

Transformers are Universal in Context Learners

Gabriel Peyré



**Takashi
Furuya**



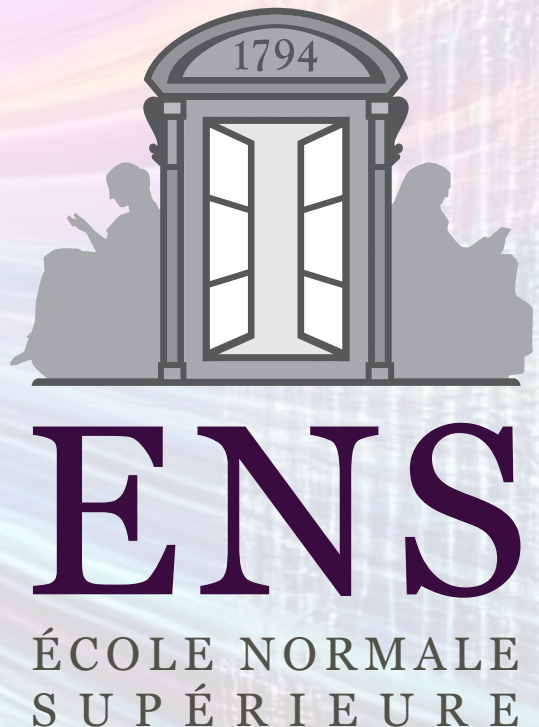
**Maarten de
Hoop**



**Valérie
Castin**



**Pierre
Ablin**



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.

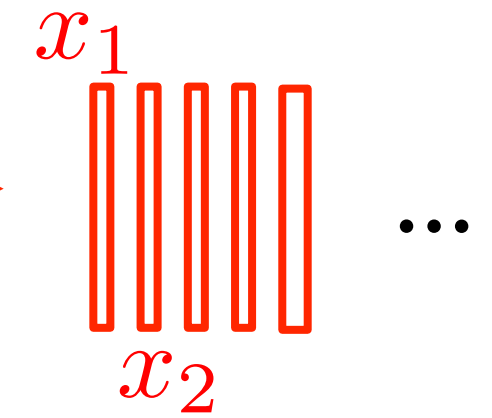
Tokenize →

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.

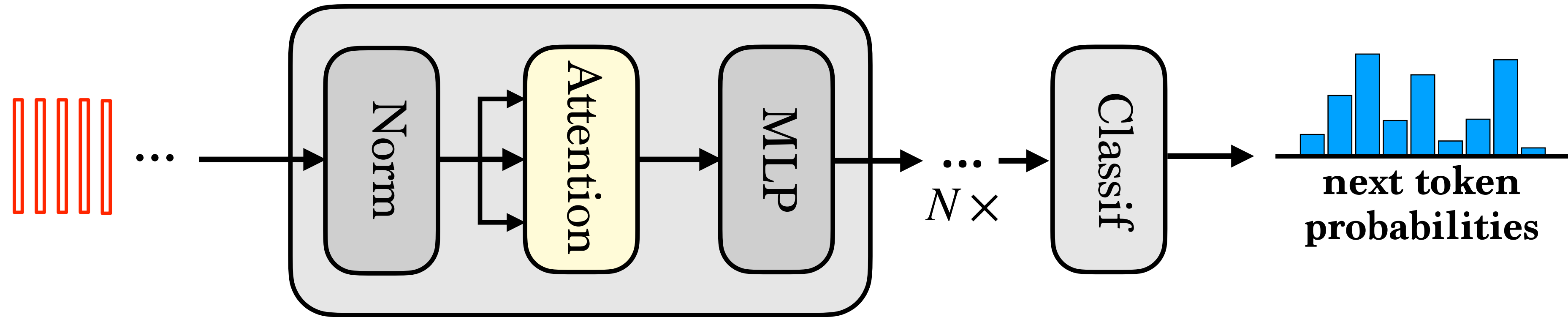
Token encoding

Positional encoding

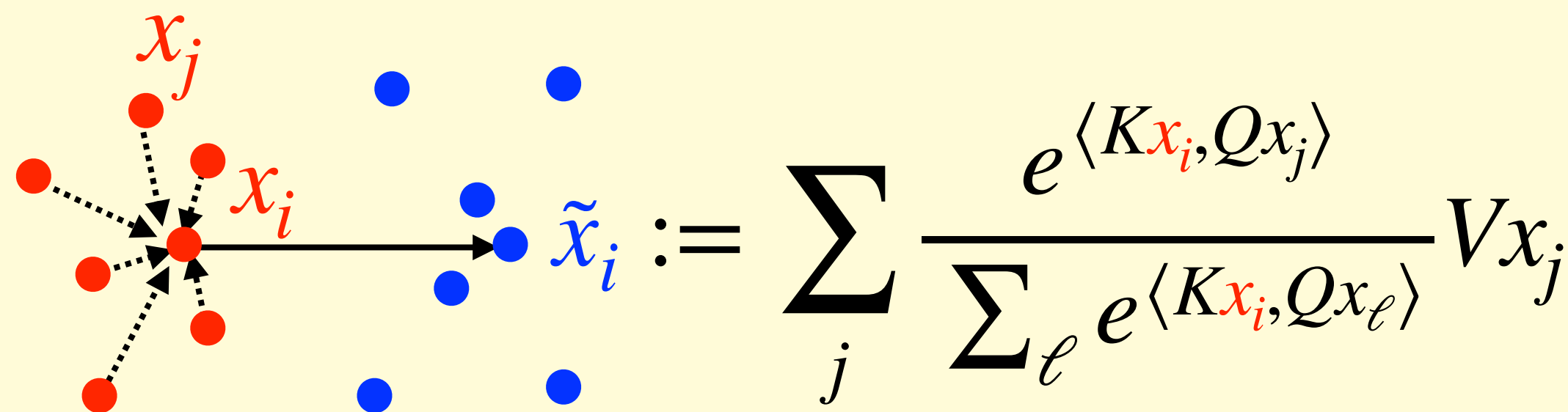
+



Points cloud
 $\{x_i\}_i$



(Unmasked) Attention layer



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.

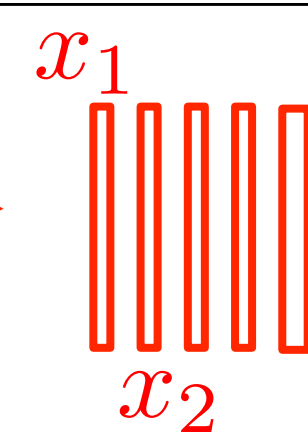
Tokenize

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.

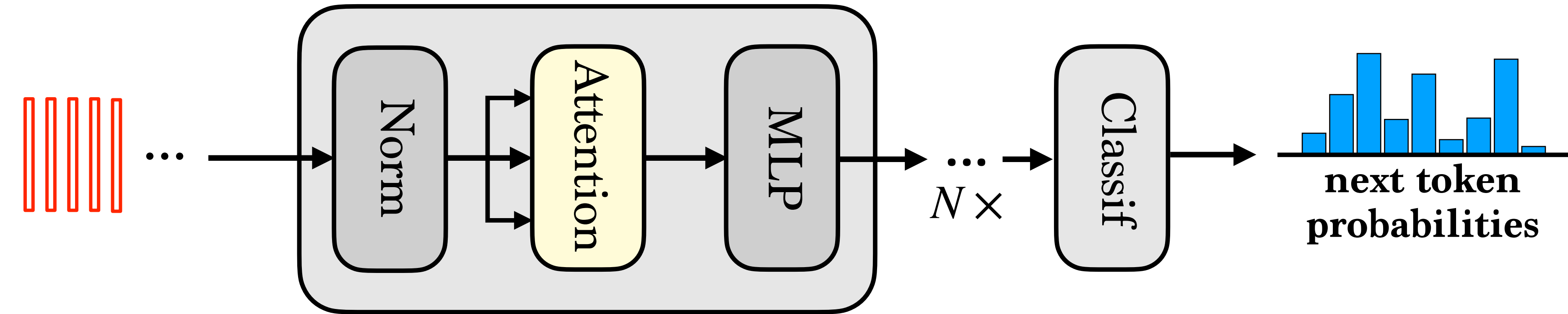
Token encoding

Positional encoding

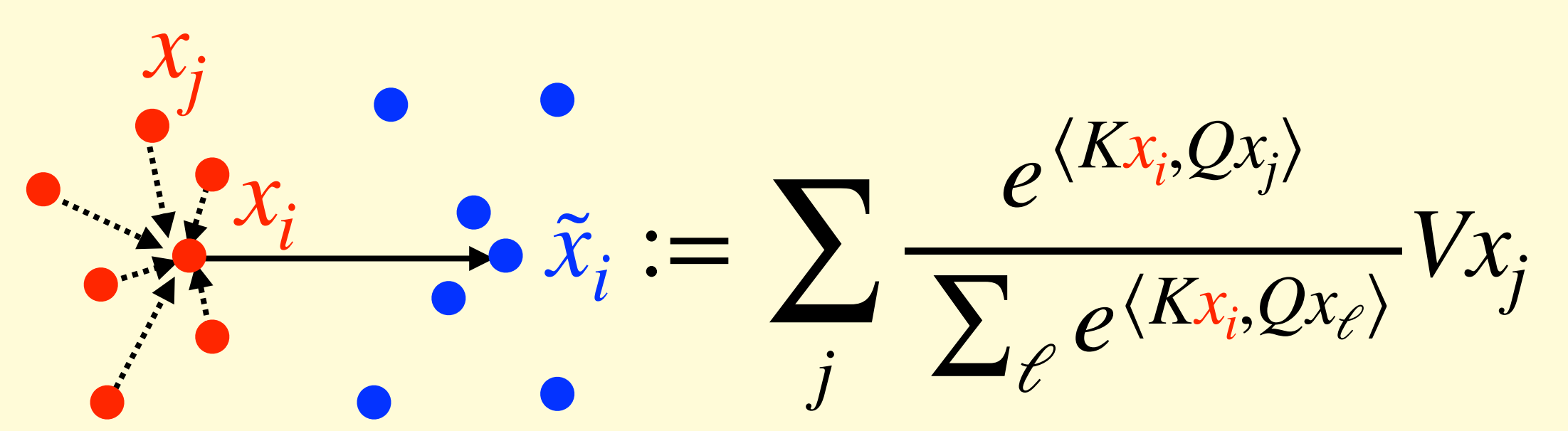
+



Points cloud
 $\{x_i\}_i$



(Unmasked) Attention layer

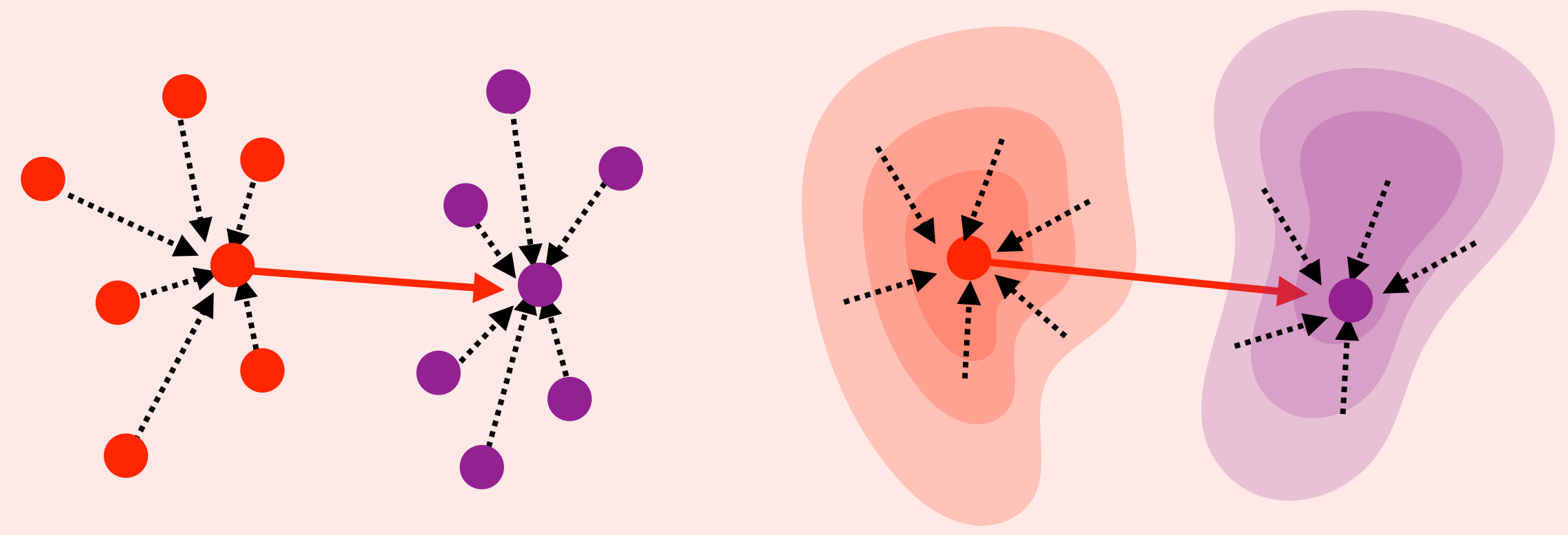


$$\tilde{x}_i := \sum_j \frac{e^{\langle Kx_i, Qx_j \rangle}}{\sum_{\ell} e^{\langle Kx_i, Qx_{\ell} \rangle}} Vx_j$$

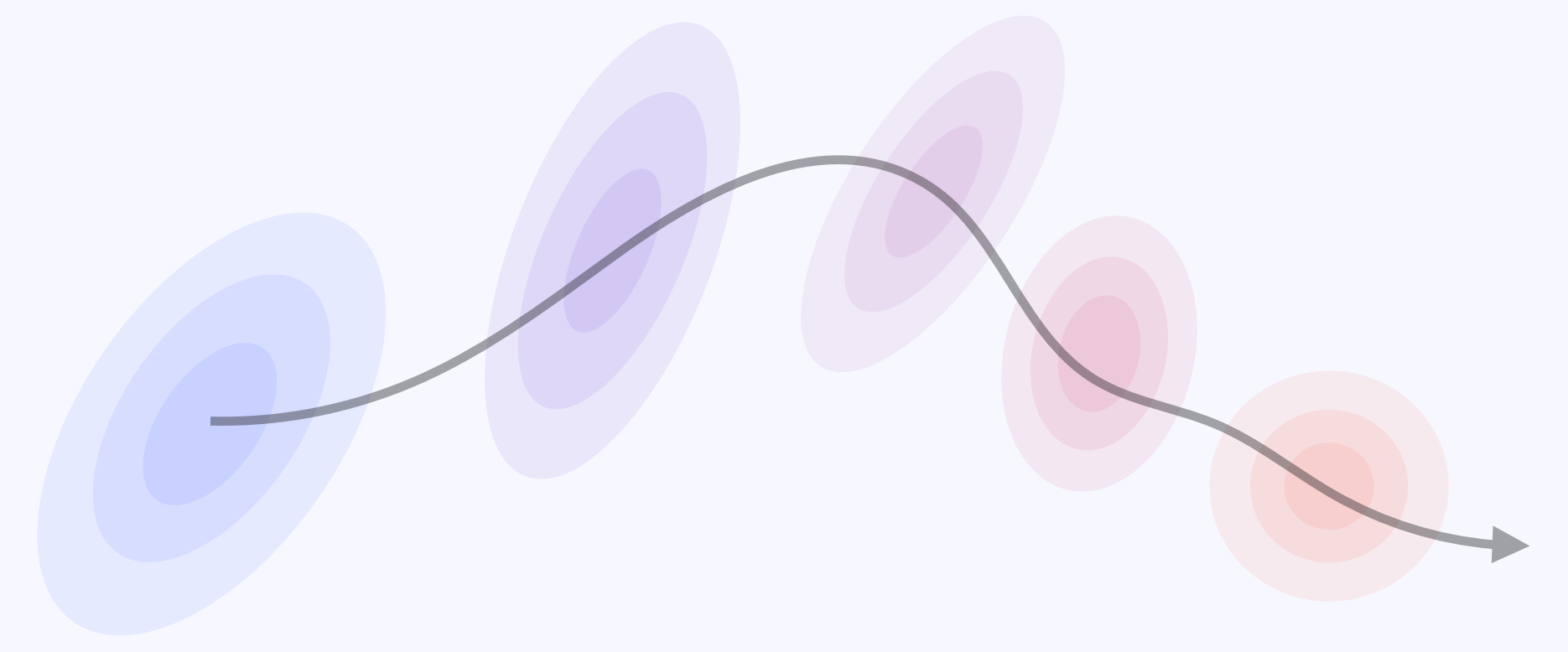
Understanding

- Arbitrary number of tokens
- Arbitrary number of layers
- Expressivity

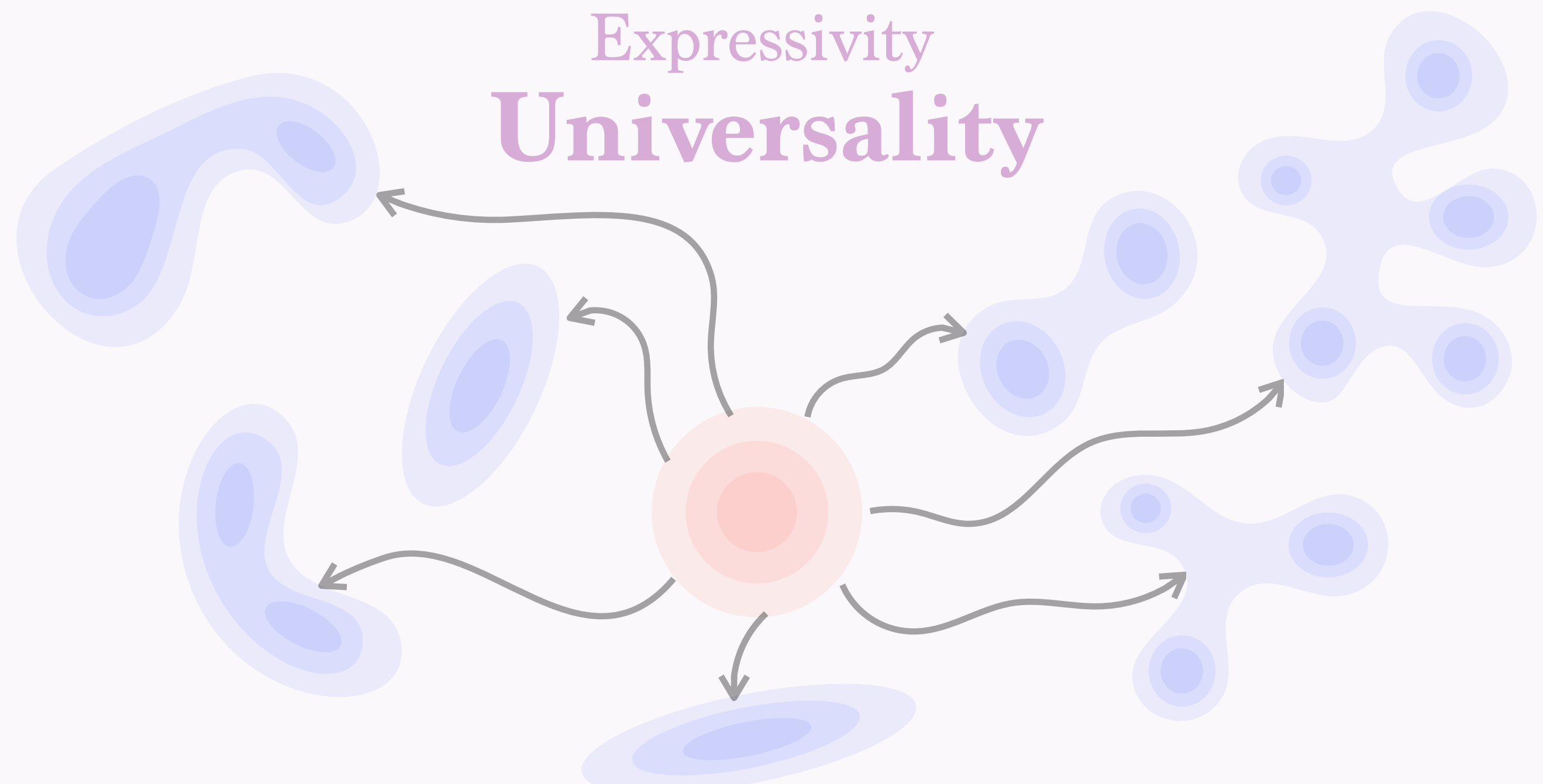
Arbitrary number of tokens
In Context Mappings
over Measures



Arbitrary number of layers
Smoothness and
PDE's



Expressivity
Universality

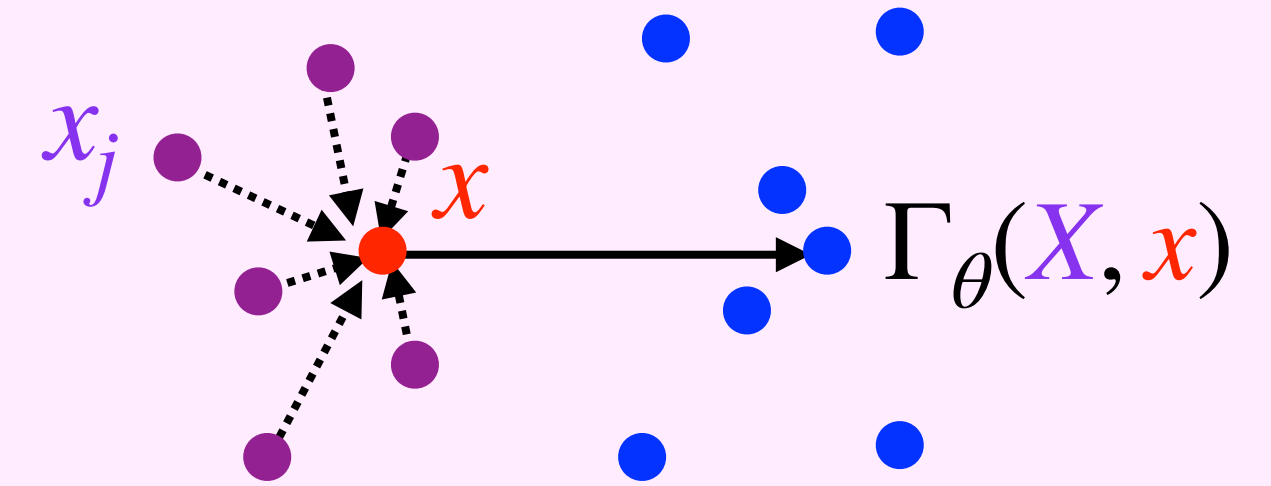


Attention as In-context Mapping

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

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

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

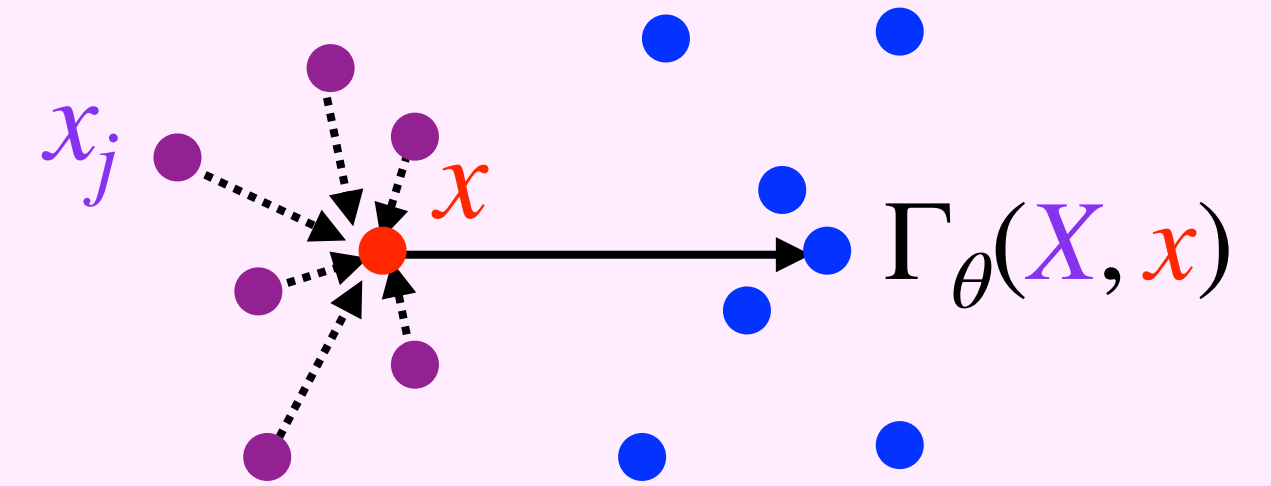


Attention as In-context Mapping

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

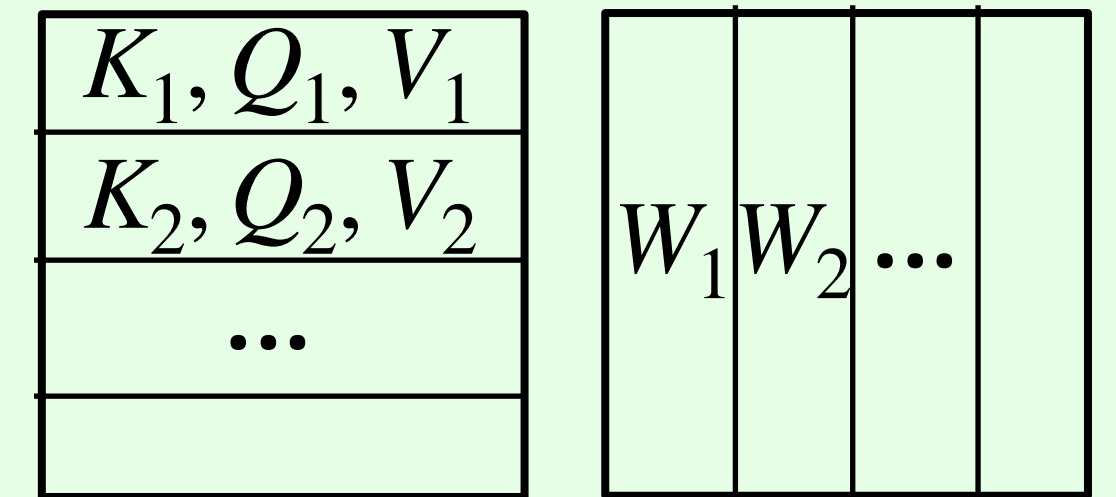
In-context mapping:
parameters $\theta := (Q, K, V)$

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



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

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

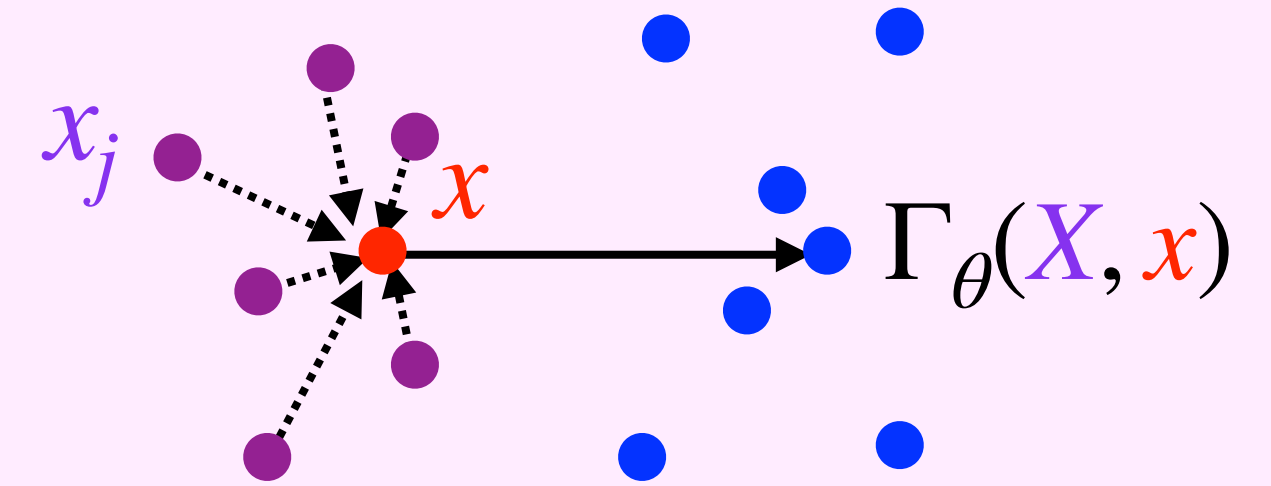


Attention as In-context Mapping

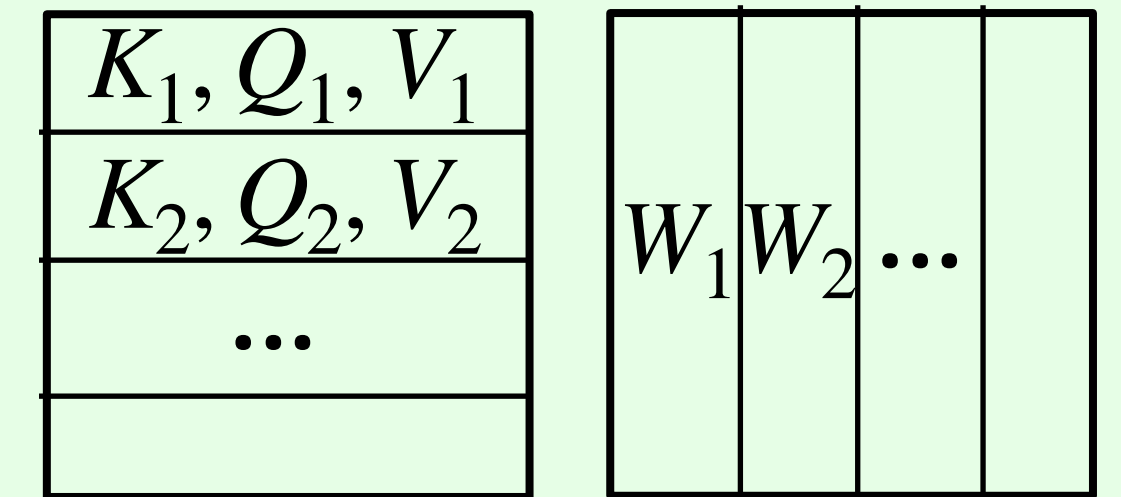
Point clouds: $X := \{x_i\}_{i=1}^n$

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

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



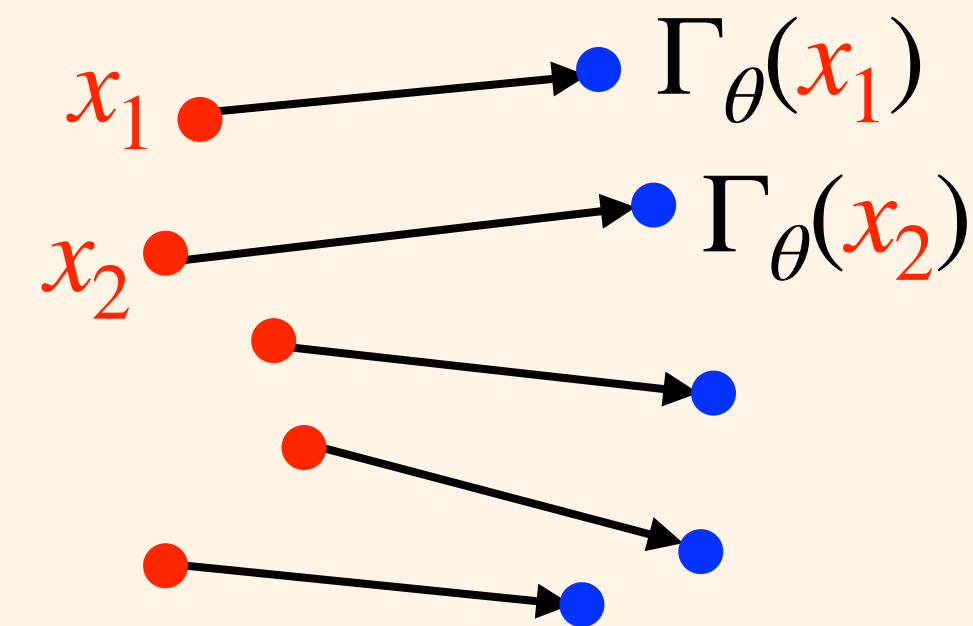
Single-head attention layer: $X \mapsto \{\Gamma_{\theta}[X](x_i)\}_{i=1}^n$



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

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

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

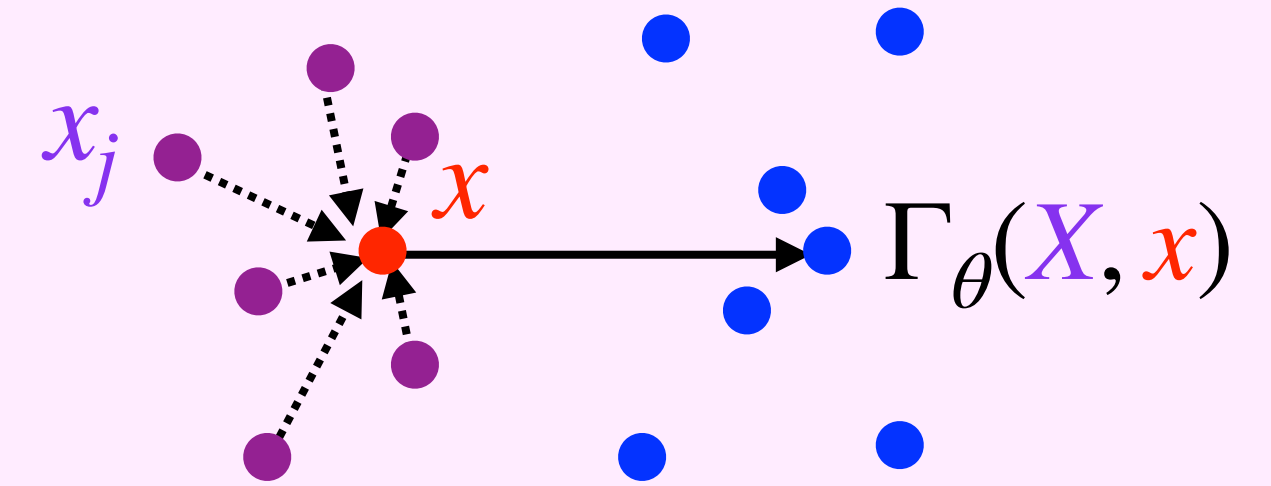


Attention as In-context Mapping

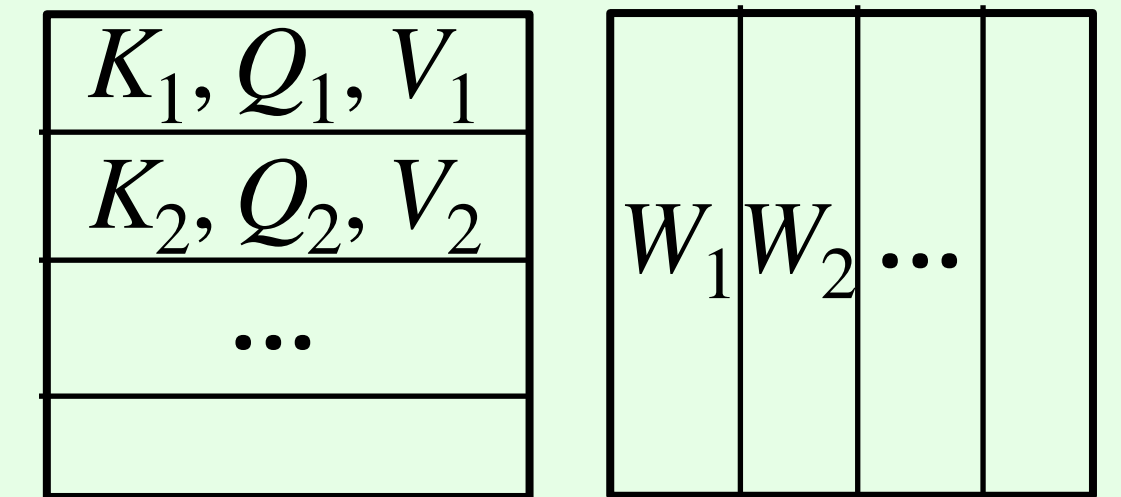
Point clouds: $X := \{x_i\}_{i=1}^n$

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

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



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

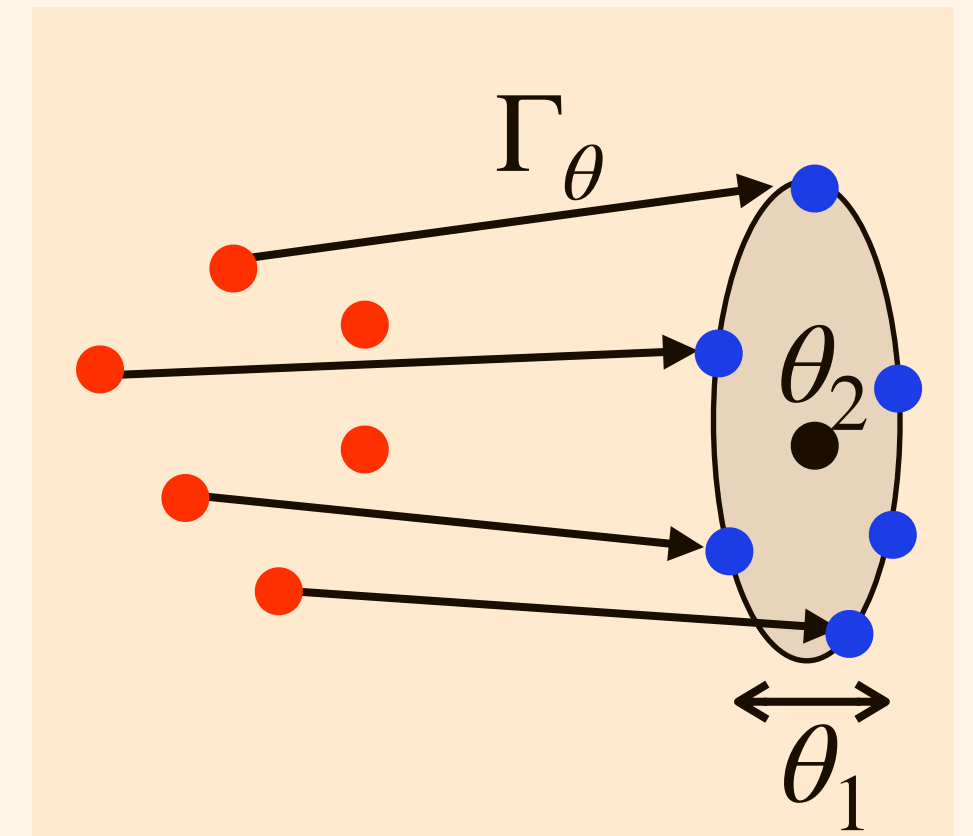
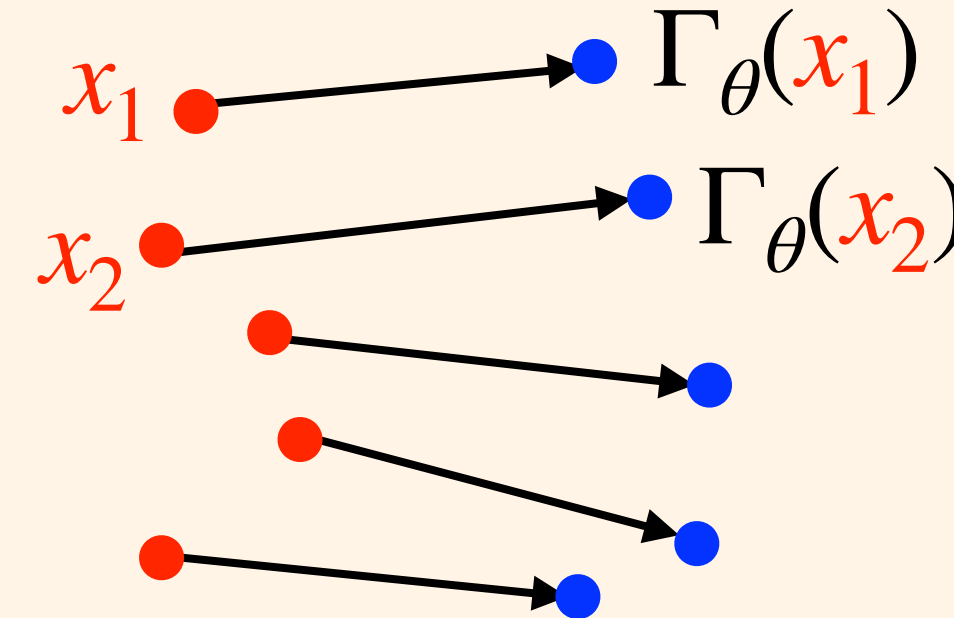


Multi-head attention layer: $X \mapsto \{\sum_{h=1}^H W_h \Gamma_{\theta_h}[X](x_i)\}_{i=1}^n$

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

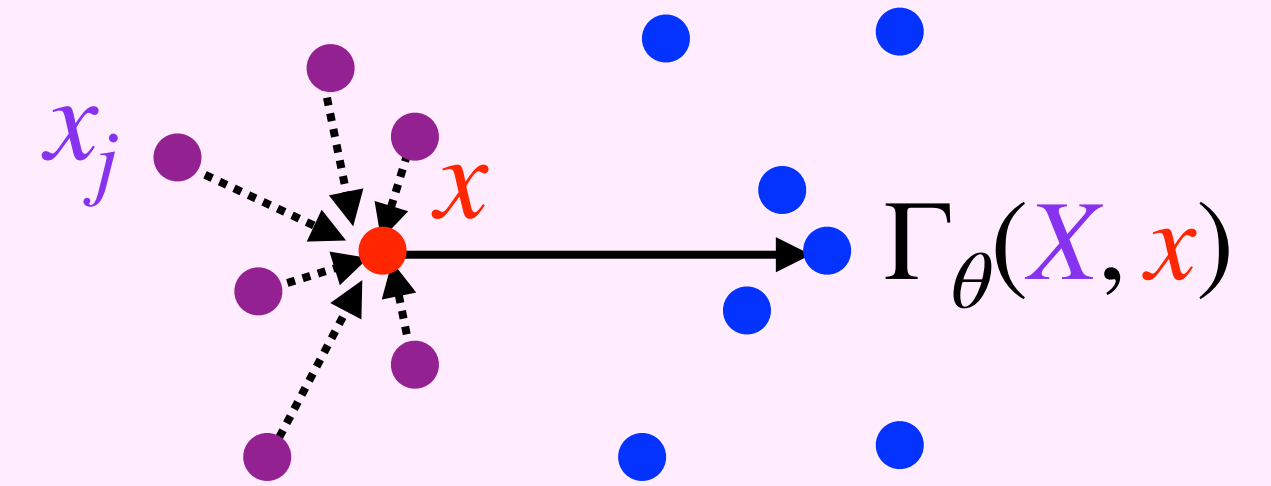


Attention as In-context Mapping

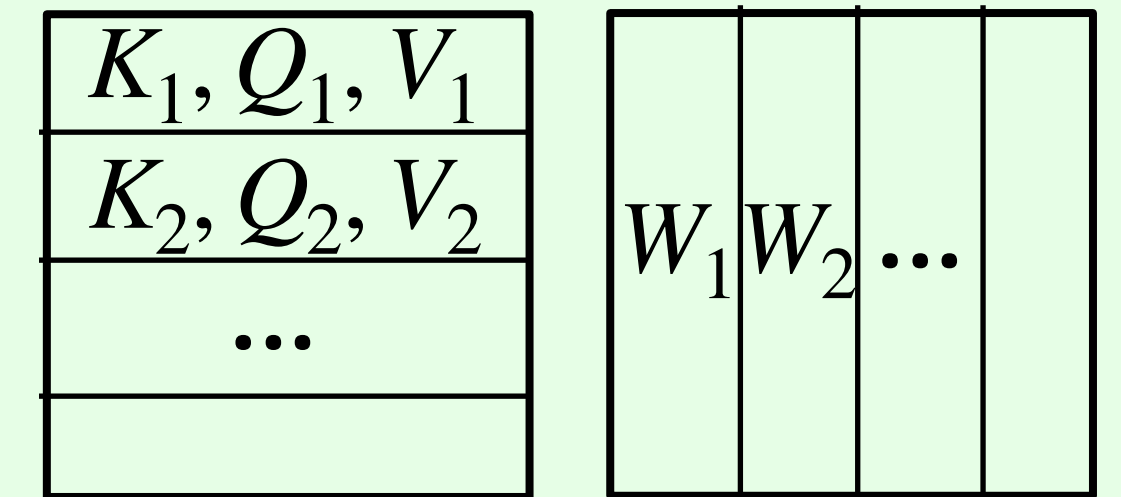
Point clouds: $X := \{x_i\}_{i=1}^n$

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

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



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

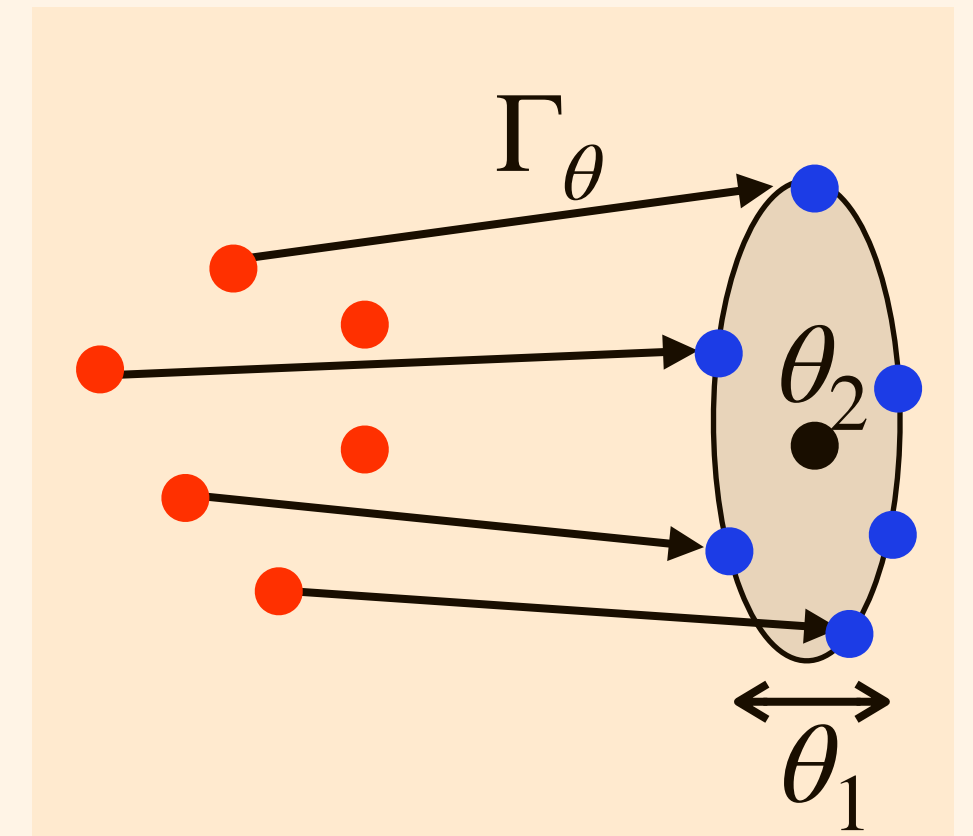
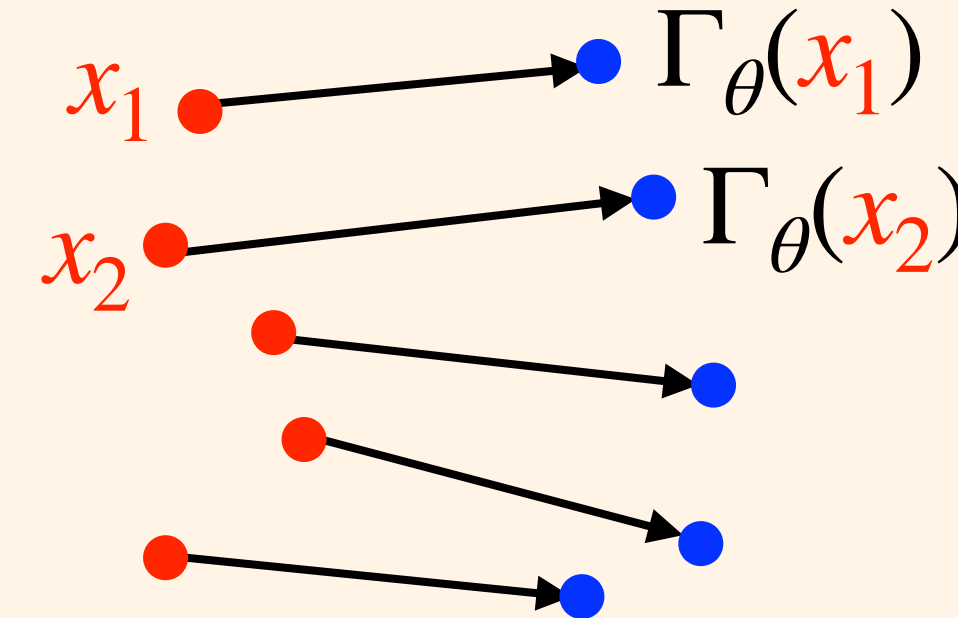


Multi-head attention layer: $X \mapsto \{\sum_{h=1}^H W_h \Gamma_{\theta_h}[X](x_i)\}_{i=1}^n$

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



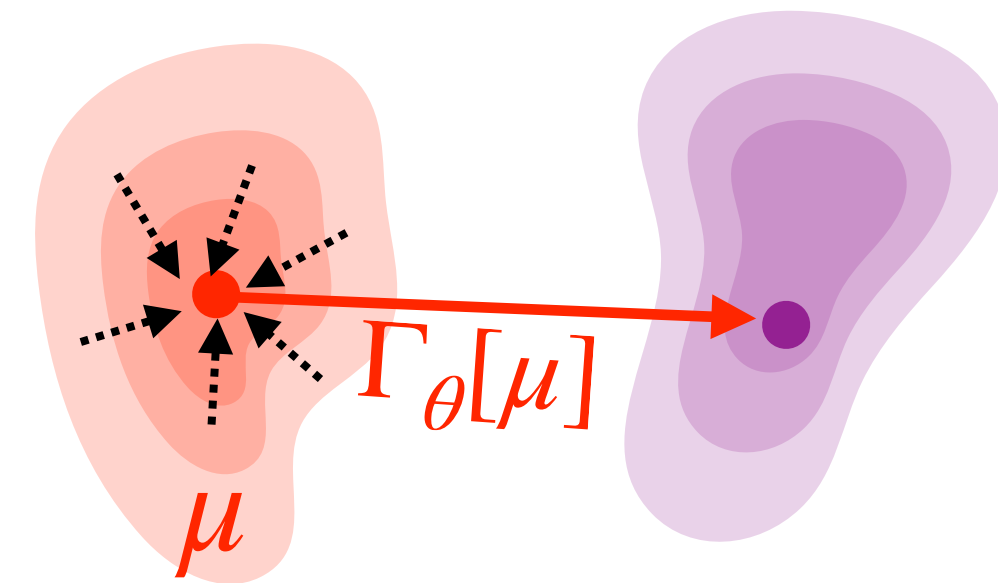
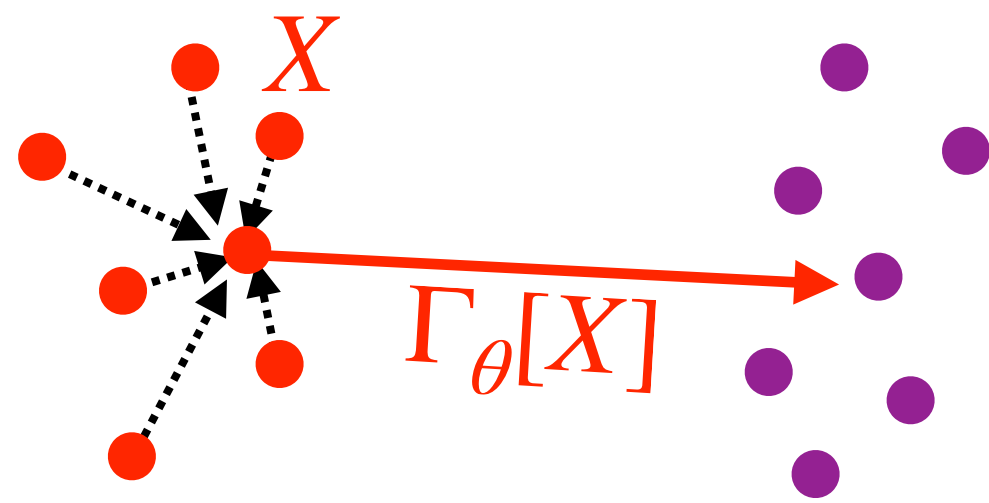
Transformer \equiv composition of in-context and context-free layers.

Attentions Operating over Measures

Number n of token is arbitrary.

(Unmasked) attention is permutation invariant.

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

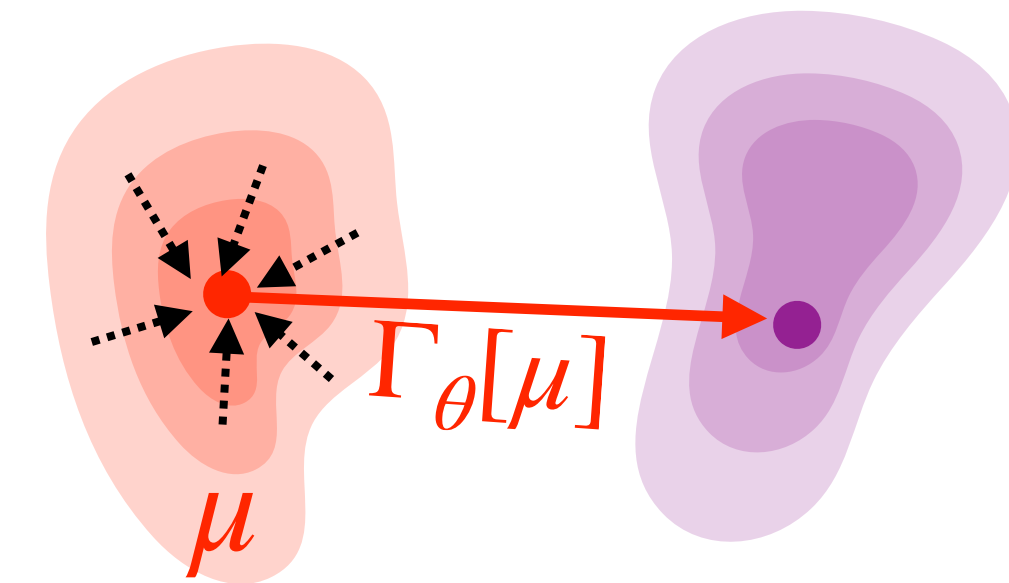
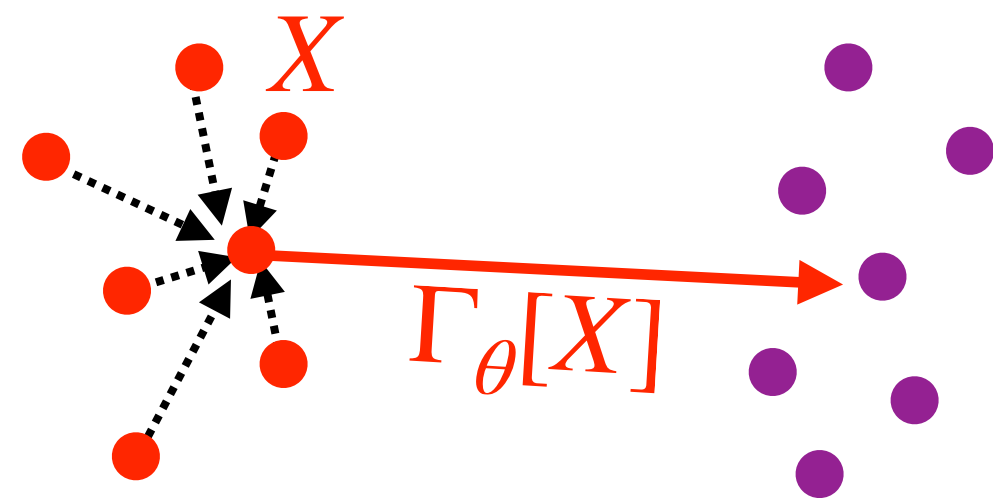


Attentions Operating over Measures

Number n of token is arbitrary.

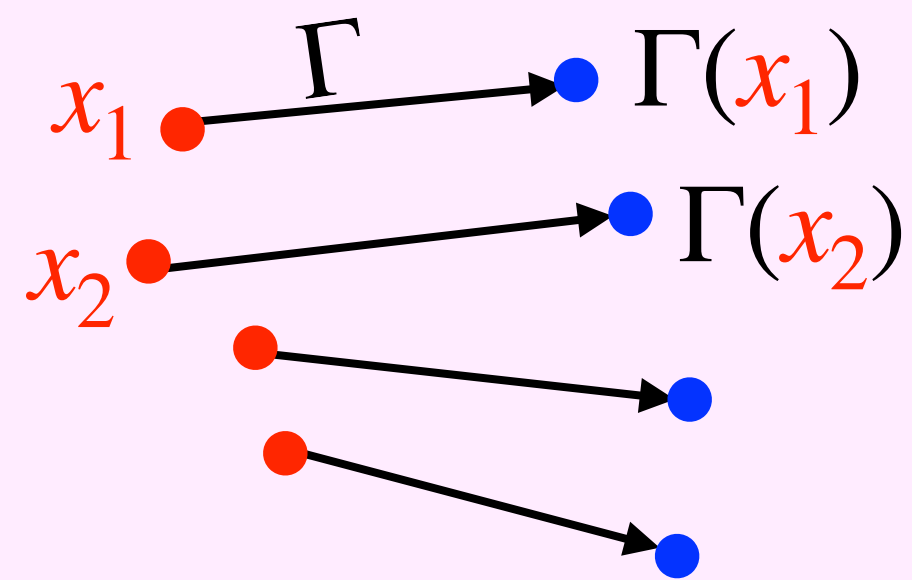
(Unmasked) attention is permutation invariant.

$$\Gamma_{\theta}[X](x) := \sum_j \frac{e^{\langle Kx, Qx_j \rangle}}{\sum_{\ell} e^{\langle Kx, Qx_{\ell} \rangle}} Vx_j \quad \xrightarrow{\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}} \quad \Gamma_{\theta}[\mu](x) := \int \frac{e^{\langle Kx, Qy \rangle}}{\int e^{\langle Kx, Qy \rangle} d\mu(y)} Vy d\mu(y)$$

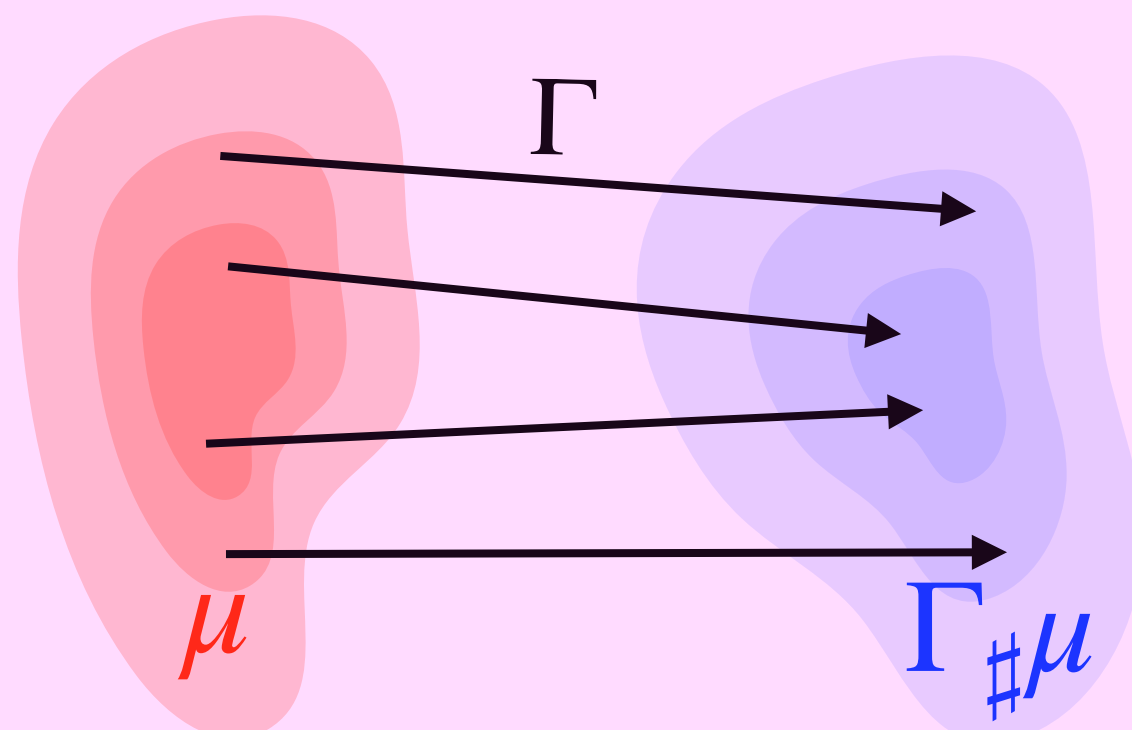


Push-forward

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



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



Attention layers

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

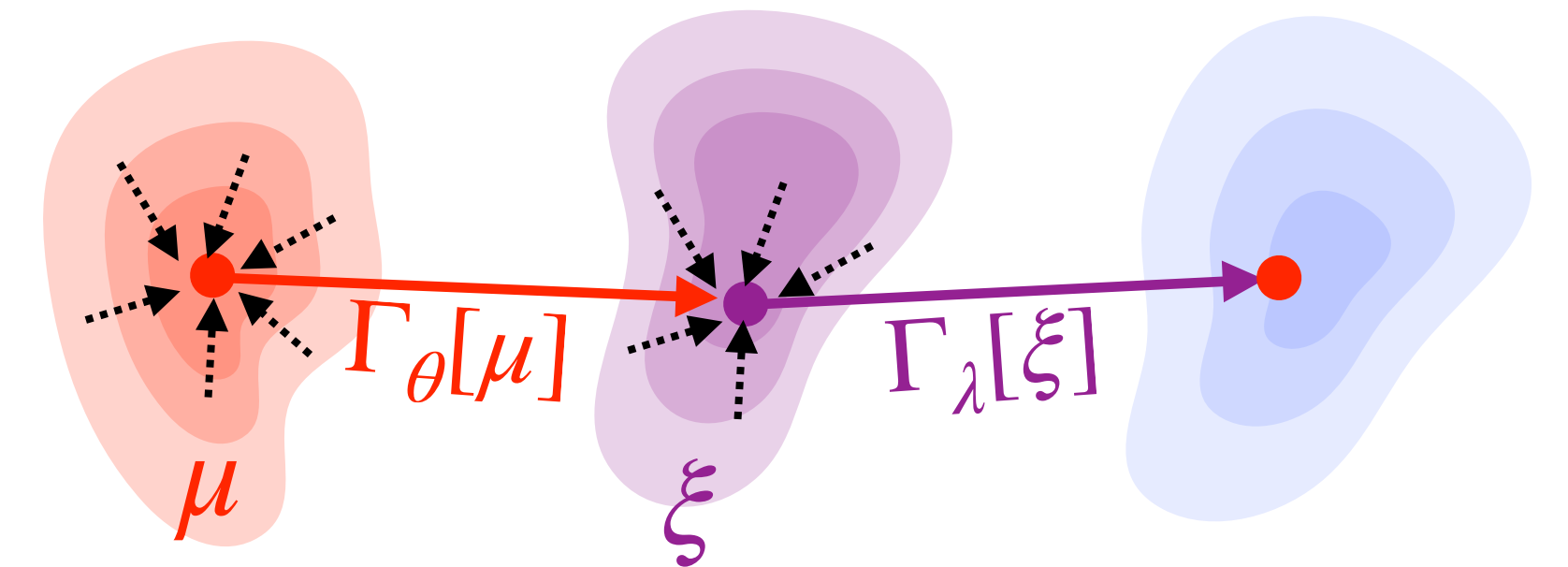
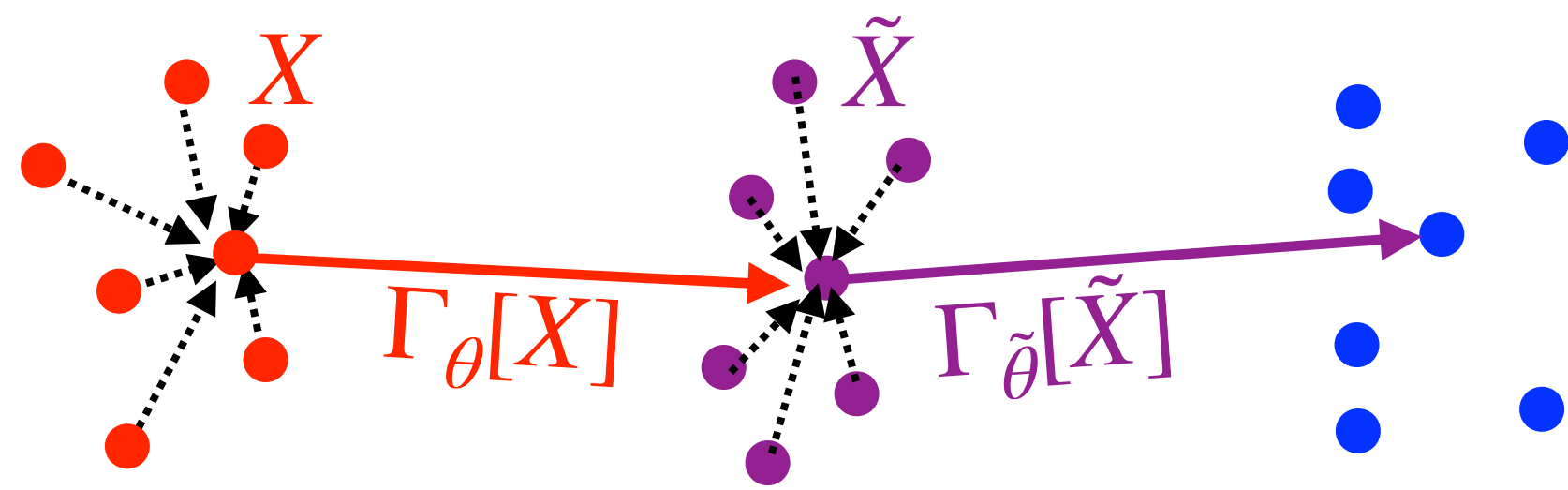
$$\mu \mapsto \Gamma_{\theta}[\mu]_{\#}\mu$$

Attentions Operating over Measures

Number n of token is arbitrary.

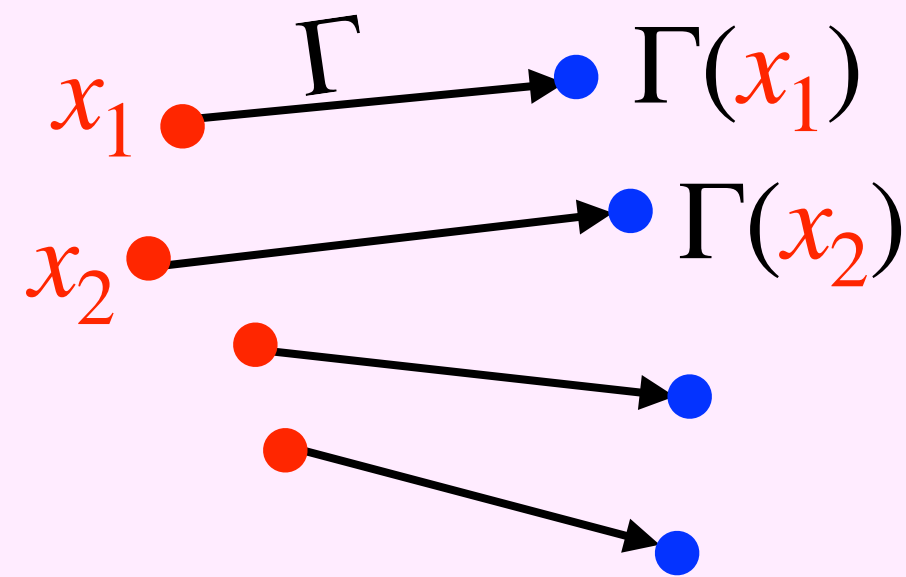
(Unmasked) attention is permutation invariant.

$$\Gamma_{\theta}[X](x) := \sum_j \frac{e^{\langle Kx, Qx_j \rangle}}{\sum_{\ell} e^{\langle Kx, Qx_{\ell} \rangle}} Vx_j \quad \xrightarrow{\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}} \quad \Gamma_{\theta}[\mu](x) := \int \frac{e^{\langle Kx, Qy \rangle}}{\int e^{\langle Kx, Qy \rangle} d\mu(y)} Vy d\mu(y)$$

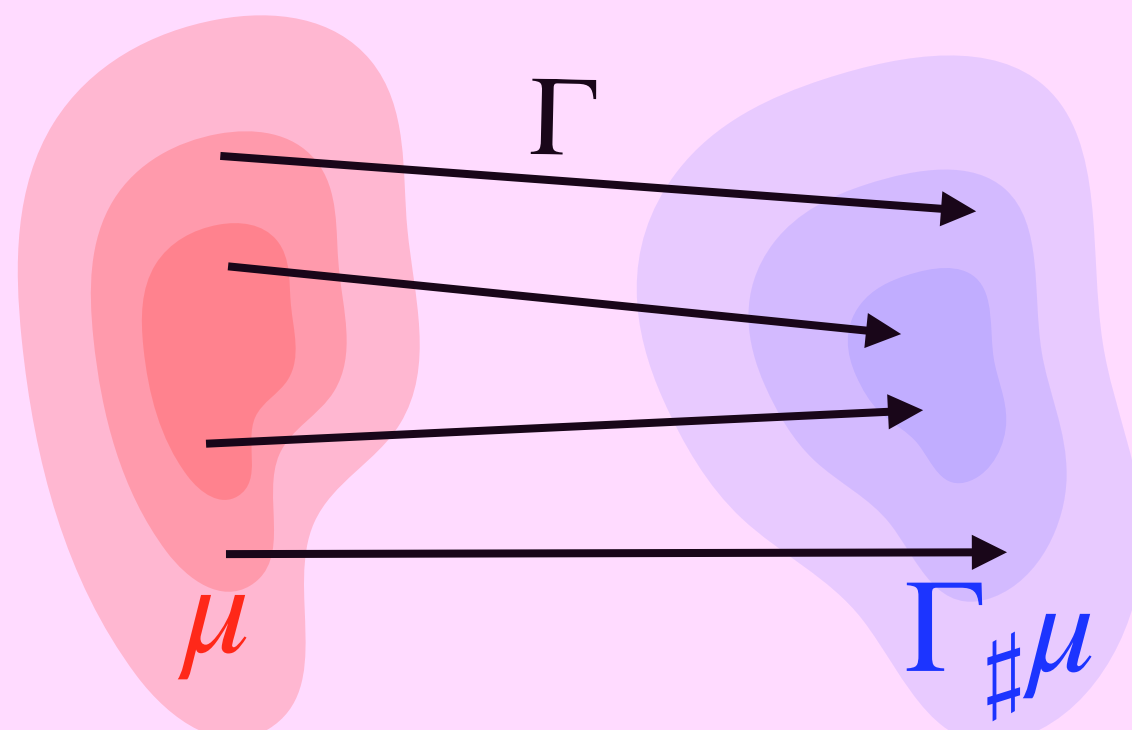


Push-forward

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



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



Attention layers

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

$$\mu \mapsto \Gamma_{\theta}[\mu]_{\#}\mu$$

Composing layers

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

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

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

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

Masked Causal Attention over Measures

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

Masked attention mapping: $\Gamma_{\theta}[X](x_i) := \sum_{j \leq i} \frac{e^{\langle Kx_i, Qx_j \rangle}}{\sum_{\ell \leq i} e^{\langle Kx_i, Qx_{\ell} \rangle}} Vx_j$

→ breaks permutation invariance.

Masked Causal Attention over Measures

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

Masked attention mapping: $\Gamma_{\theta}[X](x_i) := \sum_{j \leq i} \frac{e^{\langle Kx_i, Qx_j \rangle}}{\sum_{\ell \leq i} e^{\langle Kx_i, Qx_{\ell} \rangle}} Vx_j$

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

Masked Causal Attention over Measures

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

Masked attention mapping: $\Gamma_{\theta}[X](x_i) := \sum_{j \leq i} \frac{e^{\langle Kx_i, Qx_j \rangle}}{\sum_{\ell \leq i} e^{\langle Kx_i, Qx_{\ell} \rangle}} Vx_j$

→ breaks permutation invariance.

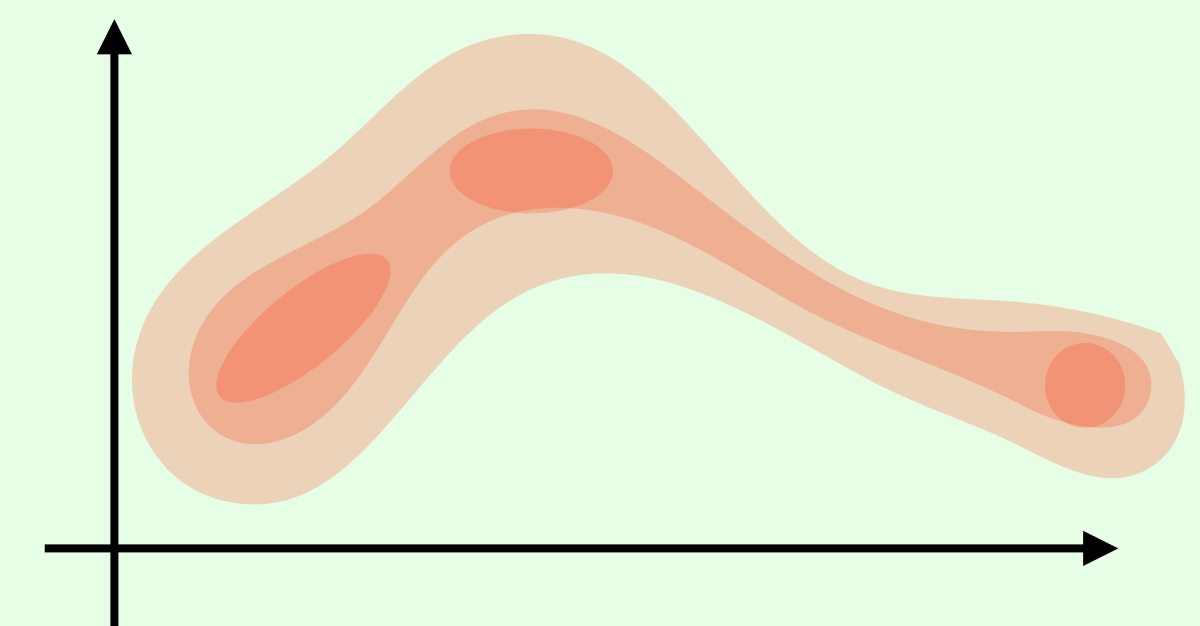
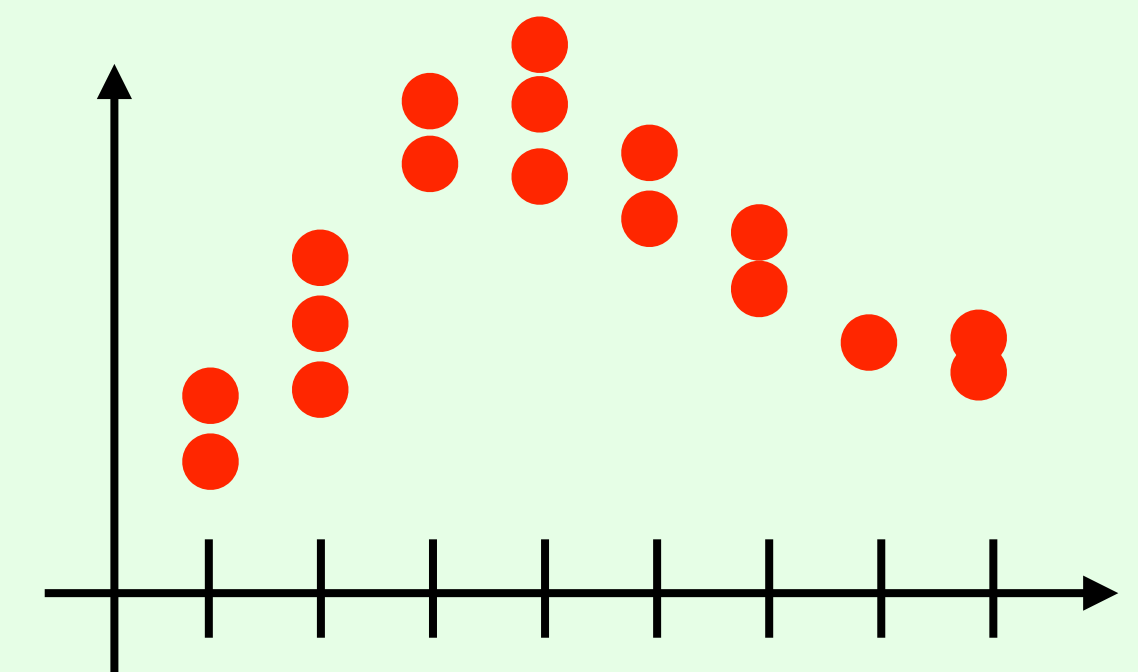
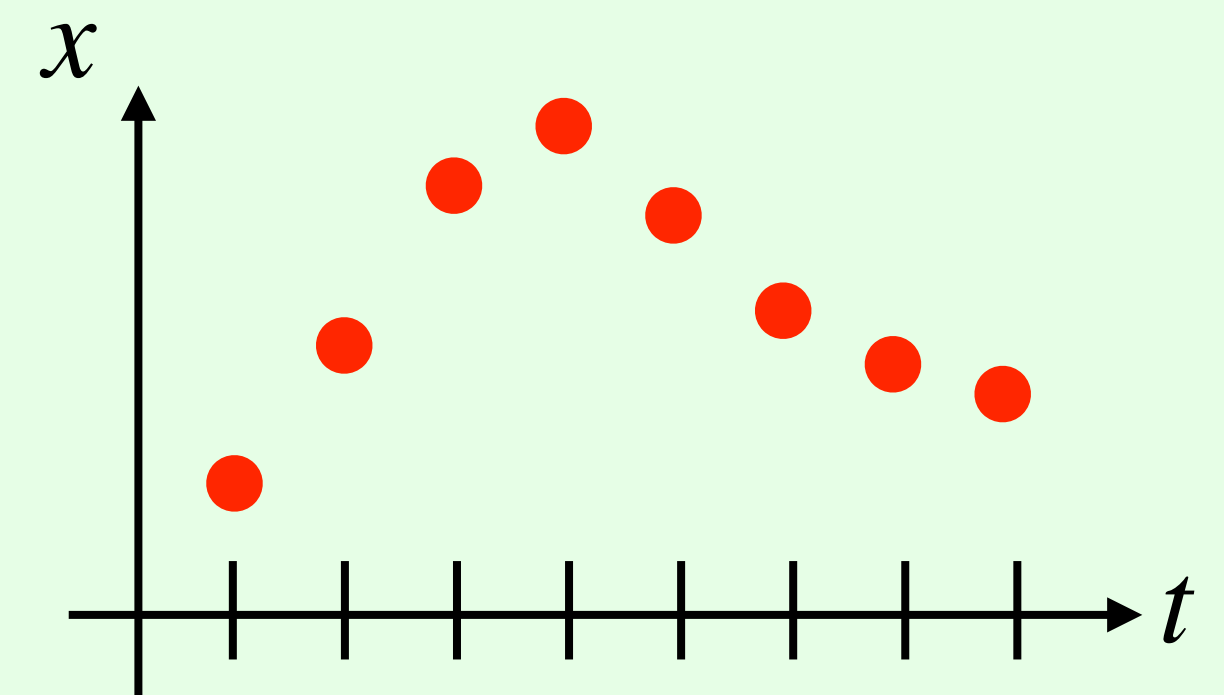
Training: next token prediction (simplified...) $\min_{\theta} \sum_X \sum_{i=1}^{n-1} \ell(\Gamma_{\theta}[X](x_i), x_{i+1})$

Testing: generative model (simplified...) $X \mapsto (x_1, \dots, x_i, \Gamma[X](x_i))$

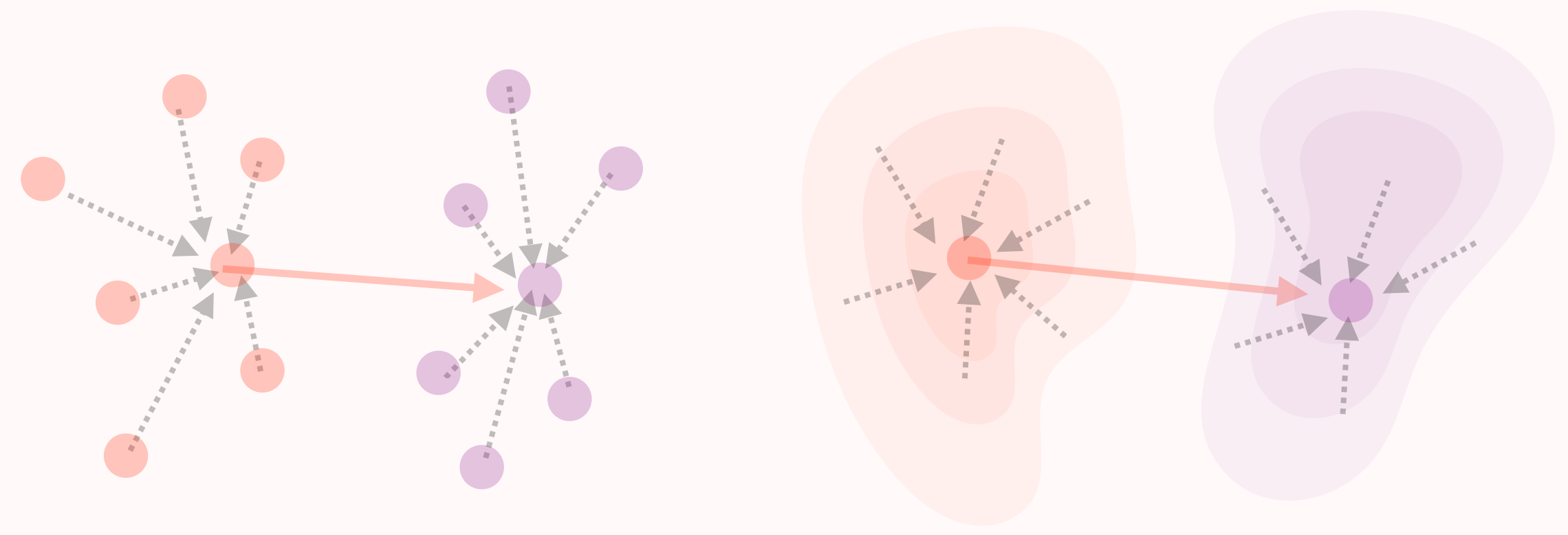
Space-time lifting:

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

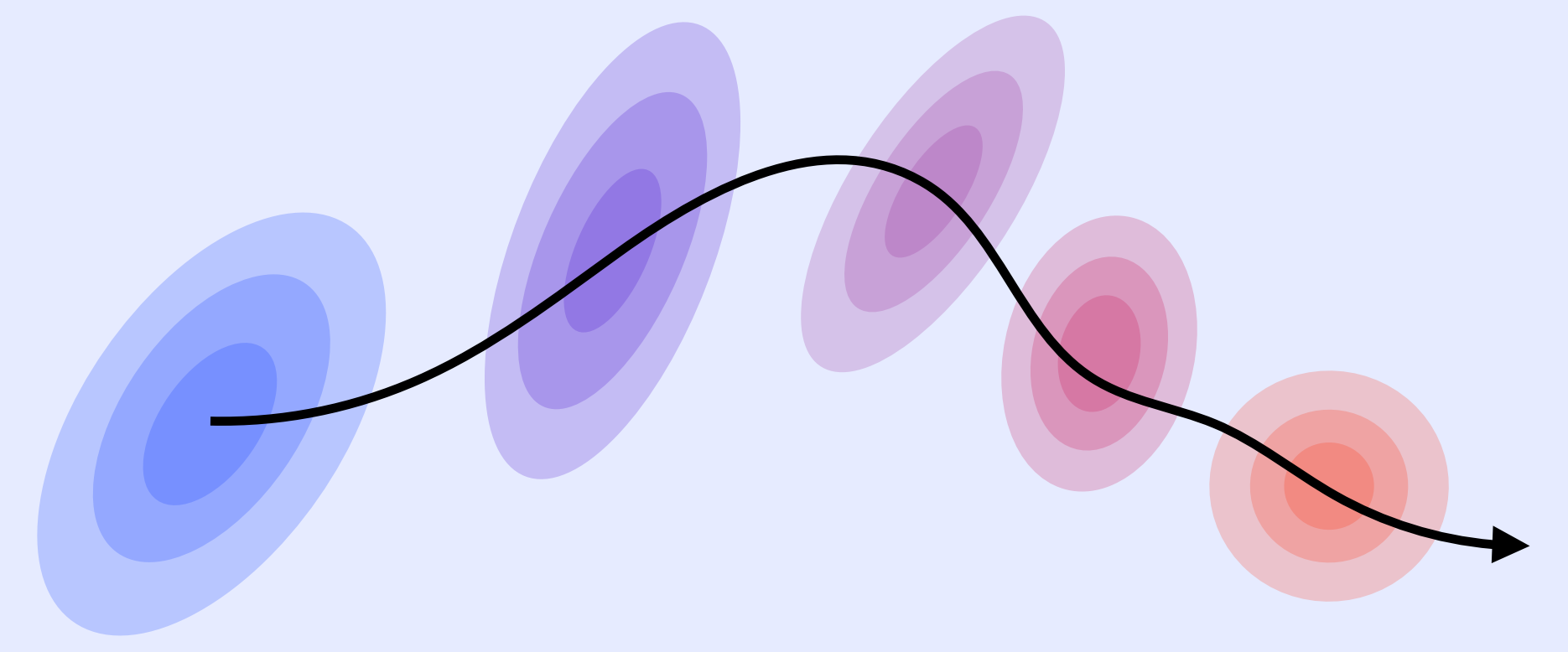
$$\Gamma_{\theta}[\mu](x, t) := \int \frac{1_{s \leq t} e^{\langle Kx, Qy \rangle}}{\int 1_{s' \leq t} e^{\langle Kx, Qy' \rangle} d\mu(y', s')} Vy d\mu(y, s)$$



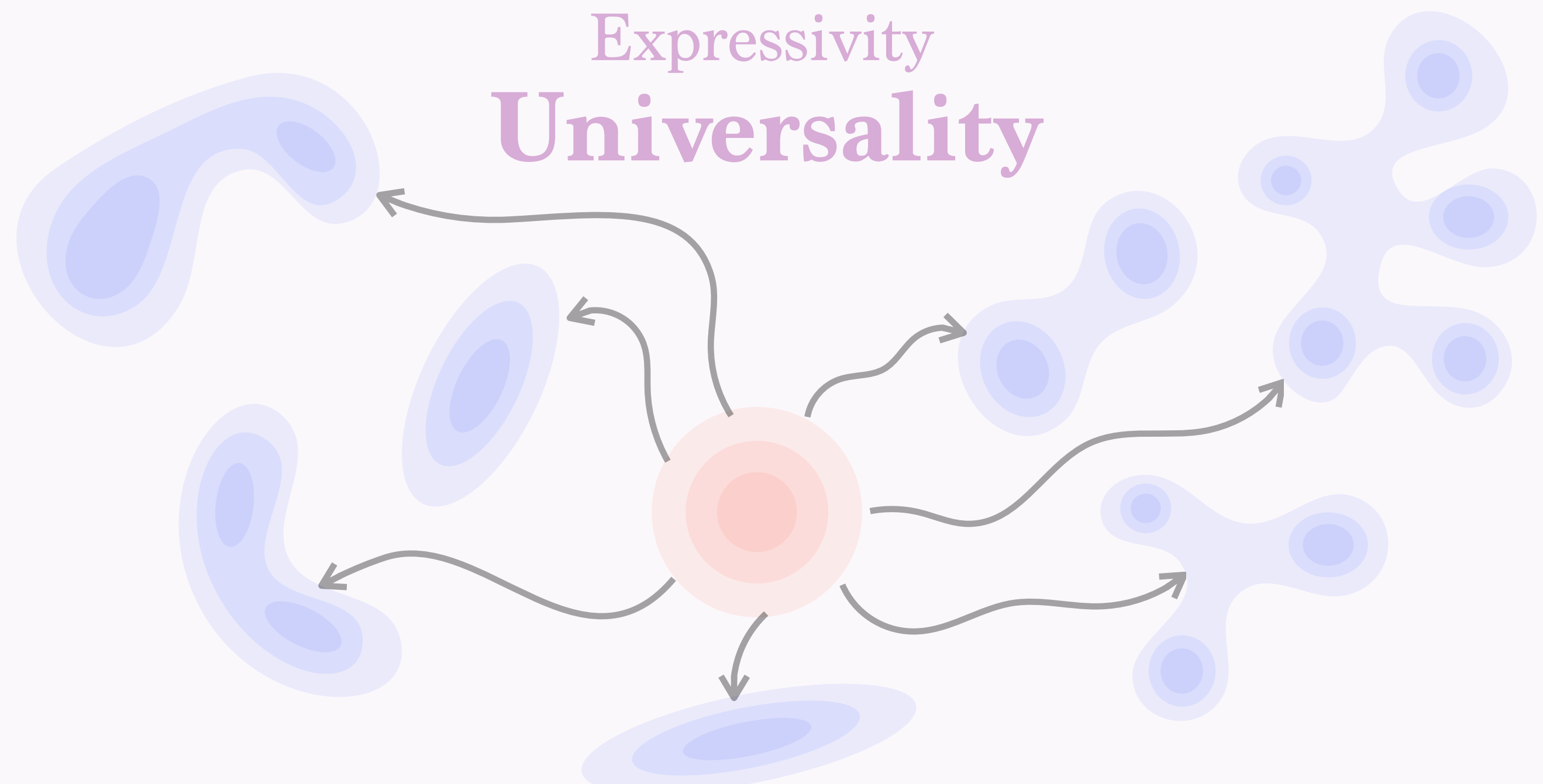
Arbitrary number of layers
In Context Mappings
over Measures



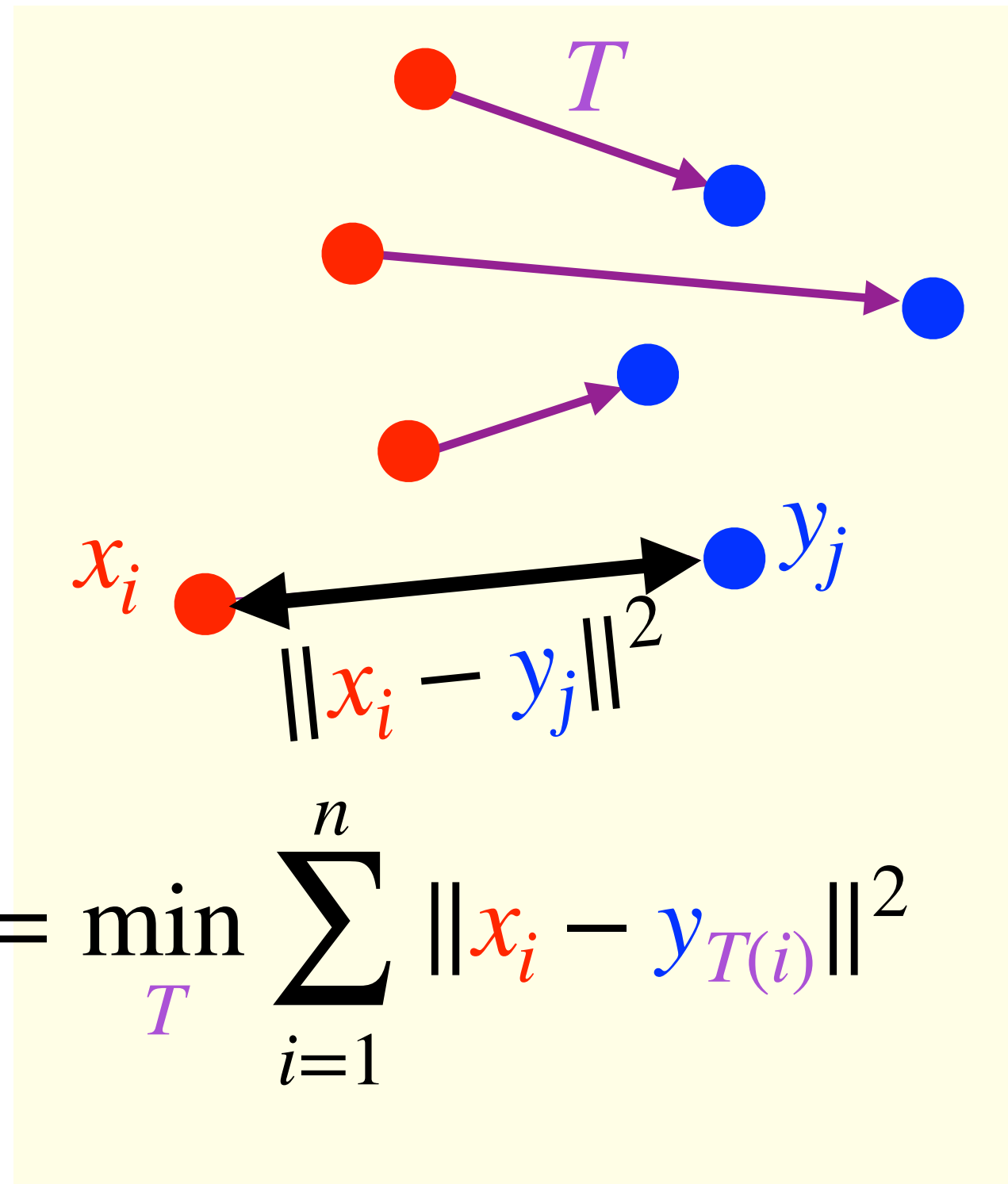
Arbitrary number of layers
Smoothness and
PDE's



Expressivity
Universality



Optimal Transport (Wasserstein) Distance

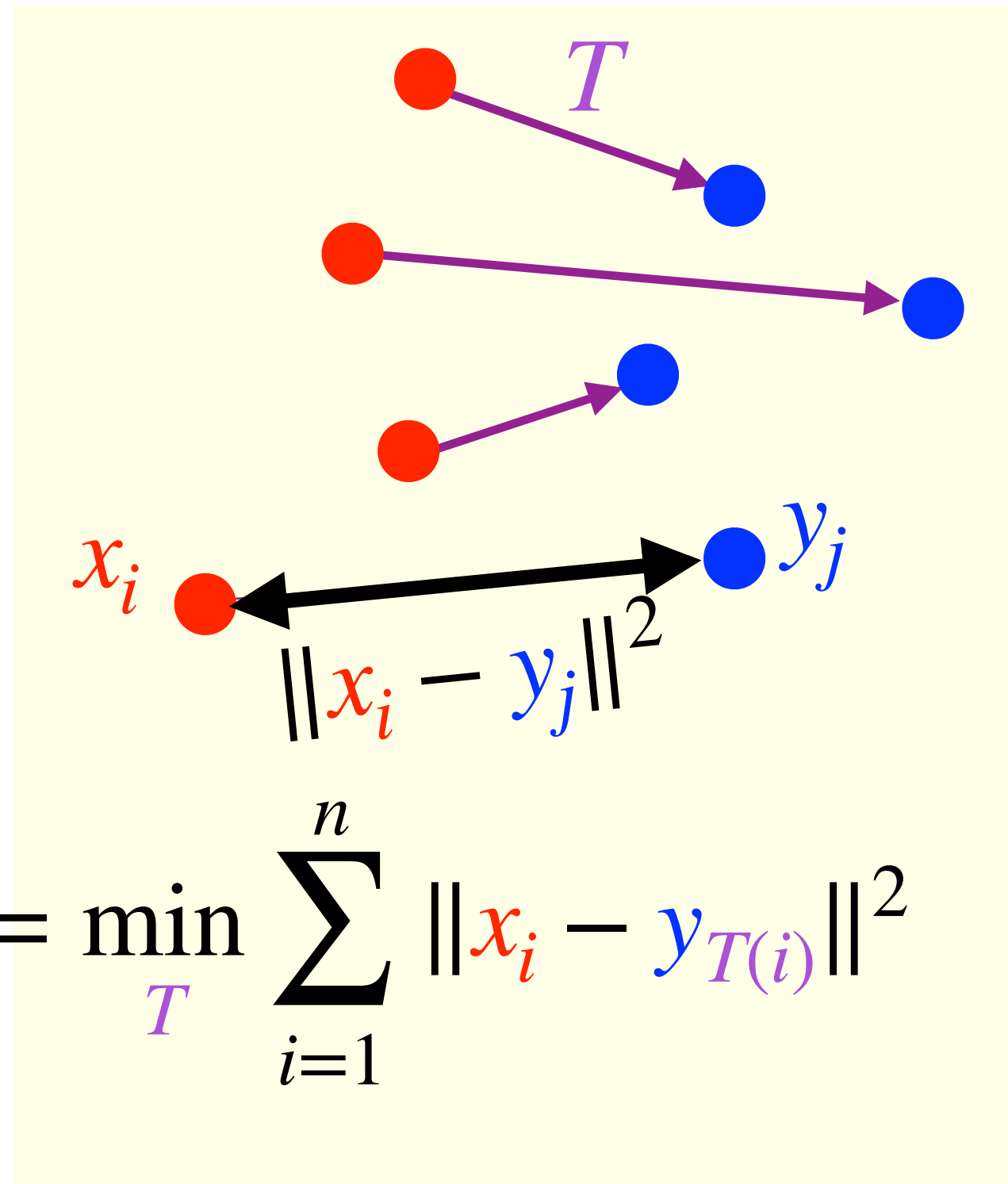


$$W_2(\mu, \nu)^2 := \min_T \sum_{i=1}^n \|x_i - y_{T(i)}\|^2$$



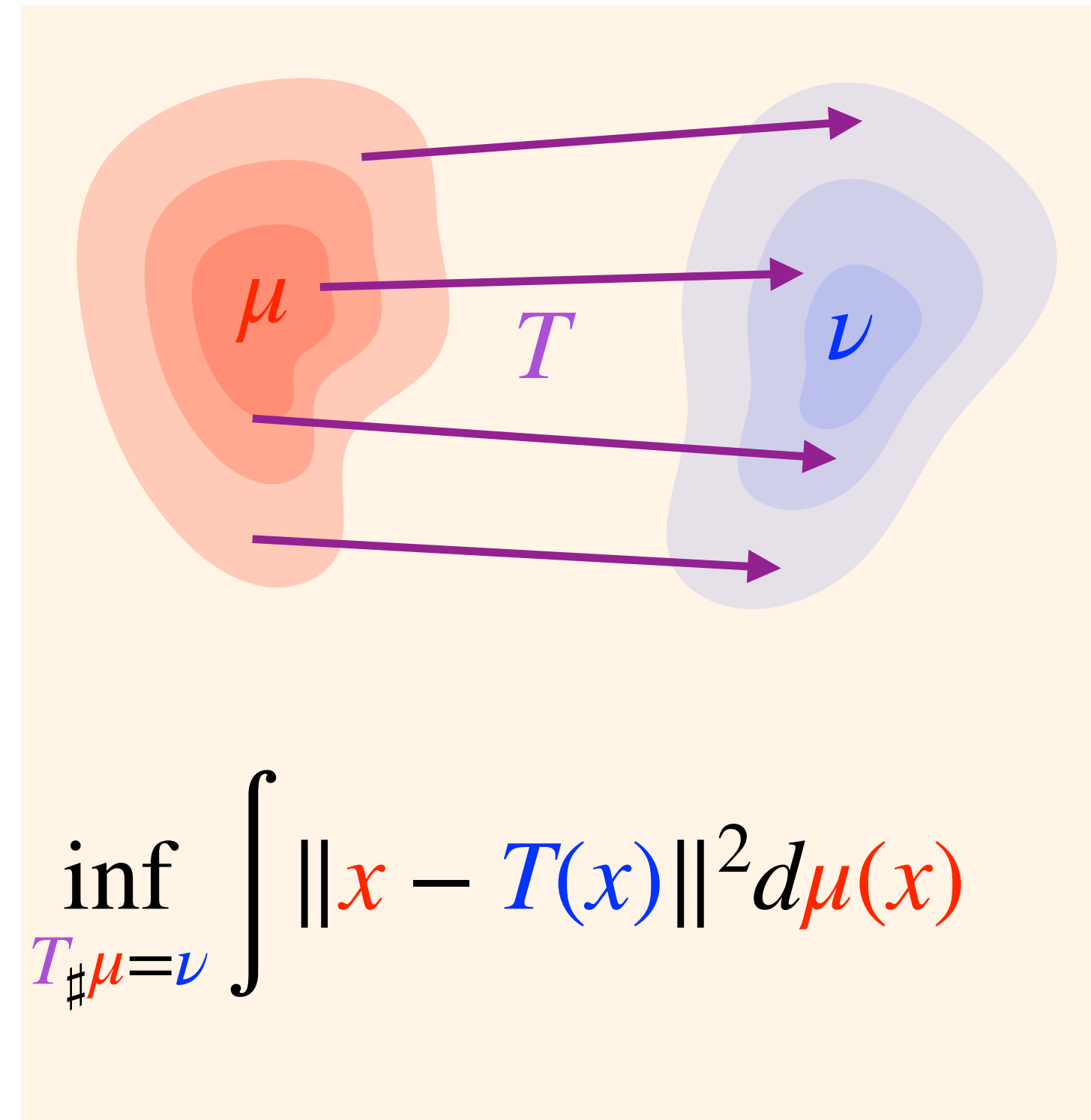
Monge 1784

Optimal Transport (Wasserstein) Distance



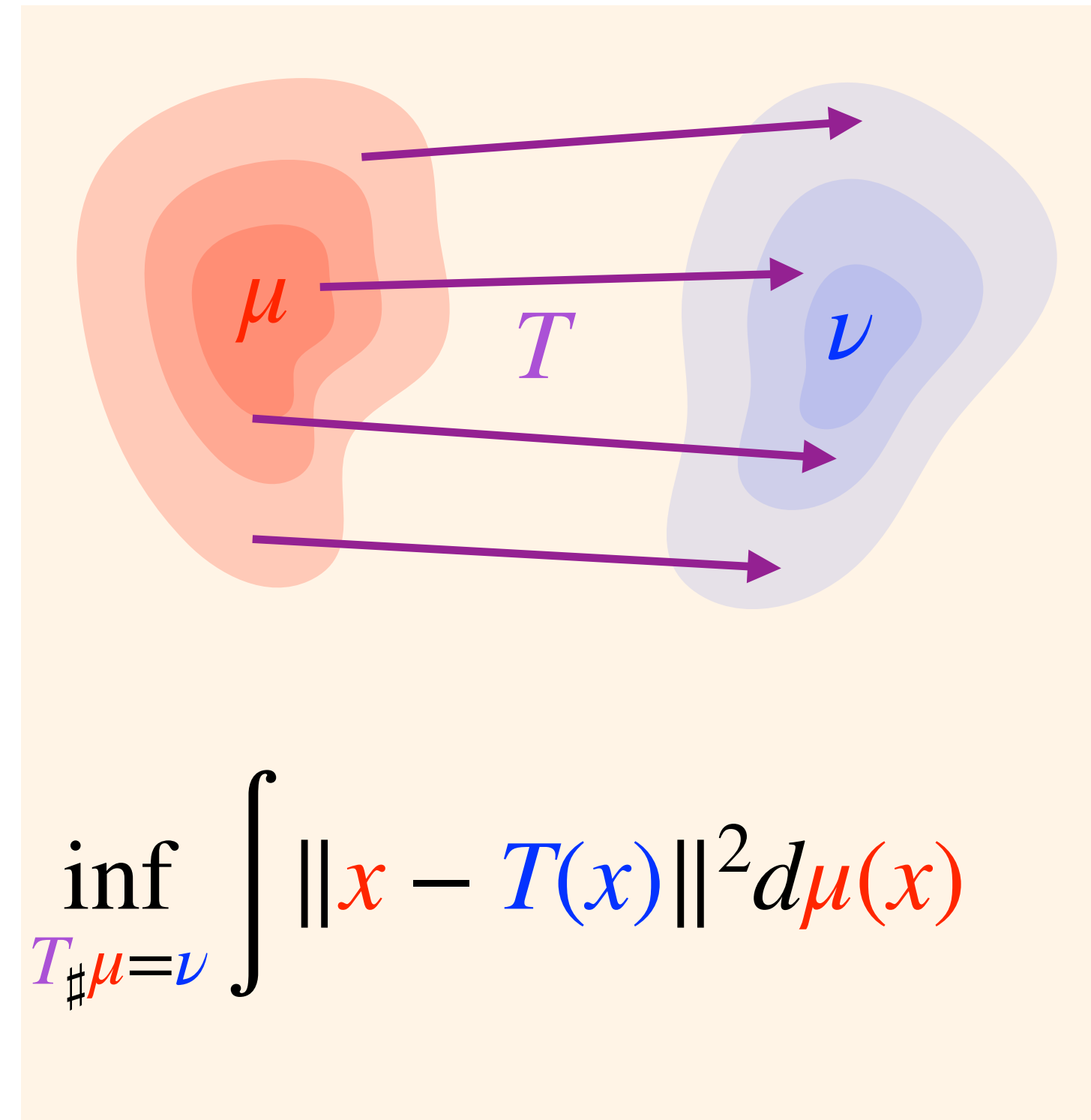
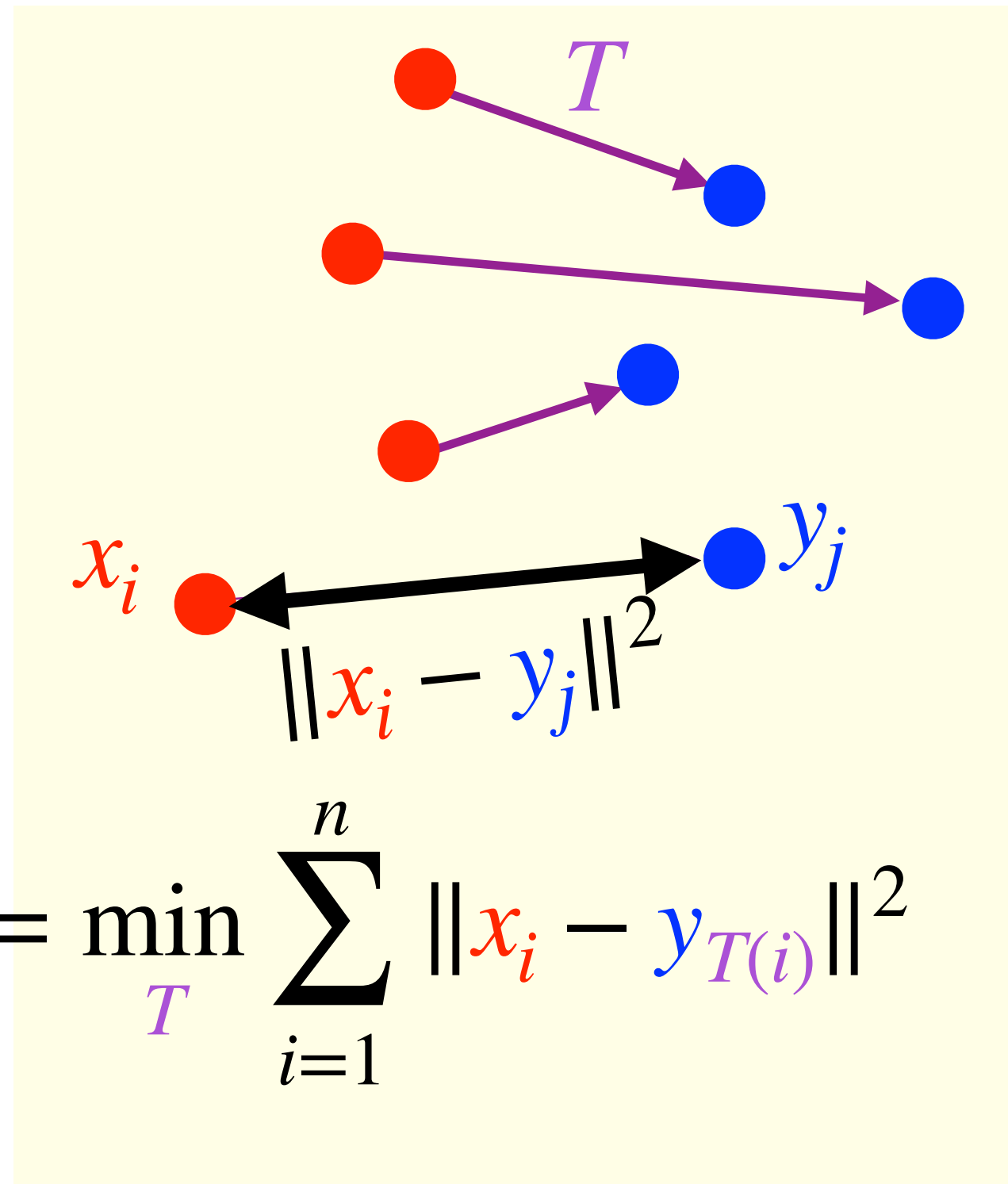
$$W_2(\mu, \nu)^2 := \min_T \sum_{i=1}^n \|x_i - y_{T(i)}\|^2$$

$$= \inf_{T_{\#}\mu = \nu} \int \|x - T(x)\|^2 d\mu(x)$$



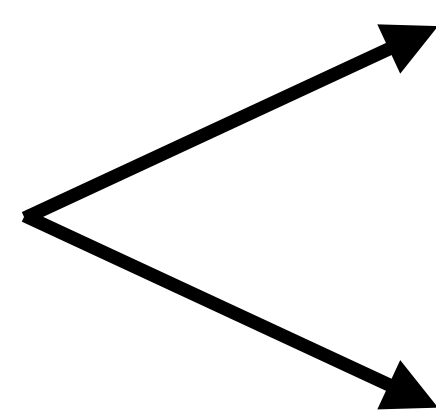
Monge 1784

Optimal Transport (Wasserstein) Distance



$$W_2(\mu, \nu)^2 := \min_T \sum_{i=1}^n \|x_i - y_{T(i)}\|^2 = \inf_{T_{\#}\mu = \nu} \int \|x - T(x)\|^2 d\mu(x)$$

General measures:



Kantorovitch relaxation

or

Approximation by discrete measures



Monge 1784



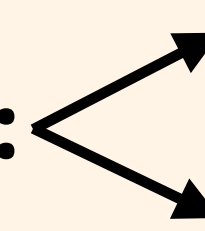
Kantorovitch 1942

How Smooth is Attention?

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

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

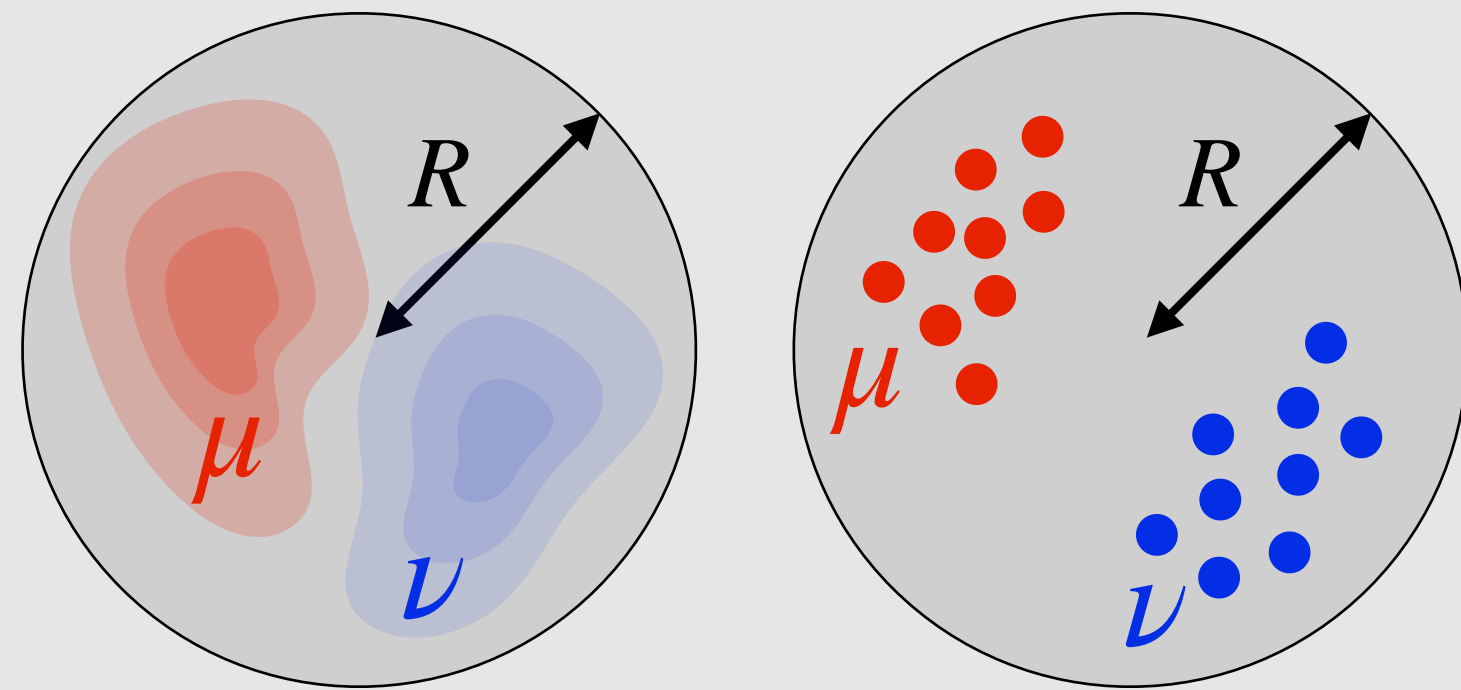
Lipschitz regularity: $W_2(\Gamma_\theta[\mu]_{\#}\mu, \Gamma_\theta[\nu]_{\#}\nu) \leq C_\theta W_2(\mu, \nu)$

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

How Smooth is Attention?

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

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



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

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

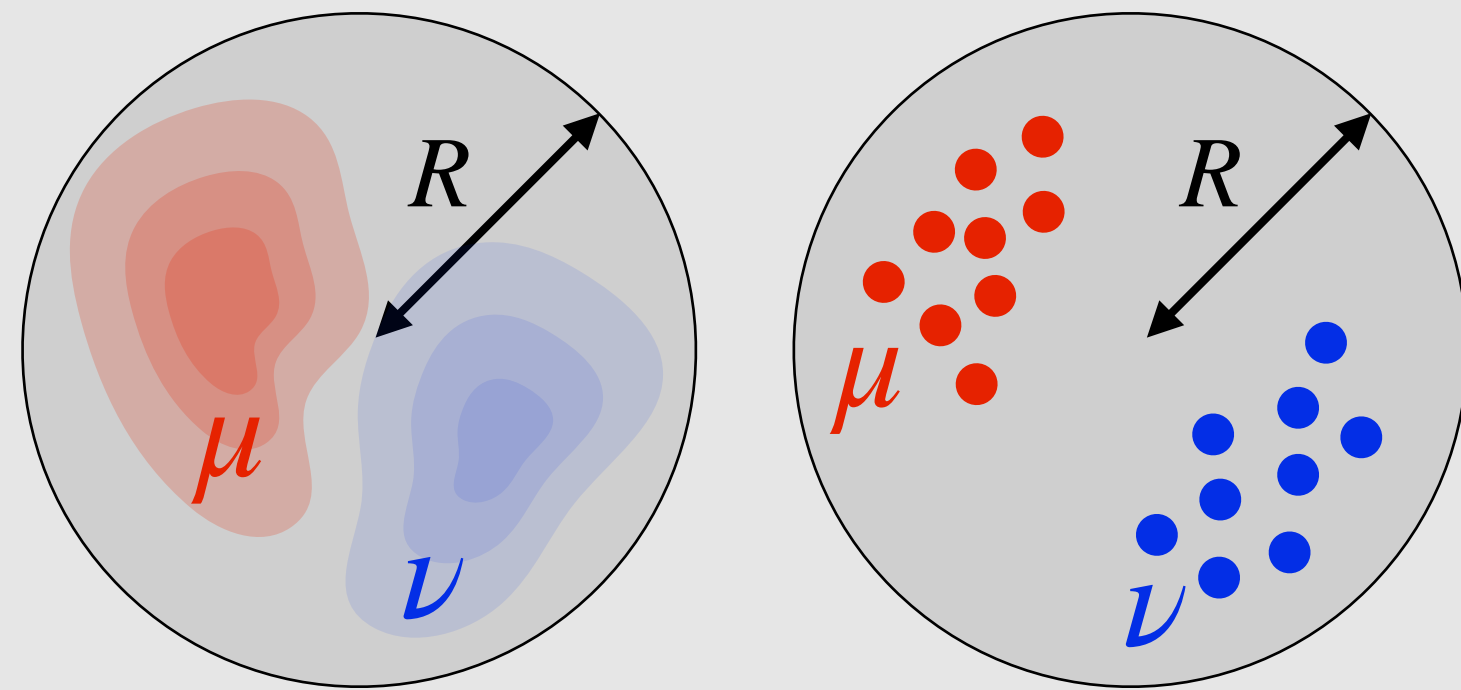
If furthermore $\mu = \frac{1}{n} \sum_i \delta_{x_i}, \nu = \frac{1}{n} \sum_i \delta_{y_i}$

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

How Smooth is Attention?

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

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



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]

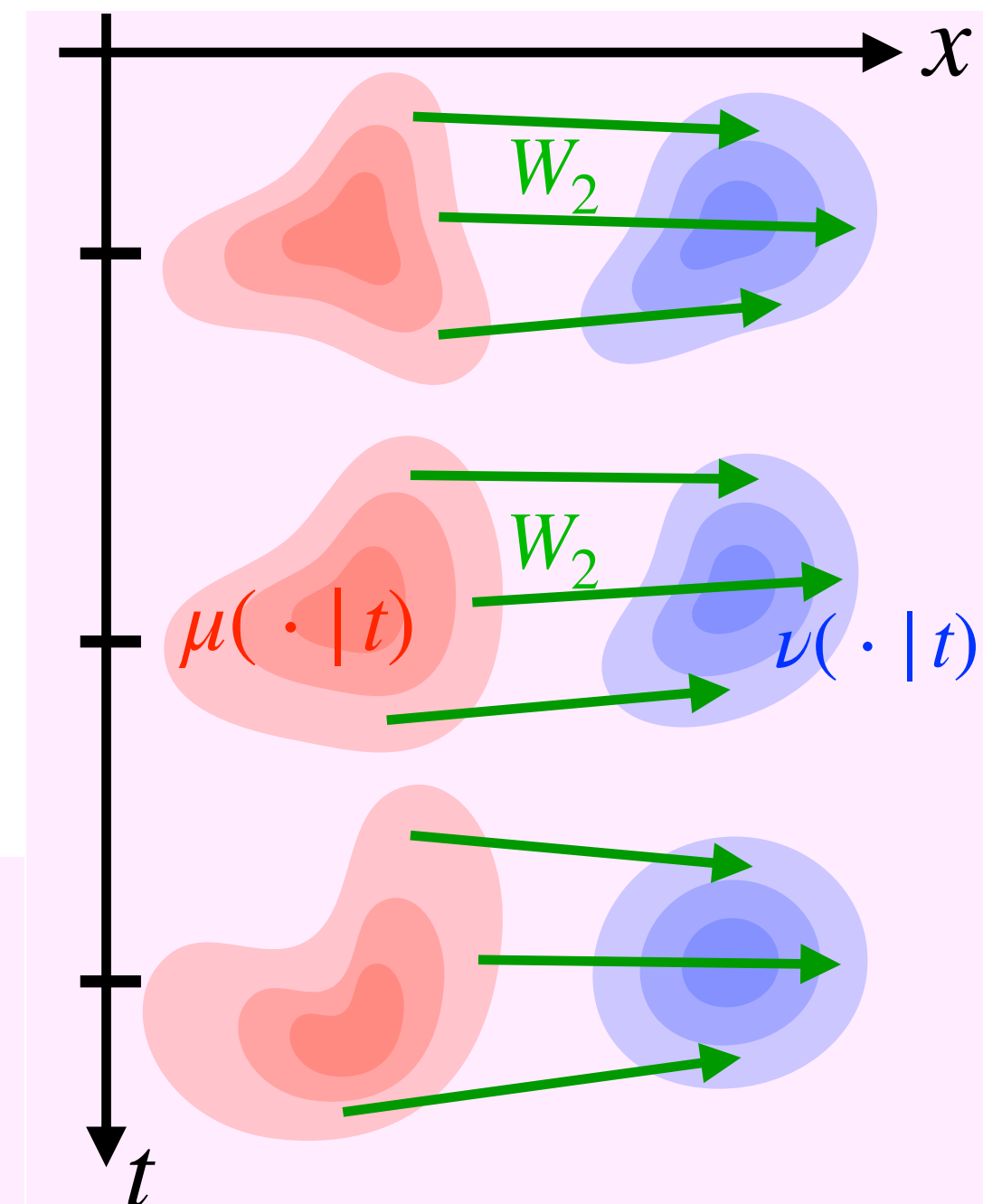
If $\text{supp}(\mu), \text{supp}(\nu) \subset B(0, R)$,

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

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

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

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

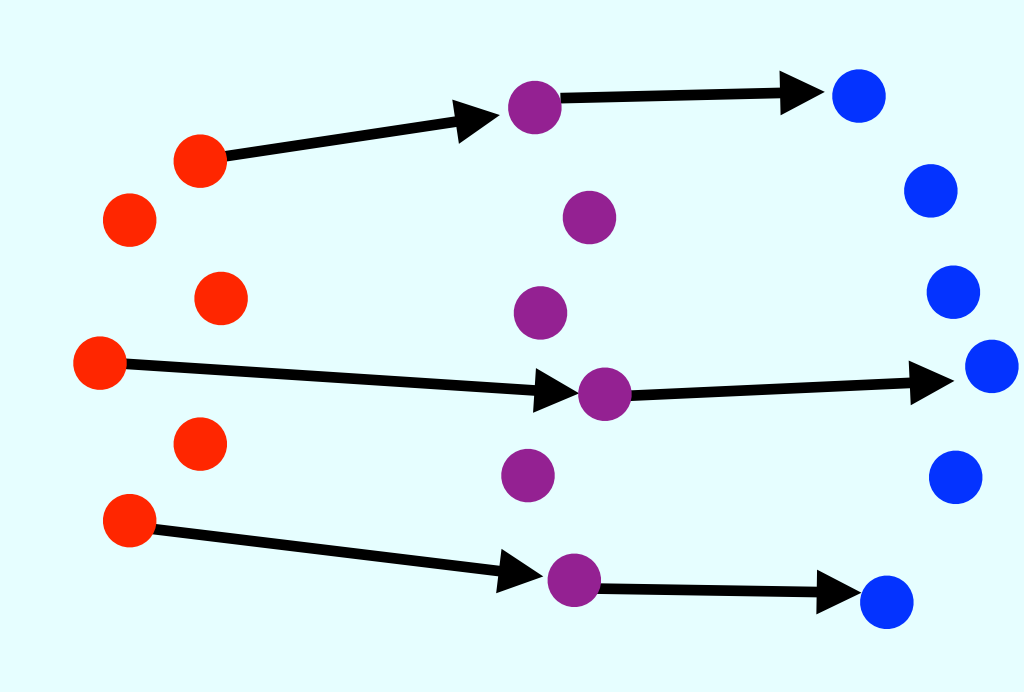


Infinite Depth as a Neural PDE

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

$$\theta = (Q, K, V)$$

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

$$x_i(t+1) = x_i(t) + \frac{1}{T} \Gamma_{\theta(t)}[\mu(t)](x_i(t))$$


The diagram illustrates the update rule for $x_i(t+1)$. It shows three particles (red dots) moving from left to right through a medium (purple dots) towards a target (blue dots). Arrows indicate the direction of movement.

Infinite Depth as a Neural PDE

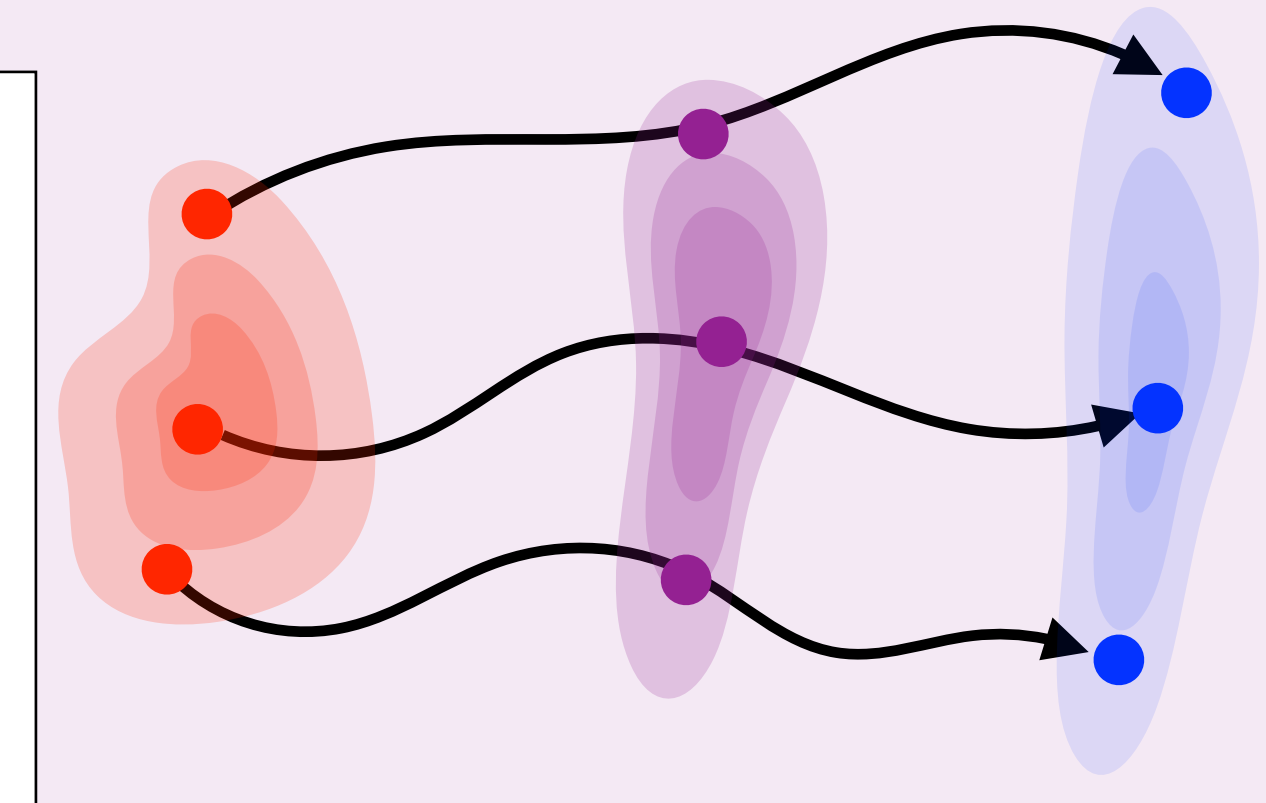
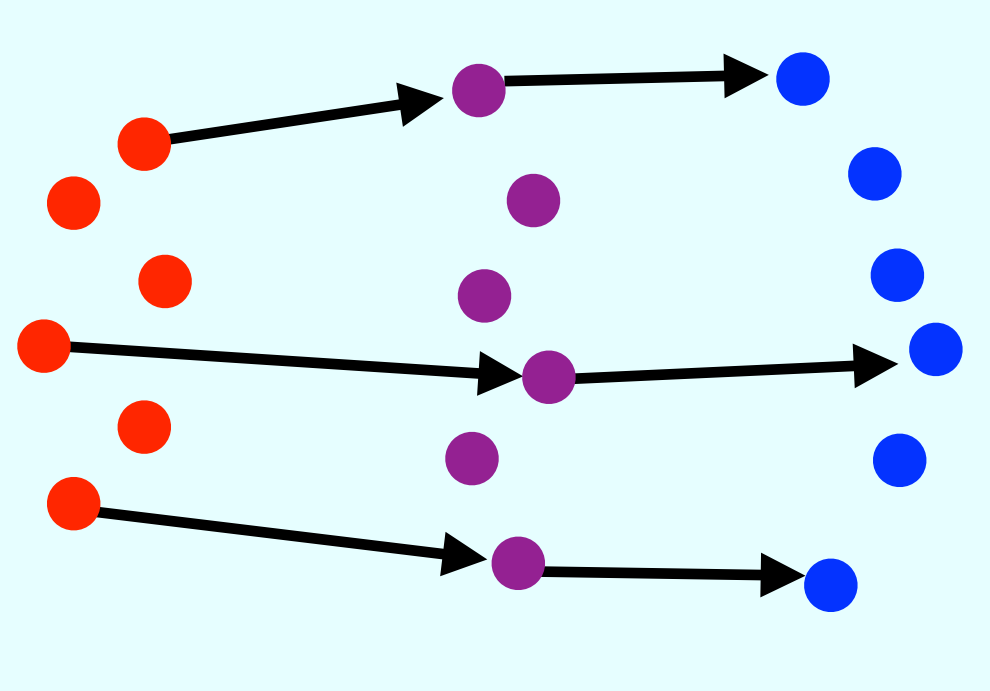
$$\Gamma_{\theta}[\mu](x) := \int \frac{e^{\langle Kx, Qy \rangle}}{\int e^{\langle Kx, Qy' \rangle} d\mu(y')} Vy d\mu(y) \quad \theta = (Q, K, V) \quad \mu(t) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$$

$$x_i(t+1) = x_i(t) + \frac{1}{T} \Gamma_{\theta(t)}[\mu(t)](x_i(t))$$

$T \rightarrow +\infty$
Infinite depth

$$\frac{dx_i}{dt}(t) = \Gamma_{\theta(t)}[\mu(t)](x_i(t))$$

Coupled
EDOs



Infinite Depth as a Neural PDE

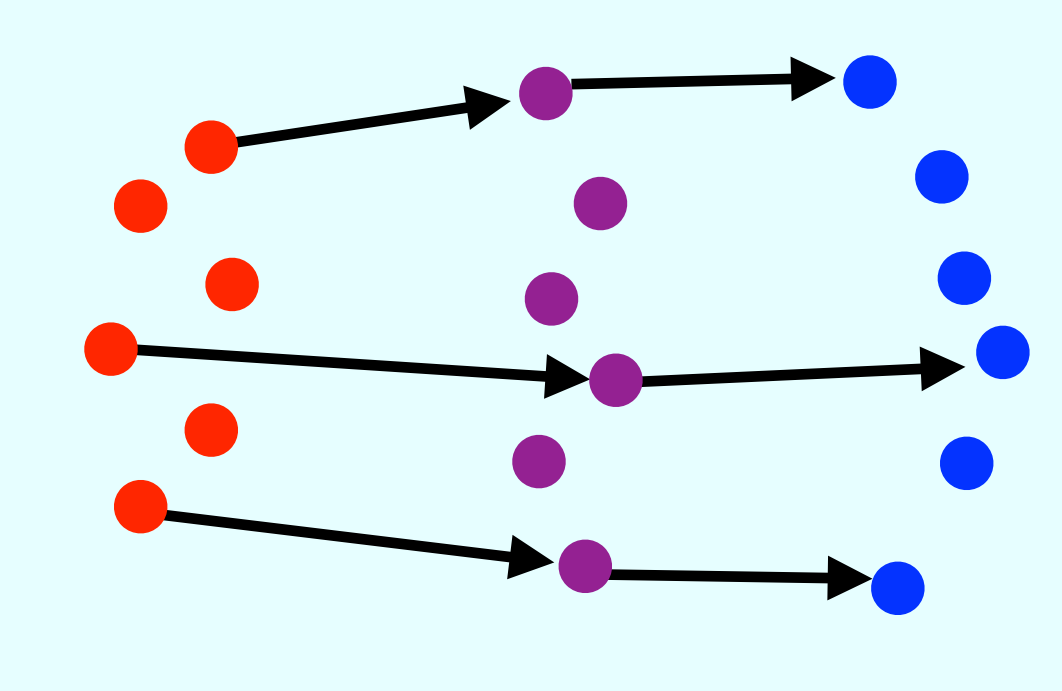
$$\Gamma_{\theta}[\mu](x) := \int \frac{e^{\langle Kx, Qy \rangle}}{\int e^{\langle Kx, Qy' \rangle} d\mu(y')} Vy d\mu(y) \quad \theta = (Q, K, V) \quad \mu(t) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$$

$$x_i(t+1) = x_i(t) + \frac{1}{T} \Gamma_{\theta(t)}[\mu(t)](x_i(t))$$

$T \rightarrow +\infty$
Infinite depth

$$\frac{dx_i}{dt}(t) = \Gamma_{\theta(t)}[\mu(t)](x_i(t))$$

Coupled
EDOs

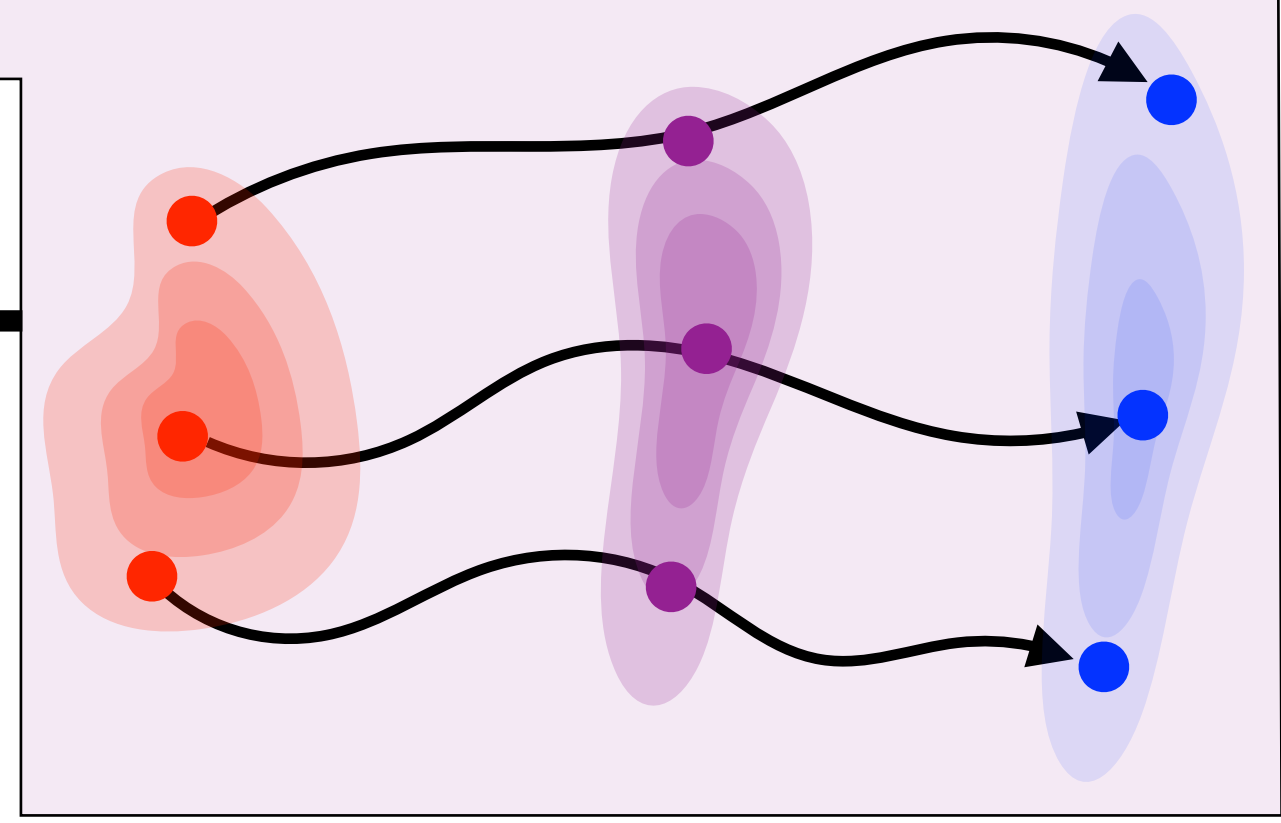


$$\frac{d\mu}{dt} + \text{div}(\mu \Gamma_{\theta}[\mu]) = 0 \quad \text{Non-linear PDE}$$

→ Not a Wasserstein flow :(
[Sander, Ablin, Blondel, Peyré, 2022]
[Geshkovski, Letrouit, Polyanskiy, Rigollet 2023]



Mean
field



Infinite Depth as a Neural PDE

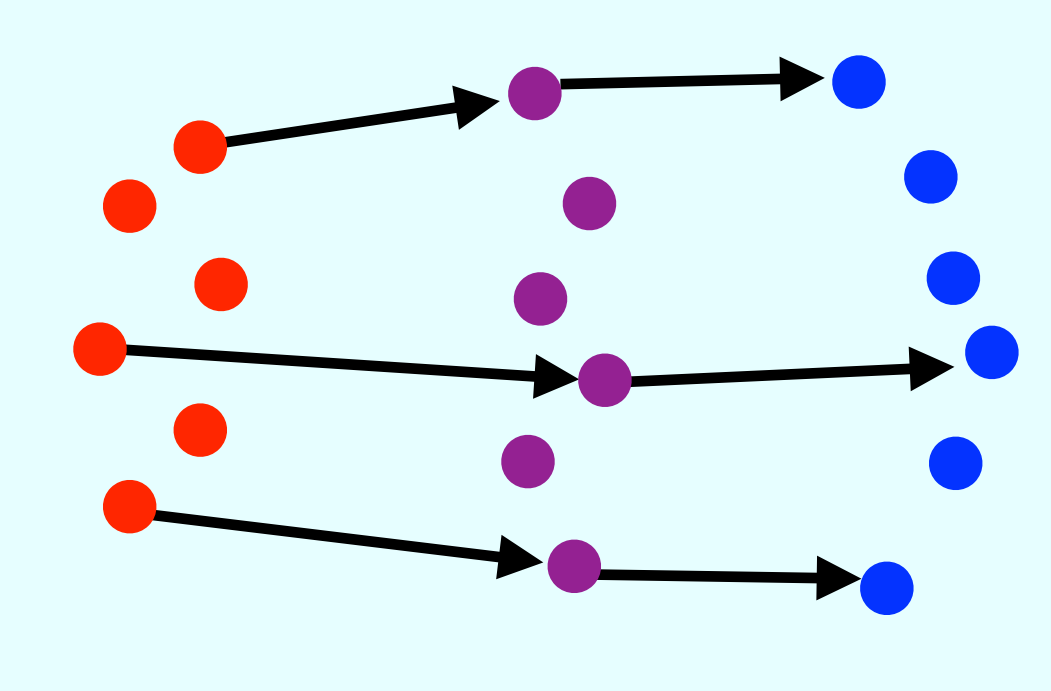
$$\Gamma_{\theta}[\mu](x) := \int \frac{e^{\langle Kx, Qy \rangle}}{\int e^{\langle Kx, Qy' \rangle} d\mu(y')} Vy d\mu(y) \quad \theta = (Q, K, V) \quad \mu(t) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$$

$$x_i(t+1) = x_i(t) + \frac{1}{T} \Gamma_{\theta(t)}[\mu(t)](x_i(t))$$

$T \rightarrow +\infty$
Infinite depth

$$\frac{dx_i}{dt}(t) = \Gamma_{\theta(t)}[\mu(t)](x_i(t))$$

Coupled EDOs

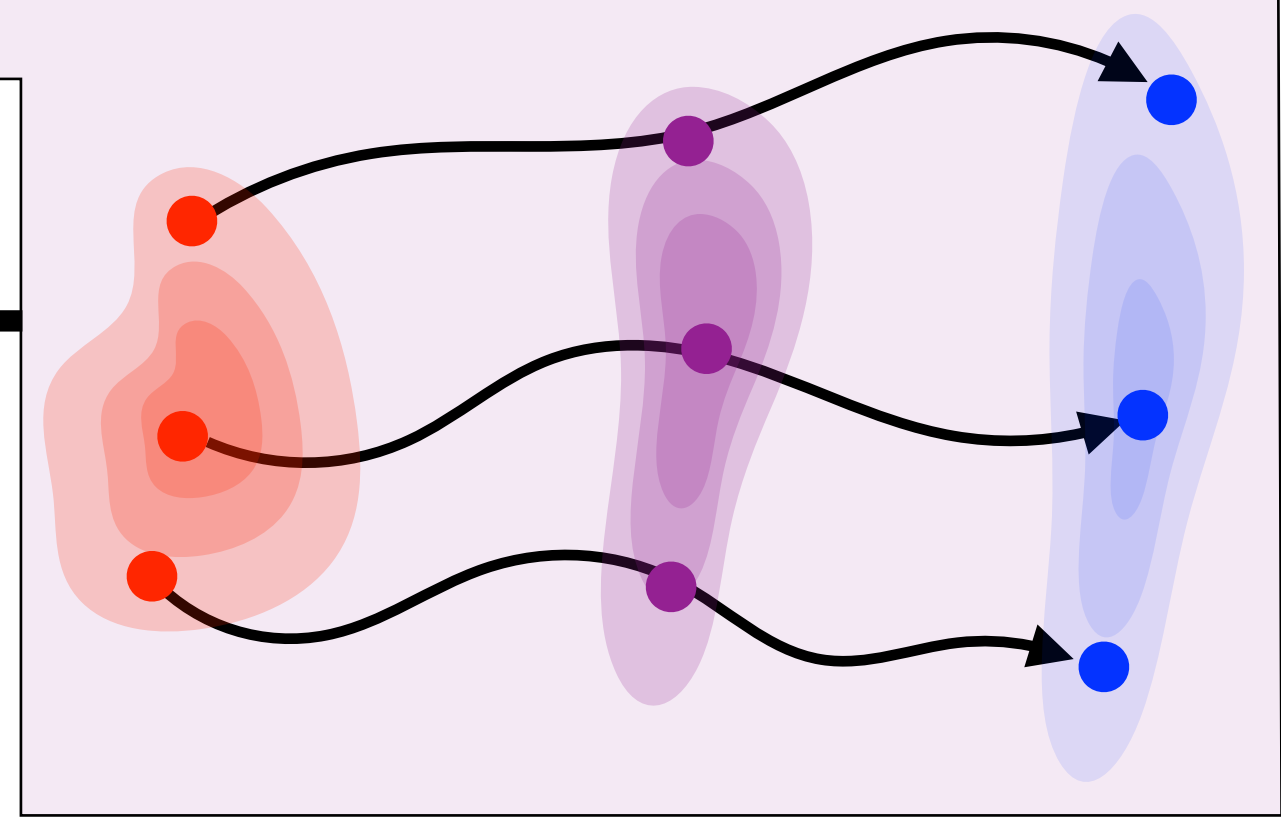


$$\frac{d\mu}{dt} + \text{div}(\mu \Gamma_{\theta}[\mu]) = 0 \quad \text{Non-linear PDE}$$

→ Not a Wasserstein flow :(
[Sander, Ablin, Blondel, Peyré, 2022]
[Geshkovski, Letrouit, Polyanskiy, Rigollet 2023]



Mean field



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

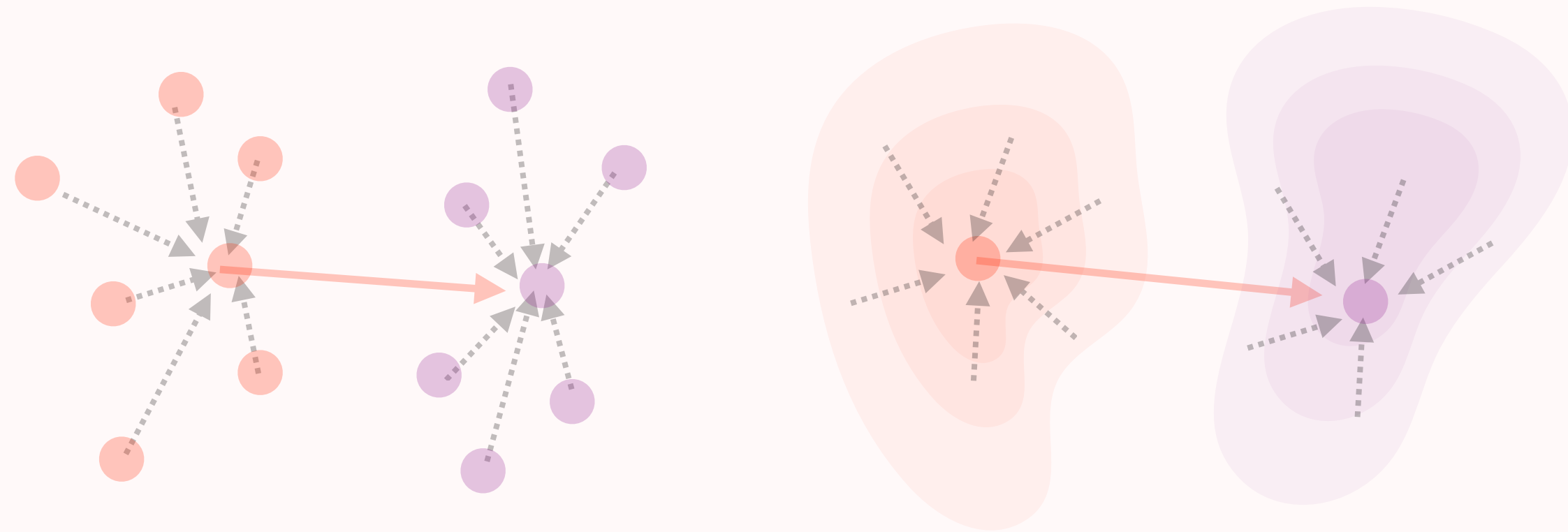
Training: $\min_{\theta} \sum_k \ell(T_{\theta}[\mu^k](x^k), y^k)$

Context Previous Next

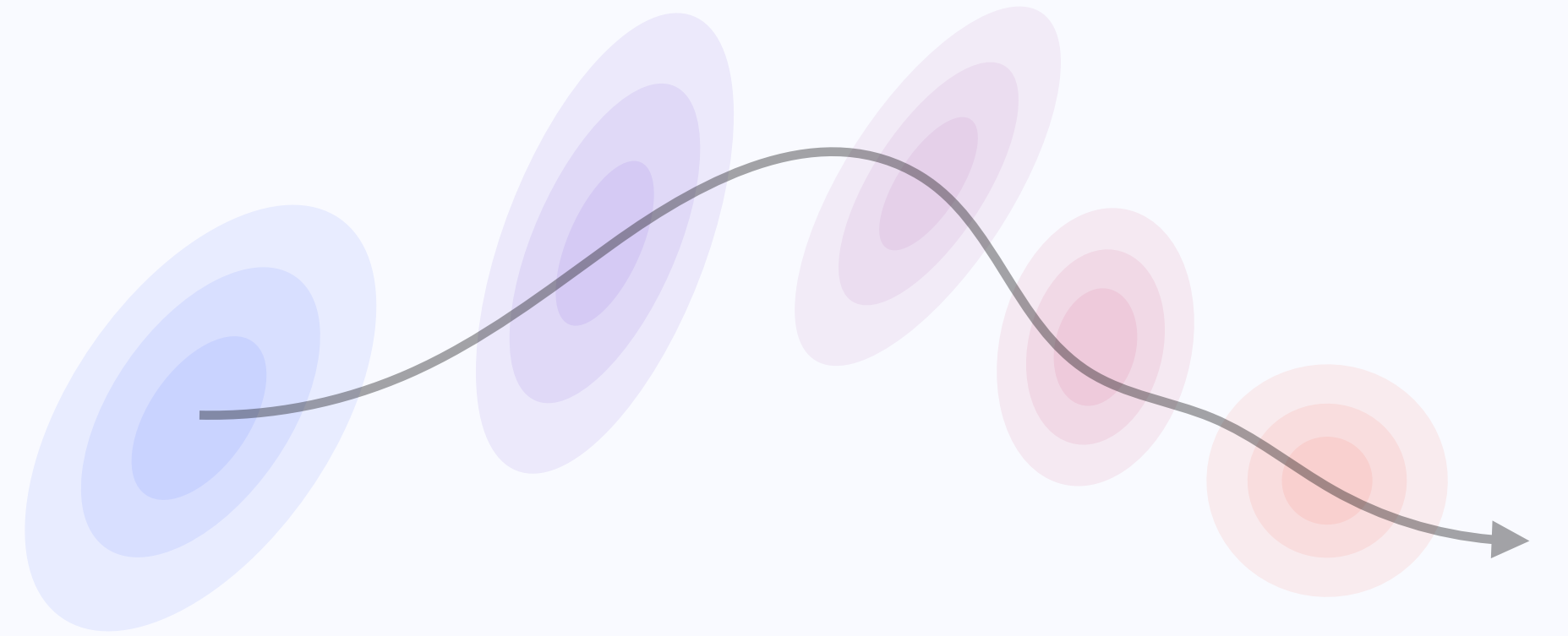
« Theorem » convergence to the global minimum if

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

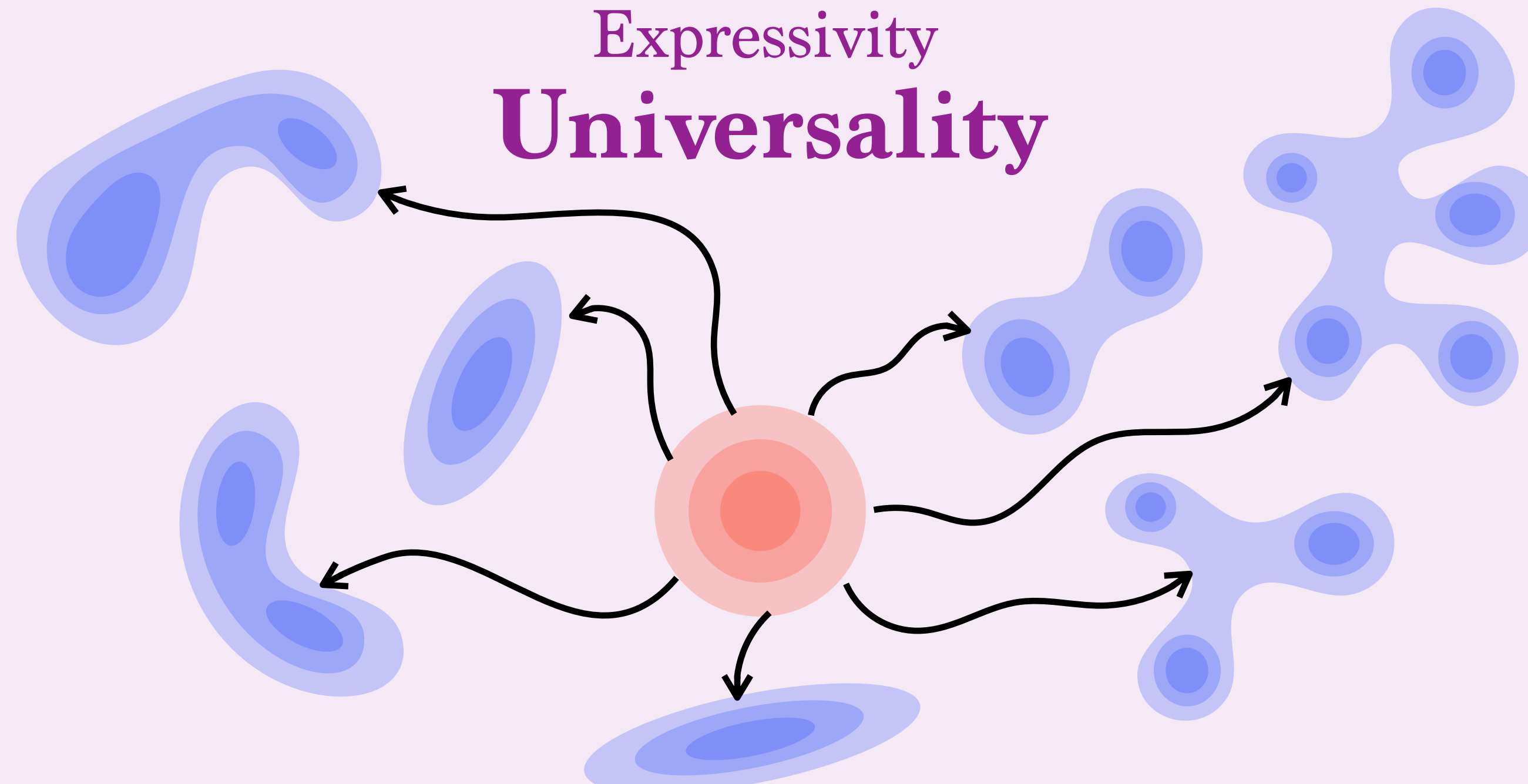
Arbitrary number of layers
**In Context Mappings
over Measures**



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



Expressivity
Universality



Universality

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

Theorem [Furuya, de Hoop, Peyré]:

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

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

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

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

Novelties:

fixed dimensions,
arbitrary # tokens.

Masked transformers:
requires Lipschitz
in time.

Universality

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

Theorem [Furuya, de Hoop, Peyré]:

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

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

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

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 #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 • MLP with skip connexion.

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

Sketch of proof

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

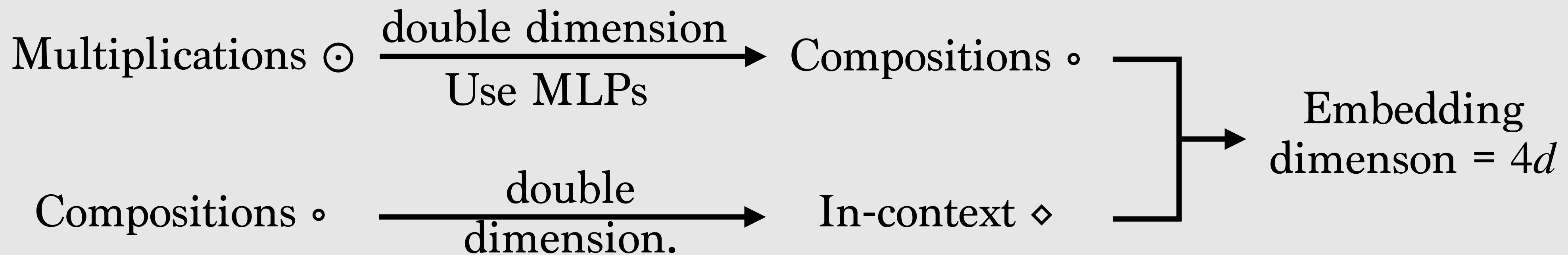
→ First component of Attention ◦ MLP with skip connexion.

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

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

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

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

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.

γ_θ are continuous.

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

$$\forall \theta, \gamma_\theta[\mu](x) = \gamma_\theta[\mu'](x') \\ \stackrel{?}{\implies} \\ (\mu, x) = (\mu', x')$$



Marshall
Stone

Karl
Weierstrass

Stone-Weierstrass
theorem

Sketch of Proof

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

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.

γ_θ are continuous.

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

$\forall \theta, \gamma_\theta[\mu](x) = \gamma_\theta[\mu'](x')$
 $\stackrel{?}{\implies}$
 $(\mu, x) = (\mu', x')$

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



Stone-Weierstrass
theorem

Sketch of Proof

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

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.

γ_θ are continuous.

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

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

$\stackrel{?}{\implies}$

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

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

$$A = c = u = 0: \quad L_1(\mu)(b) = L_1(\mu')(b)$$

$$\text{In 1-D:} \quad L_k(\mu)(b) := \int \frac{e^{by} y^k v}{\int e^{by'} d\mu(y')} d\mu(y)$$

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

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

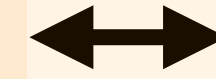
In higher dimensions: use Radon transform.

Stone-Weierstrass theorem

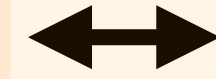
Open Problems

Smoothness: bridge the gap

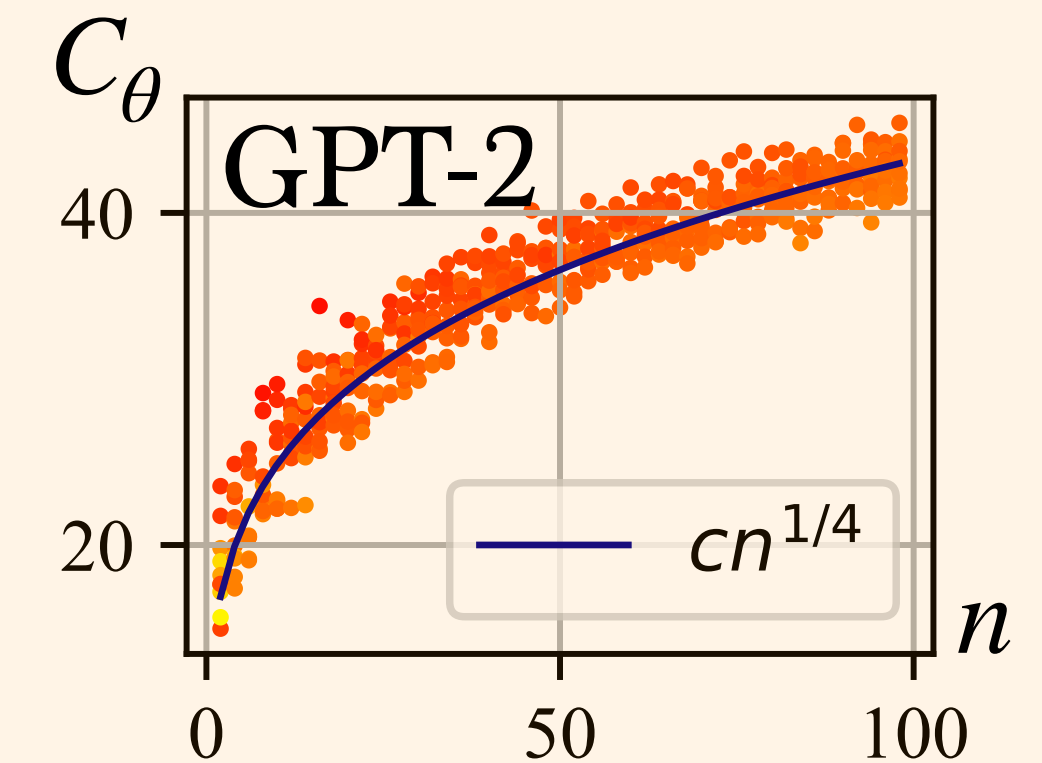
mean-field
 e^R



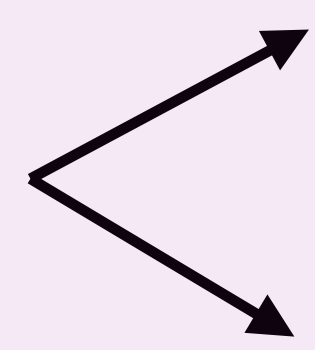
discrete
 \sqrt{n}



practice
 $n^{1/4}$



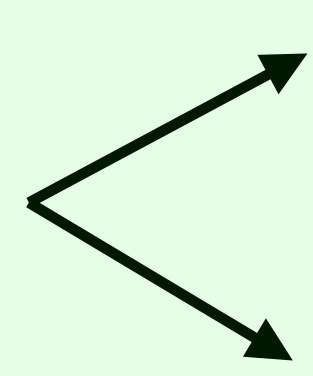
Universality:



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?