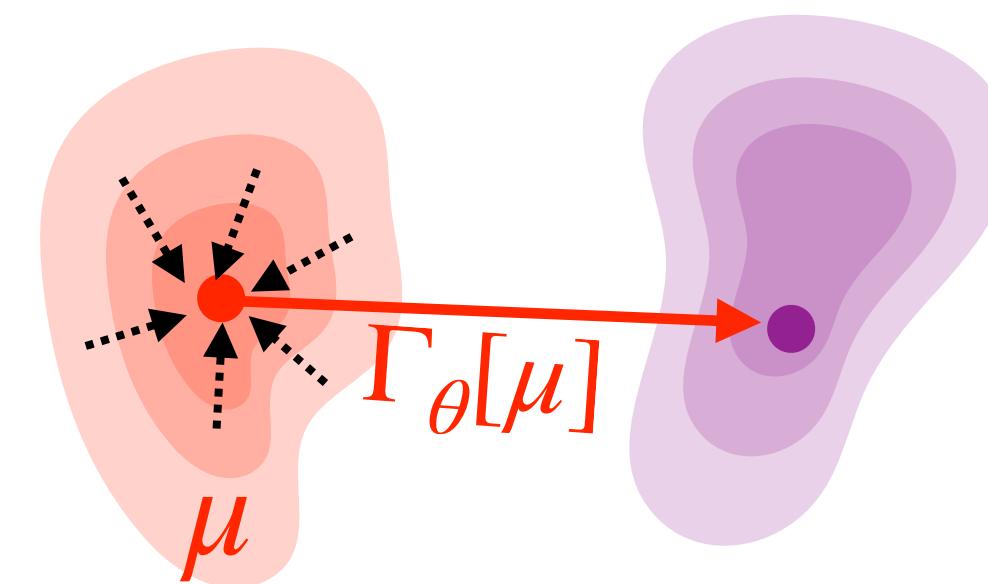
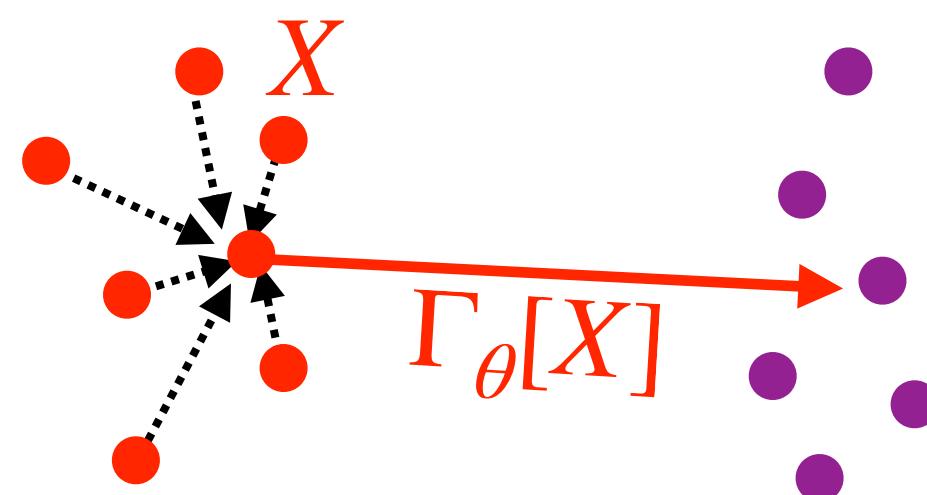


# Attentions Operating over Measures

Number  $n$  of token is arbitrary.

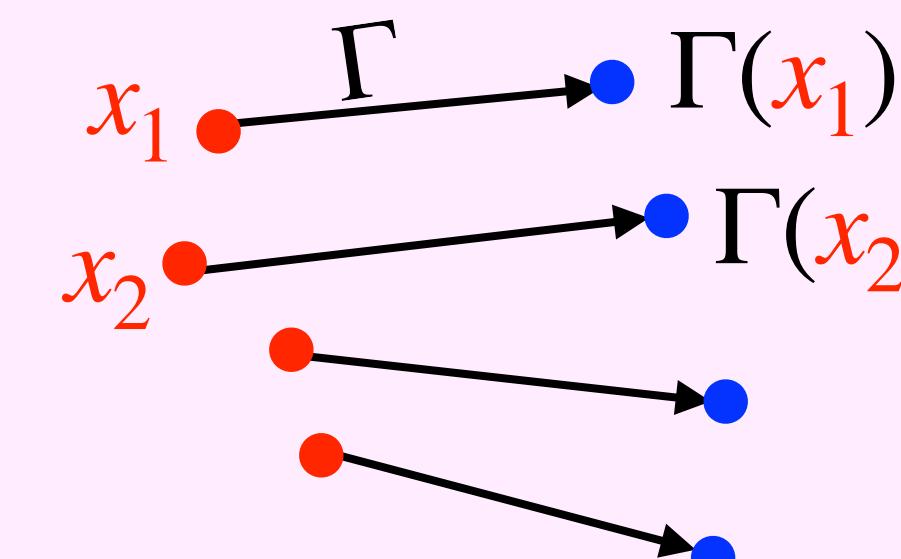
(Unmasked) attention is permutation invariant.

$$\Gamma_\theta[\textcolor{red}{X}](x) := \sum_j \frac{e^{\langle Qx, Kx_j \rangle}}{\sum_\ell e^{\langle Qx, Kx_\ell \rangle}} Vx_j \quad \boxed{\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}} \quad \Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky \rangle} d\mu(y)} Vy d\mu(y)$$



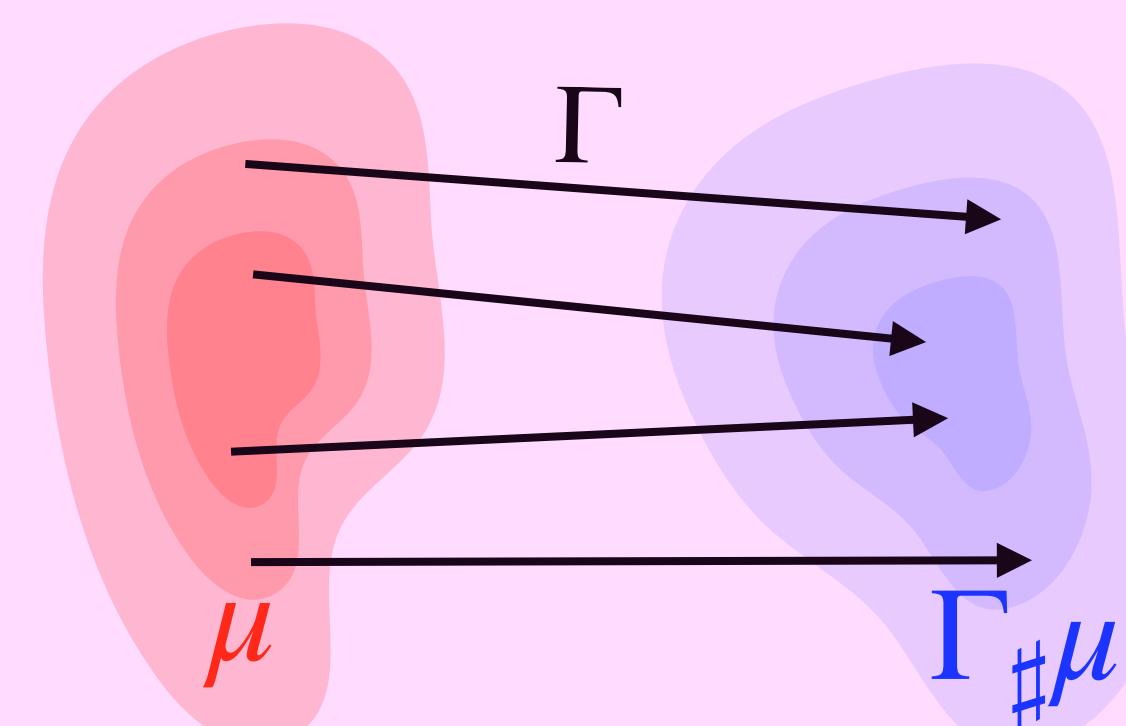
## Push-forward

$$\Gamma_\sharp \sum_i \delta_{x_i} := \sum_i \delta_{\Gamma(x_i)}$$



## Attention layers

$$X \mapsto \{\Gamma_\theta[\textcolor{violet}{X}](x_i)\}_{i=1}^n$$



$$(\Gamma_\sharp \mu)(B) := \mu(\Gamma^{-1}(B))$$

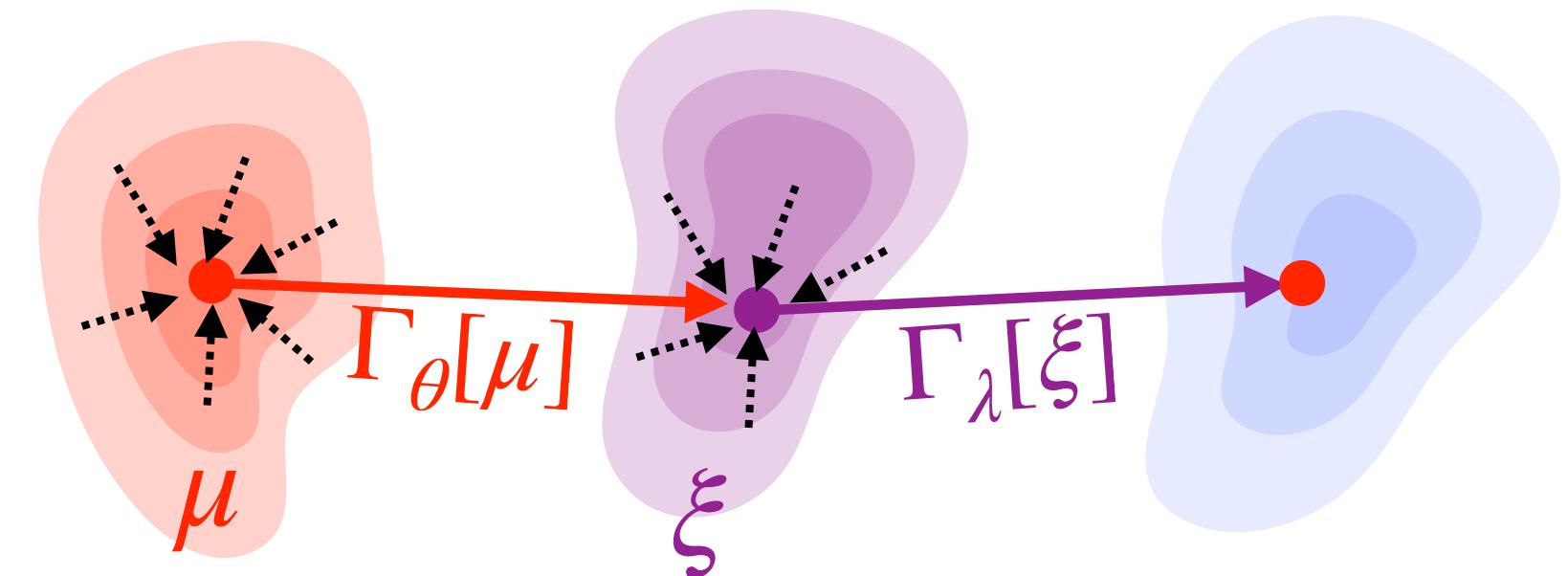
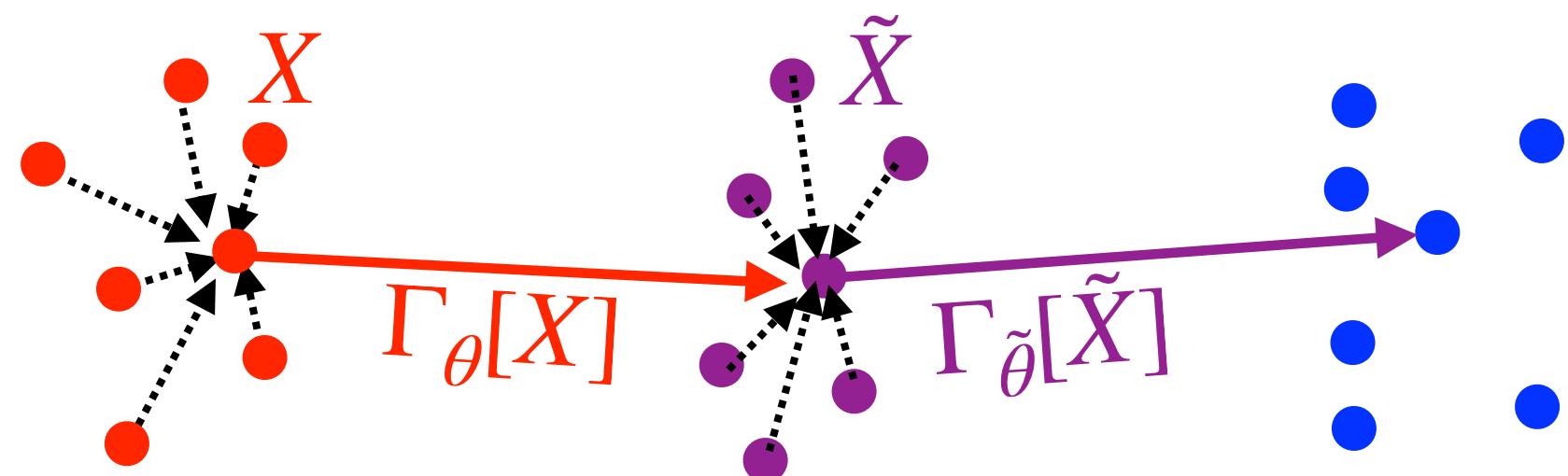
$$\mu \mapsto \Gamma_\theta[\textcolor{violet}{\mu}]_\sharp \mu$$

# Attentions Operating over Measures

Number  $n$  of token is arbitrary.

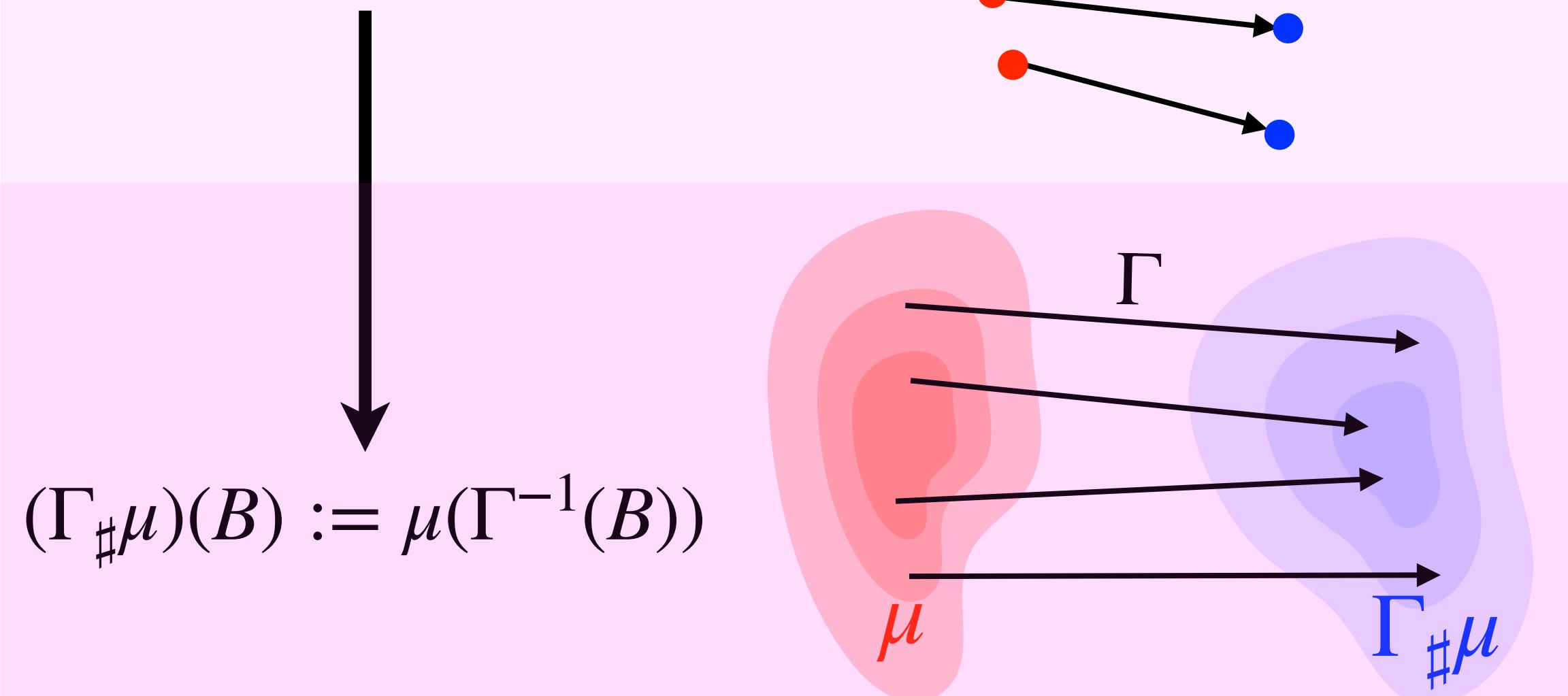
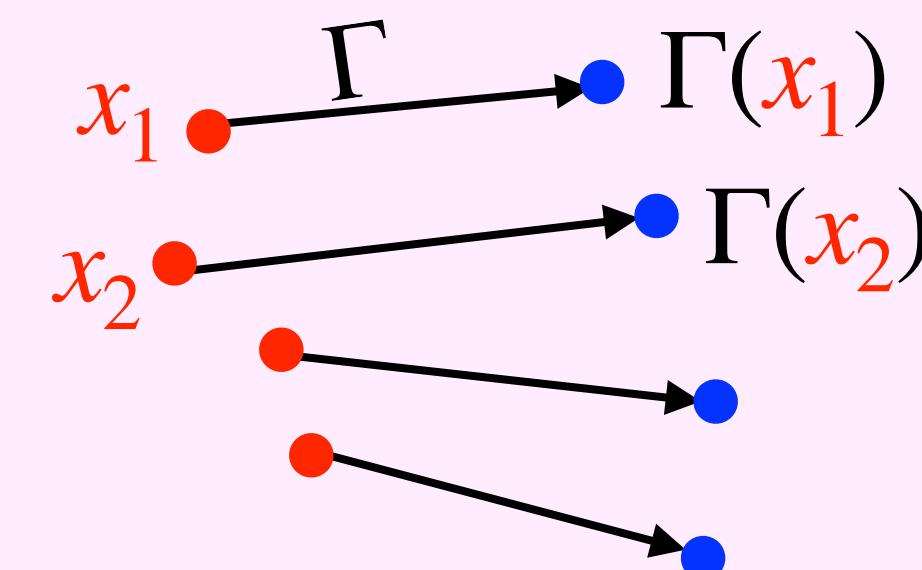
(Unmasked) attention is permutation invariant.

$$\Gamma_\theta[X](x) := \sum_j \frac{e^{\langle Qx, Kx_j \rangle}}{\sum_\ell e^{\langle Qx, Kx_\ell \rangle}} Vx_j \quad \boxed{\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}} \quad \Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky \rangle} d\mu(y)} Vy d\mu(y)$$



## Push-forward

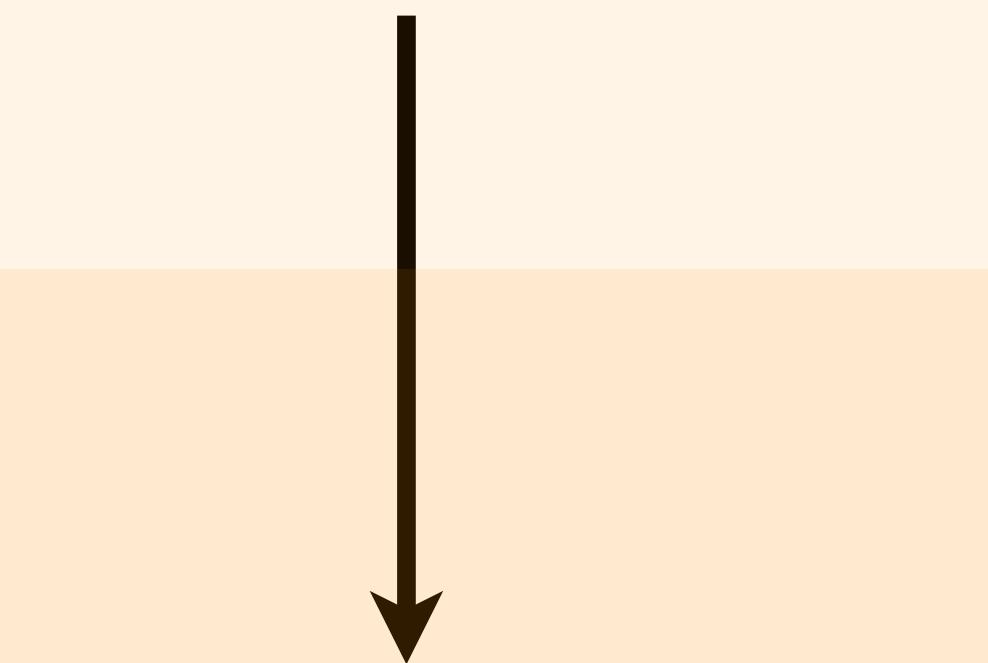
$$\Gamma_\sharp \sum_i \delta_{x_i} := \sum_i \delta_{\Gamma(x_i)}$$



$$(\Gamma_\sharp \mu)(B) := \mu(\Gamma^{-1}(B))$$

## Attention layers

$$X \mapsto \{\Gamma_\theta[X](x_i)\}_{i=1}^n$$

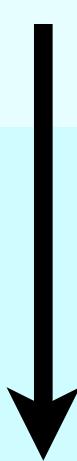


$$\mu \mapsto \Gamma_\theta[\mu]_\sharp \mu$$

## Composing layers

$$(\Gamma_\lambda \diamond \Gamma_\theta)[X] := \Gamma_\lambda[Y] \circ \Gamma_\theta[X]$$

$$\text{where } Y := (\Gamma_\theta[X](x_i))_i$$



$$(\Gamma_\lambda \diamond \Gamma_\theta)[\mu] := \Gamma_\lambda[\xi] \circ \Gamma_\theta[\mu]$$

$$\text{where } \xi := \Gamma_\theta[\mu]_\sharp \mu$$

# Masked Causal Attention over Measures

For NLP: architectures must be **causal** for next token prediction & generative modeling.

*Masked attention mapping:*     $\Gamma_\theta[\mathbf{X}](\mathbf{x}_i) := \sum_{j \leq i} \frac{e^{\langle Q\mathbf{x}_i, K\mathbf{x}_j \rangle}}{\sum_{\ell \leq i} e^{\langle Q\mathbf{x}_i, K\mathbf{x}_\ell \rangle}} V\mathbf{x}_j$

→ breaks permutation invariance.

# Masked Causal Attention over Measures

For NLP: architectures must be **causal** for next token prediction & generative modeling.

*Masked attention mapping:*

$$\Gamma_\theta[\mathbf{X}](\mathbf{x}_i) := \sum_{j \leq i} \frac{e^{\langle Q\mathbf{x}_i, K\mathbf{x}_j \rangle}}{\sum_{\ell \leq i} e^{\langle Q\mathbf{x}_i, K\mathbf{x}_\ell \rangle}} V\mathbf{x}_j$$

→ breaks permutation invariance.

*Training:* next token prediction  
(*simplified...*)

$$\min_{\theta} \sum_{\mathbf{X}} \sum_{i=1}^{n-1} \ell(\Gamma_\theta[\mathbf{X}](x_i), x_{i+1})$$

*Testing:* generative model  
(*simplified...*)

$$\mathbf{X} \mapsto (x_1, \dots, x_i, \Gamma[\mathbf{X}](x_i))$$

# Masked Causal Attention over Measures

For NLP: architectures must be **causal** for next token prediction & generative modeling.

Masked attention mapping:

$$\Gamma_\theta[X](\mathbf{x}_i) := \sum_{j \leq i} \frac{e^{\langle Q\mathbf{x}_i, K\mathbf{x}_j \rangle}}{\sum_{\ell \leq i} e^{\langle Q\mathbf{x}_i, K\mathbf{x}_\ell \rangle}} V\mathbf{x}_j$$

→ breaks permutation invariance.

Training: next token prediction  
(simplified...)

$$\min_{\theta} \sum_X \sum_{i=1}^{n-1} \ell(\Gamma_\theta[X](x_i), x_{i+1})$$

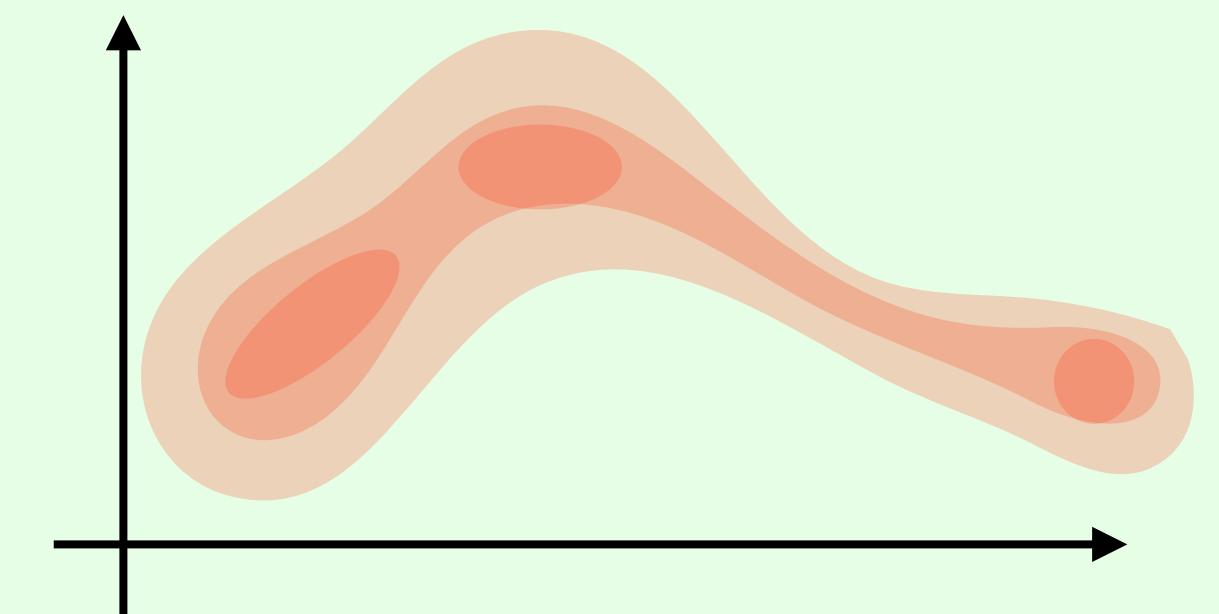
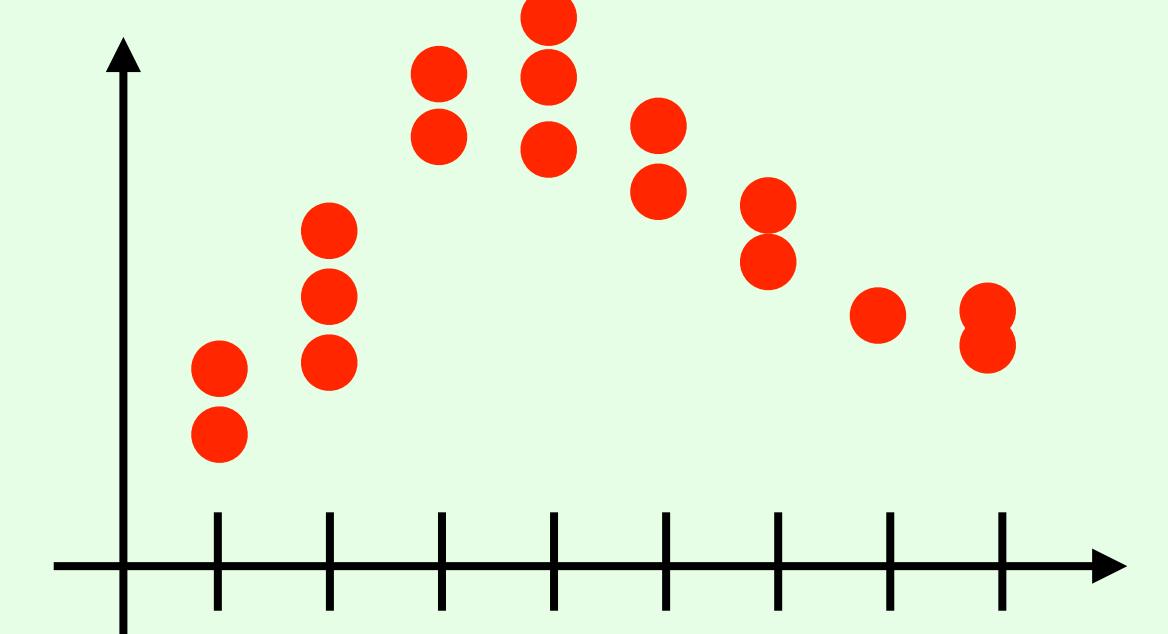
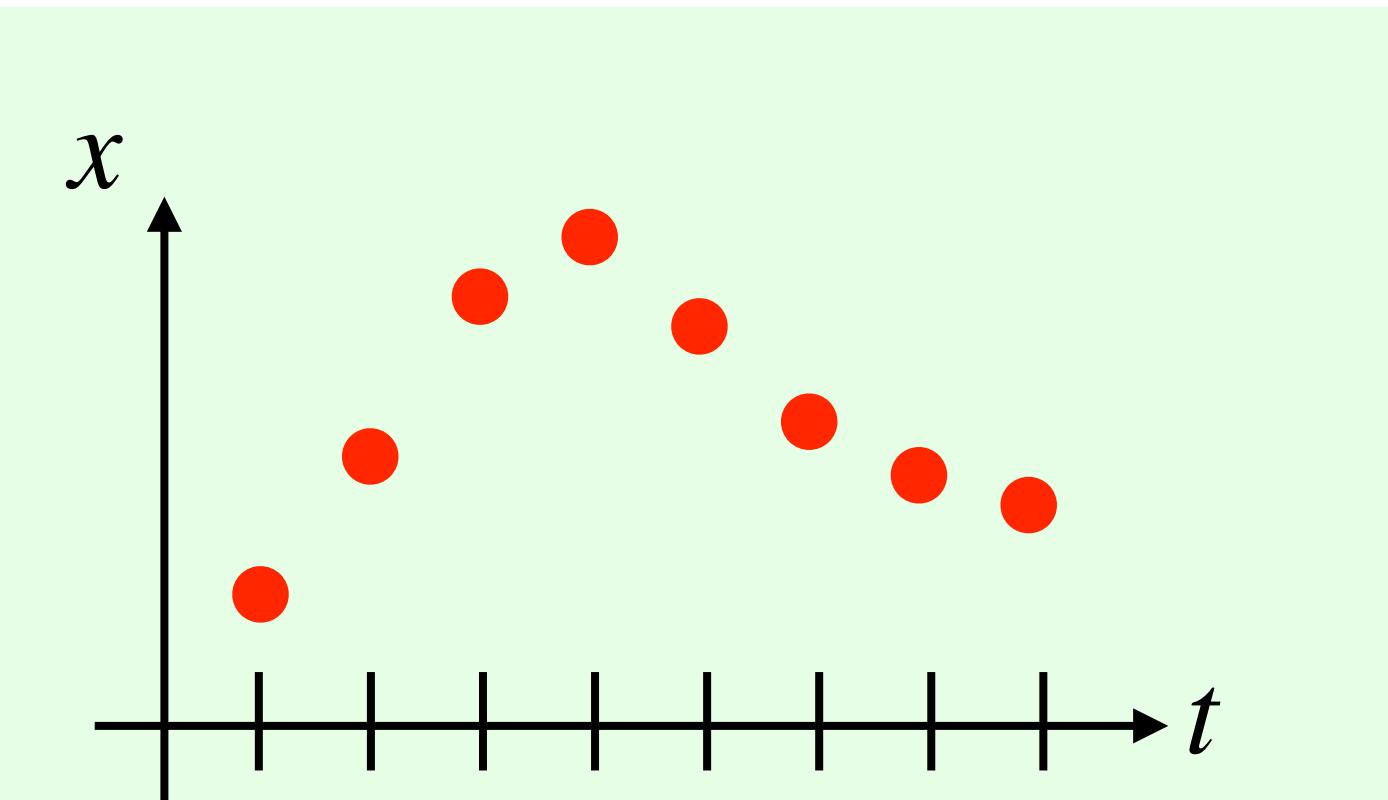
Testing: generative model  
(simplified...)

$$X \mapsto (x_1, \dots, x_i, \Gamma[X](x_i))$$

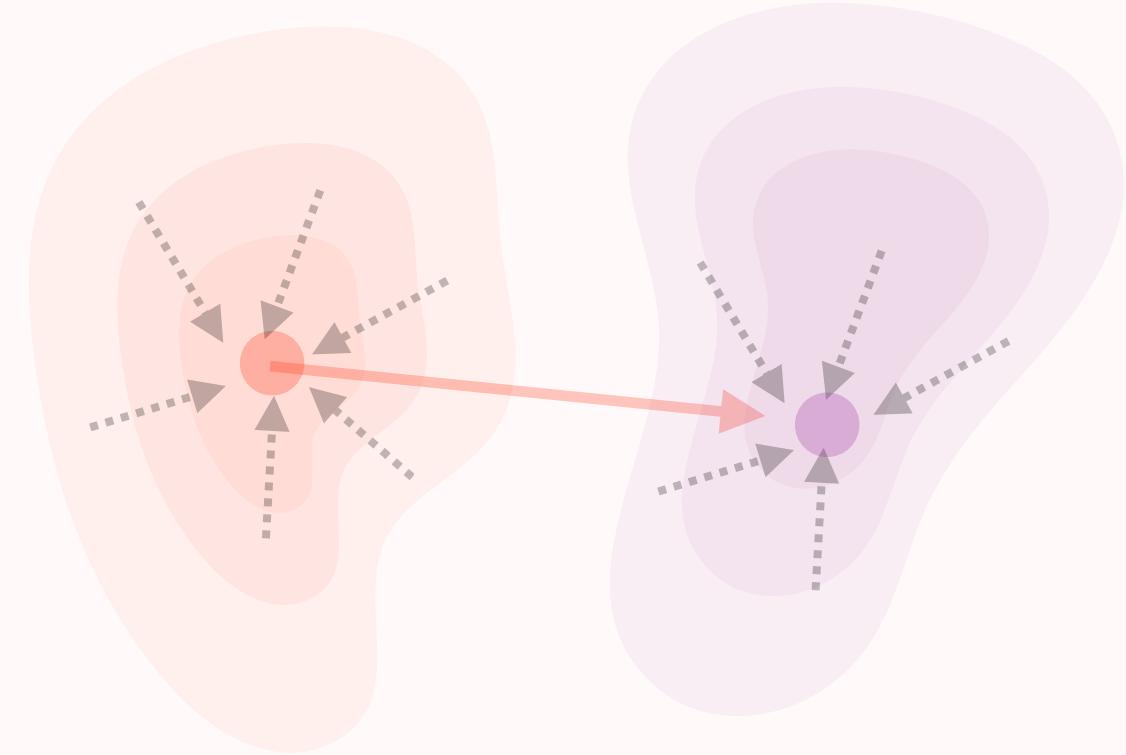
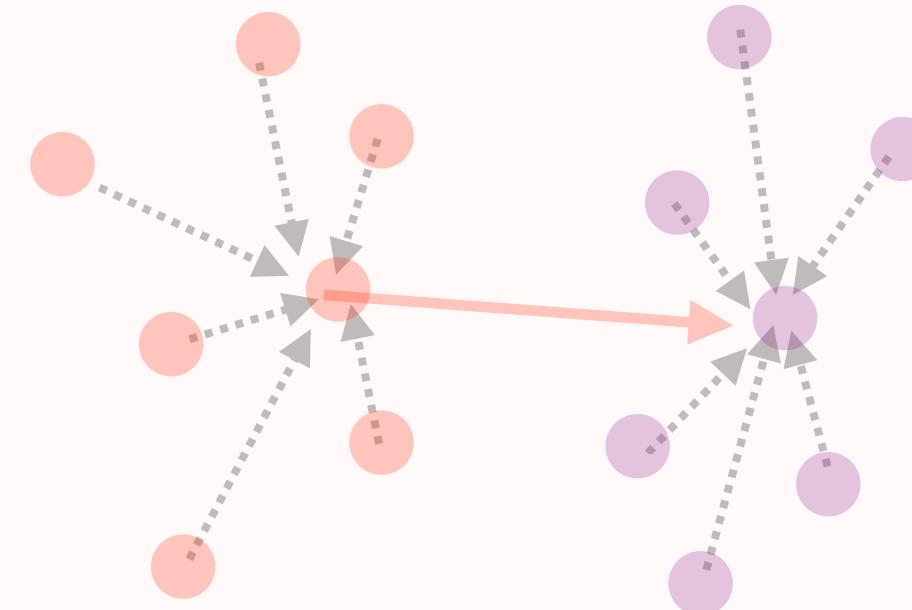
Space-time lifting:

$$\mu = \frac{1}{n} \sum_{i=1}^n \delta_{(x_i, t_i)}$$

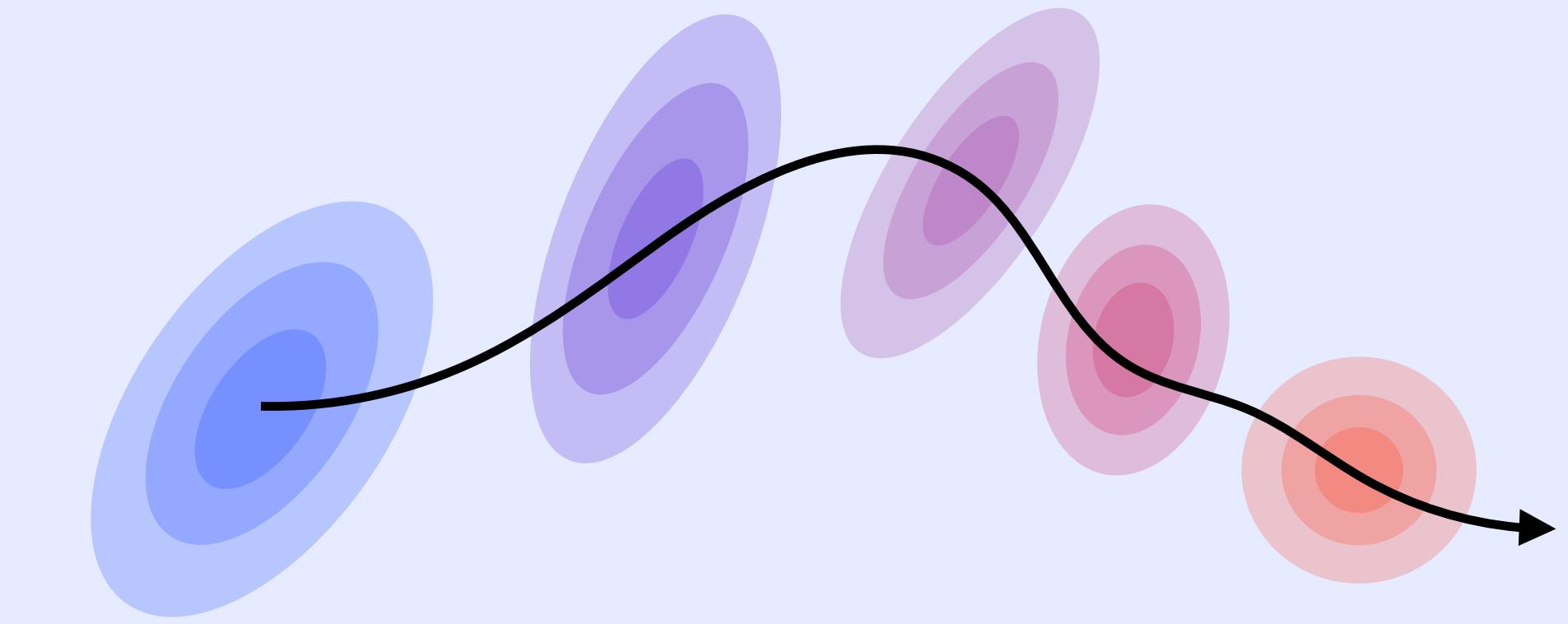
$$\Gamma_\theta[\mu](\mathbf{x}, \mathbf{t}) := \int \frac{1_{s \leq t} e^{\langle Q\mathbf{x}, K\mathbf{y} \rangle}}{\int 1_{s' \leq t} e^{\langle Q\mathbf{x}, K\mathbf{y}' \rangle} d\mu(y', s')} V\mathbf{y} d\mu(y, s)$$



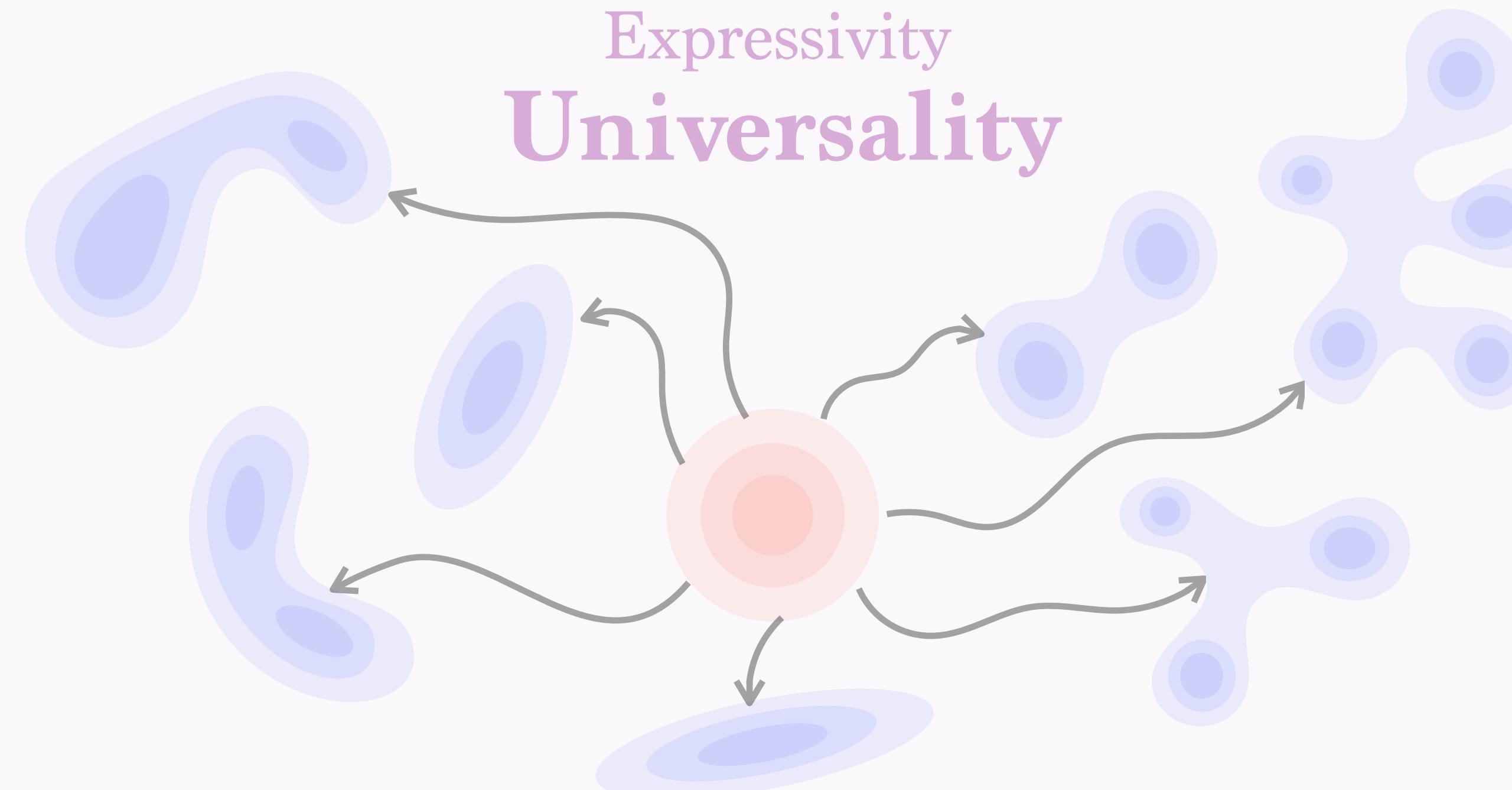
Arbitrary number of layers  
**In Context Mappings  
over Measures**



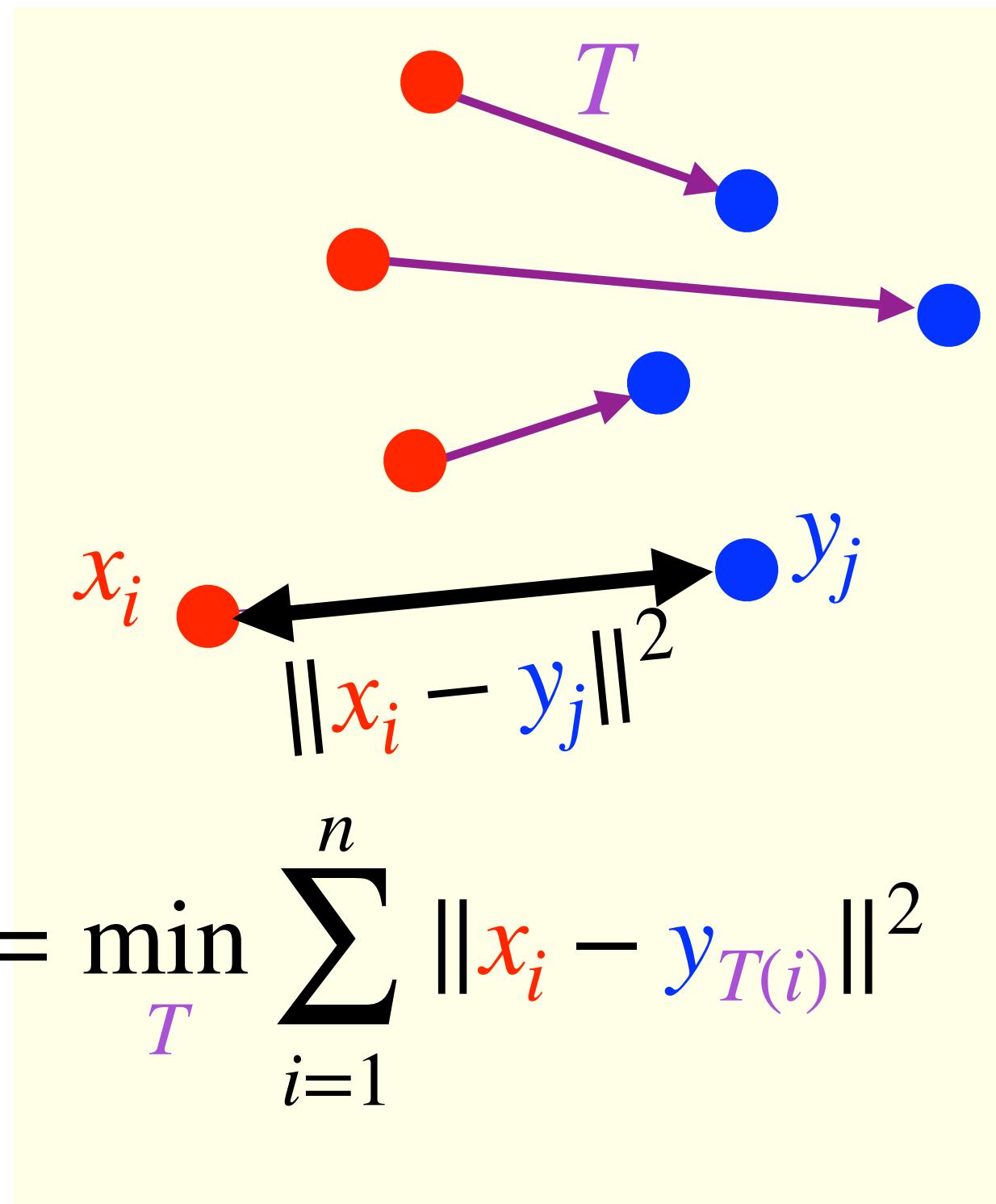
Arbitrary number of layers  
**Smoothness and  
PDE's**



Expressivity  
**Universality**

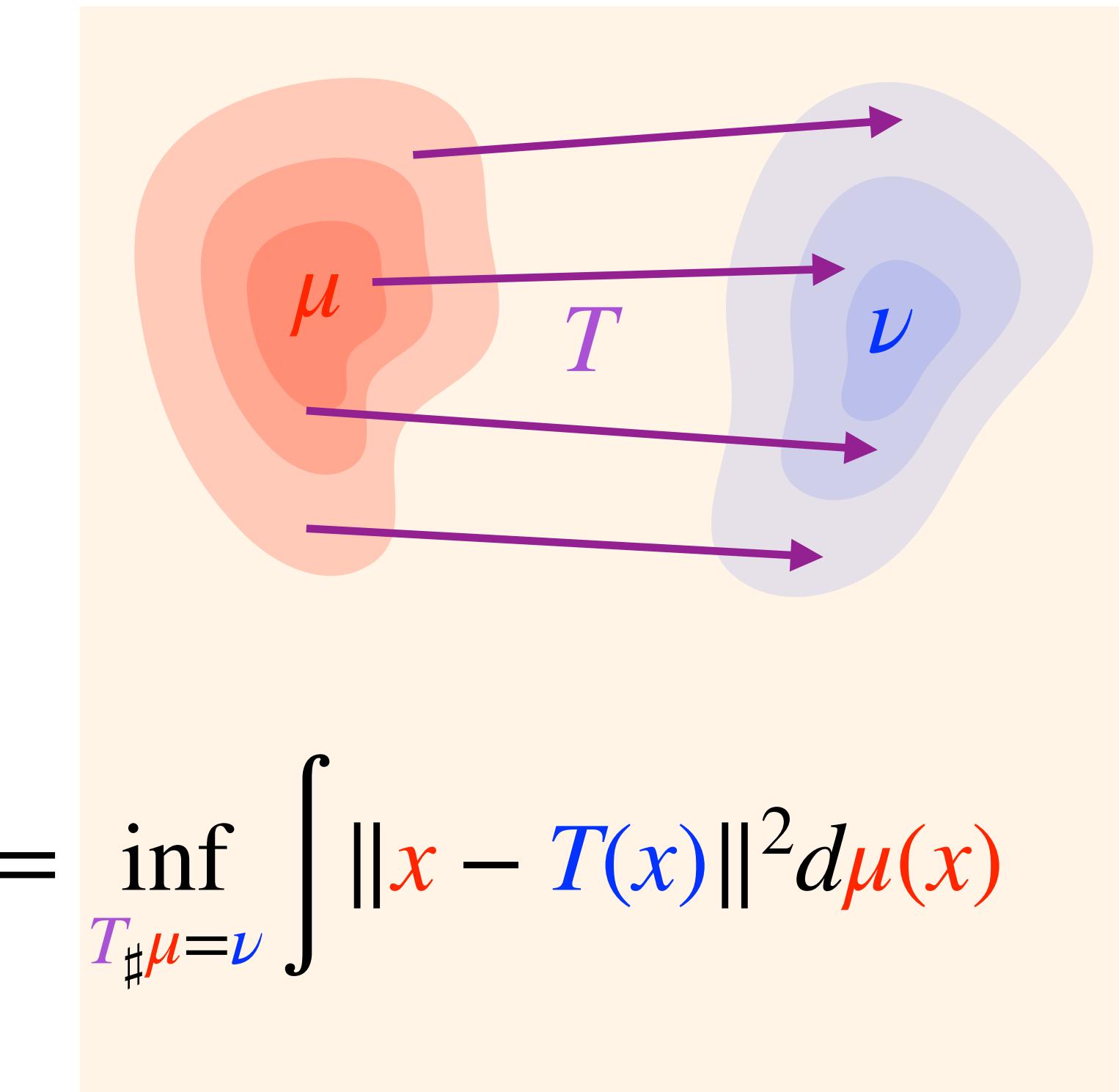
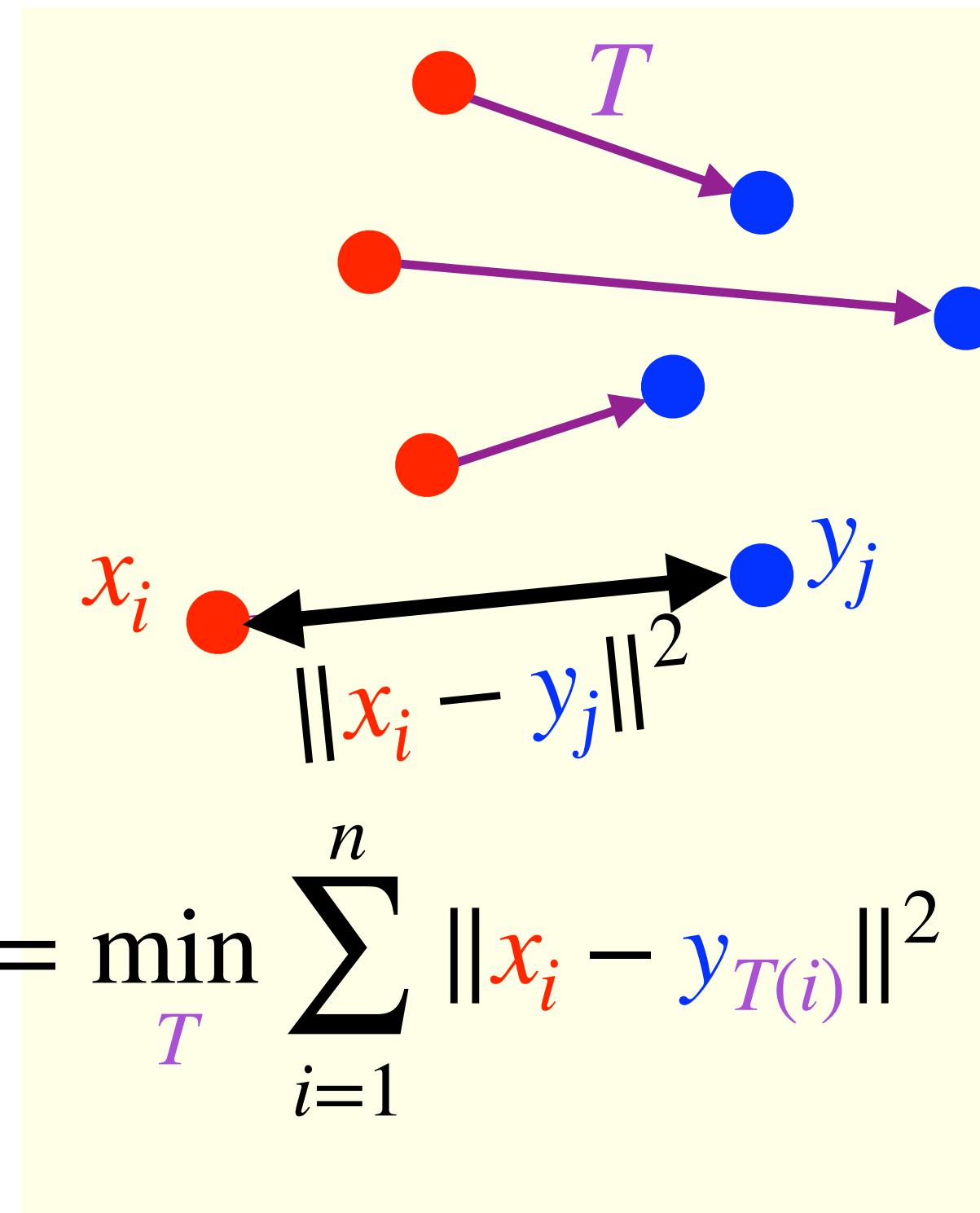


# Optimal Transport (Wasserstein) Distance



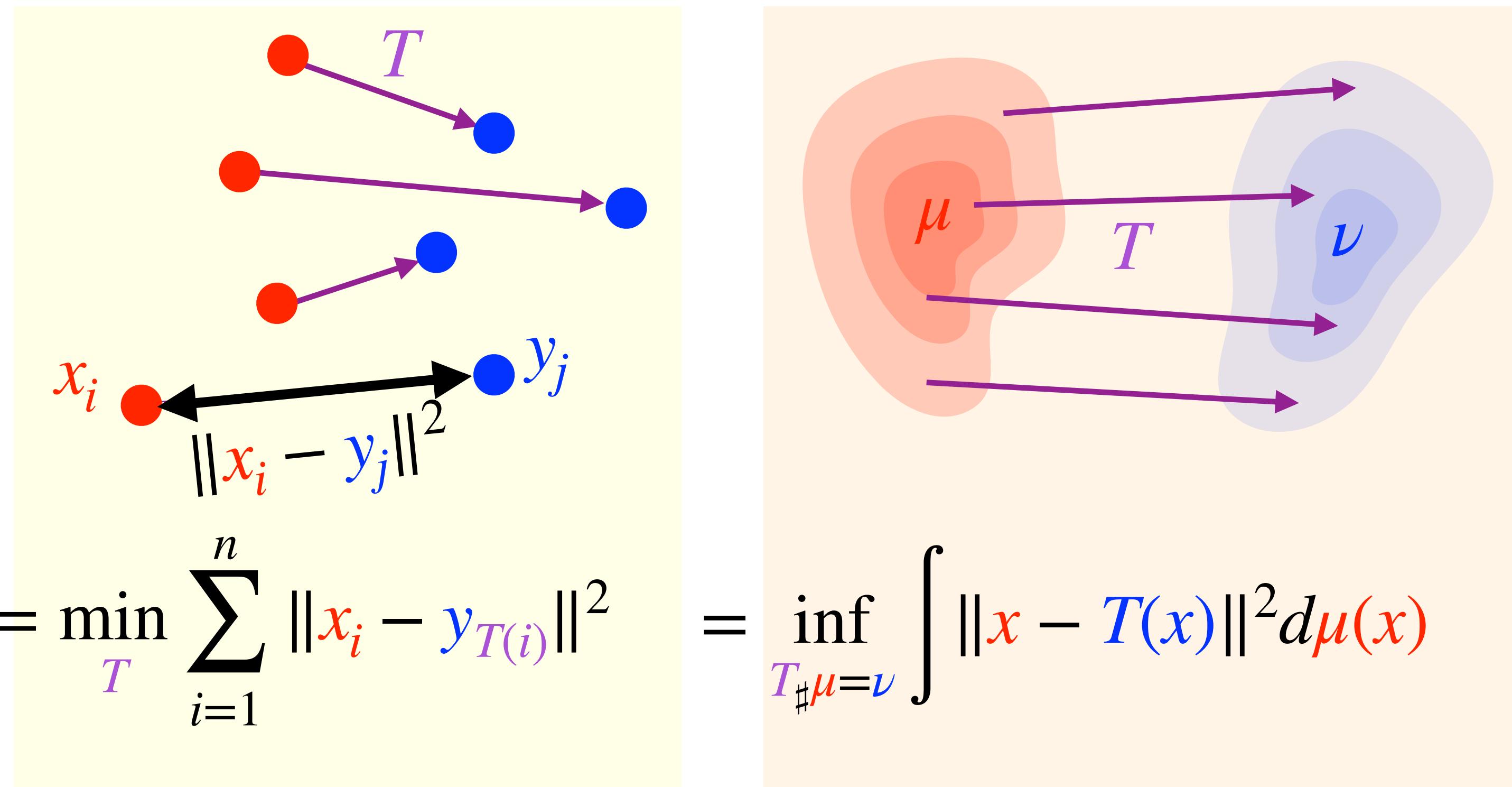
Monge 1784

# Optimal Transport (Wasserstein) Distance

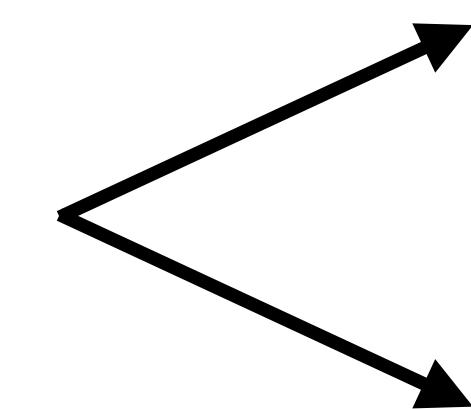


Monge 1784

# Optimal Transport (Wasserstein) Distance



General measures:



Kantorovitch relaxation

or

Approximation by discrete measures



Monge 1784



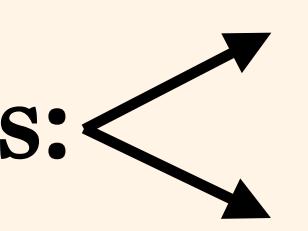
Kantorovitch 1942

# How Smooth is Attention?

Attention layer:  $\mu \mapsto \Gamma_\theta[\mu] \# \mu$

$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky' \rangle} d\mu(y')} Vy \, d\mu(y)$$

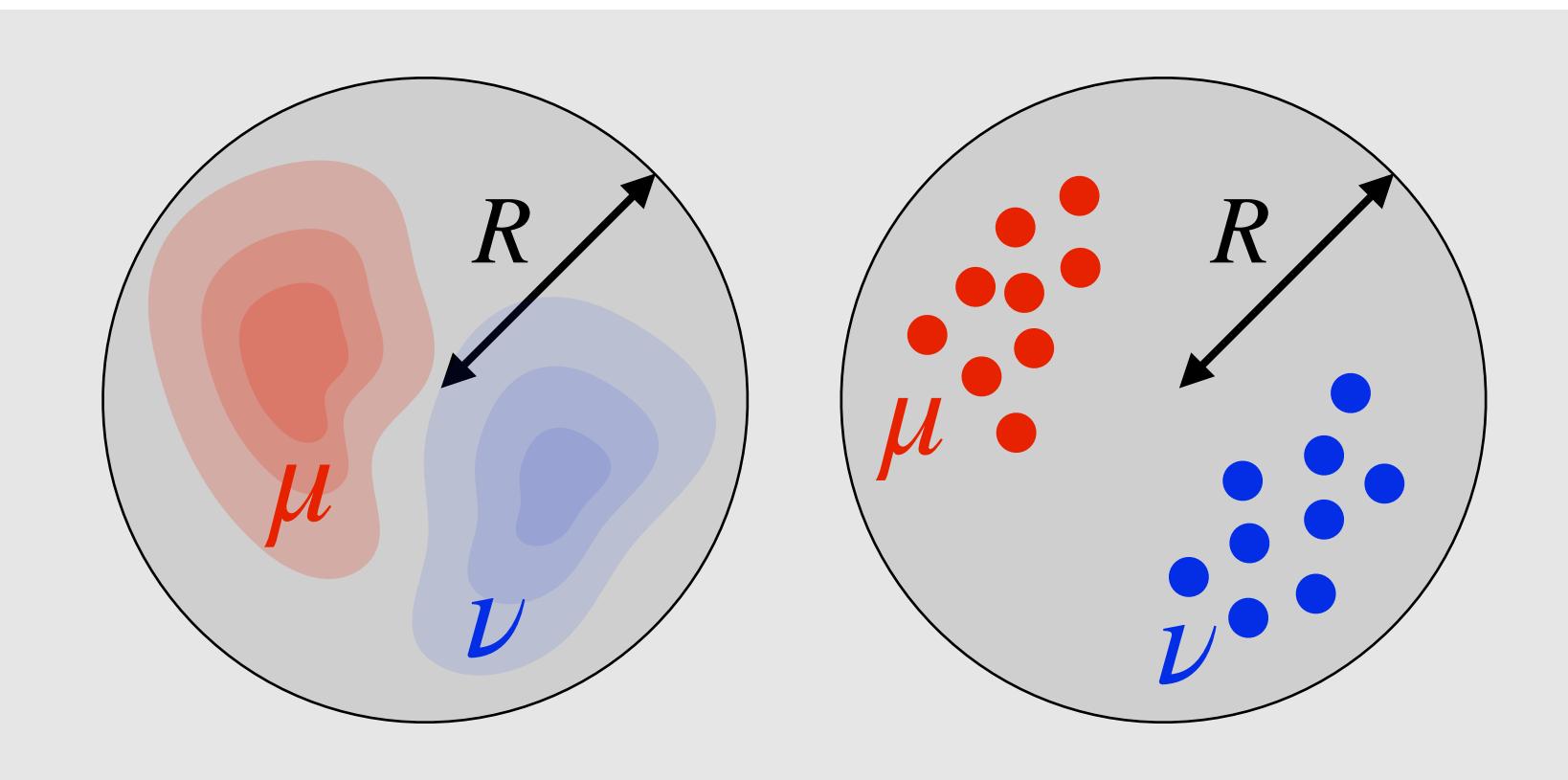
Lipschitz regularity:  $W_2(\Gamma_\theta[\mu] \# \mu, \Gamma_\theta[\nu] \# \nu) \leq C_\theta W_2(\mu, \nu)$

**Applications:**  Understanding robustness to attacks.  
Well-posedness of very deep transformers.

# How Smooth is Attention?

Attention layer:  $\mu \mapsto \Gamma_\theta[\mu] \# \mu$

$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky' \rangle} d\mu(y')} Vy \, d\mu(y)$$



Lipschitz regularity:  $W_2(\Gamma_\theta[\mu] \# \mu, \Gamma_\theta[\nu] \# \nu) \leq C_\theta W_2(\mu, \nu)$

**Applications:**

Understanding robustness to attacks.

Well-posedness of very deep transformers.

*Theorem:* [Castin, Peyré, Ablin]

If  $\text{supp}(\mu), \text{supp}(\nu) \subset B(0, R)$ ,

$$C_\theta \leq \|V\| (1 + 3\|K^\top Q\| R^2) e^{2\|K^\top Q\| R^2}$$

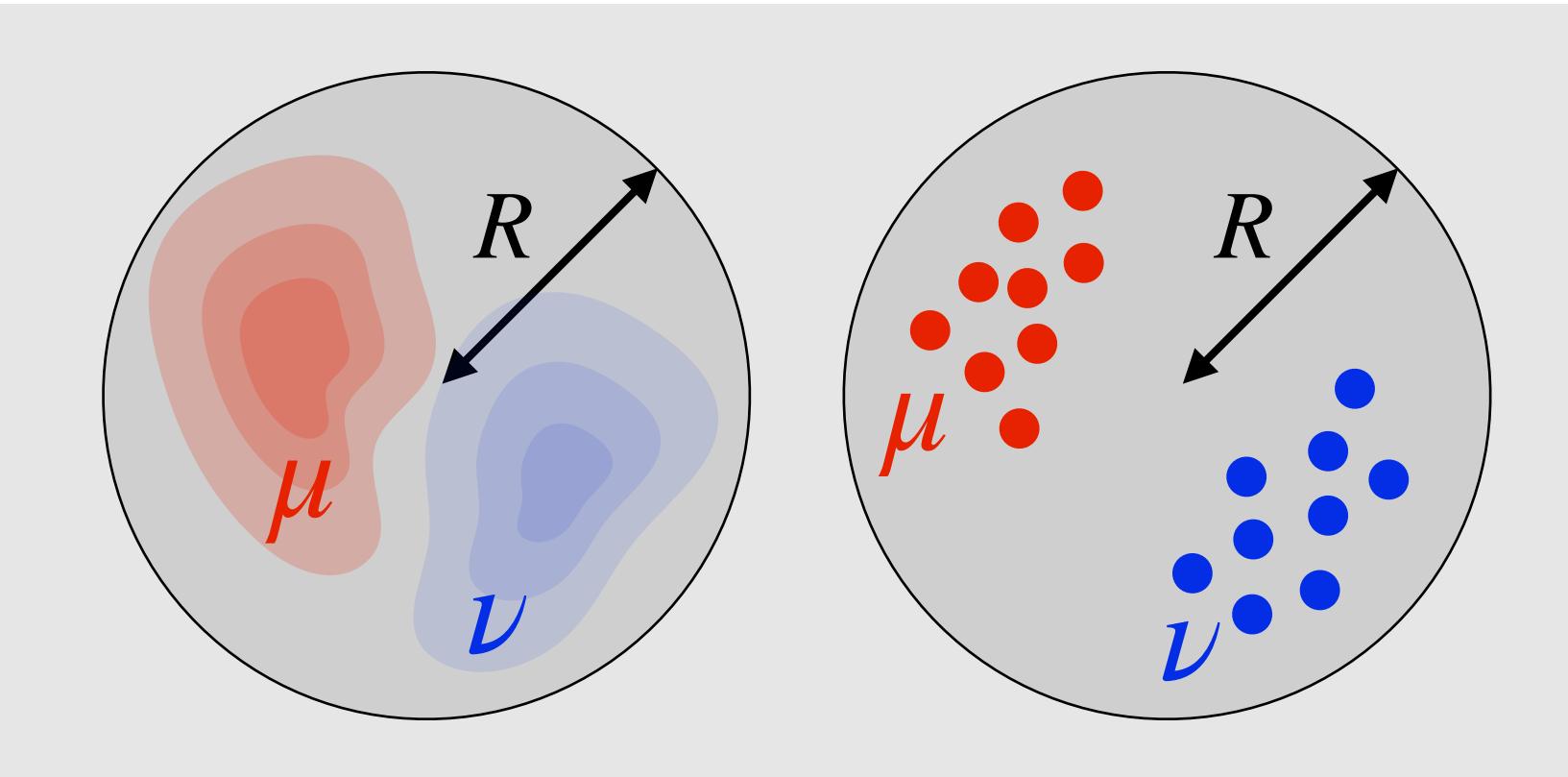
If furthermore  $\mu = \frac{1}{n} \sum_i \delta_{x_i}, \nu = \frac{1}{n} \sum_i \delta_{y_i}$

$$C_\theta \leq \|V\| \|K^\top Q\| R^2 \sqrt{12n + 3}$$

# How Smooth is Attention?

Attention layer:  $\mu \mapsto \Gamma_\theta[\mu] \# \mu$

$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky' \rangle} d\mu(y')} Vy \, d\mu(y)$$



Lipschitz regularity:  $W_2(\Gamma_\theta[\mu] \# \mu, \Gamma_\theta[\nu] \# \nu) \leq C_\theta W_2(\mu, \nu)$

**Applications:**

- Understanding robustness to attacks.
- Well-posedness of very deep transformers.

Theorem: [Castin, Peyré, Ablin]

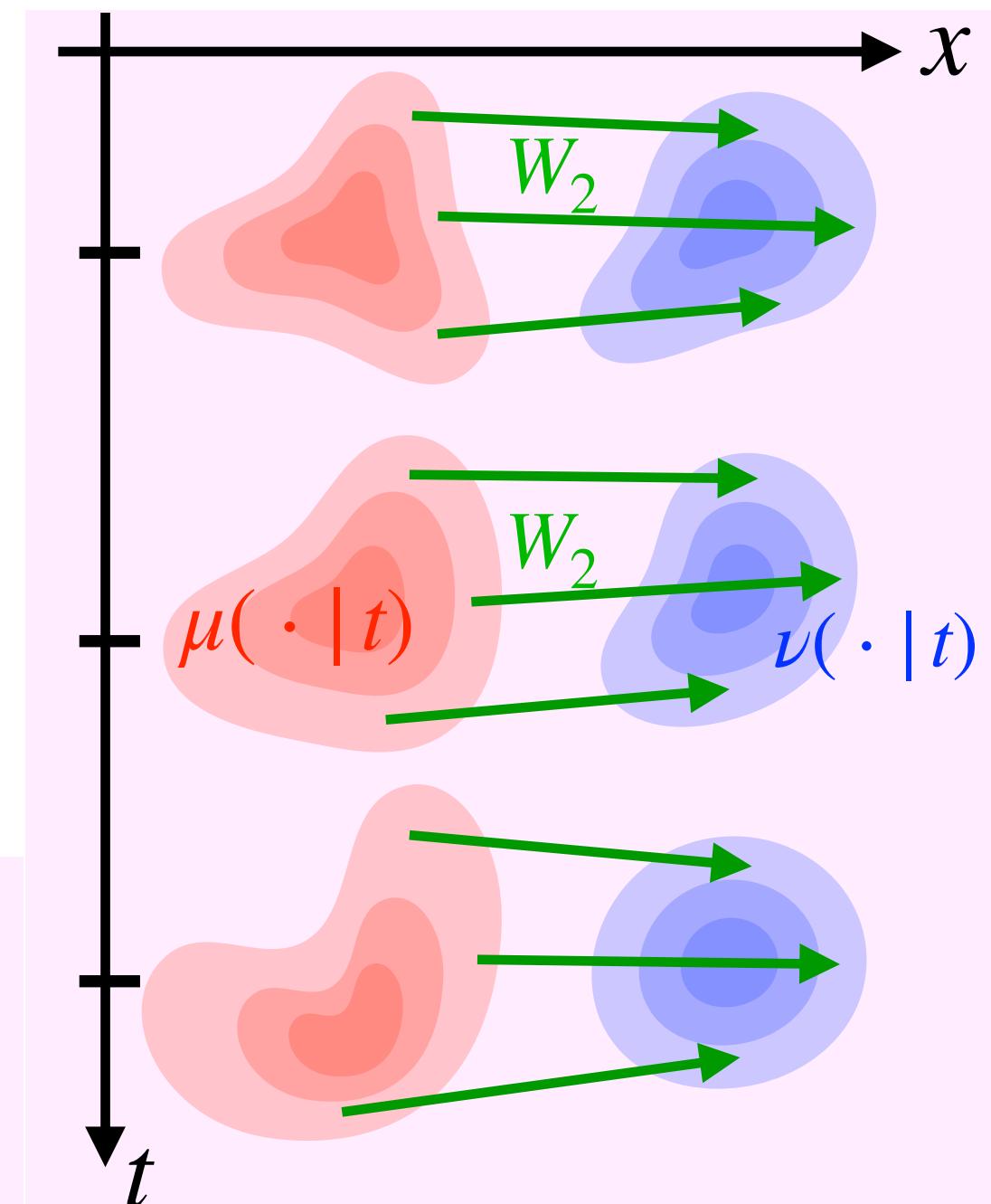
If  $\text{supp}(\mu), \text{supp}(\nu) \subset B(0, R)$ ,

$$C_\theta \leq \|V\|(1 + 3\|K^\top Q\|R^2)e^{2\|K^\top Q\|R^2}$$

If furthermore  $\mu = \frac{1}{n} \sum_i \delta_{x_i}, \nu = \frac{1}{n} \sum_i \delta_{y_i}$

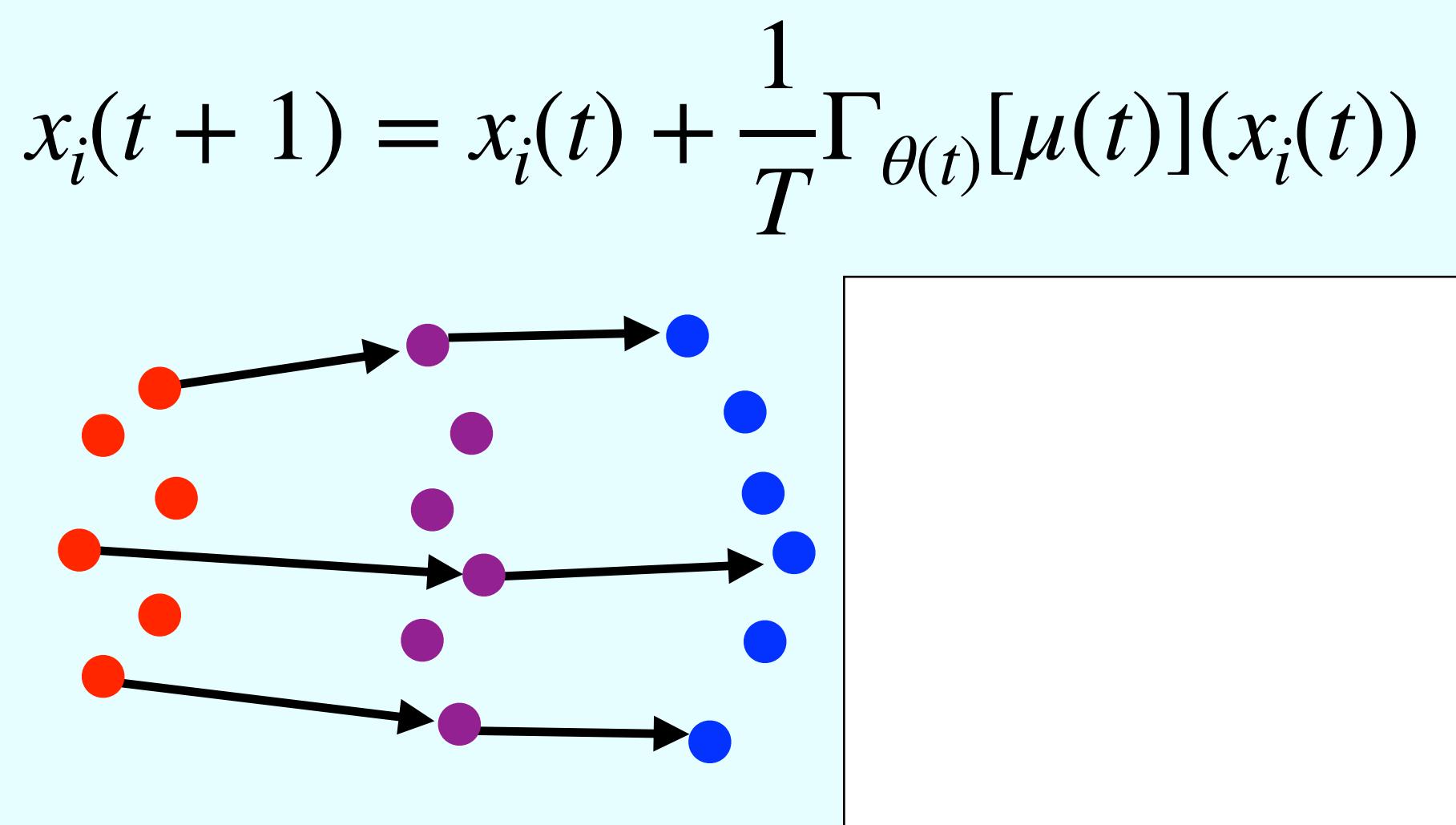
$$C_\theta \leq \|V\| \|K^\top Q\| R^2 \sqrt{12n + 3}$$

Extension to masked attention: use  $W_2^{\text{cond}}(\mu, \nu)^2 := \int_0^1 W_2^2(\mu(\cdot | t), \nu(\cdot | t)) d\mu_{[0,1]}(t)$



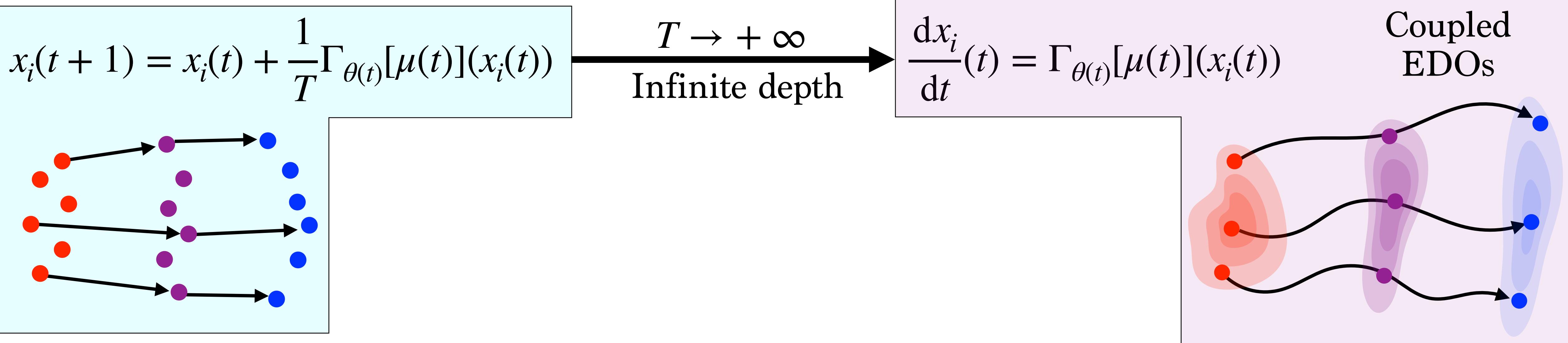
# Infinite Depth as a Neural PDE

$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky' \rangle} d\mu(y')} V y \, d\mu(y) \quad \theta = (Q, K, V) \quad \mu(t) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$$



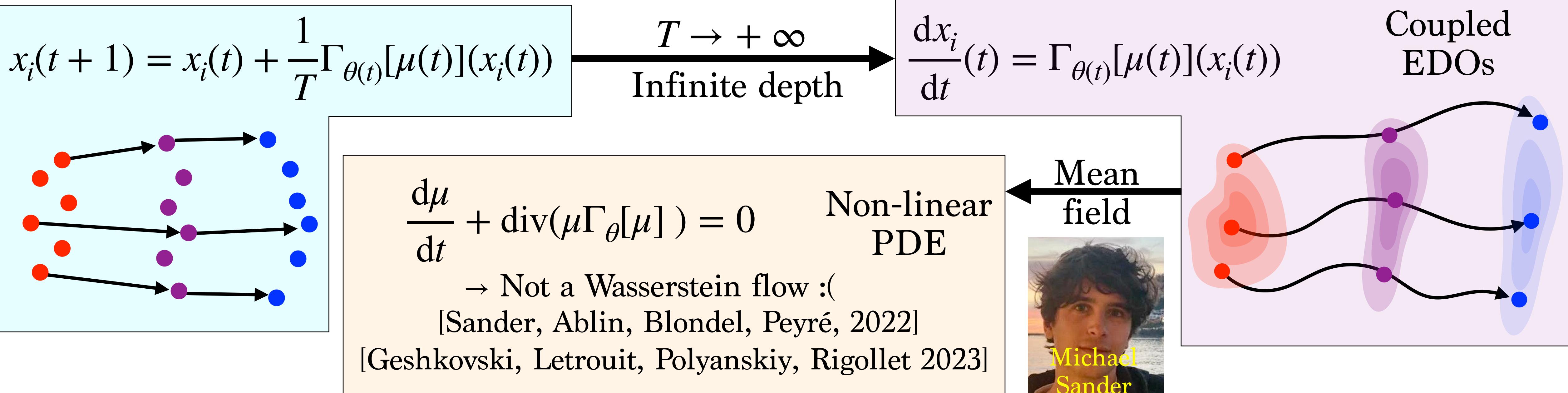
# Infinite Depth as a Neural PDE

$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky' \rangle} d\mu(y')} V y \, d\mu(y) \quad \theta = (Q, K, V) \quad \mu(t) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$$



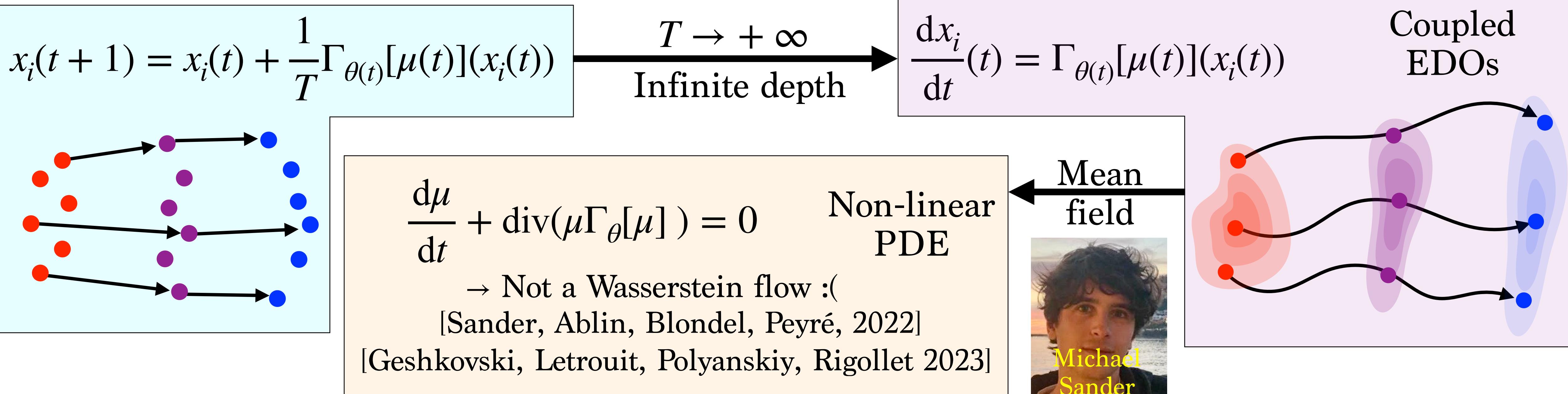
# Infinite Depth as a Neural PDE

$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky' \rangle} d\mu(y')} Vy \, d\mu(y) \quad \theta = (Q, K, V) \quad \mu(t) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$$



# Infinite Depth as a Neural PDE

$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky' \rangle} d\mu(y')} Vy \, d\mu(y) \quad \theta = (Q, K, V) \quad \mu(t) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$$



Transformer:  $T_\theta[\mu_0] : x(t=0) \xrightarrow[\mu(t=0) = \mu_0]{\dot{x} = \Gamma_\theta[\mu](x)} x(t=1)$

Training:

$$\min_\theta \sum_k \ell(T_\theta[\mu^k](x^k), y^k)$$

Context Previous Next

« Theorem » convergence to the global minimum if

- initial loss small enough
- enough heads
- $(\mu^k)_k$  separated

→ Talks by  
Pierre Marion and  
Raphaël Barboni

# Gaussian Case and Clustering

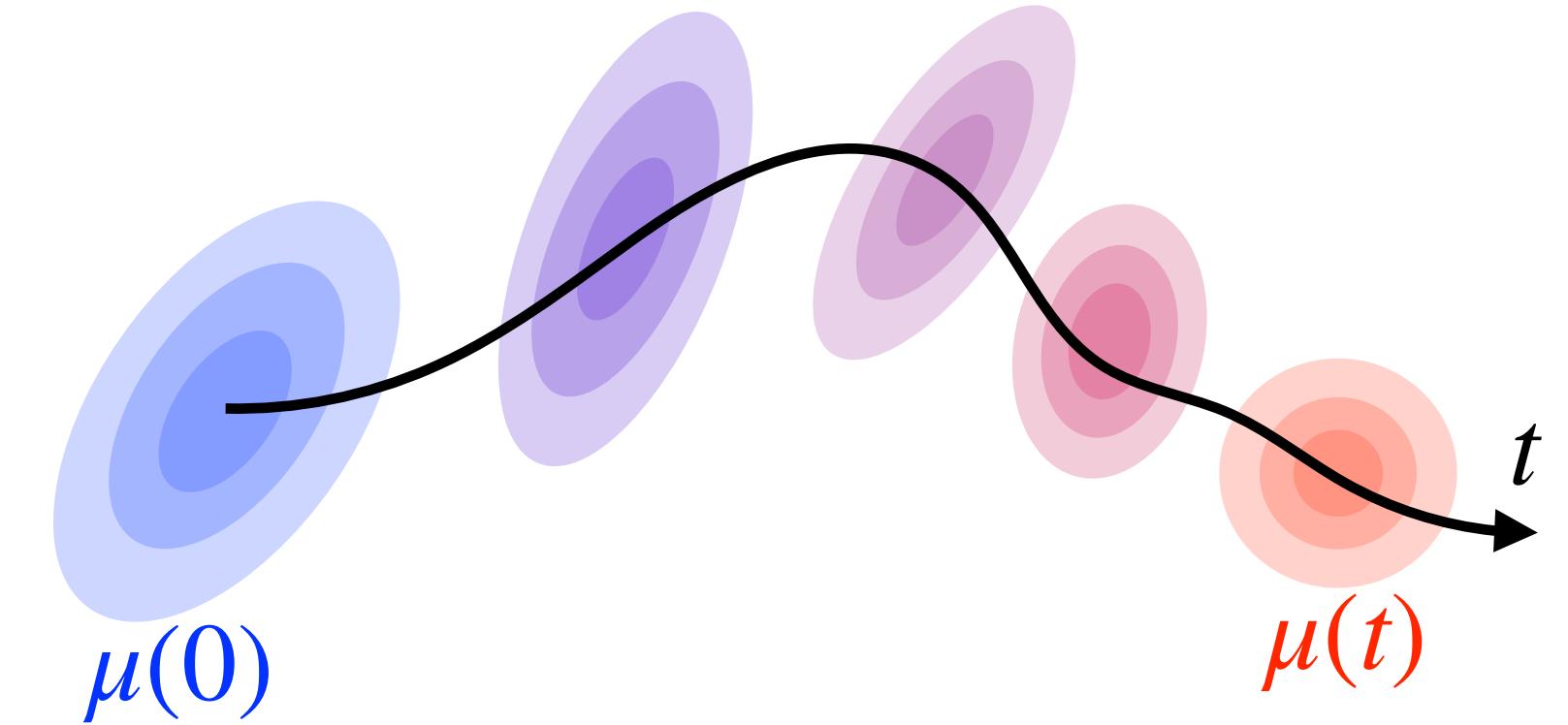
$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky' \rangle} d\mu(y')} Vy \, d\mu(y) \quad \theta(t) = (Q(t), K(t), V(t))$$

$$\frac{d\mu}{dt} + \operatorname{div}(\mu \Gamma_\theta[\mu]) = 0$$

*Theorem* [Valérie Castin]: If  $\mu(0) = \mathcal{N}(\mathbf{m}(0), \Sigma(0))$ ,

then  $\mu(s) = \mathcal{N}(\mathbf{m}(s), \Sigma(s))$

$$\begin{aligned} \dot{\mathbf{m}} &= V(\operatorname{Id} + \Sigma Q^\top K) \mathbf{m} \\ \dot{\Sigma} &= V \Sigma Q^\top K \Sigma + \Sigma K^\top Q \Sigma V^\top \end{aligned}$$



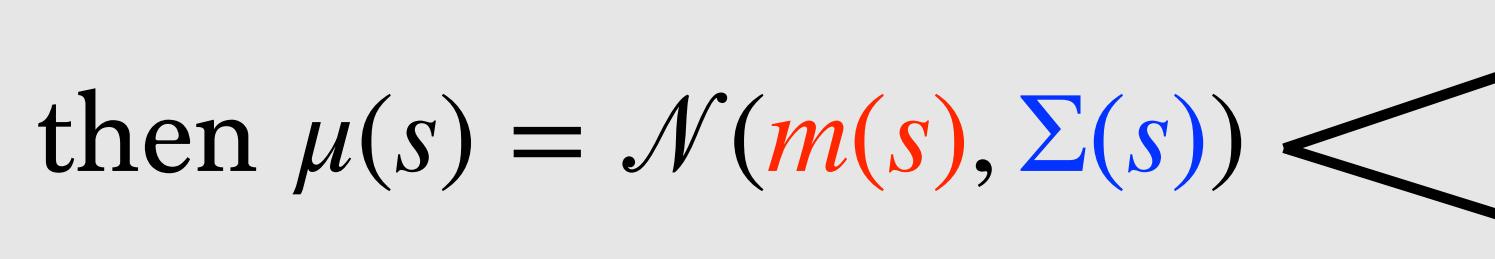
# Gaussian Case and Clustering

$$\Gamma_\theta[\mu](x) := \int \frac{e^{\langle Qx, Ky \rangle}}{\int e^{\langle Qx, Ky' \rangle} d\mu(y')} Vy d\mu(y) \quad \theta(t) = (Q(t), K(t), V(t))$$

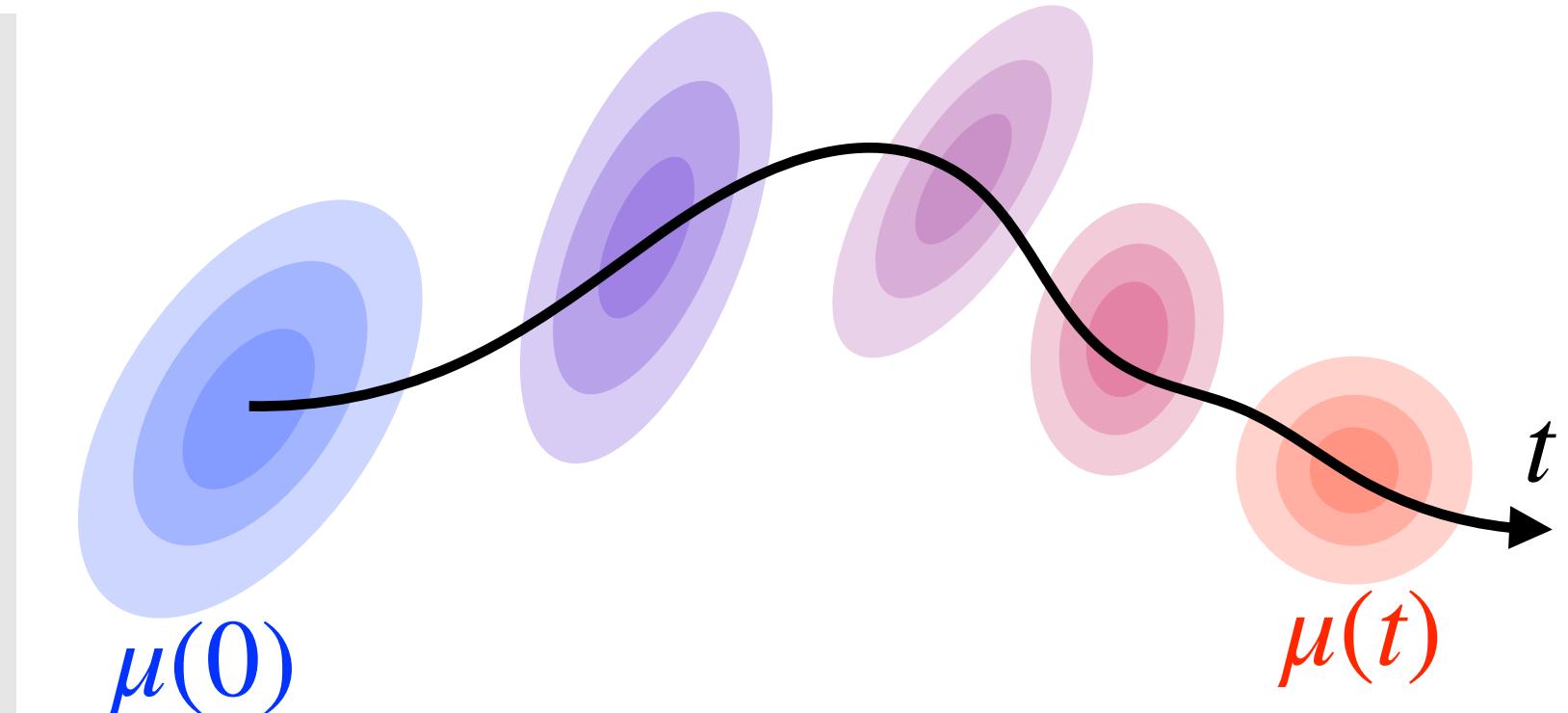
$$\frac{d\mu}{dt} + \operatorname{div}(\mu \Gamma_\theta[\mu]) = 0$$

*Theorem* [Valérie Castin]: If  $\mu(0) = \mathcal{N}(\mathbf{m}(0), \Sigma(0))$ ,

then  $\mu(s) = \mathcal{N}(\mathbf{m}(s), \Sigma(s))$

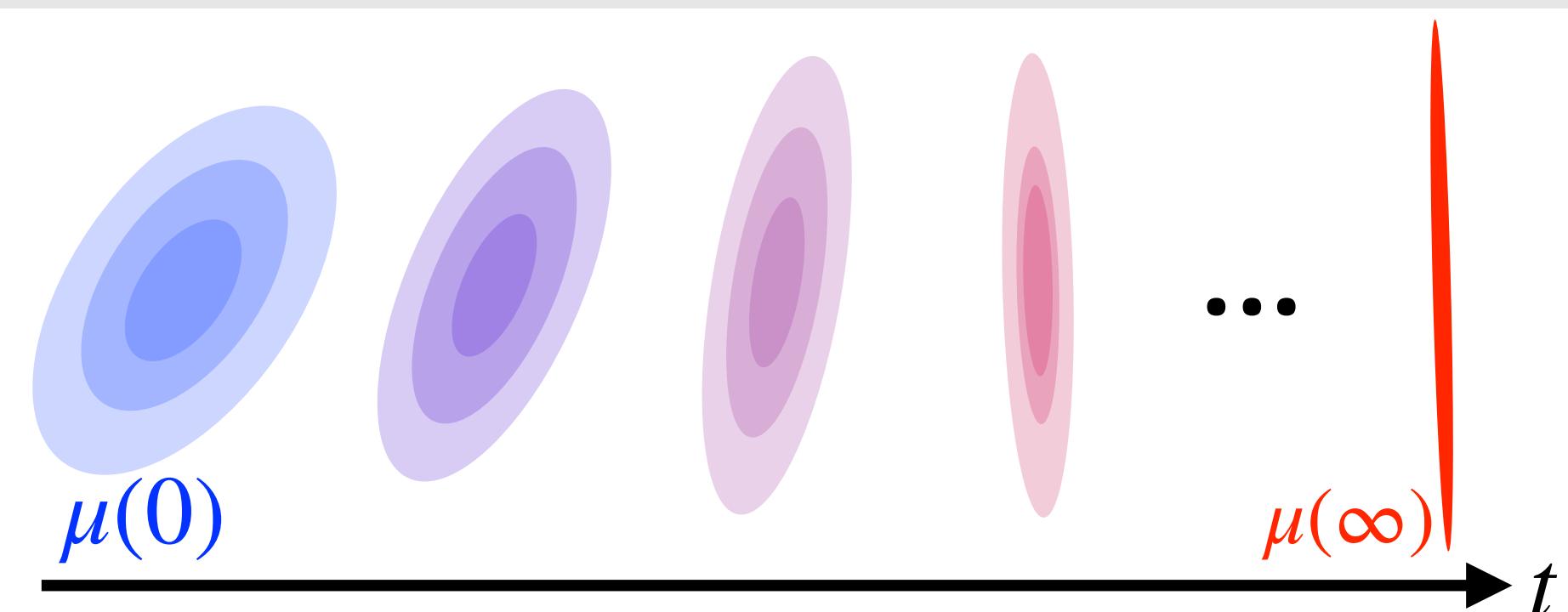


$$\begin{aligned}\dot{\mathbf{m}} &= V(\operatorname{Id} + \Sigma Q^\top K)\mathbf{m} \\ \dot{\Sigma} &= V\Sigma Q^\top K\Sigma + \Sigma K^\top Q\Sigma V^\top\end{aligned}$$



*Theorem* [Valérie Castin]:

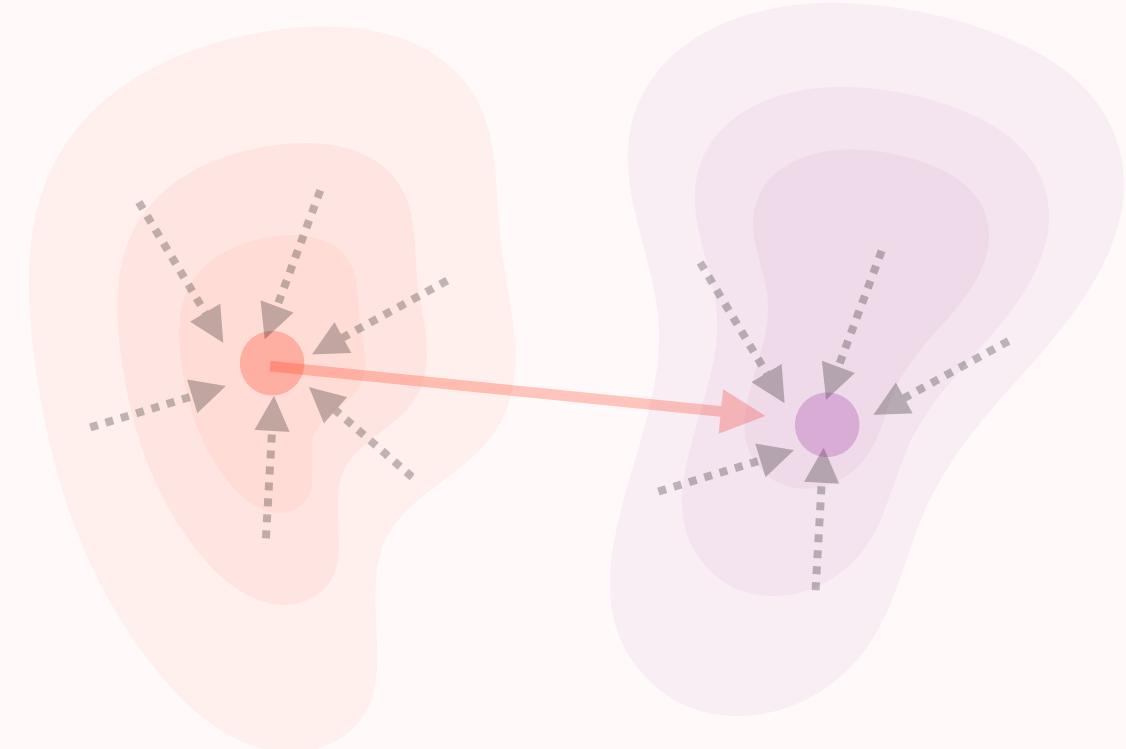
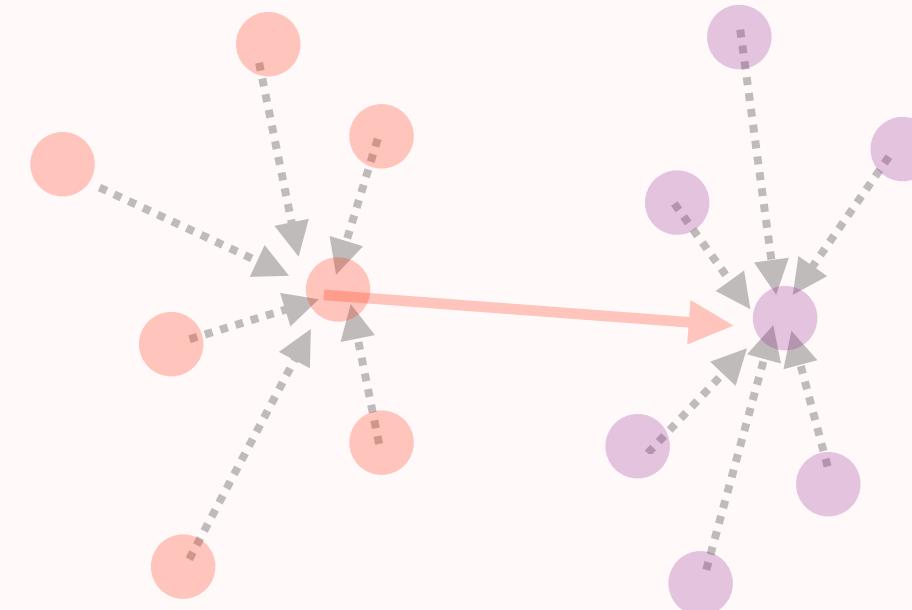
If  $V(t) = \operatorname{Id}$  and  $K(t)^\top Q(t)$  symmetric, stationary points of  $\Sigma(t)$  have rank less than  $d/2$ .



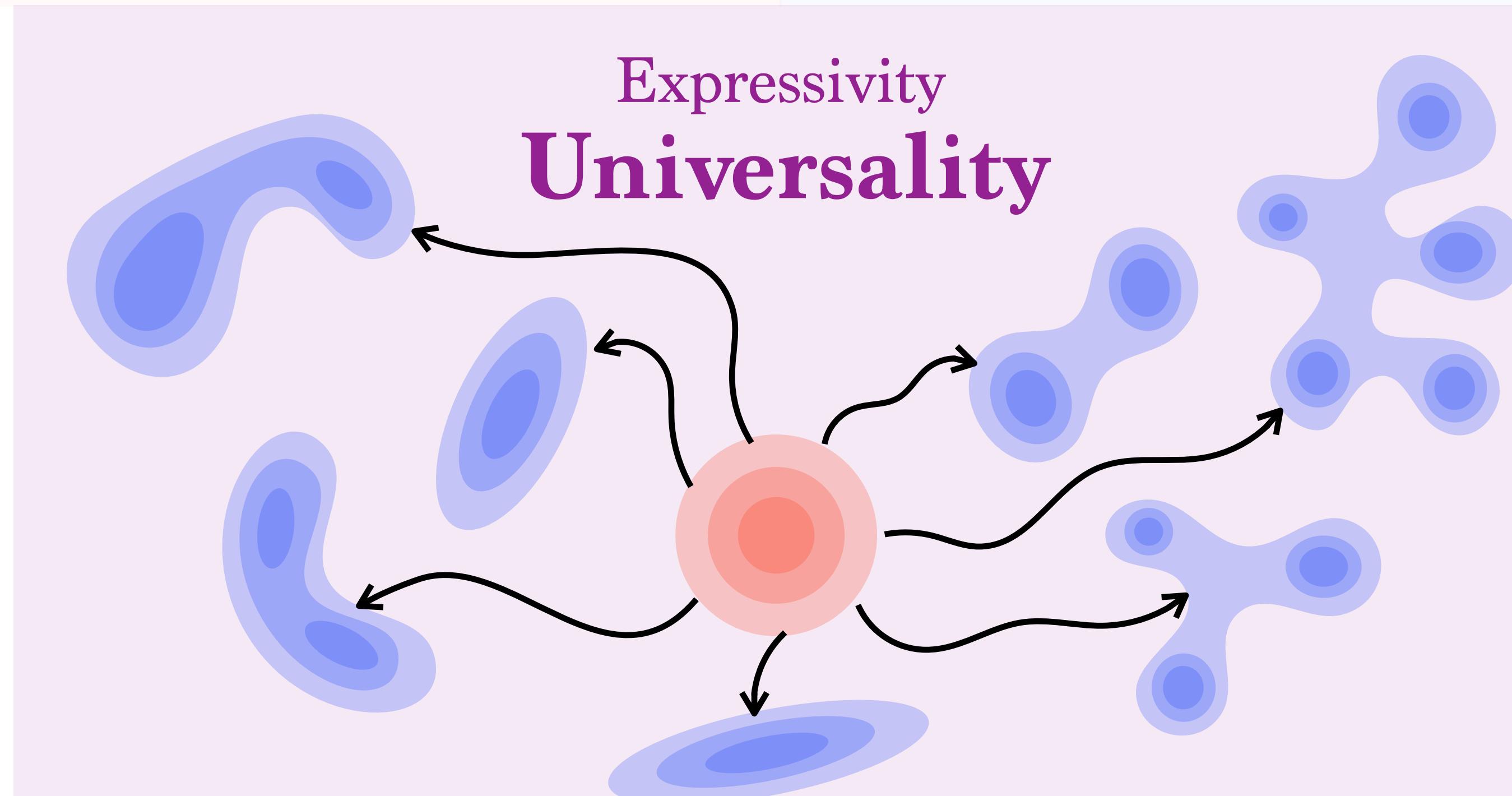
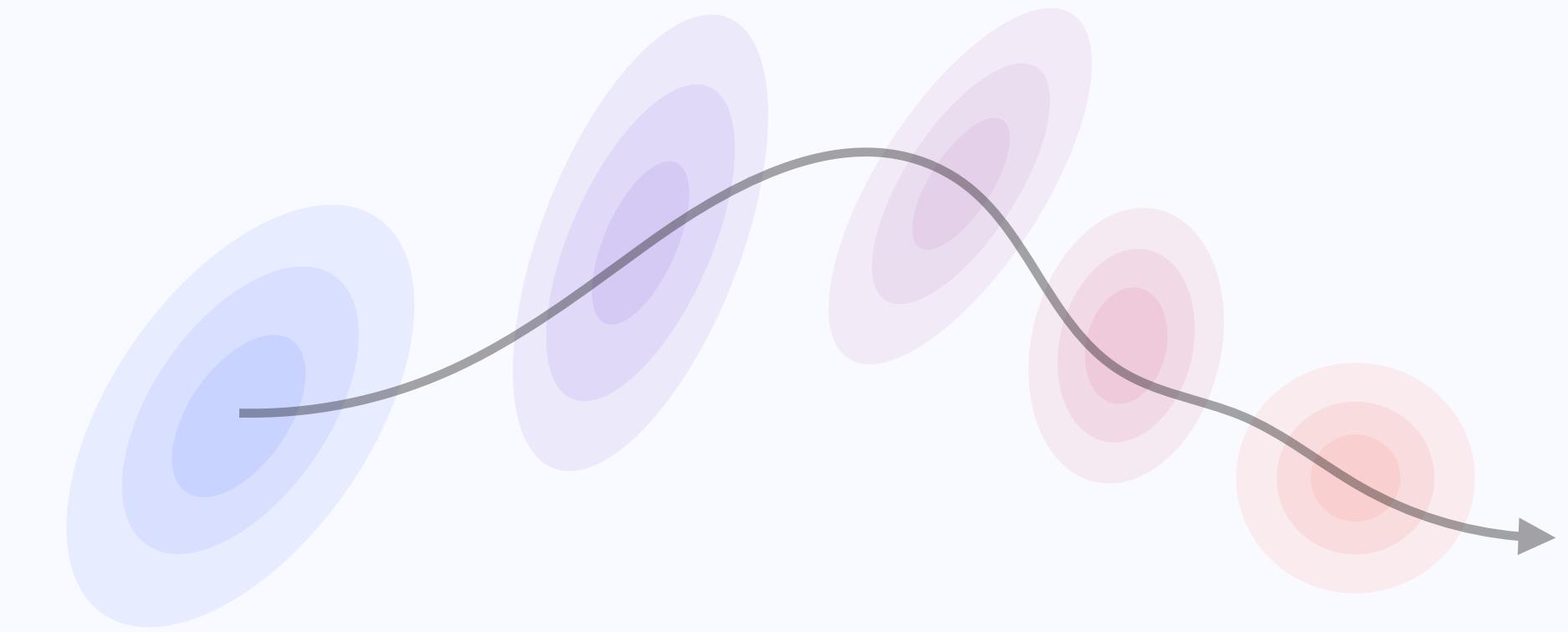
*Conjecture:* low-rank stationary covariances for any  $K, Q, V$ .

[Geshkovski, Letrouit, Polyanskiy, Rigollet 2023]  
 → The attention matrix converges to low-rank.  
 → Clustering of  $\mu$  for un-normalized attention.

Arbitrary number of layers  
**In Context Mappings  
over Measures**



Arbitrary number of layers  
**Smoothness and  
PDE's**



# Universality

$$\Gamma_\theta[\mu](x) := x + \sum_{h=1}^H \int \frac{e^{\langle Q^h x, K^h y \rangle}}{\int e^{\langle Q^h x, K^h y' \rangle} d\mu(y')} V^h y \, d\mu(y) \quad \text{or} \quad \Gamma_\theta[\mu](x) := \text{MLP}_\theta(x)$$

*Theorem* [Furuya, de Hoop, Peyré]:

Let  $\Gamma^\star : \mathcal{P}(\Omega) \times \Omega \rightarrow \mathbb{R}^d$  be  $\text{Wass}_2 \times \ell^2$ -continuous on a compact  $\Omega \subset \mathbb{R}^d$ .

For any  $\varepsilon$  there exists  $N$  and  $(\theta_1, \dots, \theta_N)$  such that

$$\forall (\mu, x) \in \mathcal{P}(\Omega) \times \Omega, |\Gamma^\star[\mu](x) - \Gamma_{\theta_N} \diamond \dots \diamond \Gamma_{\theta_1}[\mu](x)| \leq \varepsilon$$

with token dimensions  $\leq 4d$  and  $H \leq d$ .

*Novelties:*

fixed dimensions,  
arbitrary # tokens.

*Masked transformers:*  
requires Lipschitz  
in time.

# Universality

$$\Gamma_\theta[\mu](x) := x + \sum_{h=1}^H \int \frac{e^{\langle Q^h x, K^h y \rangle}}{\int e^{\langle Q^h x, K^h y' \rangle} d\mu(y')} V^h y \, d\mu(y) \quad \text{or} \quad \Gamma_\theta[\mu](x) := \text{MLP}_\theta(x)$$

*Theorem* [Furuya, de Hoop, Peyré]:

Let  $\Gamma^\star : \mathcal{P}(\Omega) \times \Omega \rightarrow \mathbb{R}^d$  be  $\text{Wass}_2 \times \ell^2$ -continuous on a compact  $\Omega \subset \mathbb{R}^d$ .

For any  $\varepsilon$  there exists  $N$  and  $(\theta_1, \dots, \theta_N)$  such that

$$\forall (\mu, x) \in \mathcal{P}(\Omega) \times \Omega, |\Gamma^\star[\mu](x) - \Gamma_{\theta_N} \diamond \dots \diamond \Gamma_{\theta_1}[\mu](x)| \leq \varepsilon$$

with token dimensions  $\leq 4d$  and  $H \leq d$ .

*Novelties:*

fixed dimensions,  
arbitrary # tokens.

*Masked transformers:*  
requires Lipschitz  
in time.

*Previous works:*

[Yun, Bhojanapalli, Singh Rawat, Reddi, Kumar, 2019]  $\rightarrow H = 2$ , dimension  $\sim \# \text{tokens}$

[Agrachev, Letrouit 2019]  $\rightarrow$  abstract genericity hypothesis (Lie algebra/control)

Discrete tokens: transformers are universal Turing machines: e.g. [Elhage et al 2021]

# Sketch of proof

1-D elementary block:  $\gamma_\theta[\mu](x) := \langle x, u \rangle + \int \frac{e^{\langle Ax+b, y \rangle}}{\int e^{\langle Ax+b, y \rangle} d\mu(y)} (\langle v, y \rangle + c) d\mu(y)$   $\theta := (A, b, c, u, v)$

→ First component of Attention  $\circ$  MLP with skip connexion.

Cylindrical algebra:  $\mathcal{A} := \text{Span} \bigcup_N \{\gamma_{\theta_1} \odot \dots \odot \gamma_{\theta_N} : (\theta_1, \dots, \theta_N)\}$   $(\gamma_1 \odot \gamma_2)[\mu](x) := \gamma_1[\mu](x)\gamma_2[\mu](x)$

# Sketch of proof

1-D elementary block:

$$\gamma_\theta[\mu](x) := \langle x, u \rangle + \int \frac{e^{\langle Ax+b, y \rangle}}{\int e^{\langle Ax+b, y \rangle} d\mu(y)} (\langle v, y \rangle + c) d\mu(y)$$

$$\theta := (A, b, c, u, v)$$

→ First component of Attention  $\circ$  MLP with skip connexion.

Cylindrical algebra:

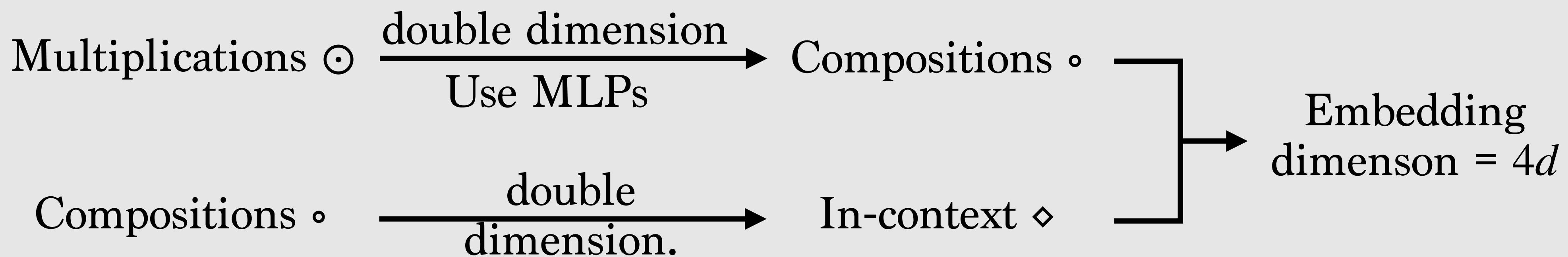
$$\mathcal{A} := \text{Span} \bigcup_N \{\gamma_{\theta_1} \odot \dots \odot \gamma_{\theta_N} : (\theta_1, \dots, \theta_N)\}$$

$$(\gamma_1 \odot \gamma_2)[\mu](x) := \gamma_1[\mu](x) \gamma_2[\mu](x)$$

*Proposition:* any map  $(\mu, x) \rightarrow (\alpha_1[\mu](x), \dots, \alpha_d[\mu](x)) \in \mathbb{R}^d$  with  $\alpha_i \in \mathcal{A}$  can be uniformly approximated by a transformer with skip connexions.

Use 1D dimension by dimension → requires  $H = d$  heads.

Proof sketch:



# Sketch of Proof

$$\gamma_\theta[\mu](x) := \langle x, u \rangle + \int \frac{e^{\langle Ax+b, y \rangle}}{\int e^{\langle Ax+b, y' \rangle} d\mu(y')} (\langle v, y \rangle + c) d\mu(y) \quad \mathcal{A} := \text{Span} \bigcup_N \{\gamma_{\theta_1} \odot \dots \odot \gamma_{\theta_N}\}$$

*Lemma:*  $\mathcal{A}$  is dense in continuous maps  $\mathcal{P}(\Omega) \times \Omega \rightarrow \mathbb{R}$  for  $\text{Wass}_2 \times \ell^2$

# Sketch of Proof

$$\gamma_\theta[\mu](x) := \langle x, u \rangle + \int \frac{e^{\langle Ax+b, y \rangle}}{\int e^{\langle Ax+b, y' \rangle} d\mu(y')} (\langle v, y \rangle + c) d\mu(y)$$
$$\mathcal{A} := \text{Span} \bigcup_N \{\gamma_{\theta_1} \odot \dots \odot \gamma_{\theta_N}\}$$

*Lemma:*  $\mathcal{A}$  is dense in continuous maps  $\mathcal{P}(\Omega) \times \Omega \rightarrow \mathbb{R}$  for  $\text{Wass}_2 \times \ell^2$

*Proof:*

$\mathcal{P}(\Omega) \times \Omega$  is compact.

Stone-Weierstrass  
theorem

$\gamma_\theta$  are continuous.

$A = b = u = v = 0, c = 1:$   
 $\gamma_\theta[\mu] = 1$

$\forall \theta, \gamma_\theta[\mu](x) = \gamma_\theta[\mu'](x')$

?

$(\mu, x) = (\mu', x')$



# Sketch of Proof

$$\gamma_\theta[\mu](x) := \langle x, u \rangle + \int \frac{e^{\langle Ax+b, y \rangle}}{\int e^{\langle Ax+b, y' \rangle} d\mu(y')} (\langle v, y \rangle + c) d\mu(y)$$
$$\mathcal{A} := \text{Span} \bigcup_N \{\gamma_{\theta_1} \odot \dots \odot \gamma_{\theta_N}\}$$

*Lemma:*  $\mathcal{A}$  is dense in continuous maps  $\mathcal{P}(\Omega) \times \Omega \rightarrow \mathbb{R}$  for  $\text{Wass}_2 \times \ell^2$

*Proof:*

$\mathcal{P}(\Omega) \times \Omega$  is compact.

Stone-Weierstrass theorem

$$A = b = u = v = 0, c = 1:$$
$$\gamma_\theta[\mu] = 1$$

$$\forall \theta, \gamma_\theta[\mu](x) = \gamma_\theta[\mu'](x')$$

?

$$(\mu, x) = (\mu', x')$$

$$\rightarrow c = v = 0: \quad \langle x, u \rangle = \langle x', u \rangle$$



Marshall  
Stone

Karl  
Weierstrass

# Sketch of Proof

$$\gamma_\theta[\mu](x) := \langle x, u \rangle + \int \frac{e^{\langle Ax+b, y \rangle}}{\int e^{\langle Ax+b, y' \rangle} d\mu(y')} (\langle v, y \rangle + c) d\mu(y)$$

$$\mathcal{A} := \text{Span} \bigcup_N \{\gamma_{\theta_1} \odot \dots \odot \gamma_{\theta_N}\}$$

*Lemma:*  $\mathcal{A}$  is dense in continuous maps  $\mathcal{P}(\Omega) \times \Omega \rightarrow \mathbb{R}$  for  $\text{Wass}_2 \times \ell^2$

*Proof:*

$\mathcal{P}(\Omega) \times \Omega$  is compact.

$\gamma_\theta$  are continuous.

$A = b = u = v = 0, c = 1:$   
 $\gamma_\theta[\mu] = 1$

$\forall \theta, \gamma_\theta[\mu](x) = \gamma_\theta[\mu'](x')$

?

$(\mu, x) = (\mu', x')$

$$\begin{array}{l} \rightarrow c = v = 0: \quad \langle x, u \rangle = \langle x', u \rangle \\ \rightarrow A = c = u = 0: \quad L_1(\mu)(b) = L_1(\mu')(b) \end{array}$$

In 1-D:

$$L_k(\mu)(b) := \int \frac{e^{by} y^k v}{\int e^{by'} d\mu(y')} d\mu(y)$$

$$L'_k = L_{k+1} - L_k L_1$$

$$L_1(\mu) = L_1(\mu') \Rightarrow \forall k, L_k(\mu) = L_k(\mu') \Rightarrow \forall k, \int y^k d\mu(y) = \int y^k d\mu'(y)$$

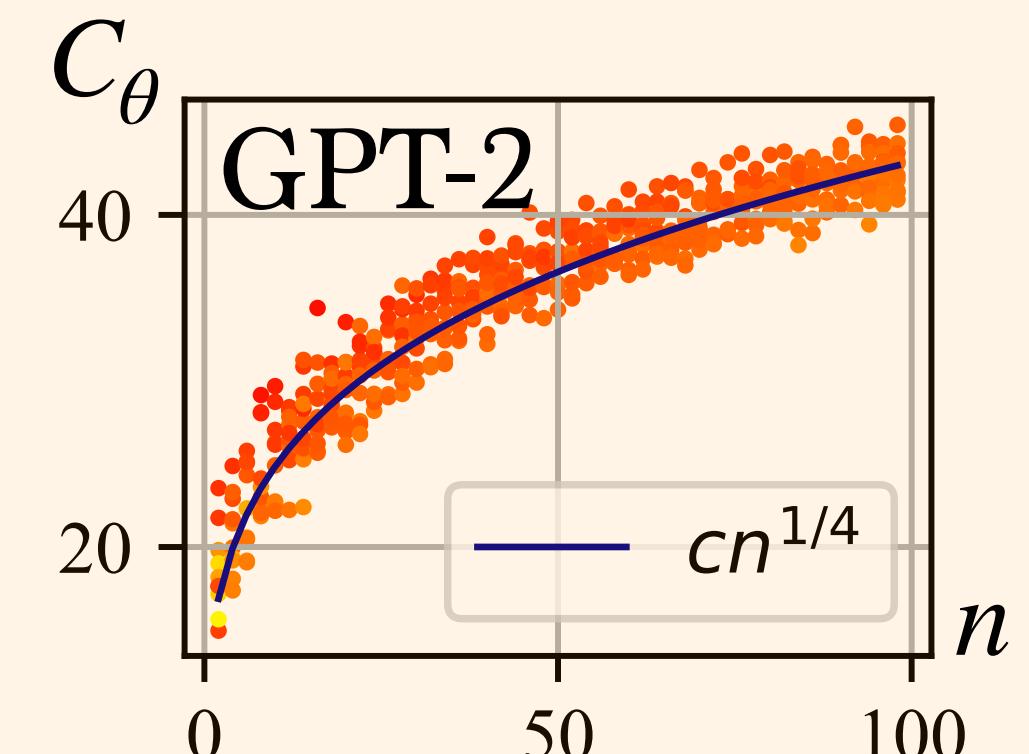
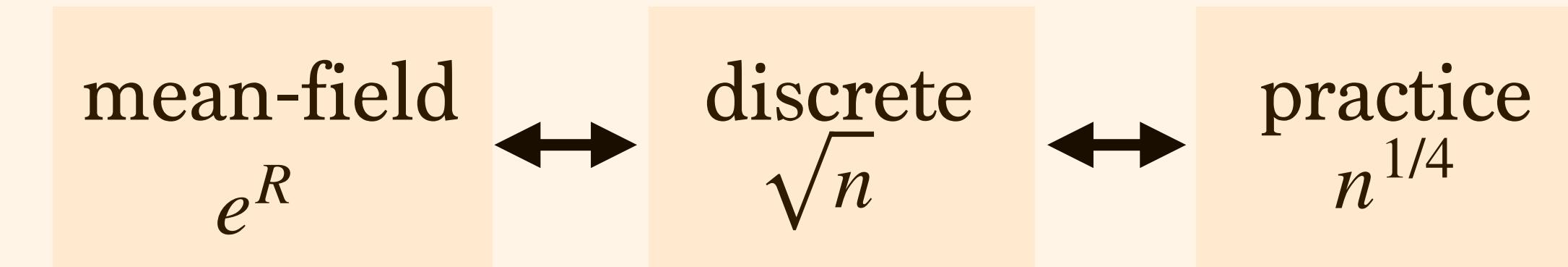
*In higher dimensions:* use Radon transform.



Stone-Weierstrass  
theorem

# Open Problems

*Smoothness:* bridge the gap



*Universality:*

- Replace scalar-valued cylindrical maps by more effective functions.
- Toward quantitative approximation bound, leverage smoothness.

*Optimisation:*

- Understand the structure of optimal  $(Q, K, V)$
- Why is Adam normalization needed for training?