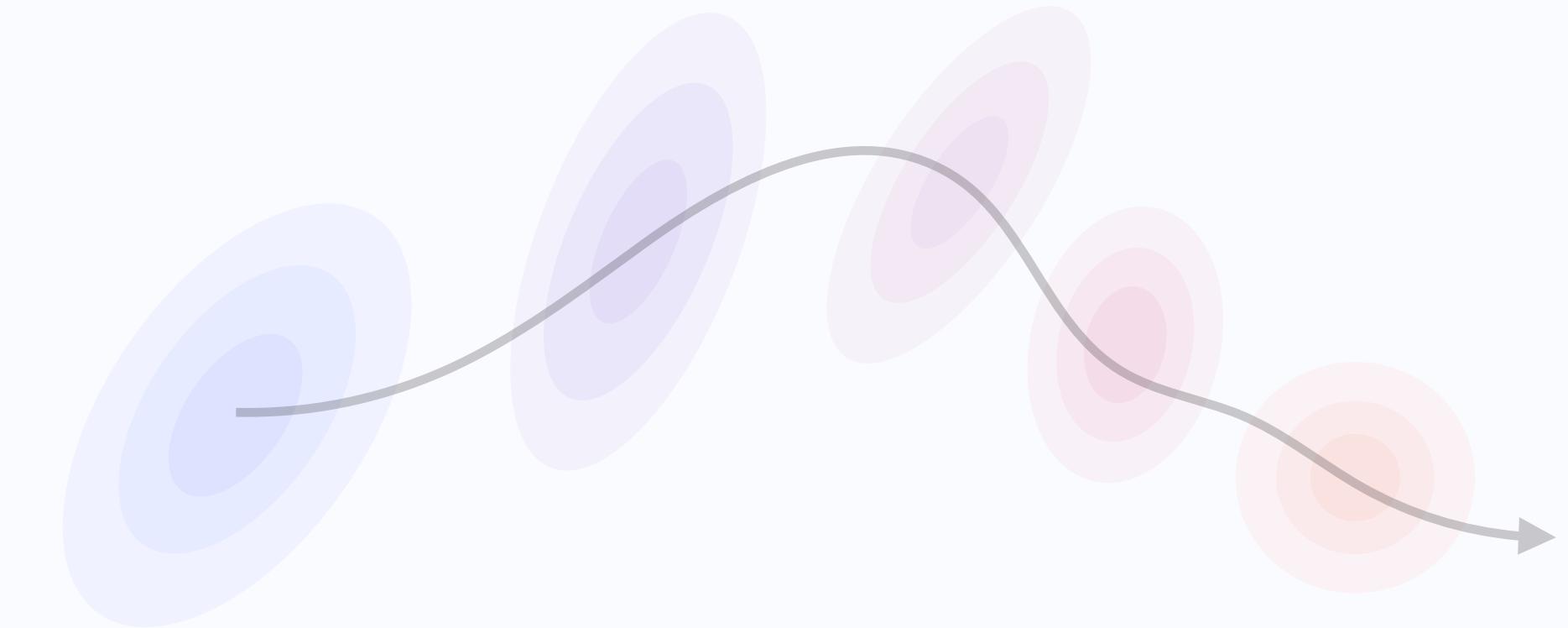
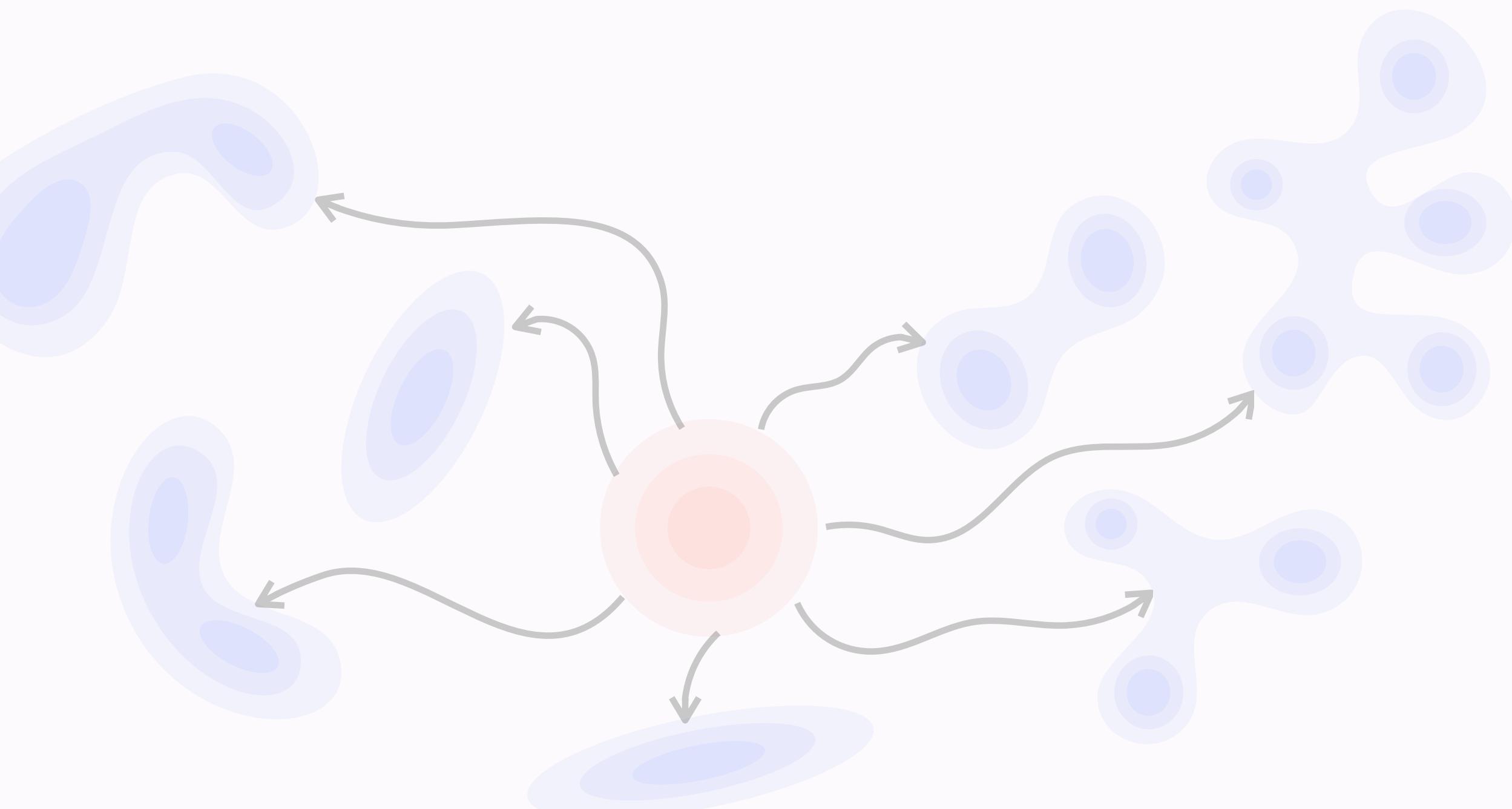
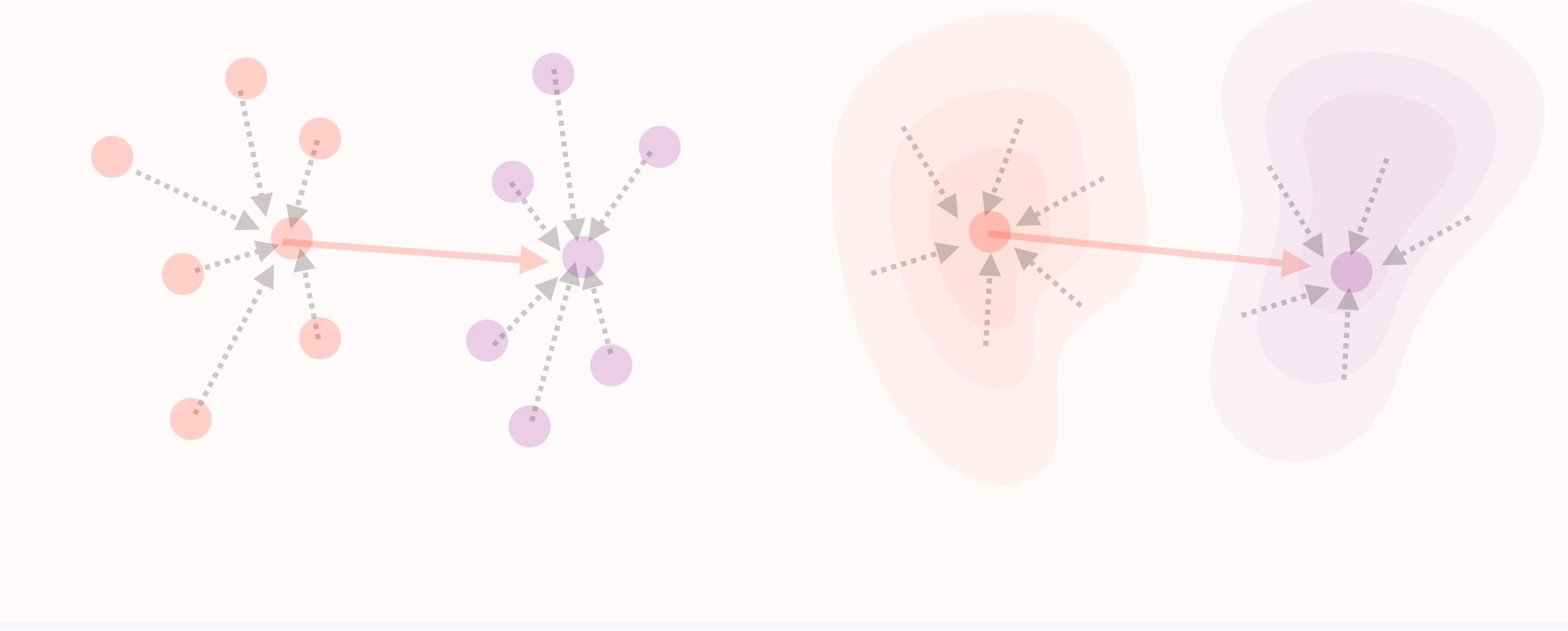
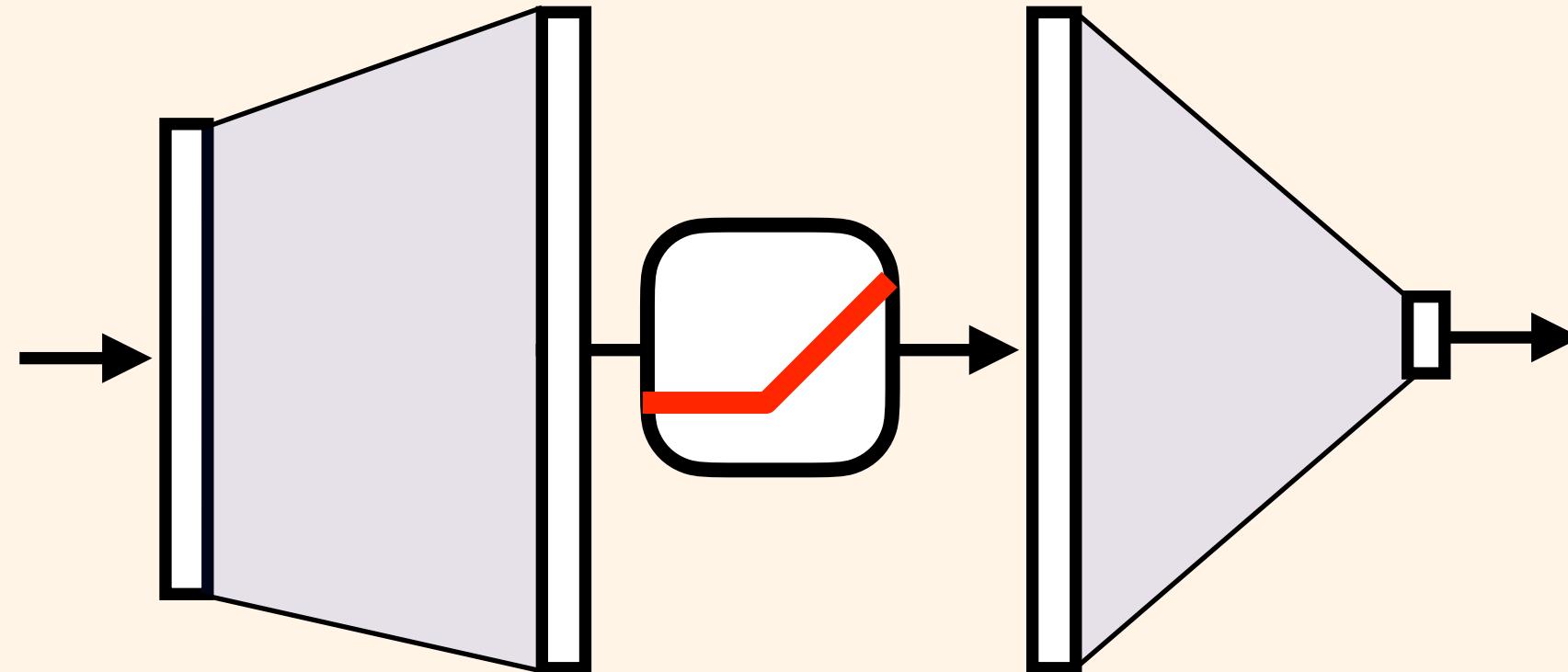


# Universal Approximation: *from Neural Networks to Transformers*

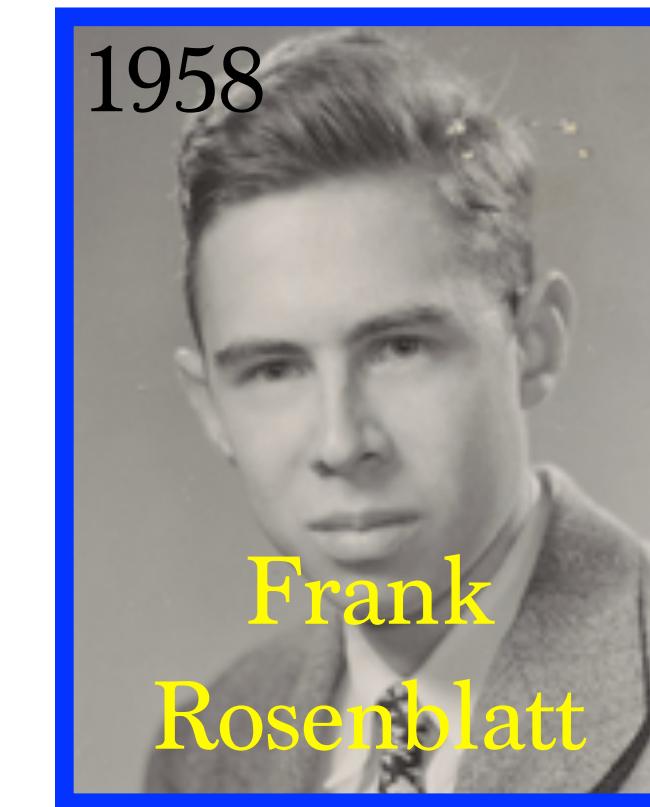
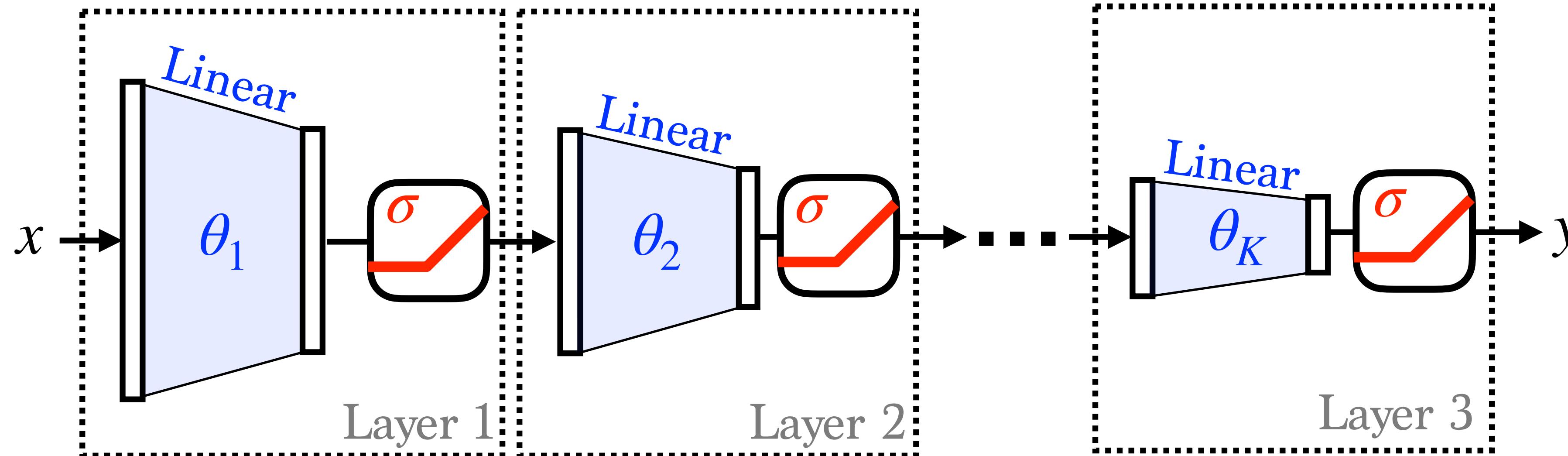
Gabriel Peyré



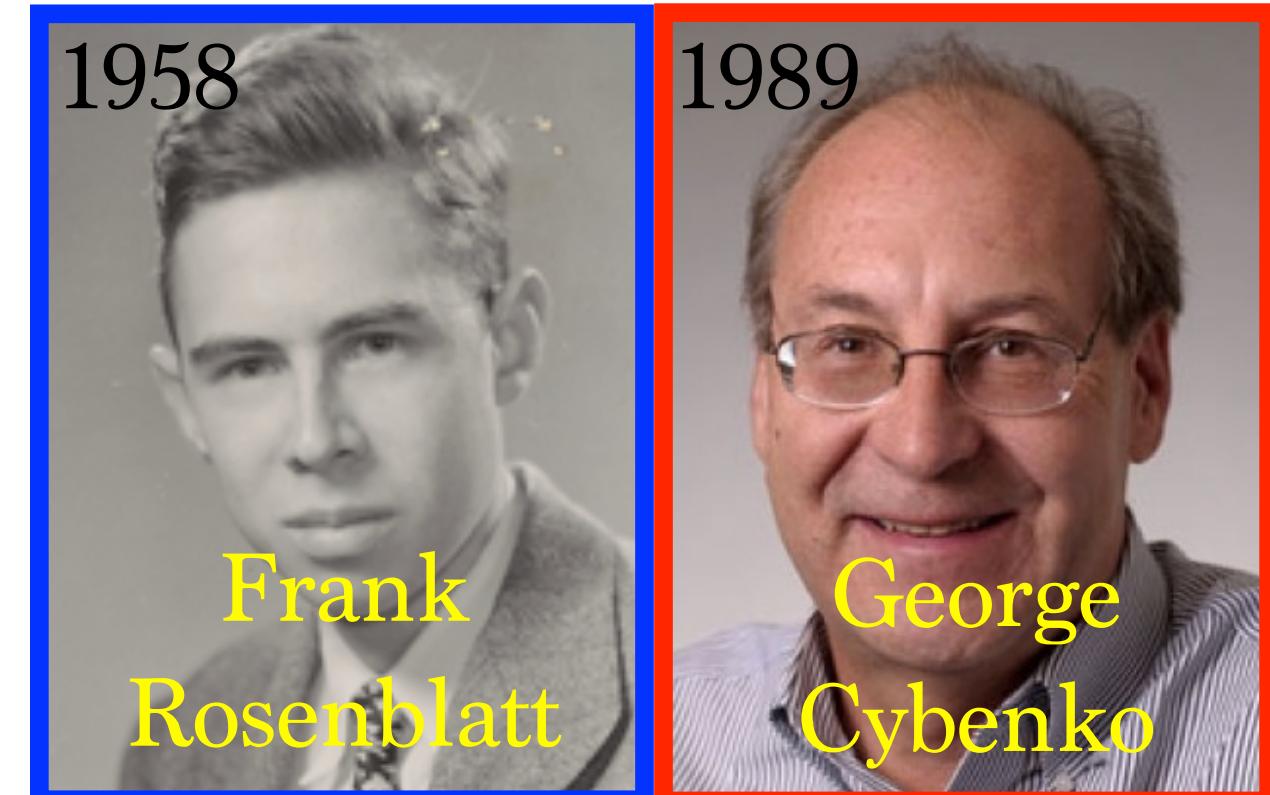
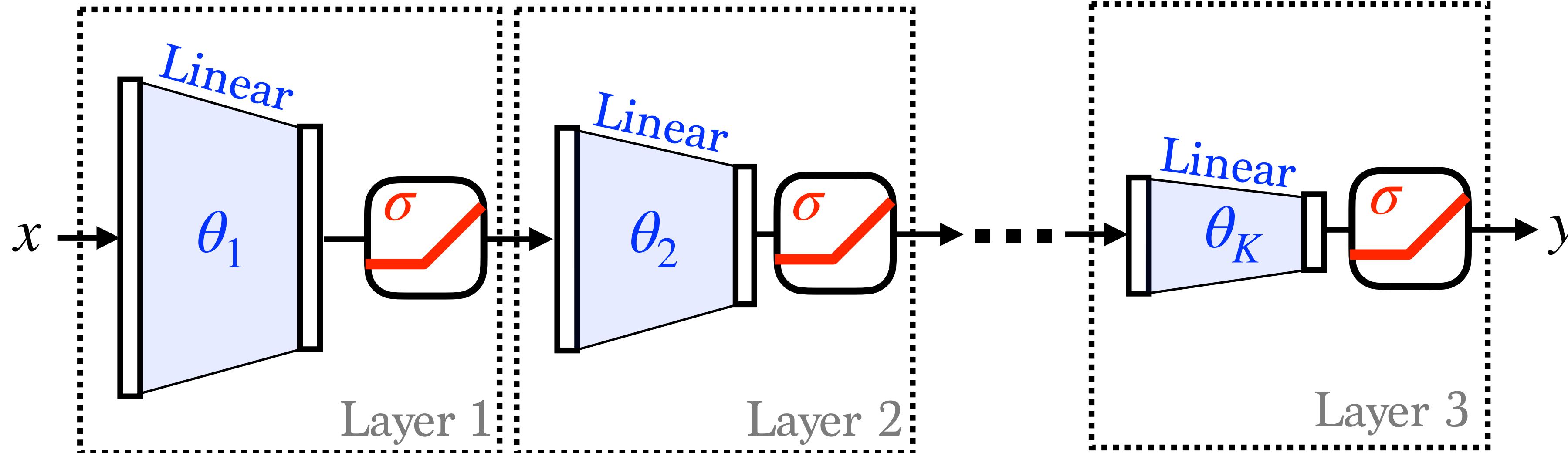
# Perceptrons: context-free universality



# Perceptrons and Universality



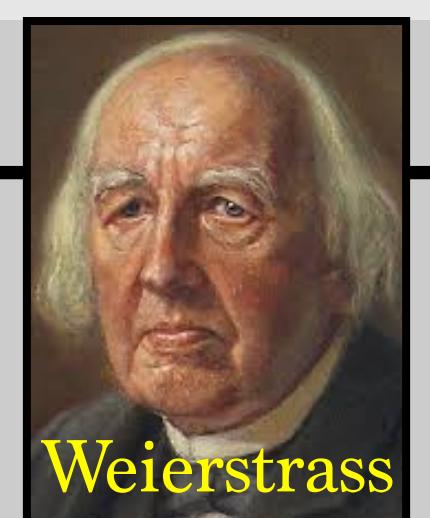
# Perceptrons and Universality



**Two layers:**  $\Gamma_\theta(x) := \sum_{i=1}^n a_i \sigma(\langle x, w_i \rangle + b_i)$

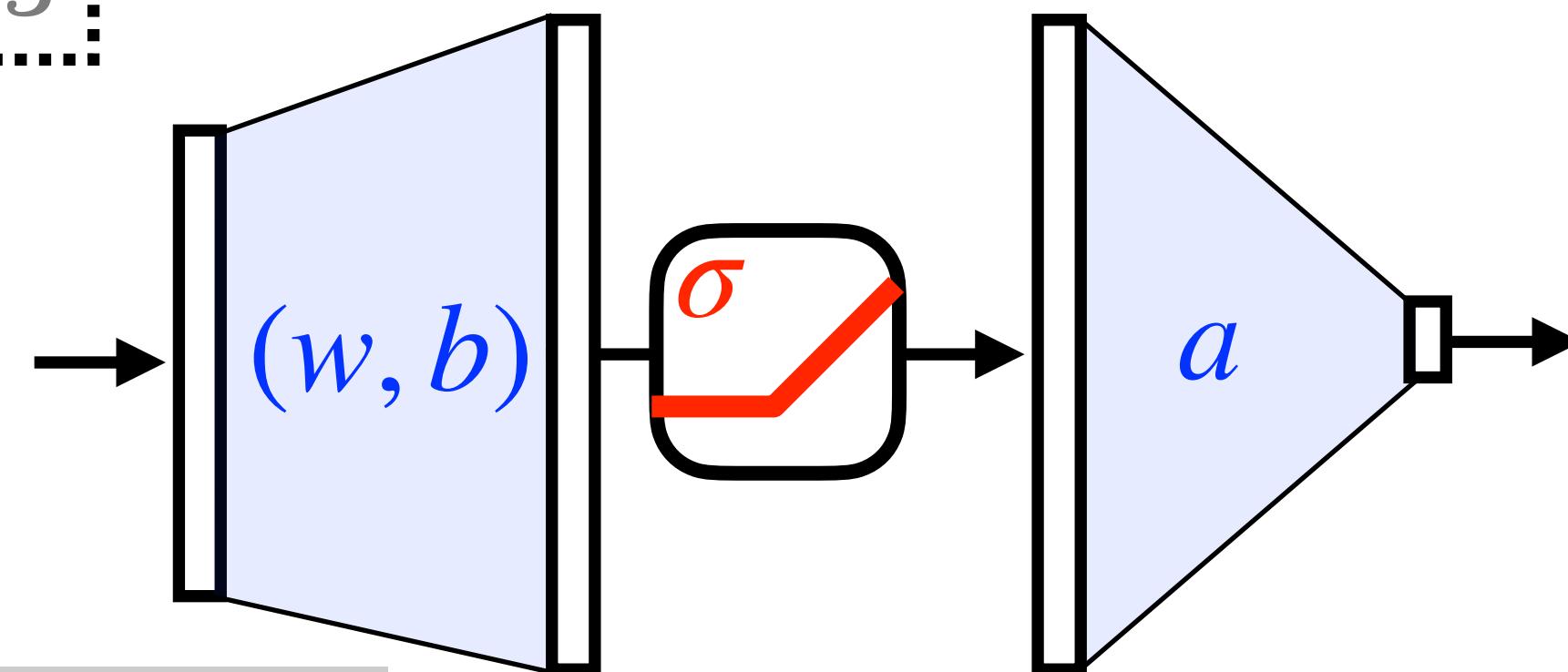
*Theorem:* for sigmoid and ReLu  $\sigma$ , on a compact  $\Omega$ ,  
 $\{\Gamma_\theta\}_\theta$  is dense in  $(\mathcal{C}(\Omega), L^\infty(\Omega))$ .

*Proof:* For  $\sigma = \sin$ ,  $\{\Gamma_\theta\}_\theta$  is an algebra

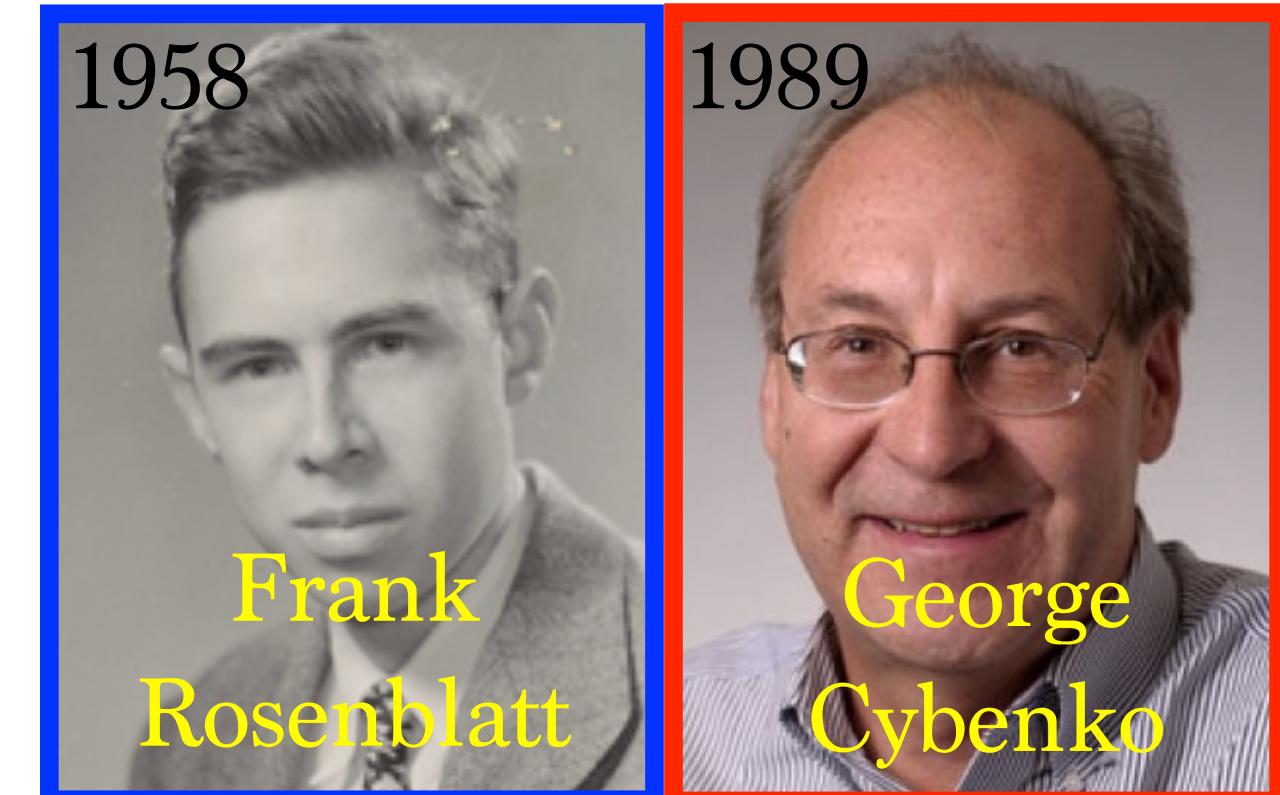
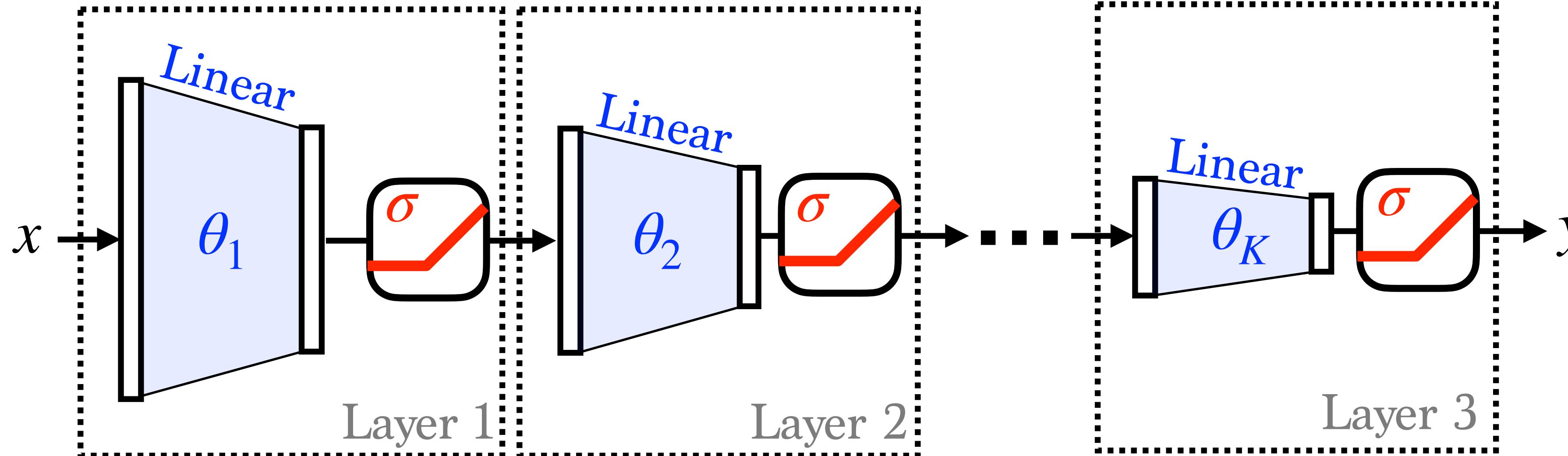


$\{\Gamma_\theta\}_\theta$  is dense

Approximate sin by generic  $\sigma$



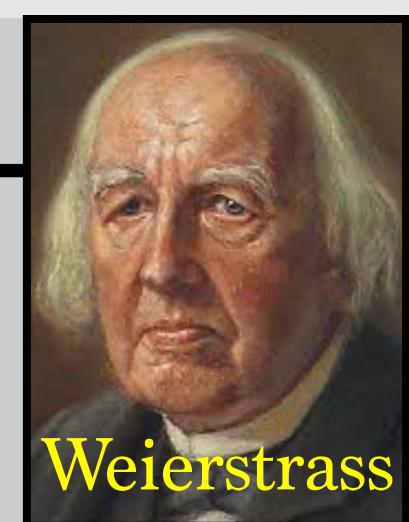
# Perceptrons and Universality



**Two layers:**  $\Gamma_\theta(x) := \sum_{i=1}^n a_i \sigma(\langle x, w_i \rangle + b_i)$

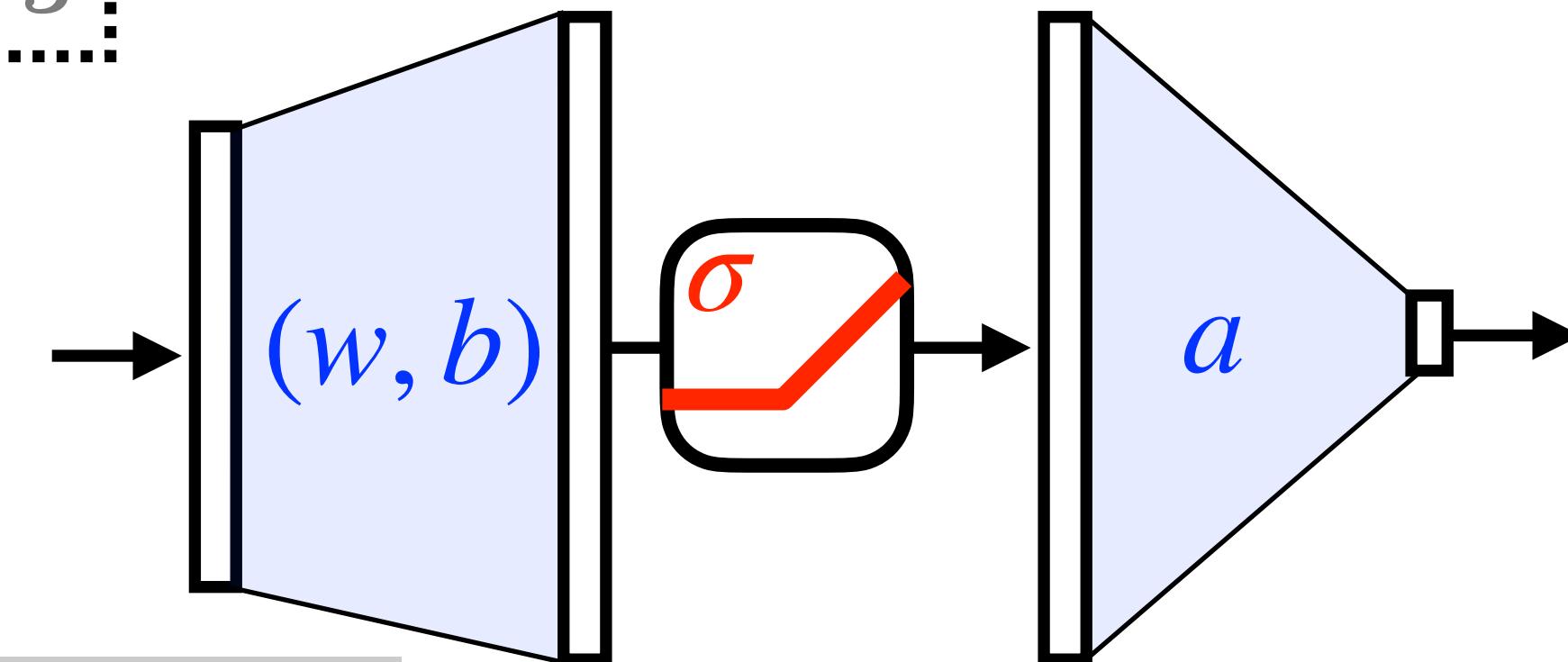
*Theorem:* for sigmoid and ReLu  $\sigma$ , on a compact  $\Omega$ ,  
 $\{\Gamma_\theta\}_\theta$  is dense in  $(\mathcal{C}(\Omega), L^\infty(\Omega))$ .

*Proof:* For  $\sigma = \sin$ ,  $\{\Gamma_\theta\}_\theta$  is an algebra



$\{\Gamma_\theta\}_\theta$  is dense

Approximate sin by generic  $\sigma$



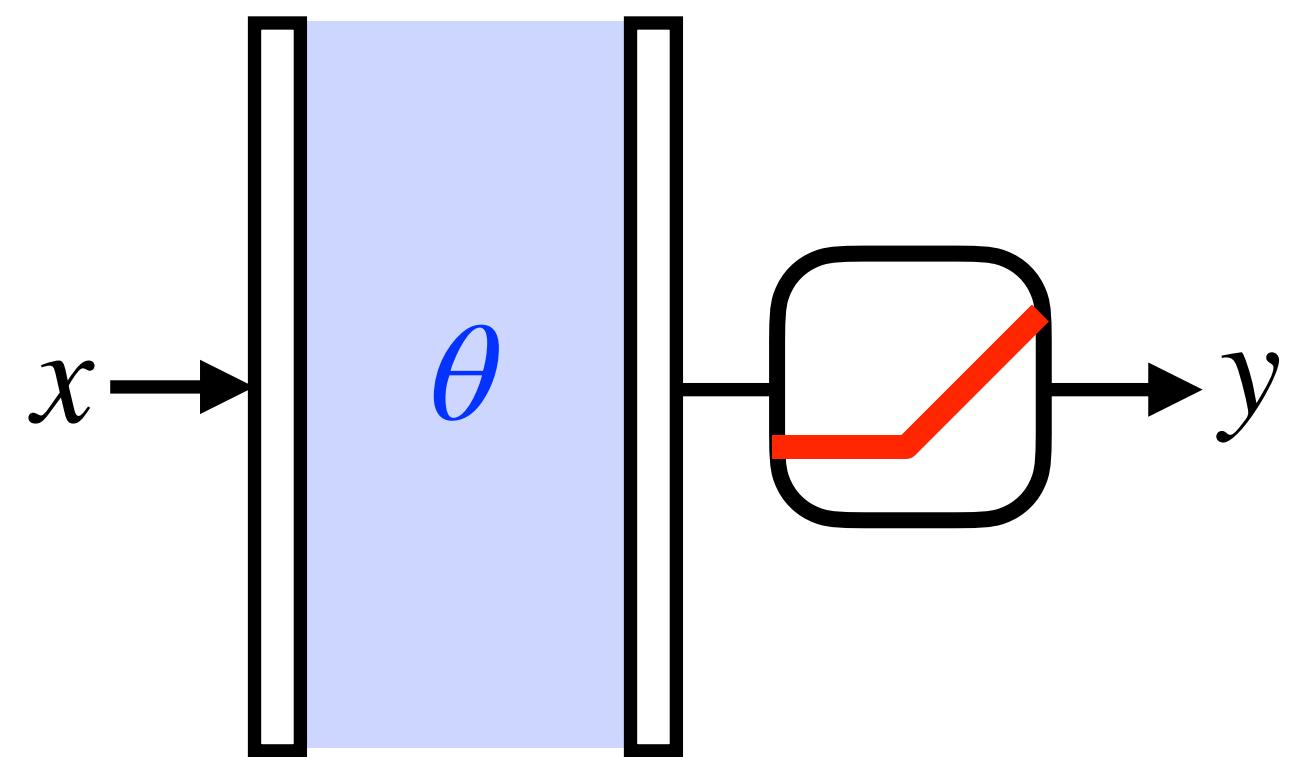
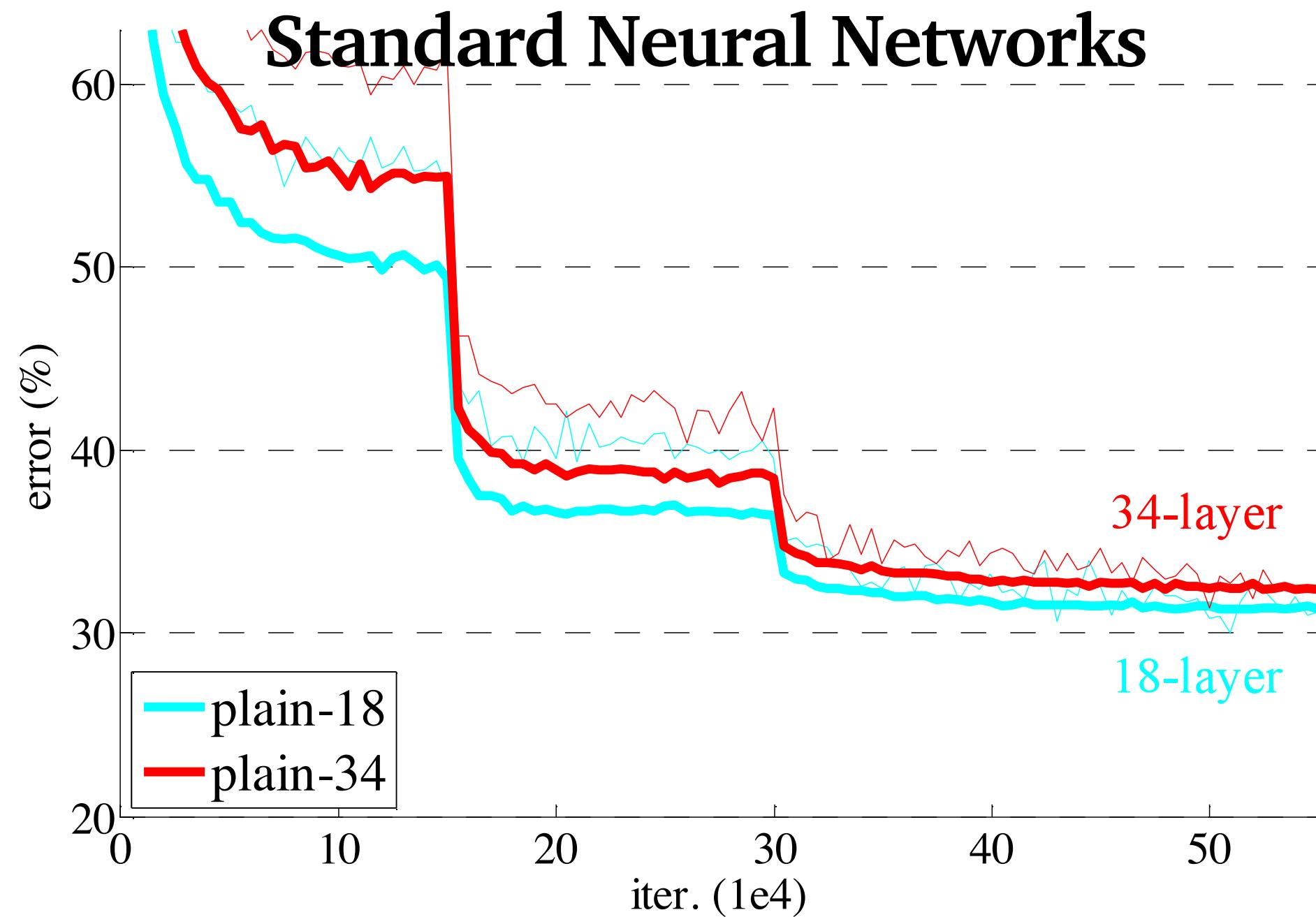
**Open problems**

→ Expressivity: role of depth?

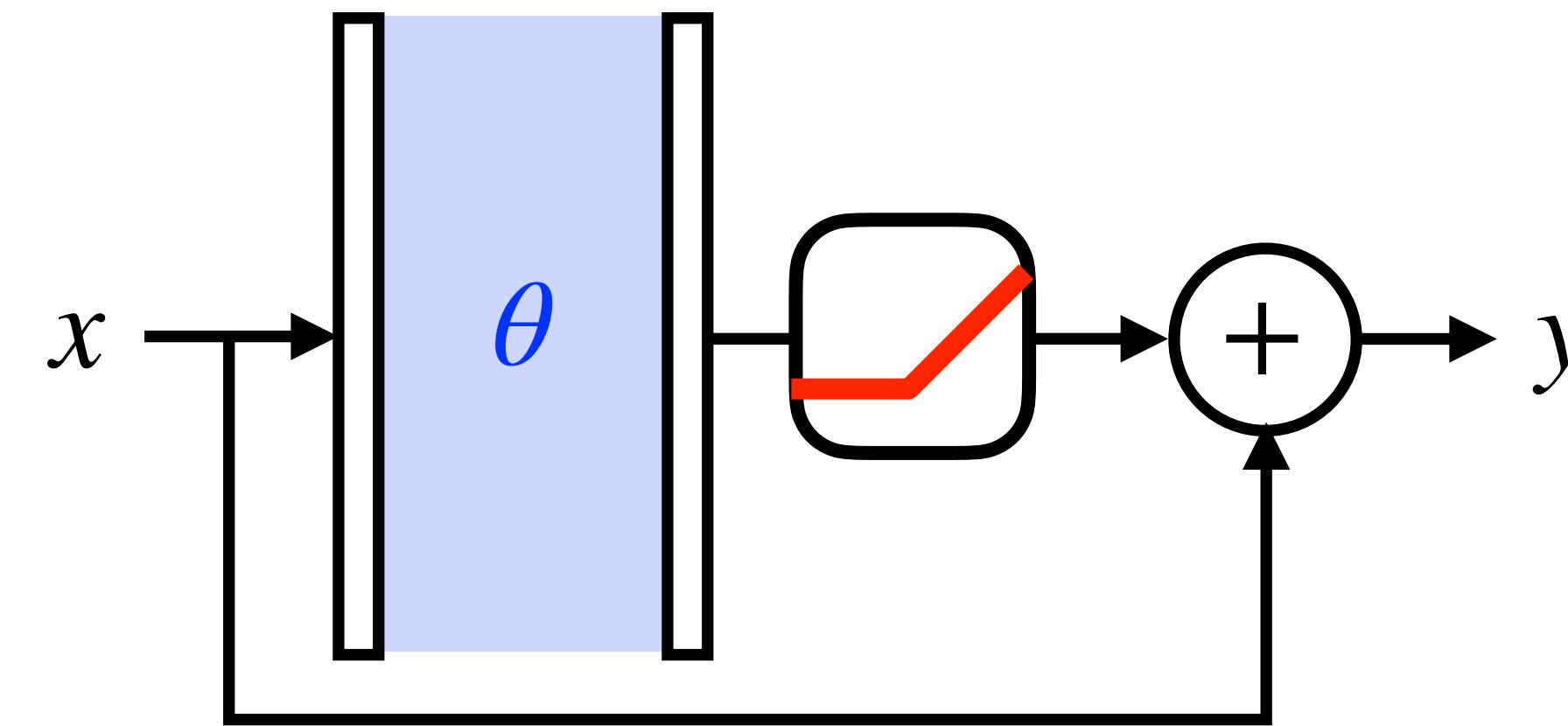
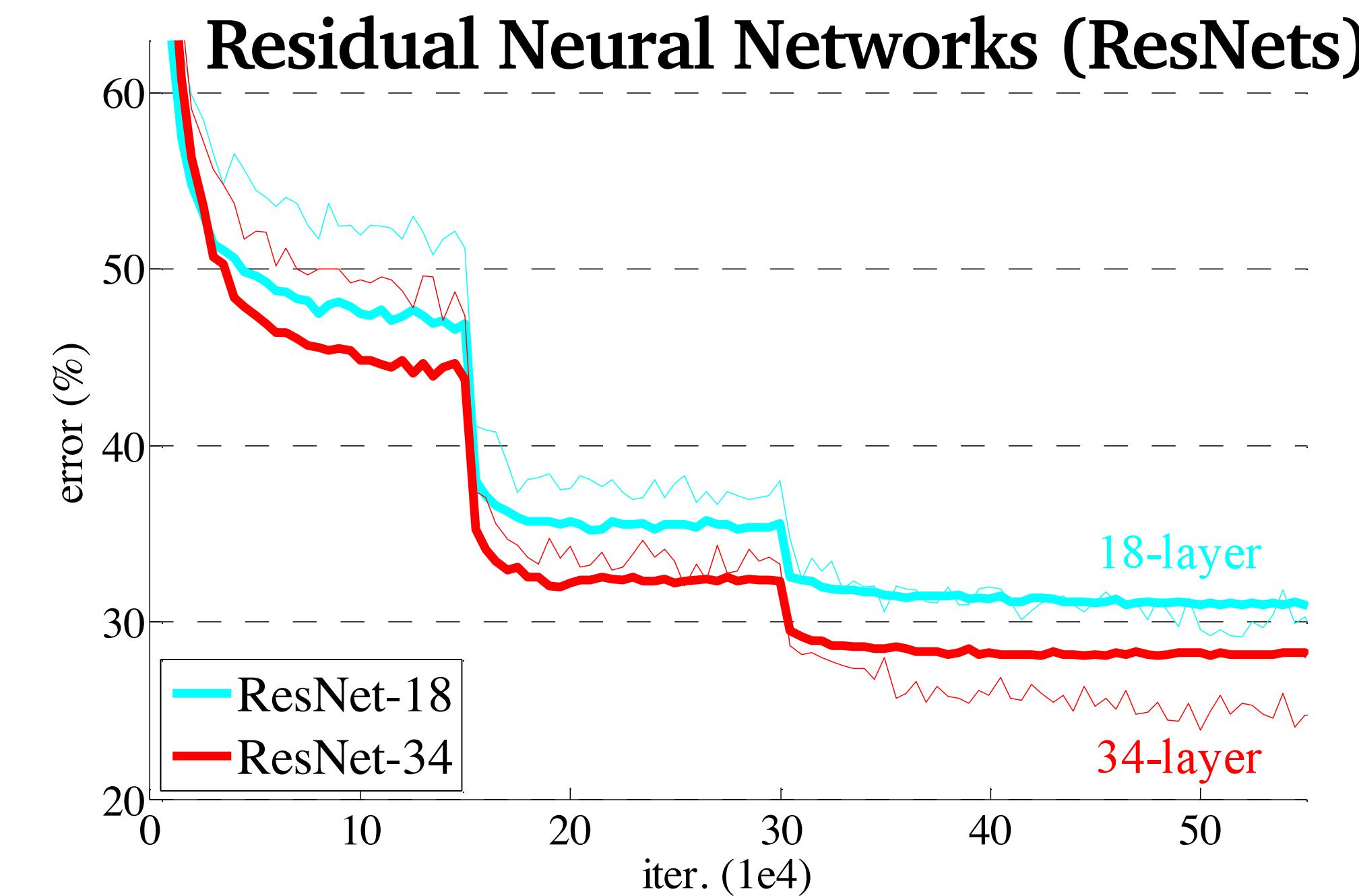
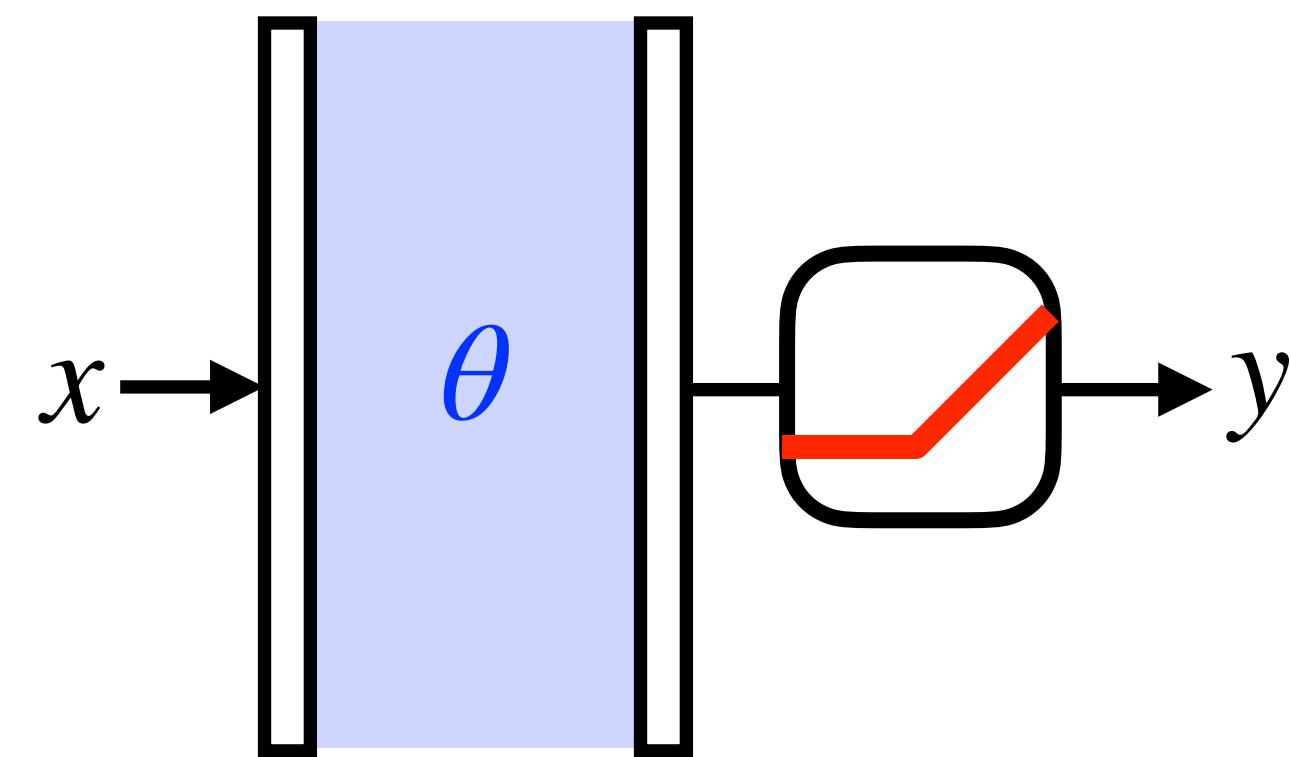
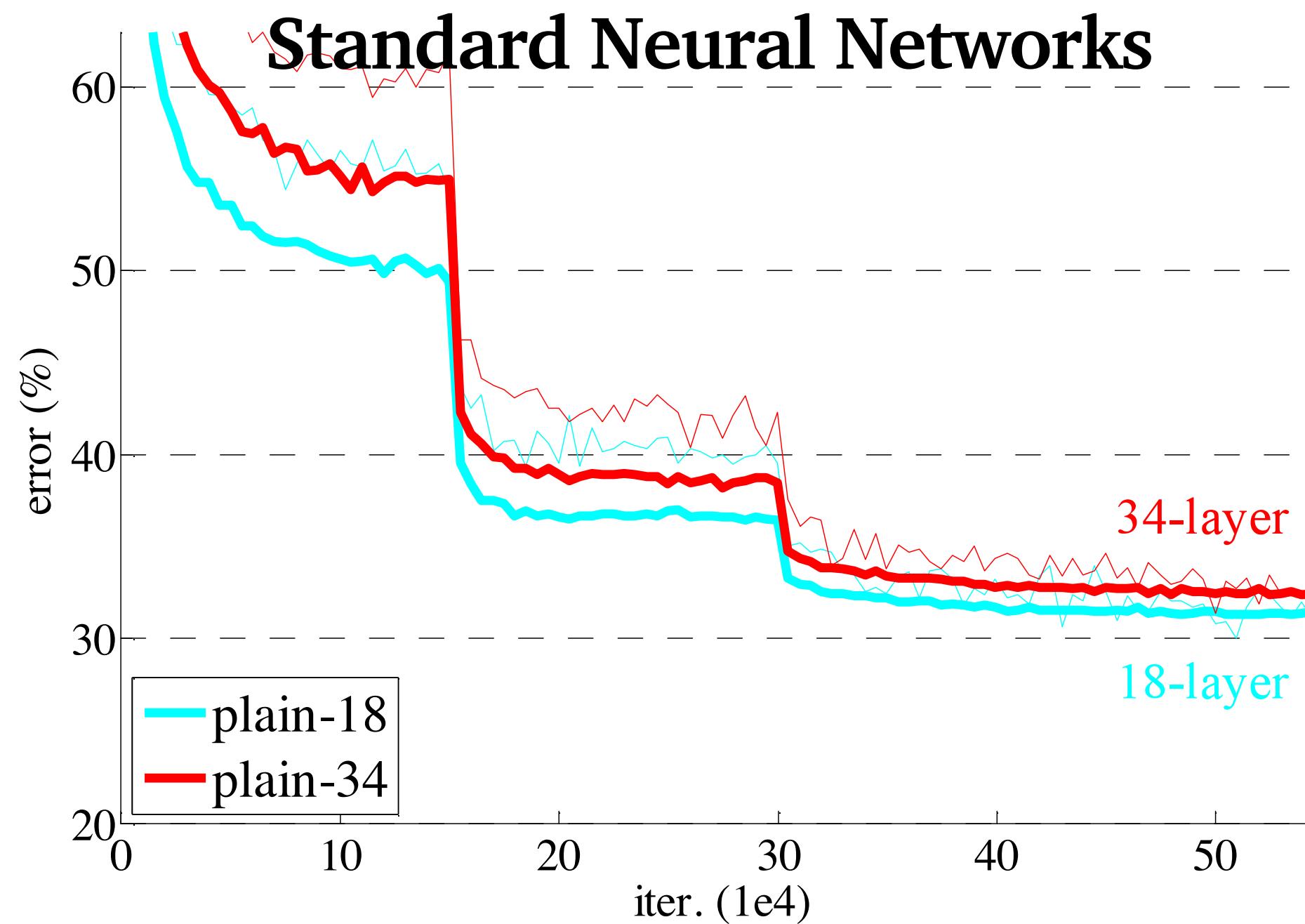
→ Optimization: convergence of gradient descent.



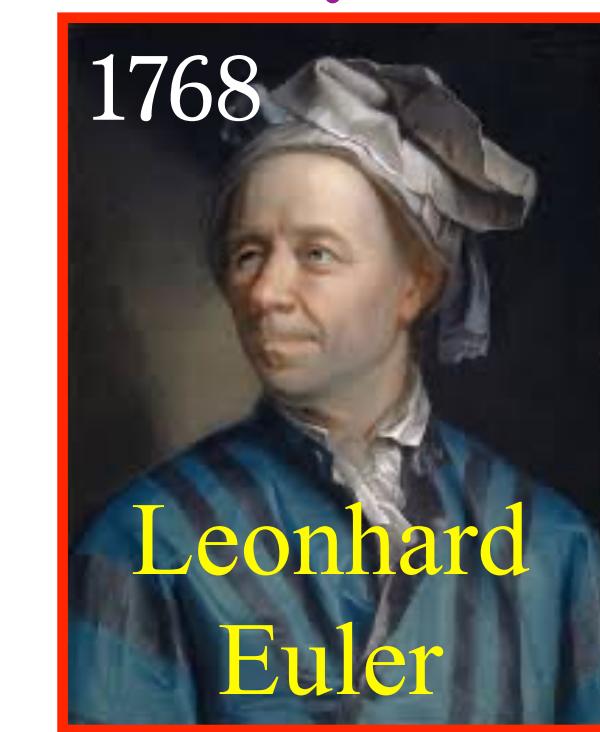
# The Deeper, the Better



# The Deeper, the Better



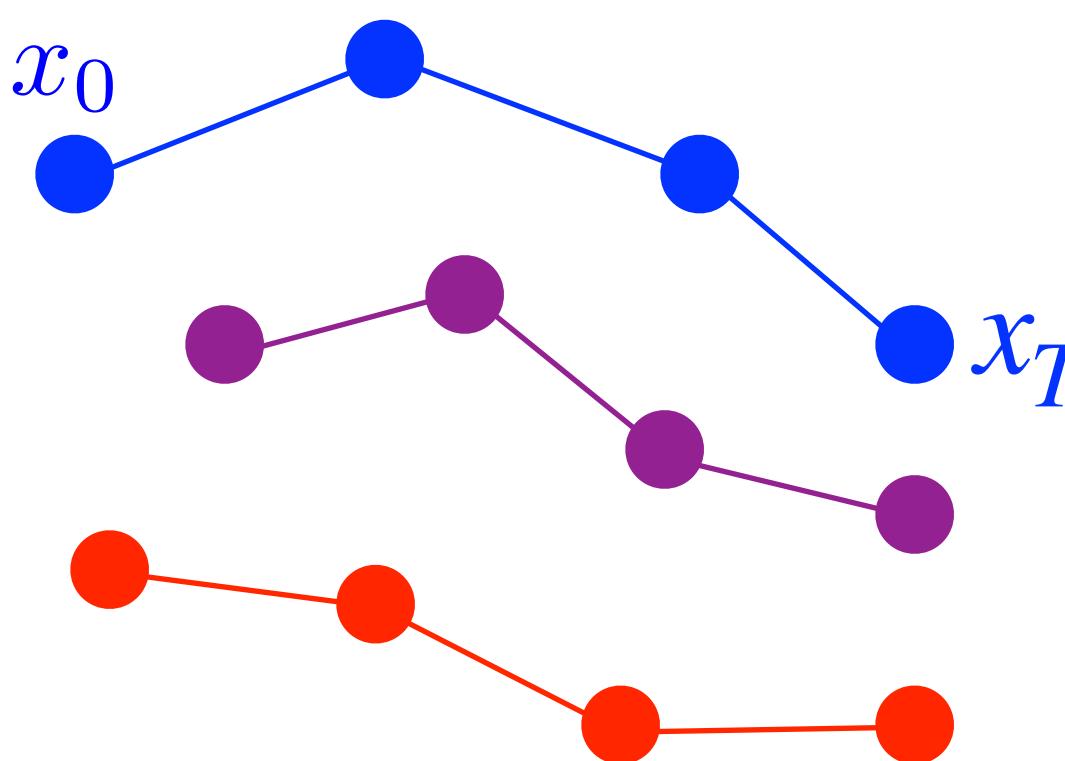
Infinite depth



Differential equations

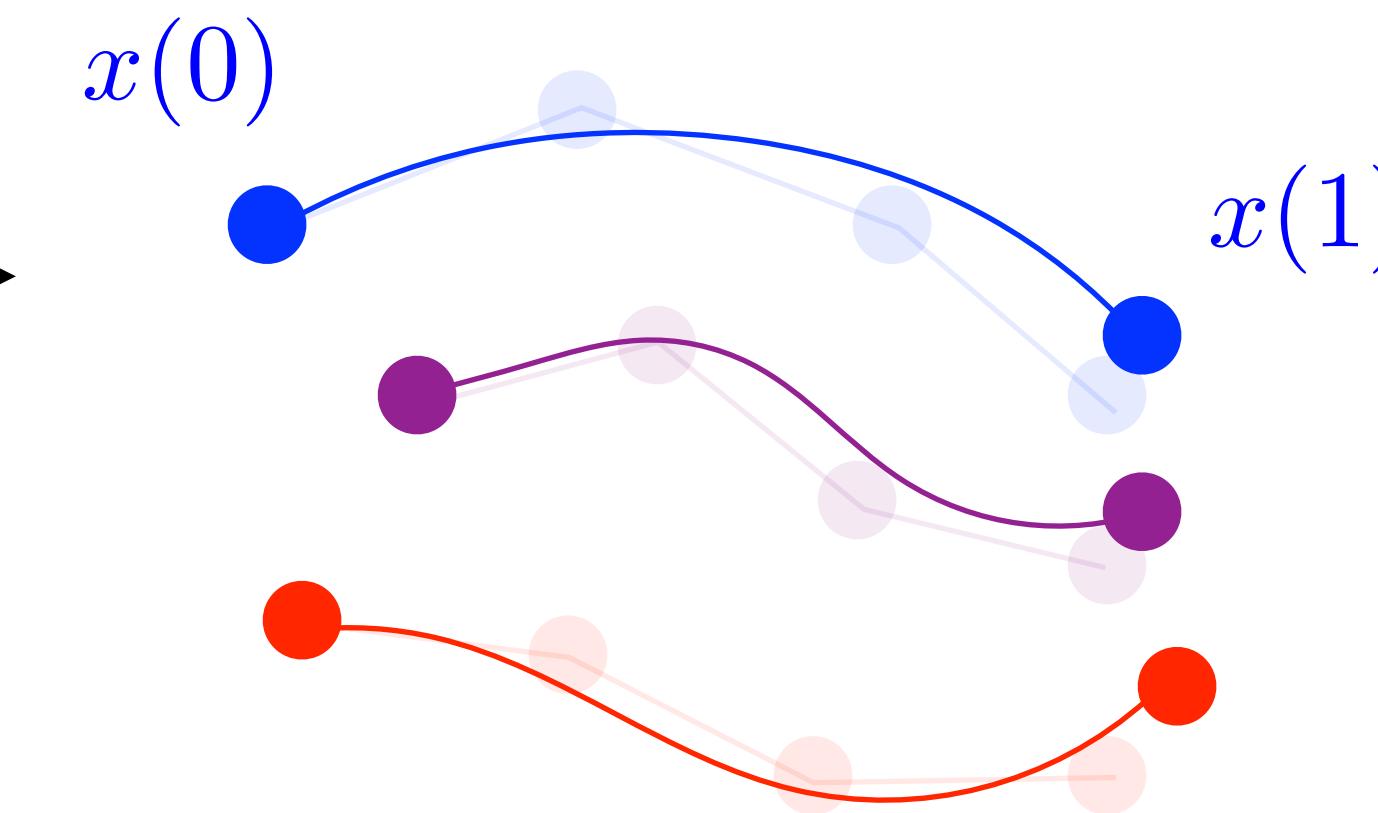
# Infinite Depth and Neural-ODEs

$T$  layers  
 $x_{t+1} = x_t + \frac{1}{T} \Gamma_{\theta_t}(x_t)$



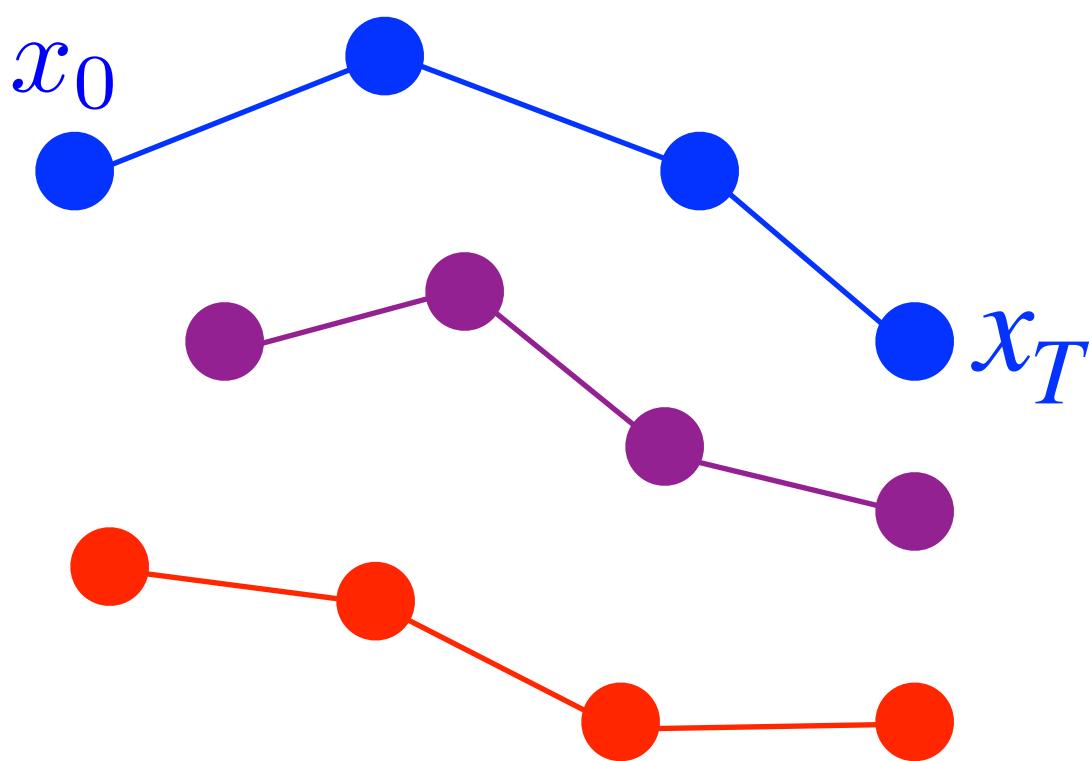
$T \rightarrow \infty$

$T = \infty$  layers  
 $\dot{x}_t = \Gamma_{\theta(t)}(x(t))$



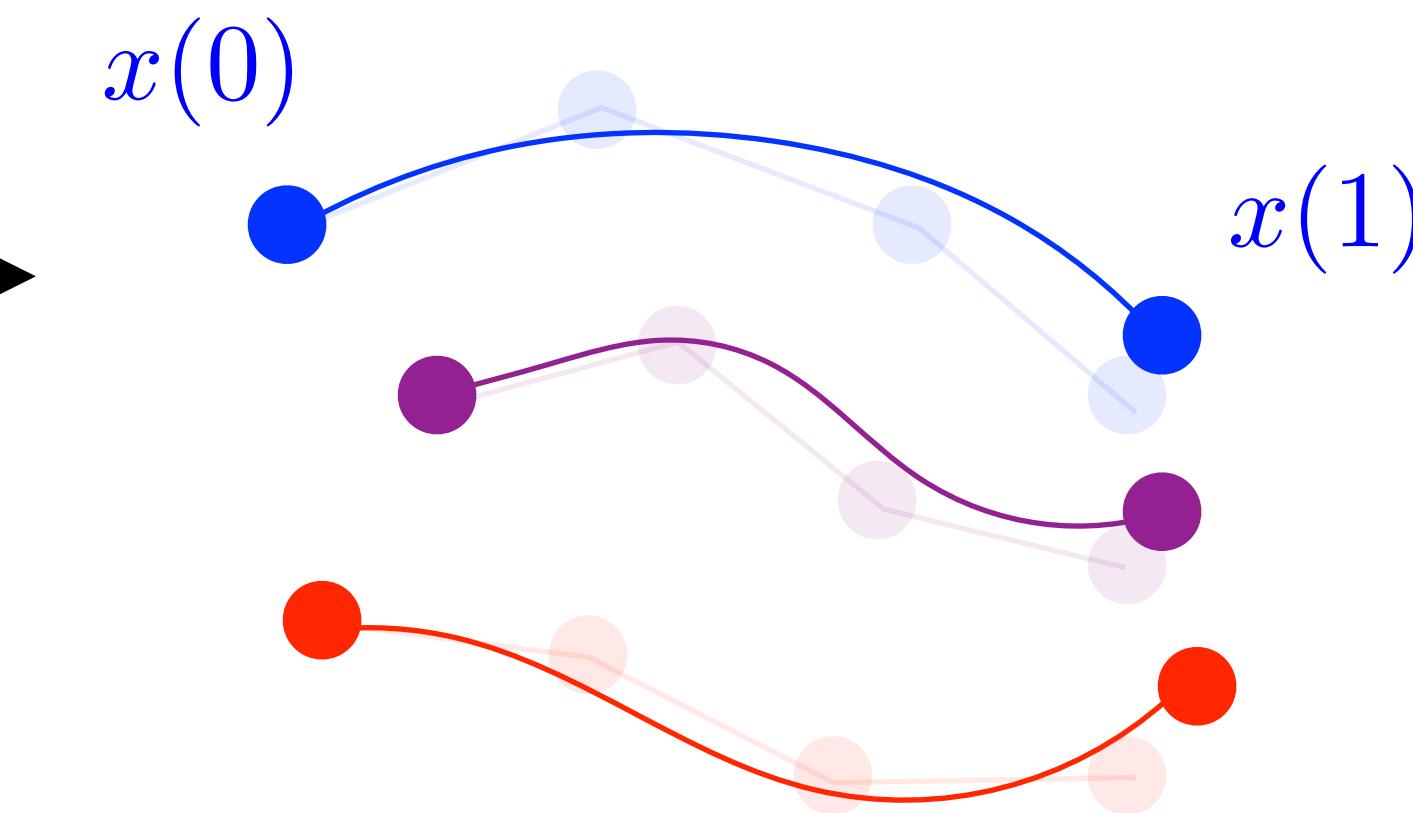
# Infinite Depth and Neural-ODEs

$T$  layers  
 $x_{t+1} = x_t + \frac{1}{T} \Gamma_{\theta_t}(x_t)$

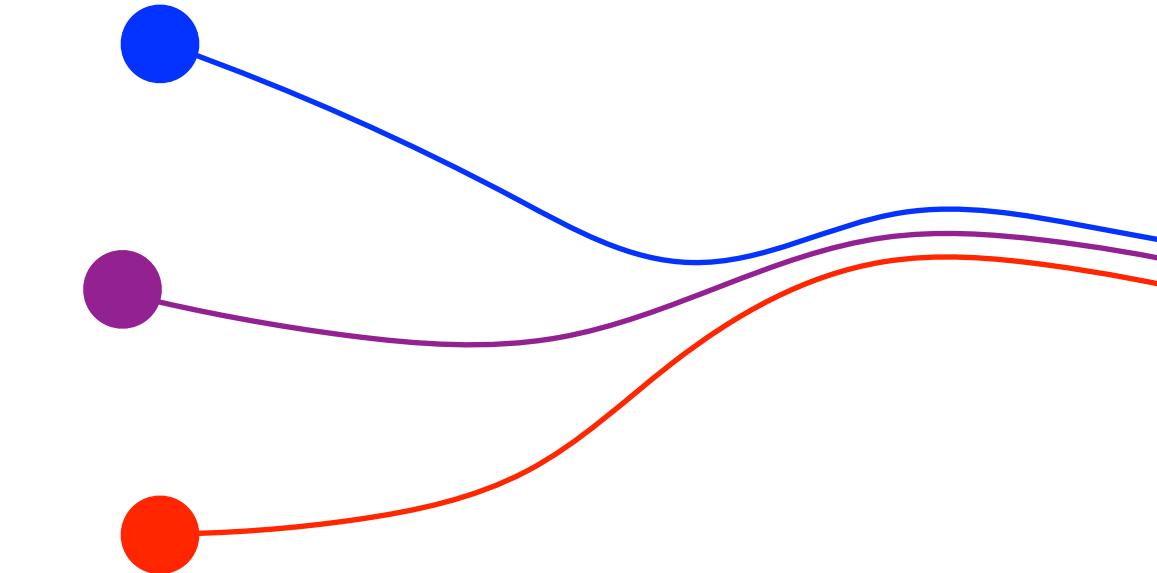
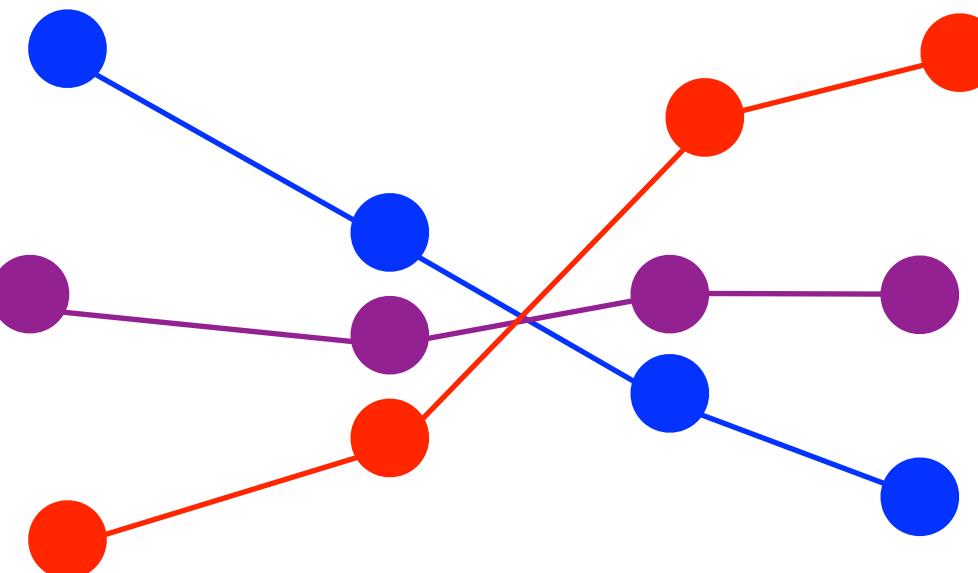


$$T \rightarrow \infty$$

$T = \infty$  layers  
 $\dot{x}_t = \Gamma_{\theta(t)}(x(t))$



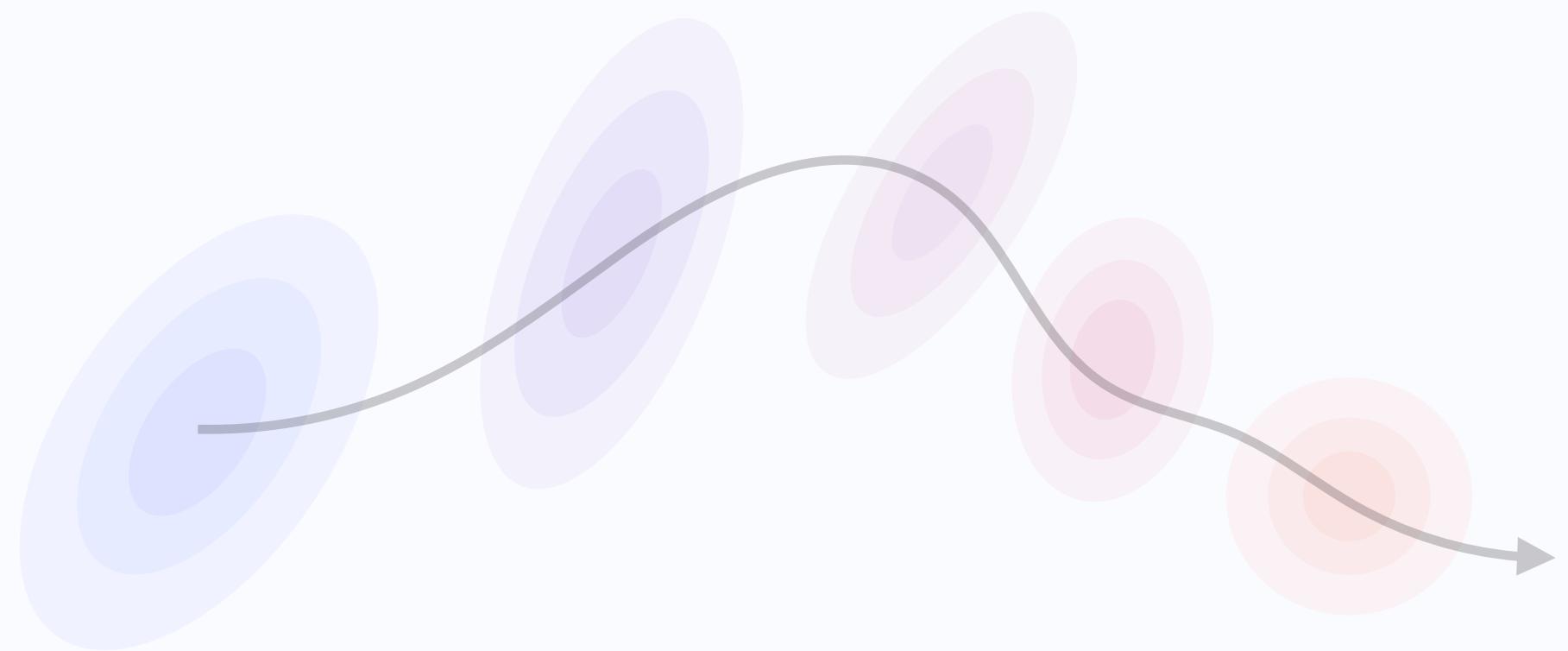
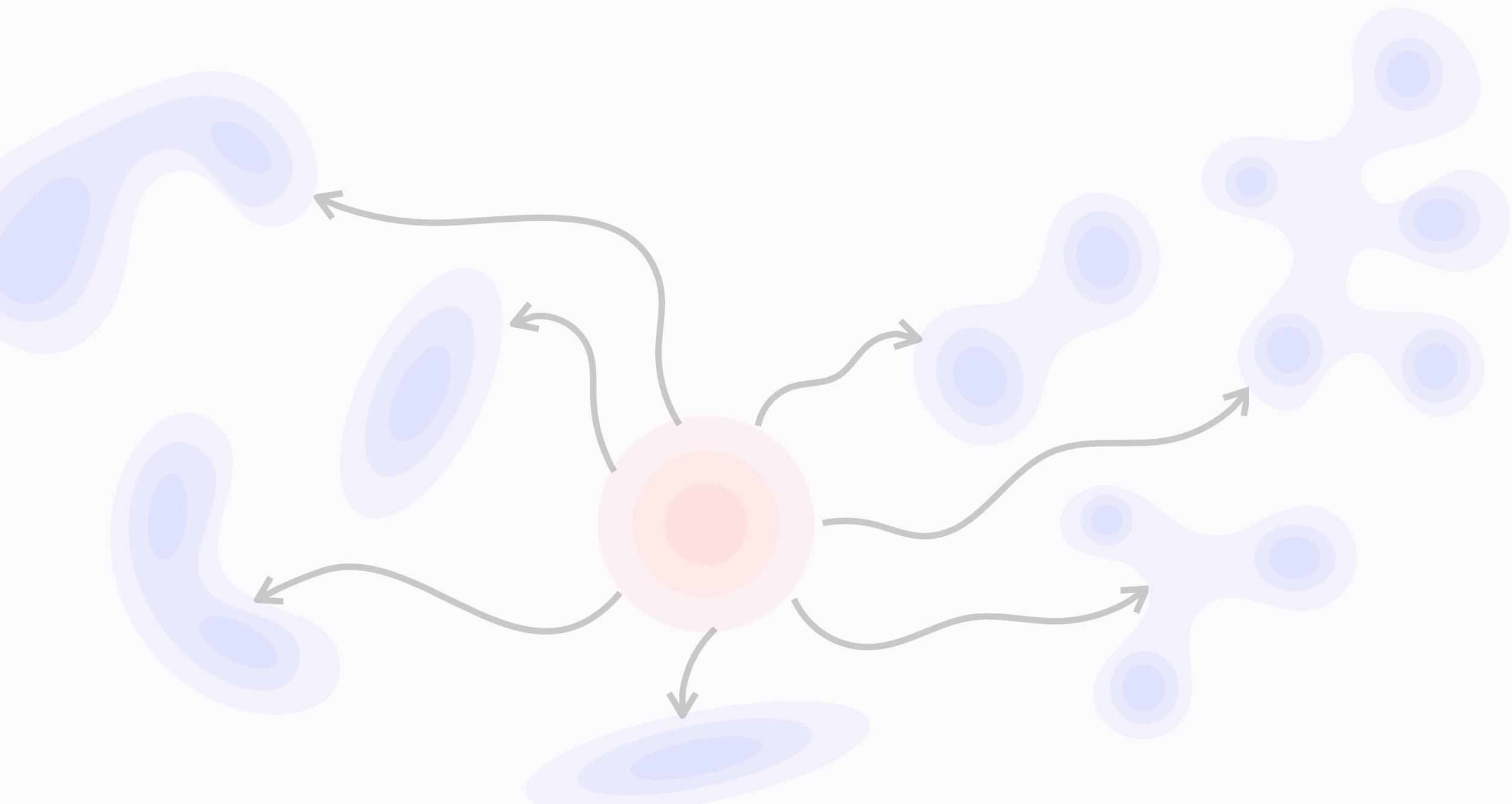
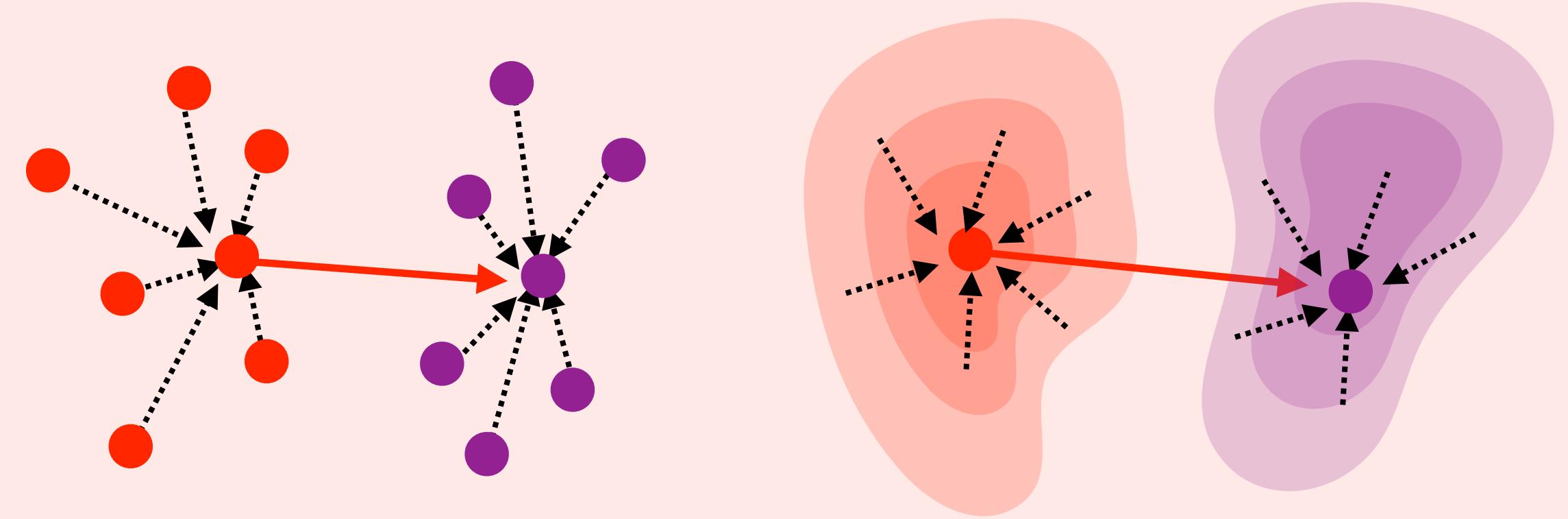
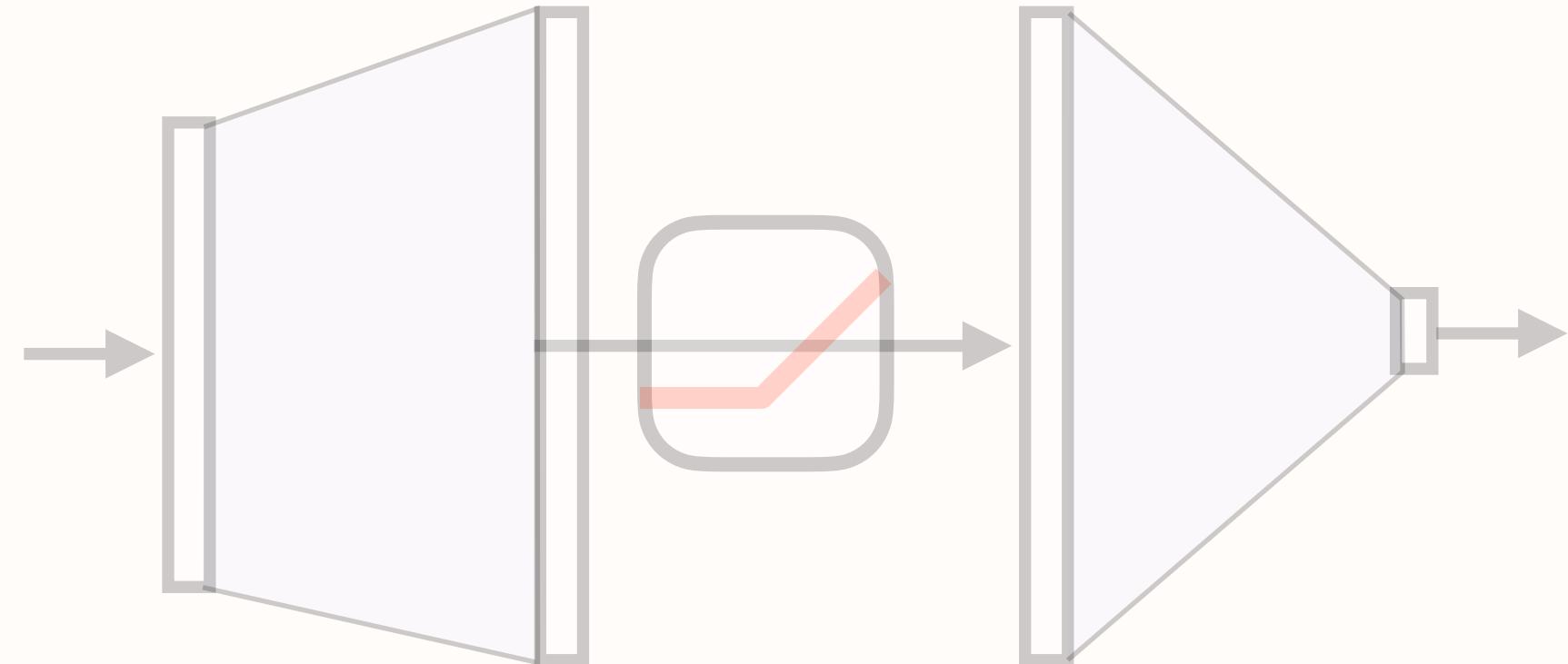
Trajectories cannot cross:  $x \rightarrow x(1)$  defines a diffeomorphism.



Open  
problems

Expressivity: universality requires embedding in  $\mathbb{R}^{d+1}$  space.  
Optimization: weights could « blows » during training.

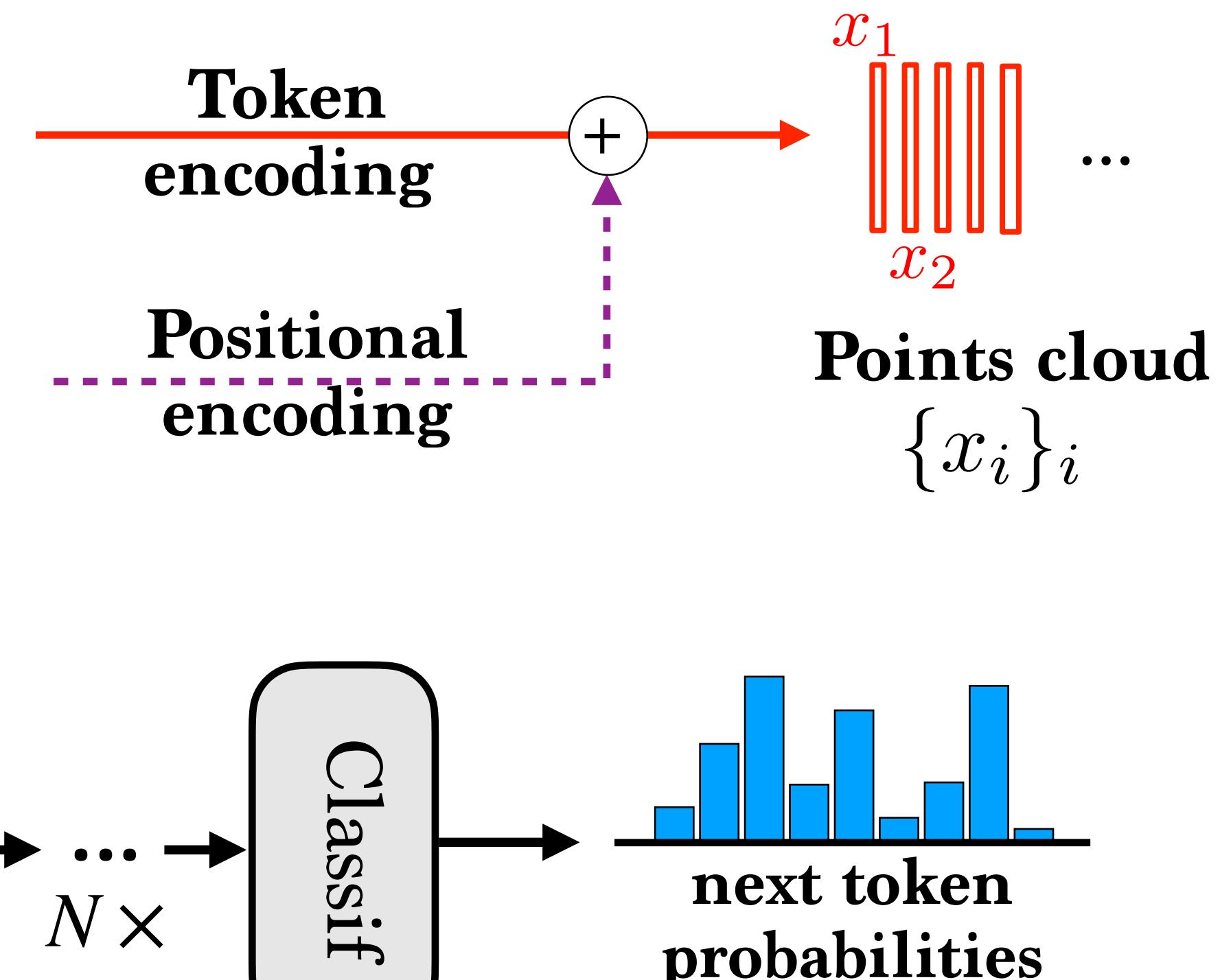
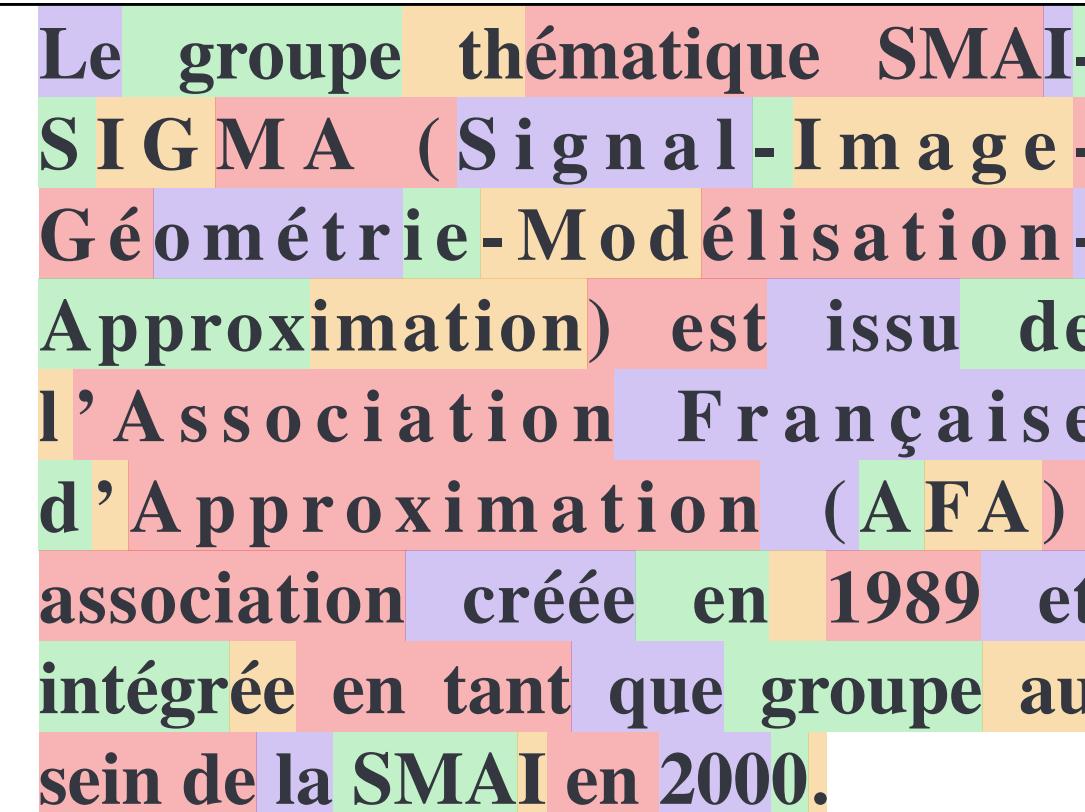
# In Context Mappings



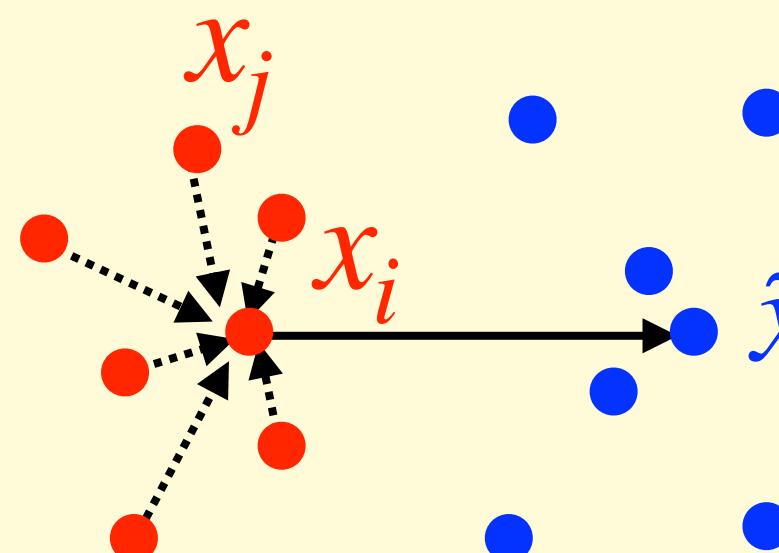
# Transformers and attention mechanism

Le groupe thématique SMAI-SIGMA (Signal-Image-Géométrie-Modélisation-Approximation) est issu de l'Association Française d'Approximation (AFA), association créée en 1989 et intégrée en tant que groupe au sein de la SMAI en 2000.

Tokenize



(Unmasked) Attention layer

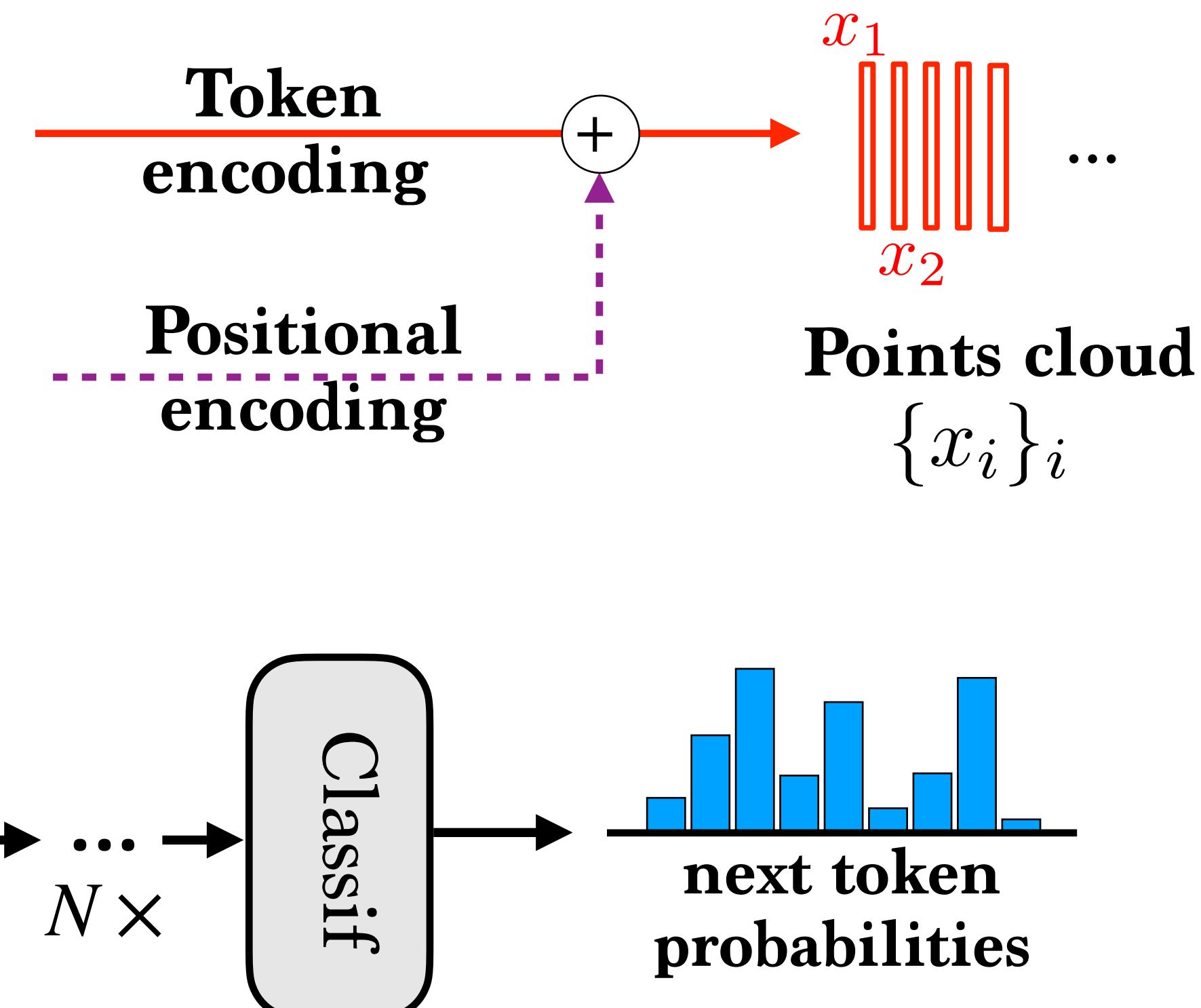
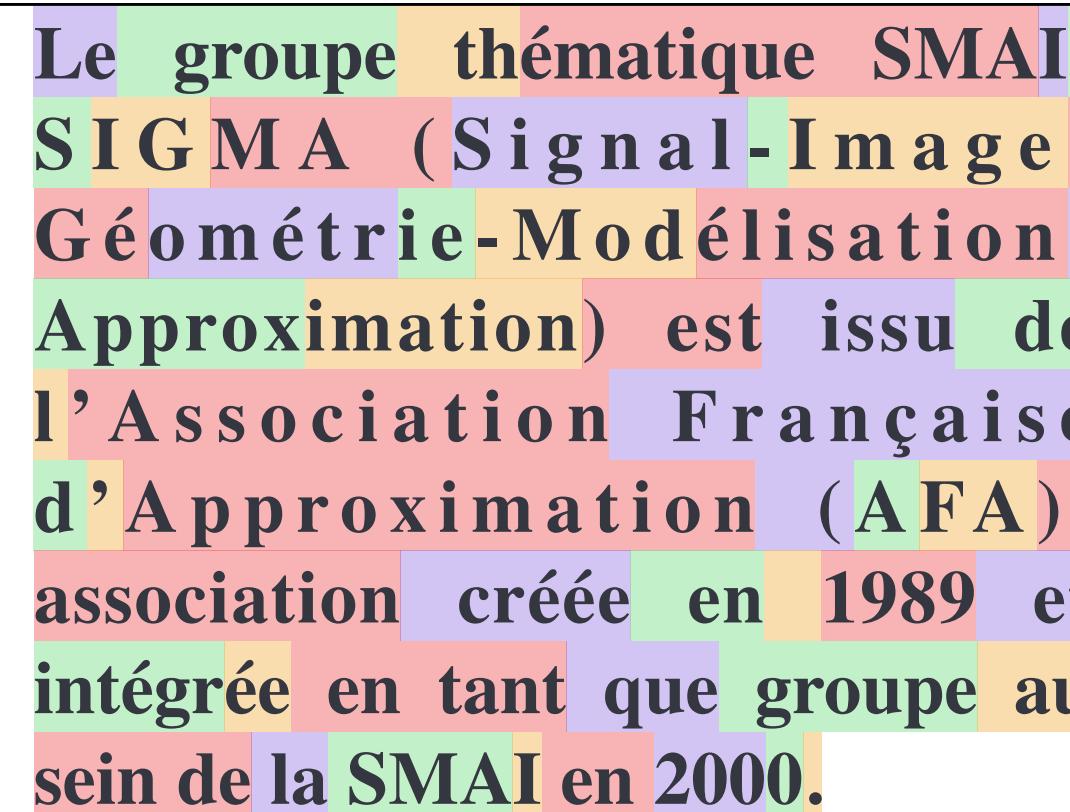


$$\tilde{x}_i := \sum_j \frac{e^{\langle Qx_i, Kx_j \rangle}}{\sum_\ell e^{\langle Qx_i, Kx_\ell \rangle}} Vx_j$$

# Transformers and attention mechanism

Le groupe thématique SMAI-SIGMA (Signal-Image-Géométrie-Modélisation-Approximation) est issu de l'Association Française d'Approximation (AFA), association créée en 1989 et intégrée en tant que groupe au sein de la SMAI en 2000.

Tokenize



(Unmasked) Attention layer

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

Understanding

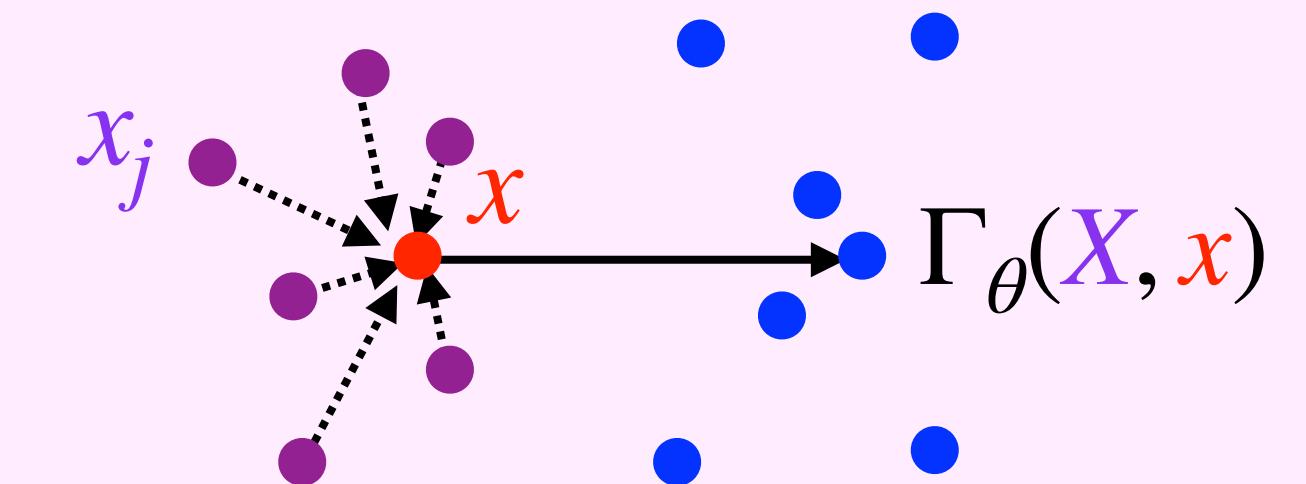
- Arbitrary number of tokens
- Arbitrary number of layers
- Expressivity

# Attention as In-context Mapping

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

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

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

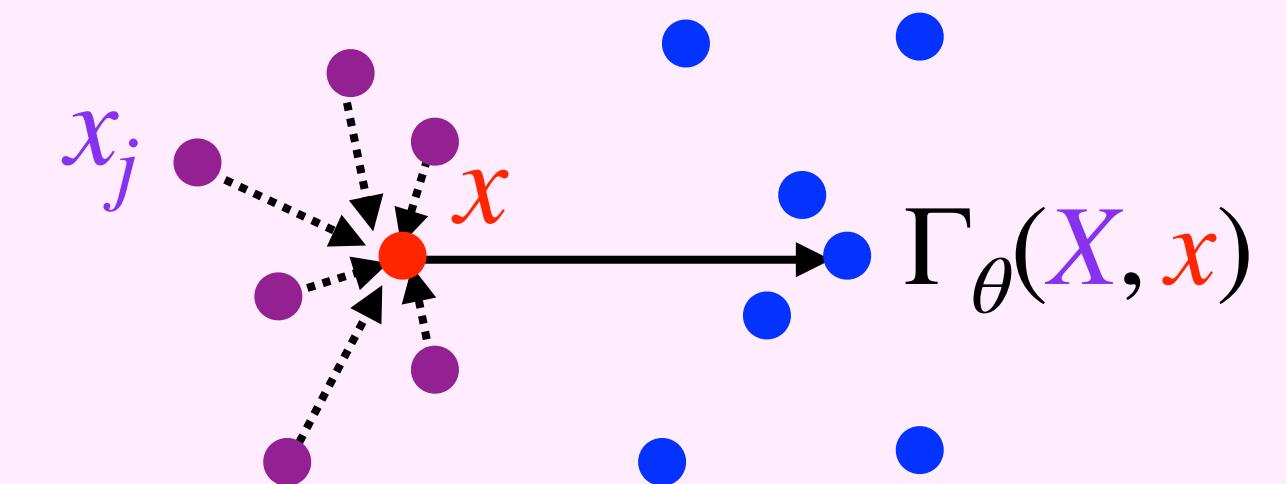


# Attention as In-context Mapping

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:**  $X \mapsto \{\sum_{h=1}^H \Gamma_{\theta_h}[\mathbf{X}](x_i)\}_{i=1}^n$

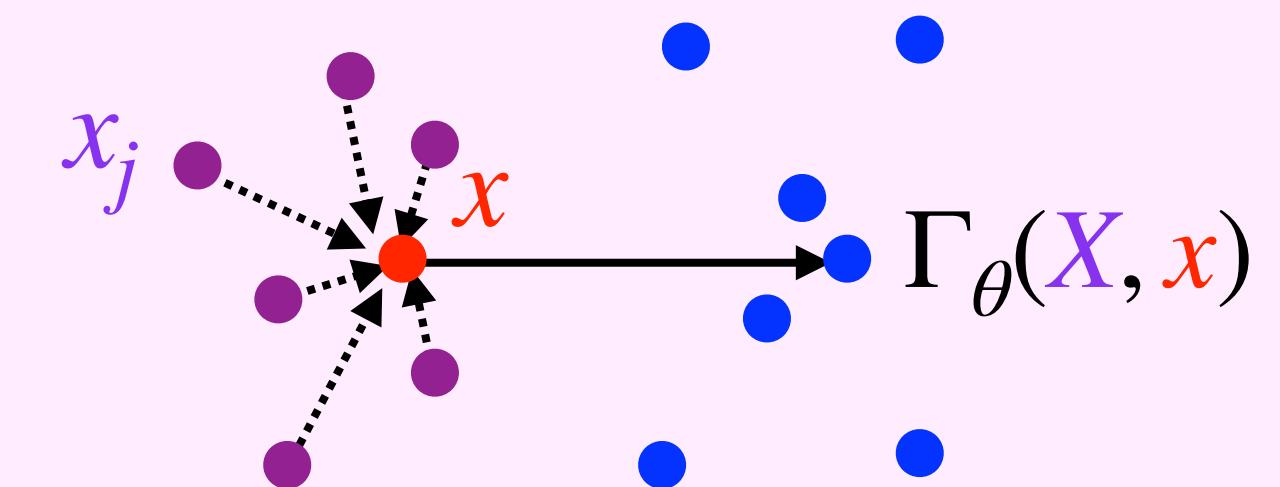
$K_1, Q_1$
$K_2, Q_2$
...
$K_H, Q_H$

# Attention as In-context Mapping

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

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

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



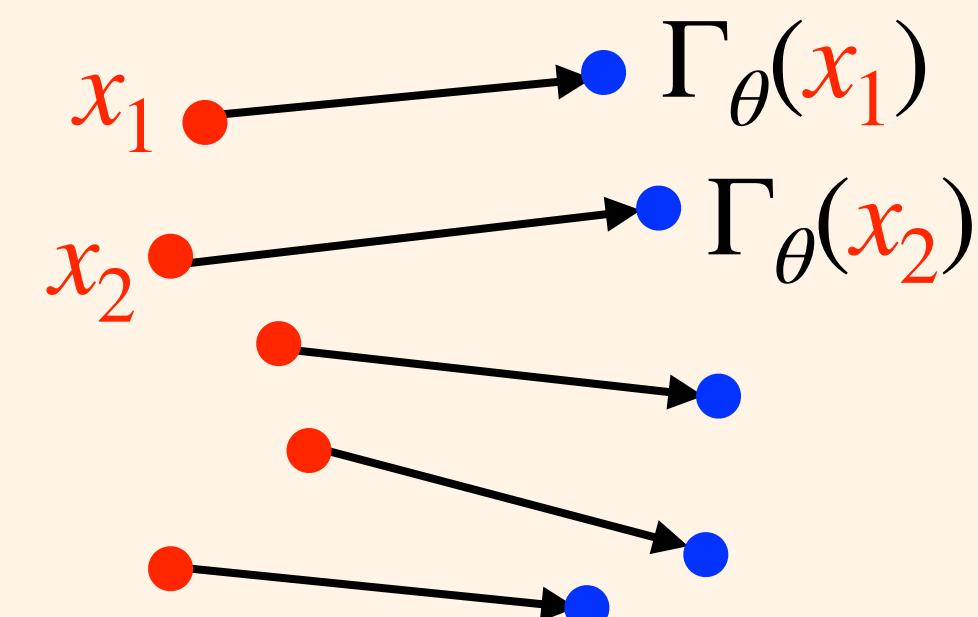
**Single-head attention layer:**  $X \mapsto \{\Gamma_\theta[\mathbf{X}](x_i)\}_{i=1}^n$

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

$K_1, Q_1$
$K_2, Q_2$
...
$K_H, Q_H$

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

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

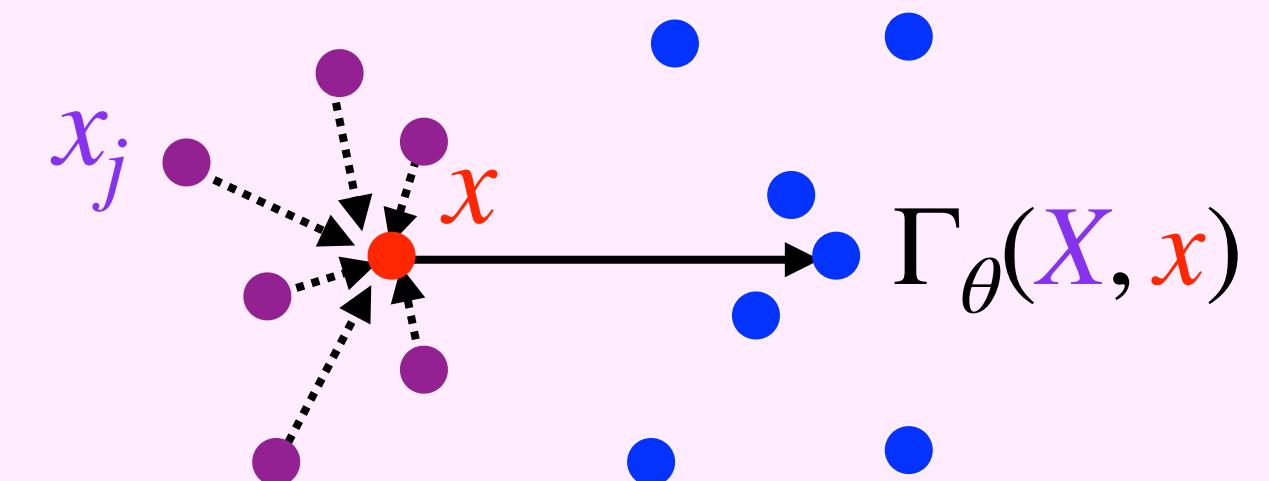


# Attention as In-context Mapping

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

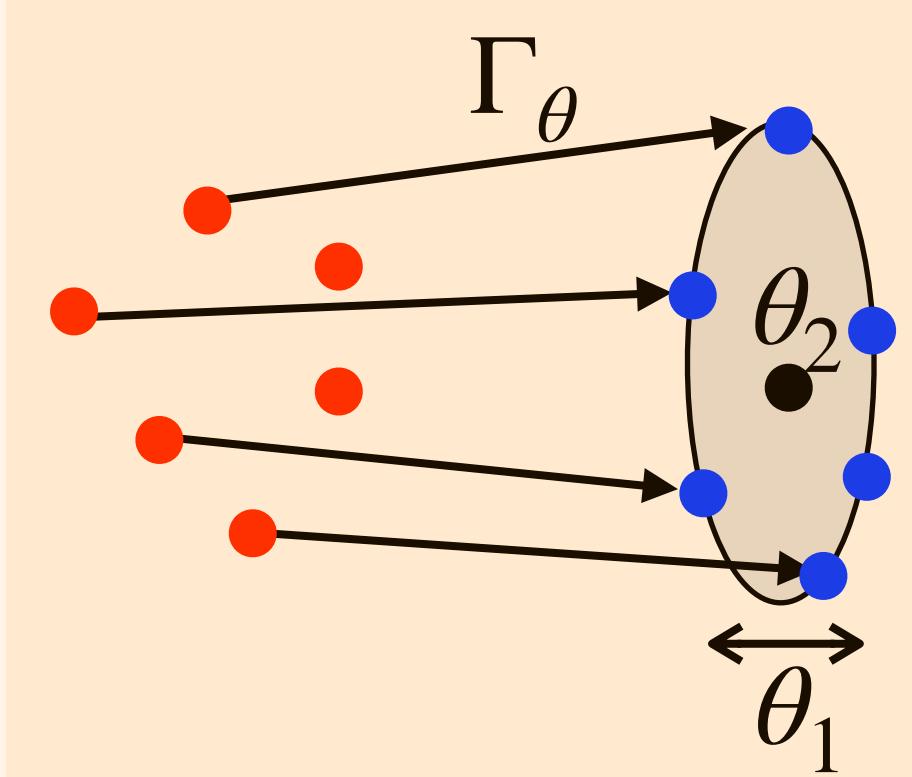
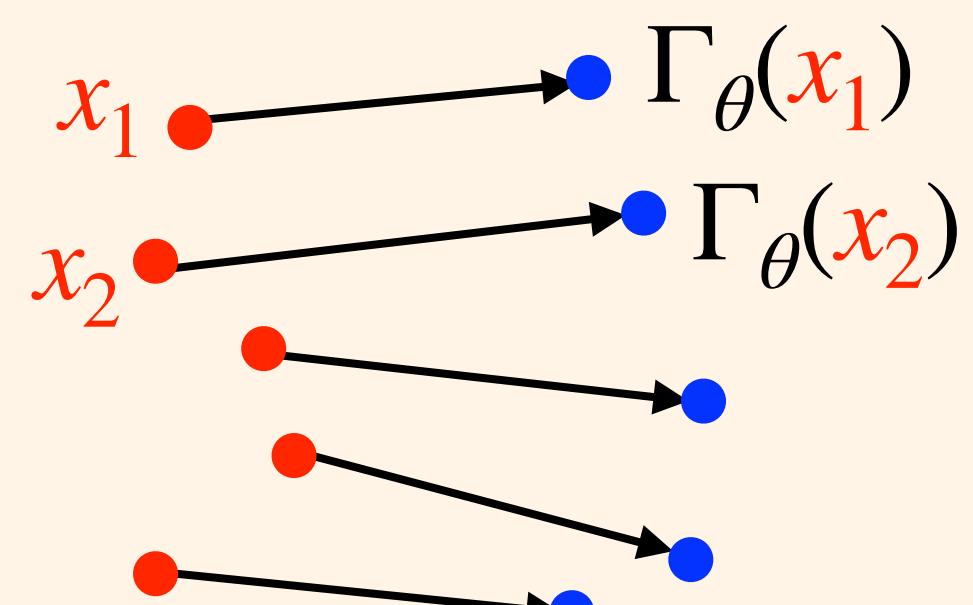
$K_1, Q_1$
$K_2, Q_2$
...
$K_H, Q_H$

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

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

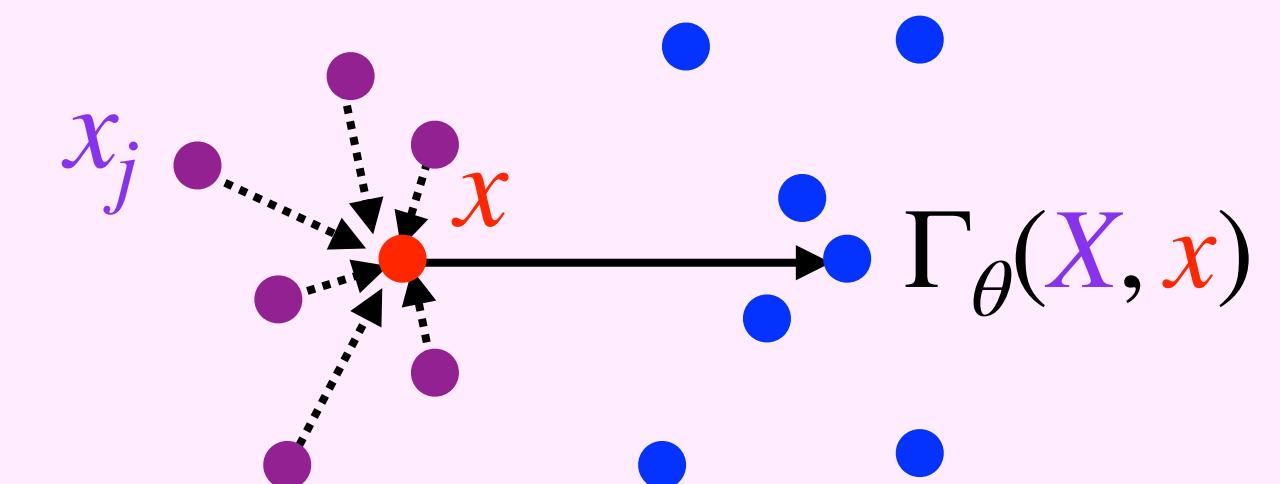


# Attention as In-context Mapping

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

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

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



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

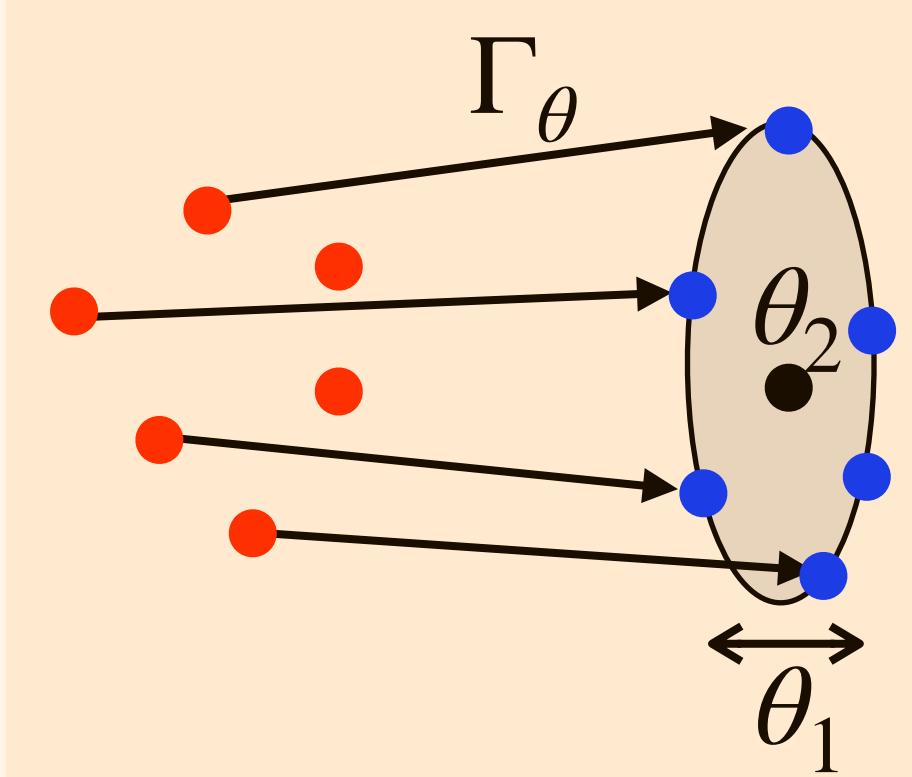
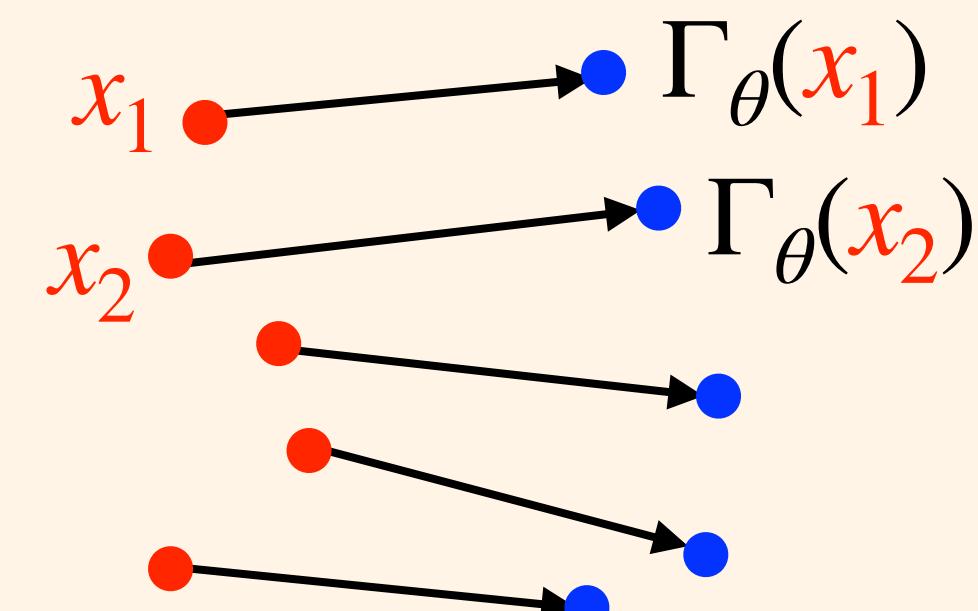
$K_1, Q_1$
$K_2, Q_2$
...
$K_H, Q_H$

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

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



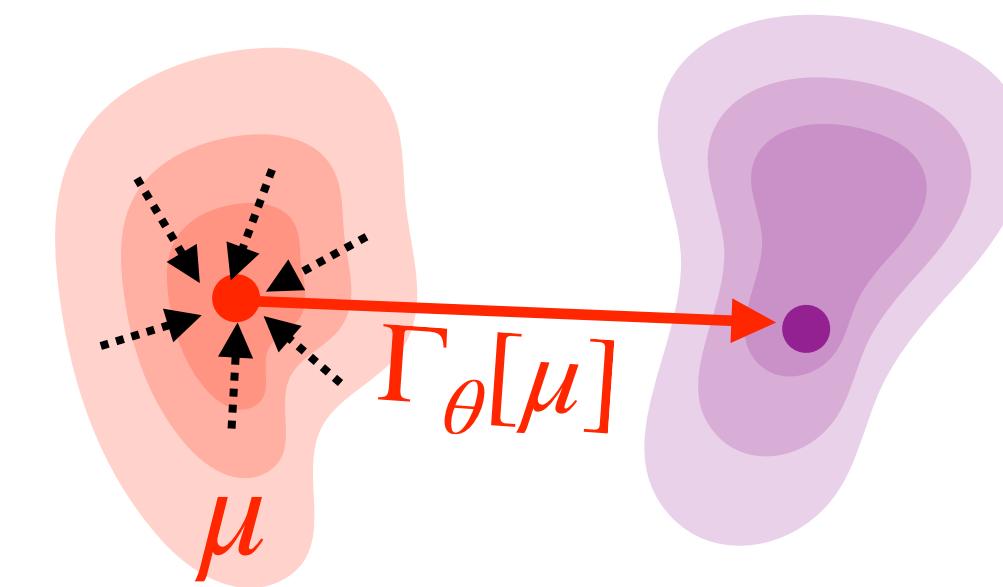
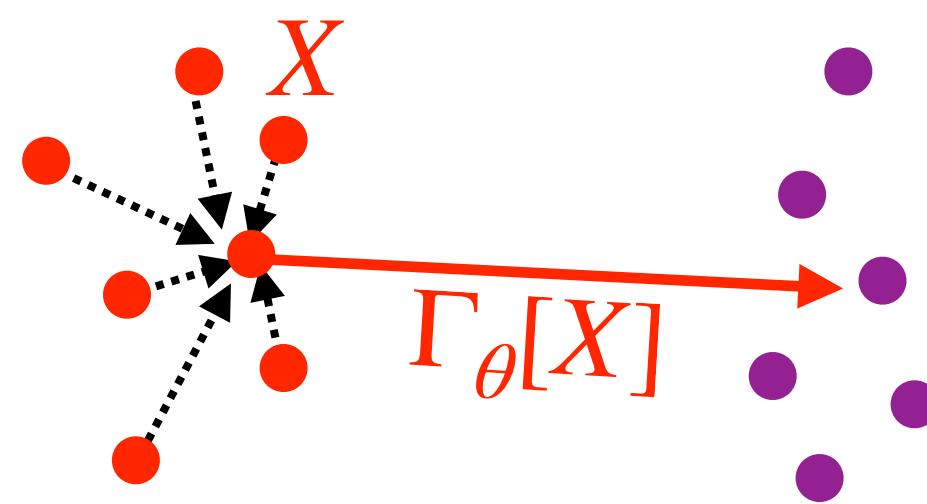
Transformer  $\equiv$  composition of in-context and context-free layers.

# Attentions Operating over Measures

Number  $n$  of token is arbitrary.

(Unmasked) attention is permutation invariant.

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

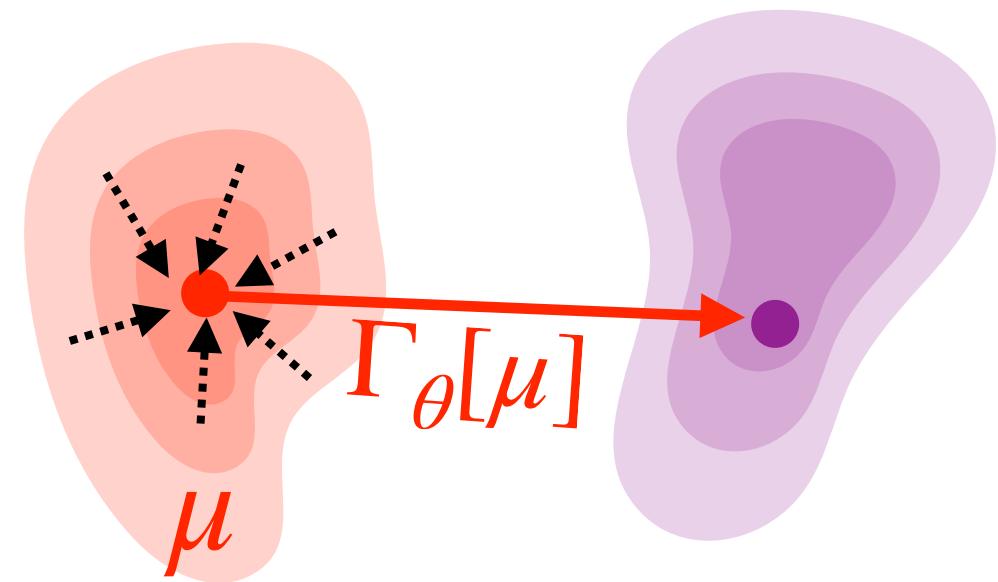
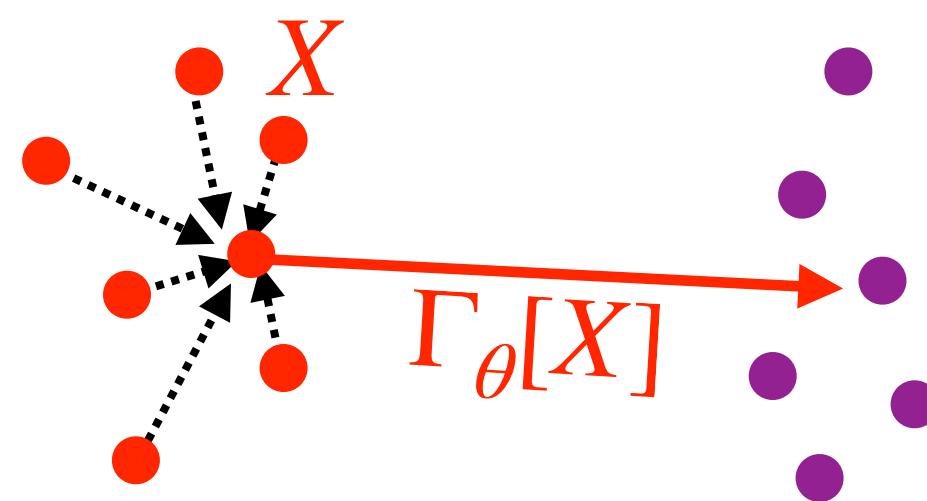


# Attentions Operating over Measures

Number  $n$  of token is arbitrary.

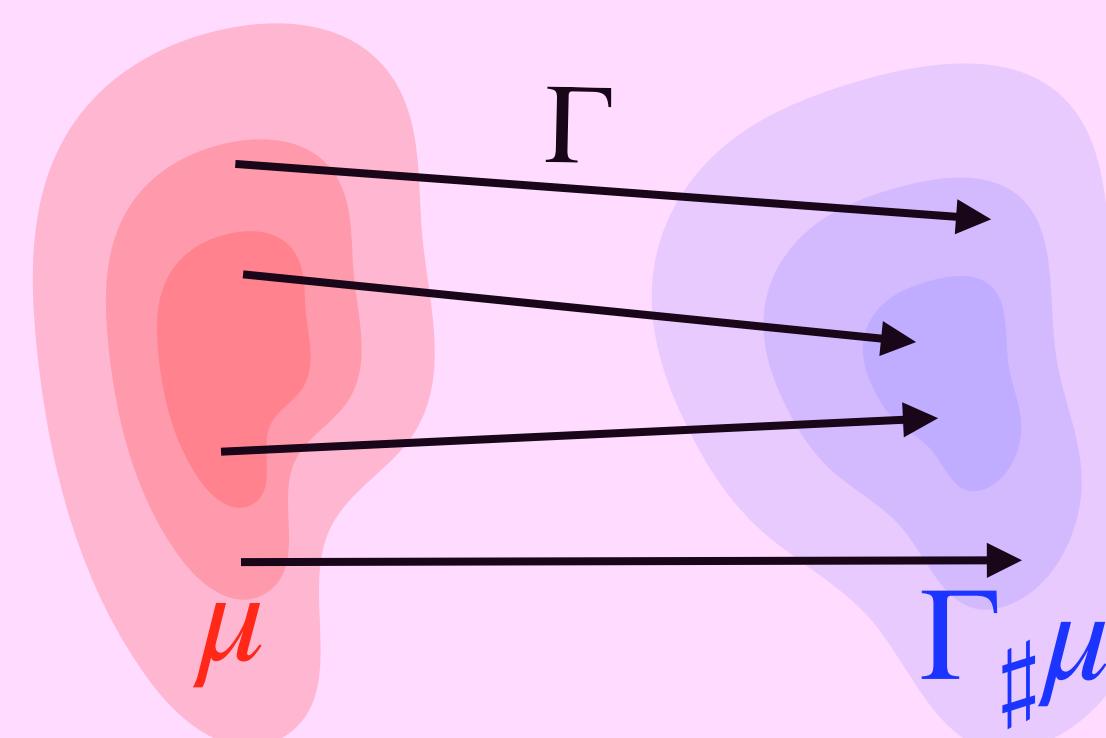
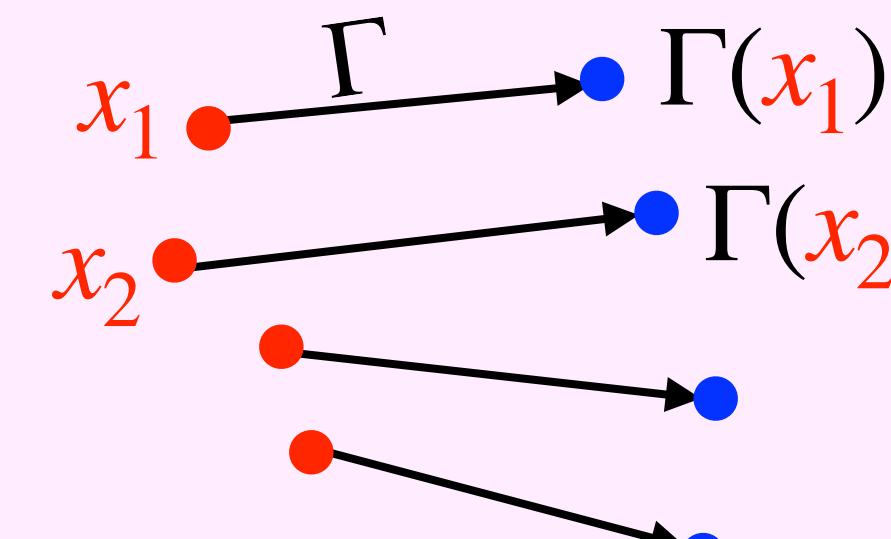
(Unmasked) attention is permutation invariant.

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



## Push-forward

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



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

## Layers

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



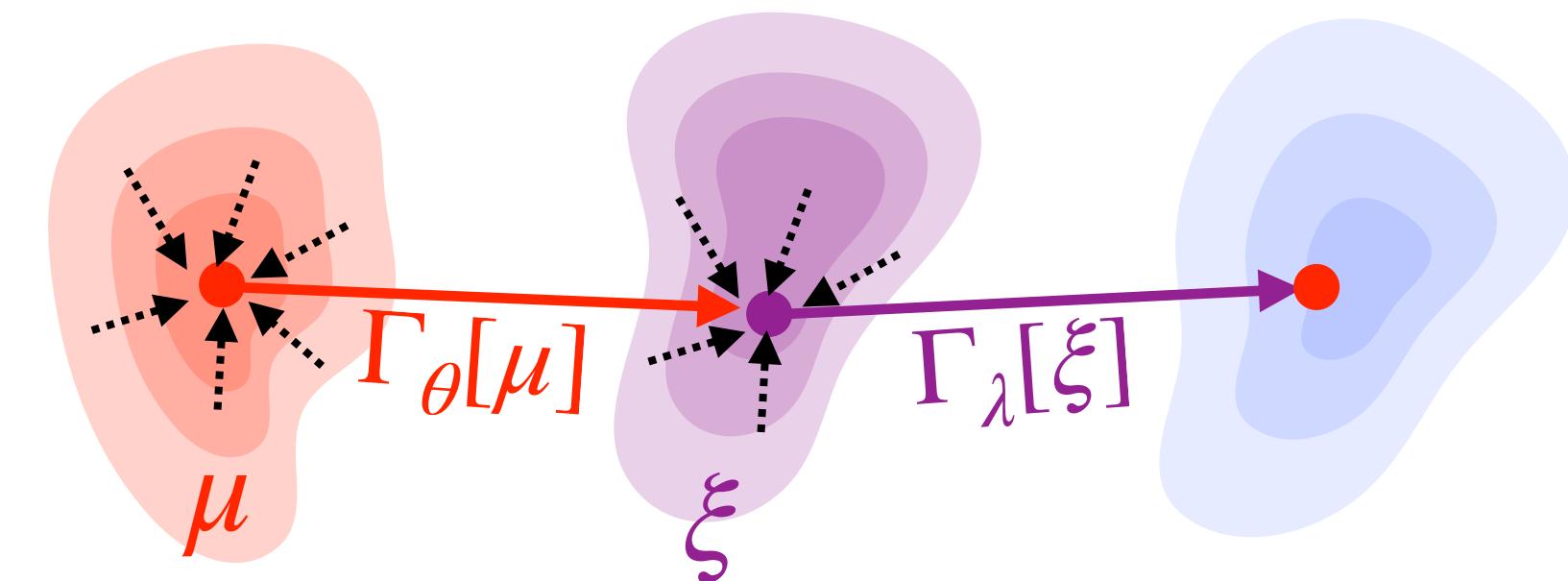
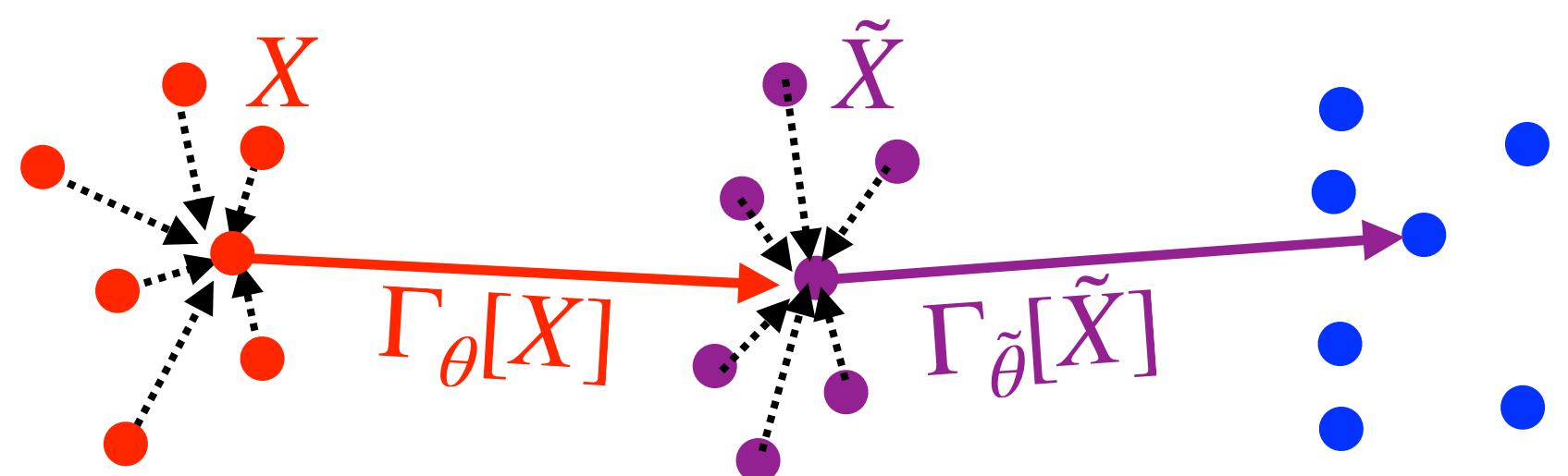
$$\mu \mapsto \Gamma[\textcolor{violet}{\mu}]_\sharp \mu$$

# Attentions Operating over Measures

Number  $n$  of token is arbitrary.

(Unmasked) attention is permutation invariant.

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

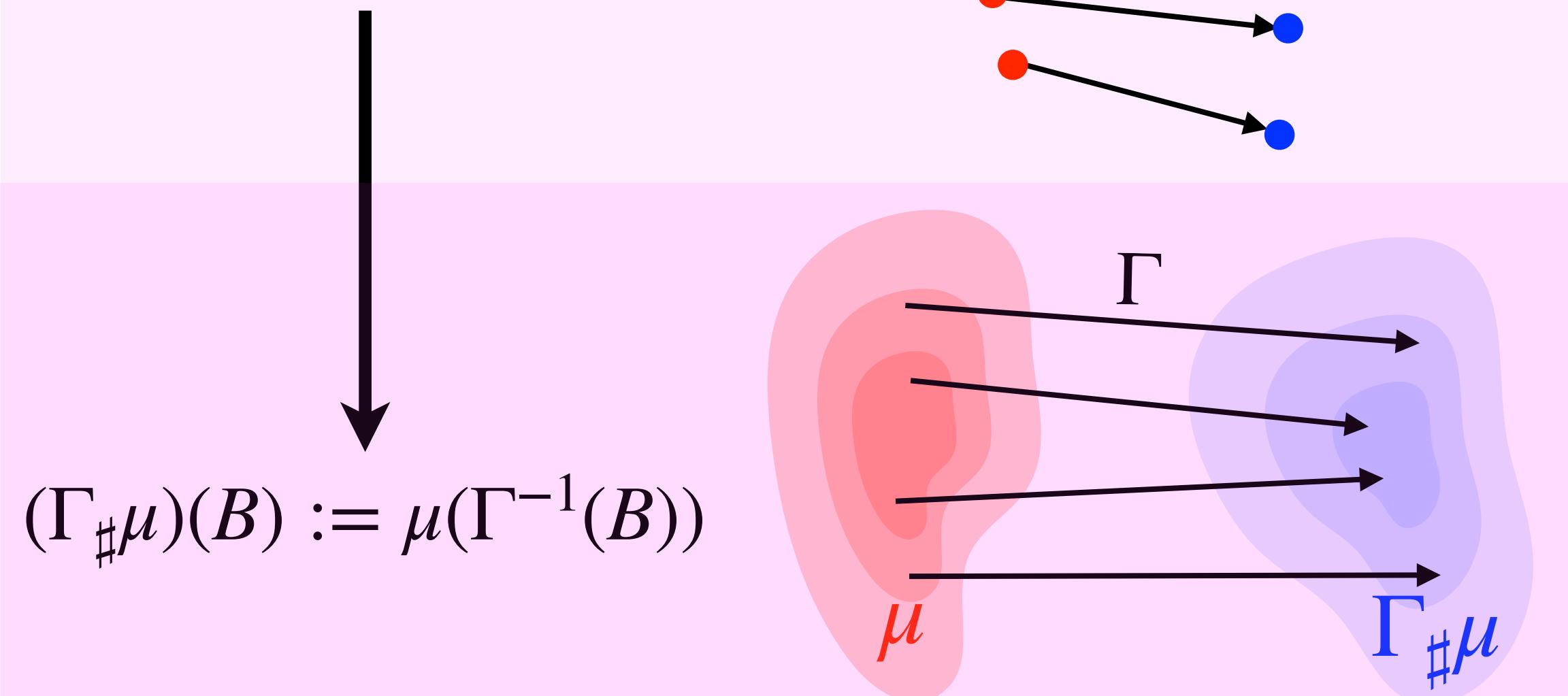


## Push-forward

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

$$x_1 \xrightarrow{\Gamma} \Gamma(x_1)$$

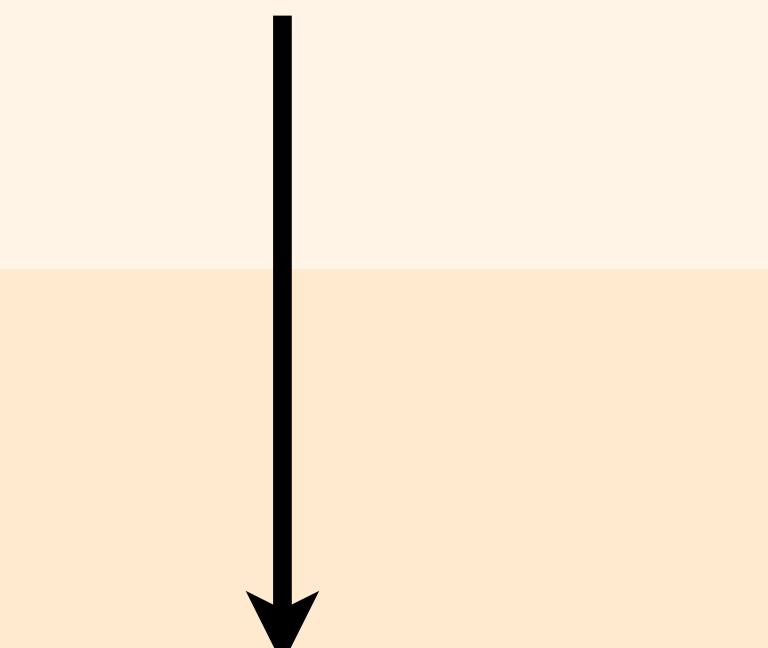
$$x_2 \xrightarrow{\Gamma} \Gamma(x_2)$$



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

## Layers

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



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

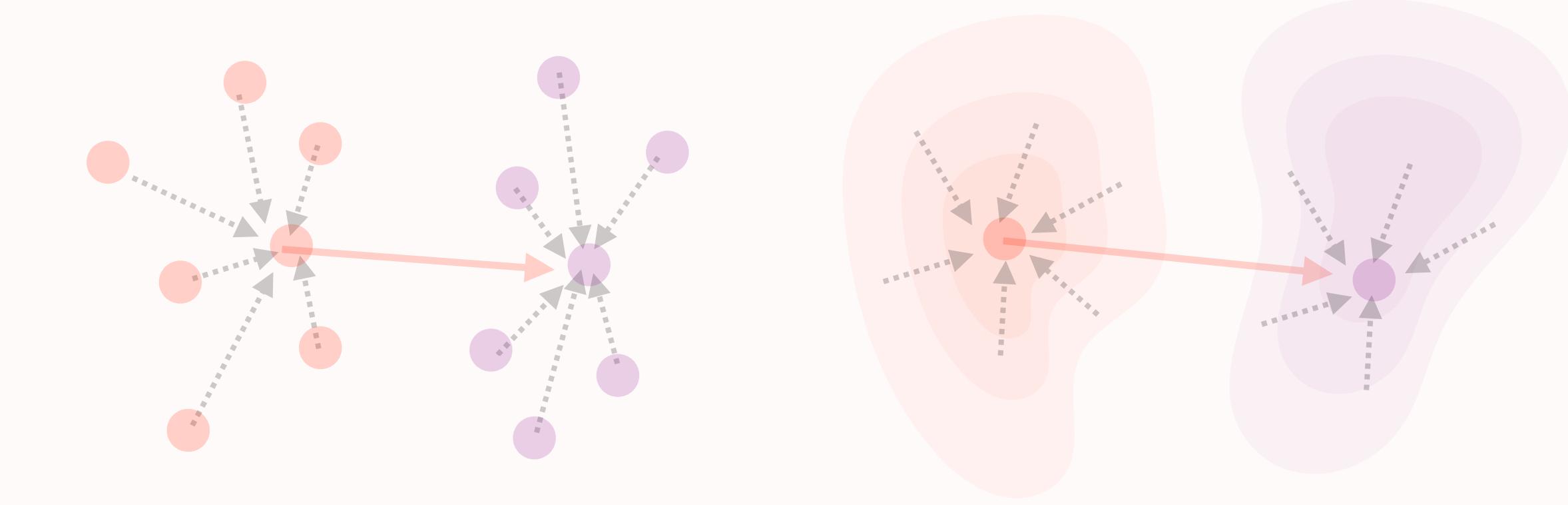
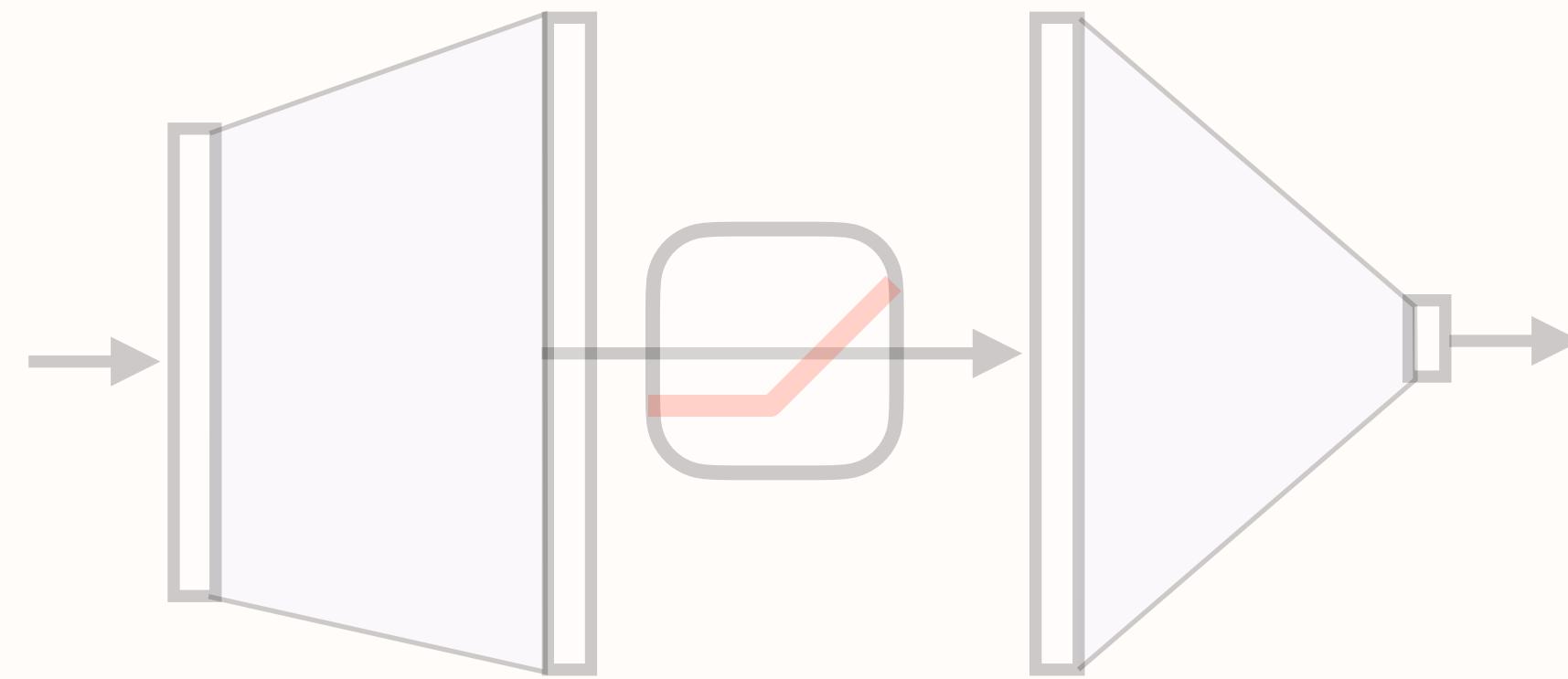
## Composing layers

$$(\Gamma' \diamond \Gamma)[X] := \Gamma'[Y] \circ \Gamma[X]$$

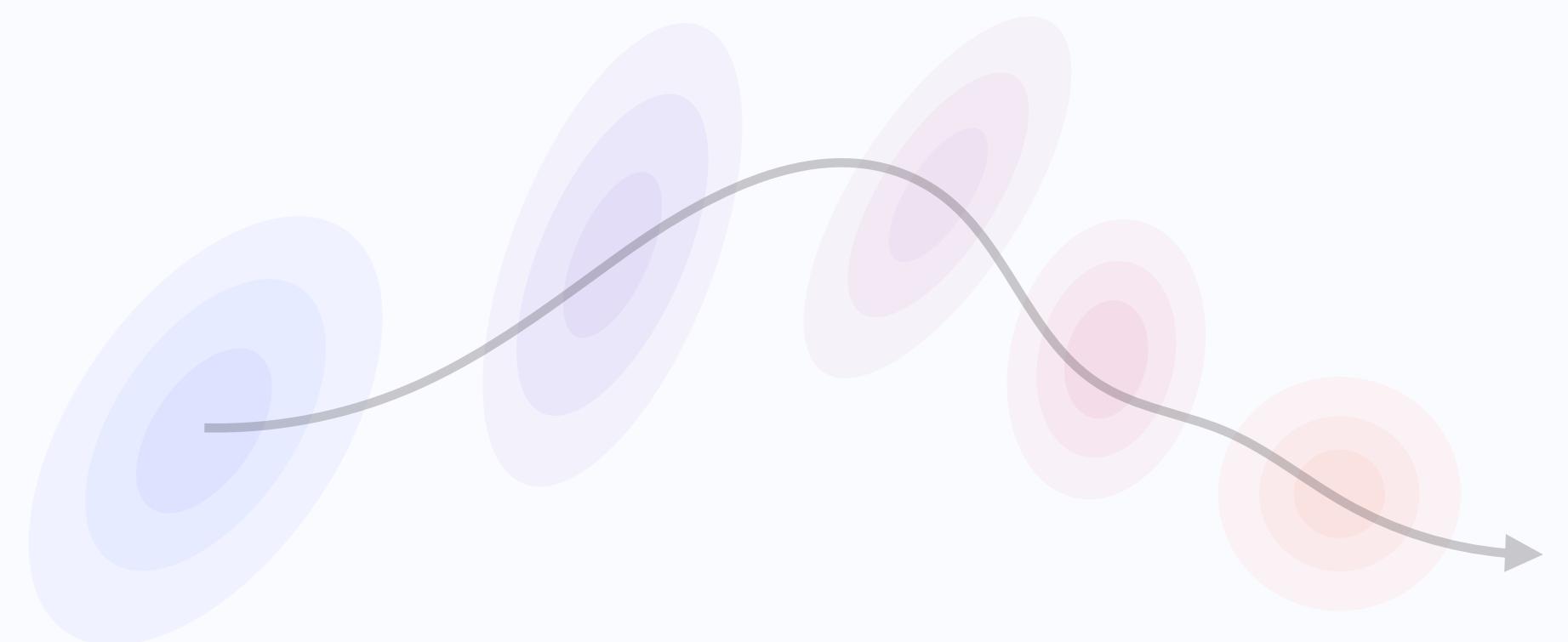
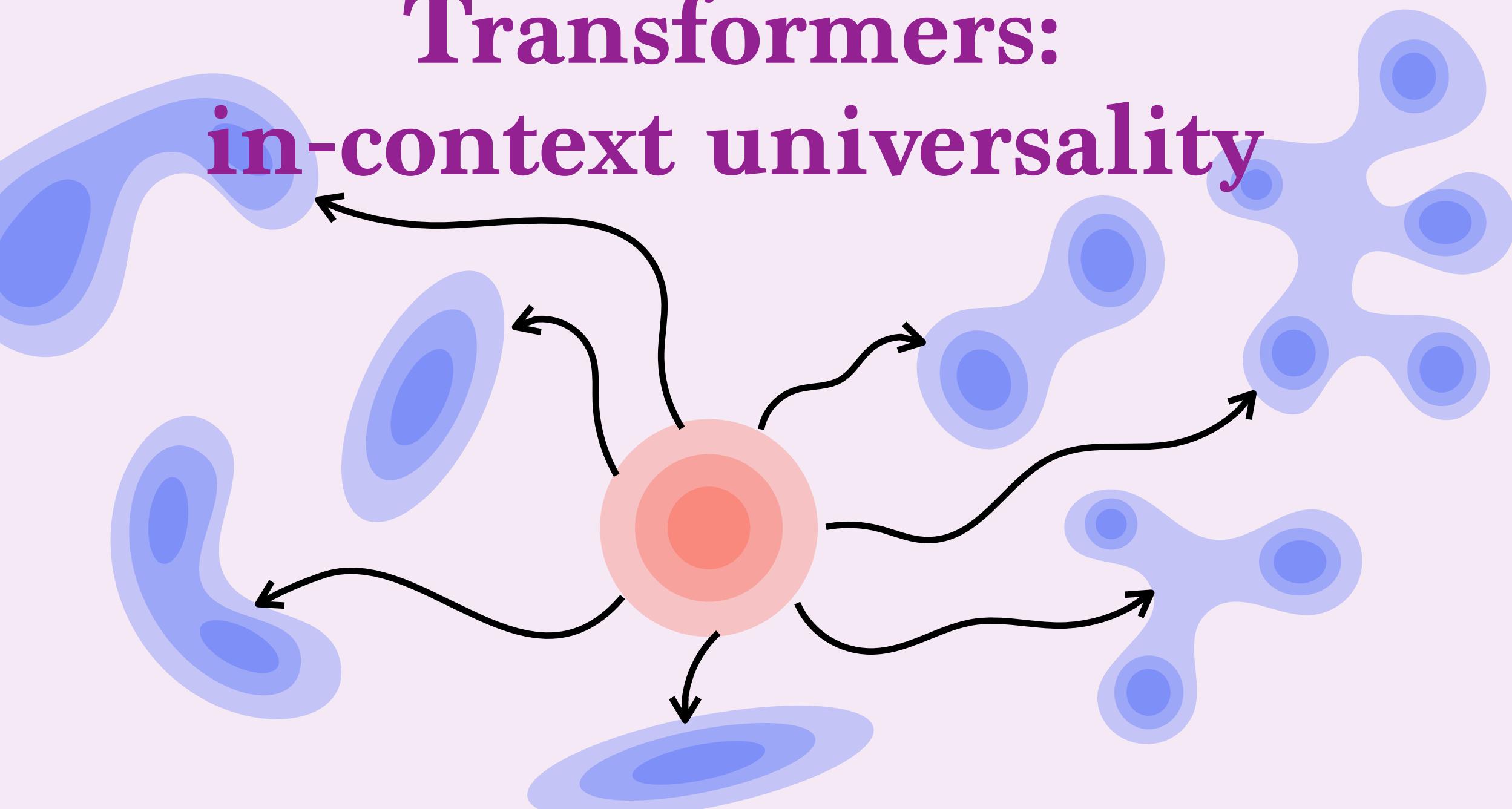
$$\text{where } Y := (\Gamma[X](x_i))_i$$

$$(\Gamma' \diamond \Gamma)[\mu] := \Gamma'[\xi] \circ \Gamma[\mu]$$

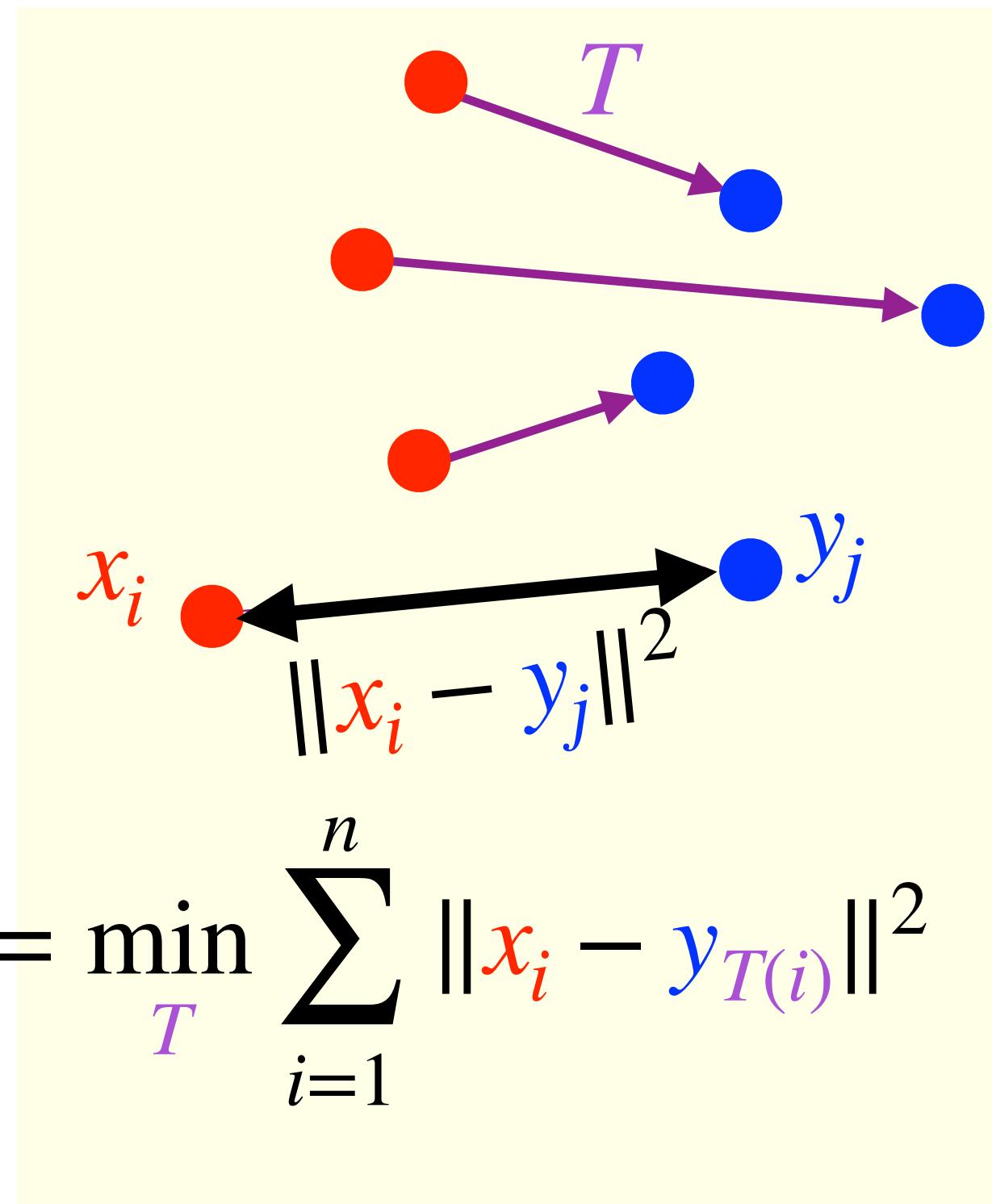
$$\text{where } \xi := \Gamma[\mu]_\sharp \mu$$



Transformers:  
in-context universality

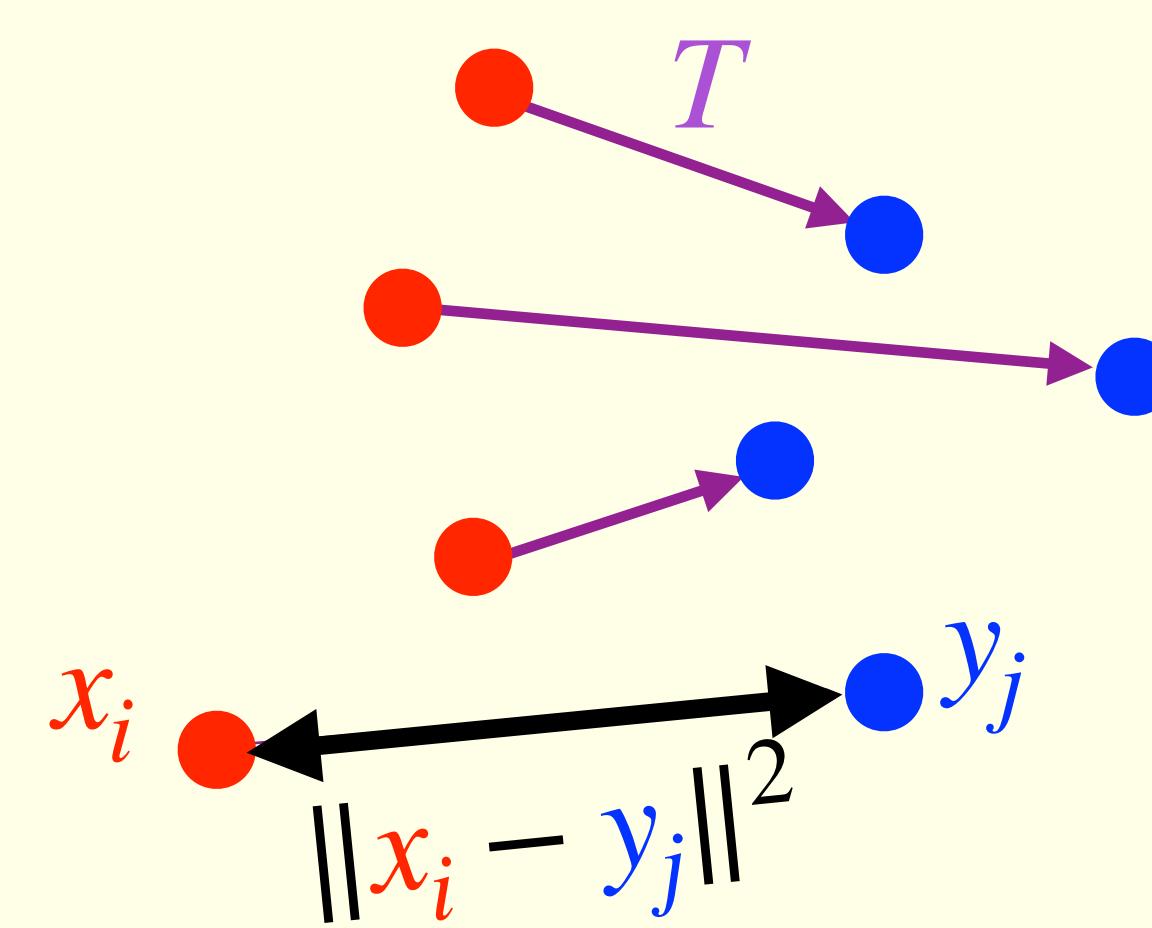


# Optimal Transport (Wasserstein) Distance

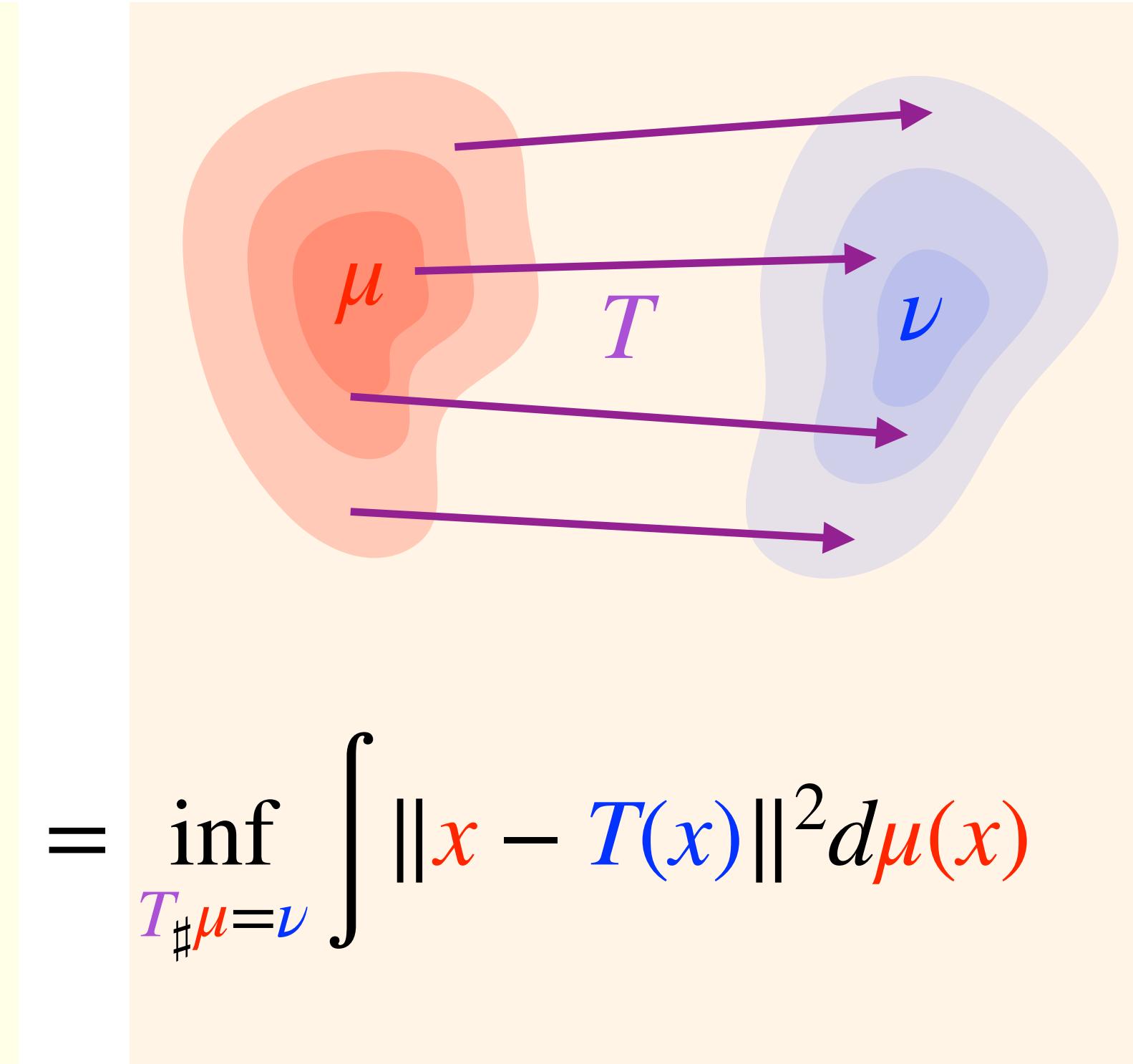


Monge 1784

# Optimal Transport (Wasserstein) Distance



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

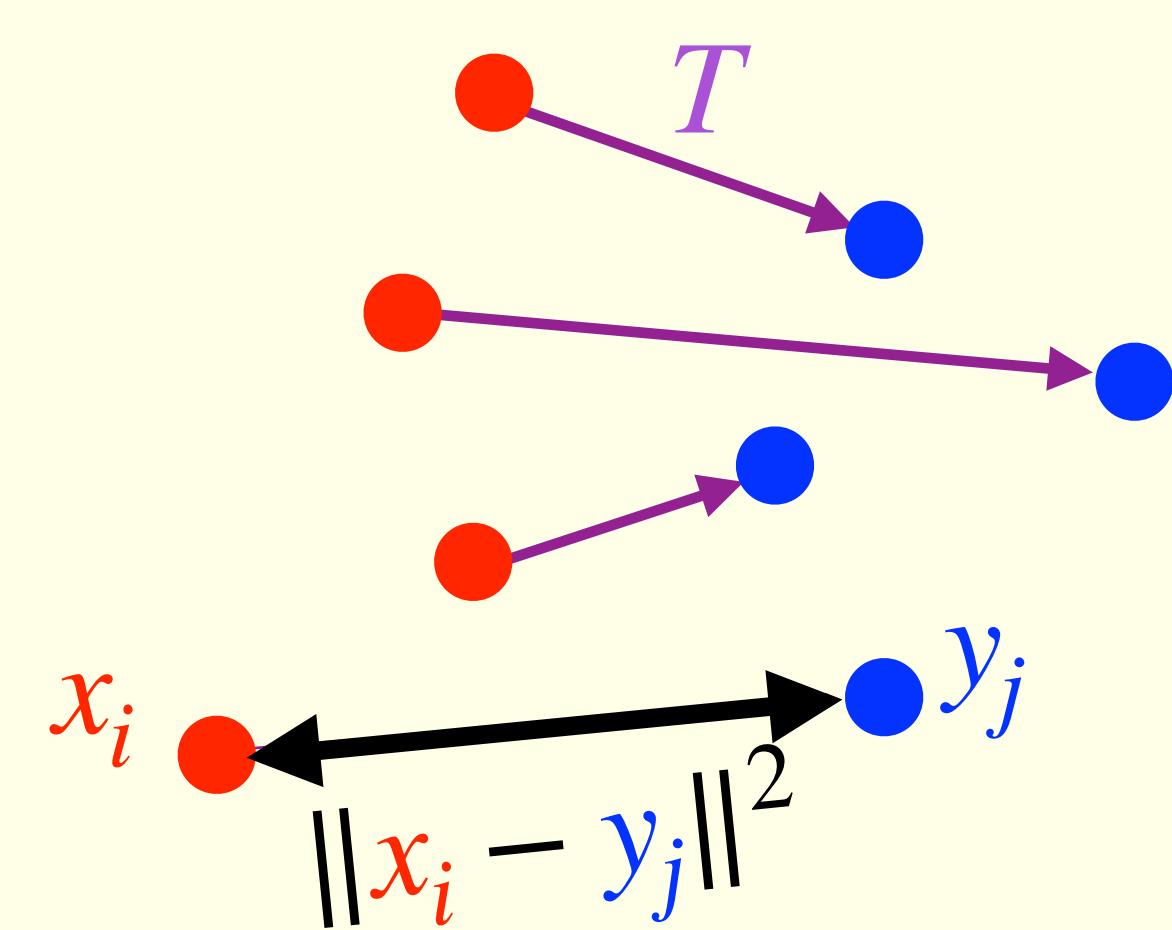


$$= \inf_{T_{\#}\mu = \nu} \int \|x - T(x)\|^2 d\mu(x)$$

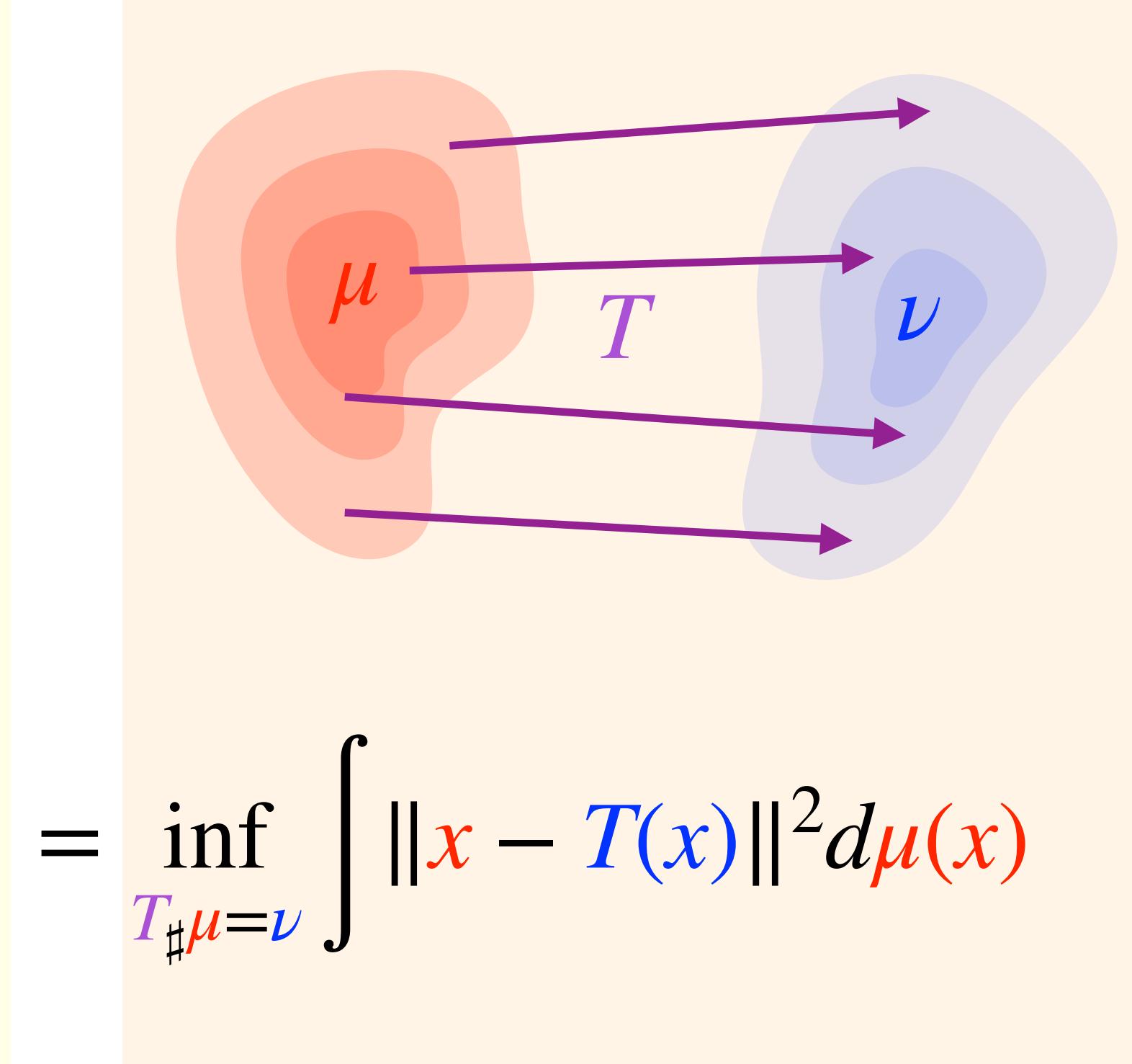


Monge 1784

# Optimal Transport (Wasserstein) Distance



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



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



Monge 1784



General measures:

Kantorovitch relaxation  
or  
Approximation by discrete measures

Kantorovitch 1942

# Universal Approximation

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

*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  $\sim \# \text{tokens}$

[Geshkovski, Rigollet, Ruiz-Balet, 2024]  $\rightarrow$  Universal interpolation.

[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:

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

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

# 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  $\circ$  MLP with skip connexion.

Cylindrical algebra:

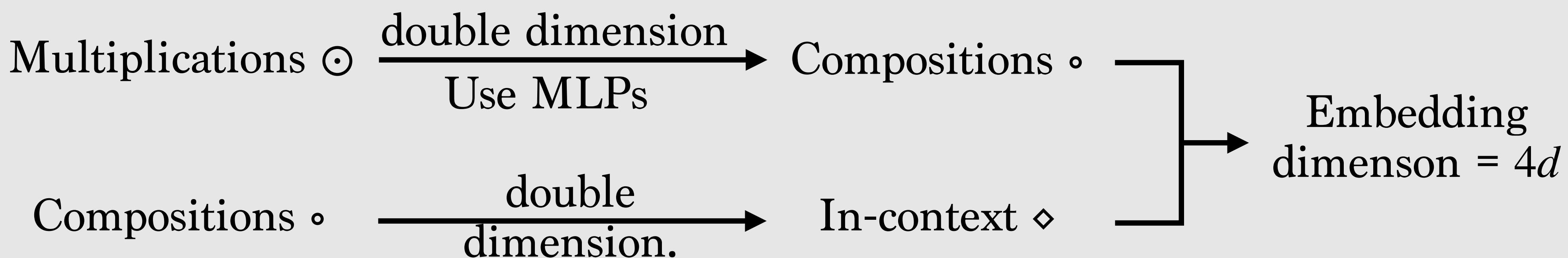
$$\mathcal{A} := \text{Span} \bigcup_N \{\gamma_{\theta_1} \odot \dots \odot \gamma_{\theta_N} : (\theta_1, \dots, \theta_N)\}$$

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

*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

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

*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

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

*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

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

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

?

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

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



# 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:

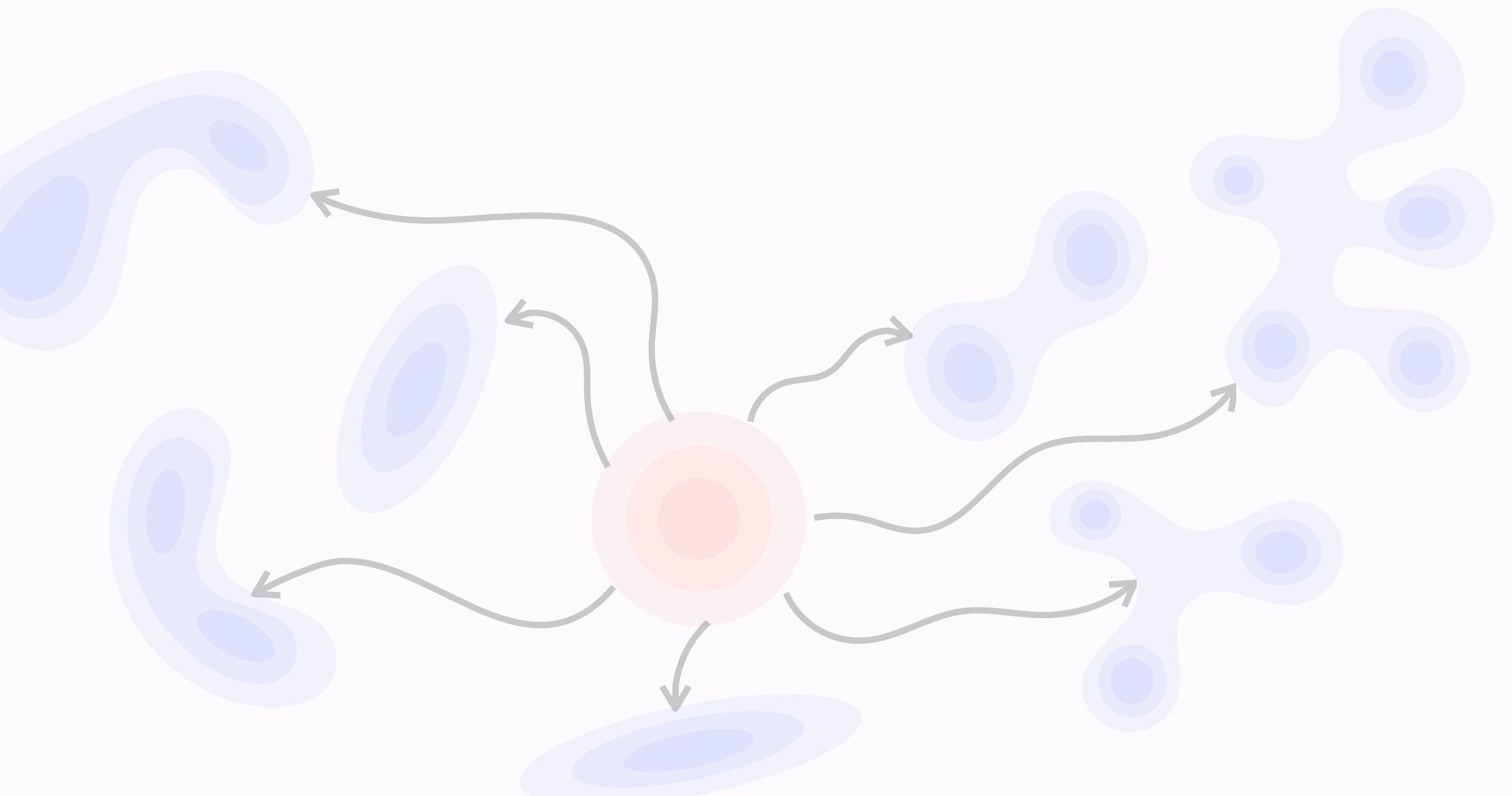
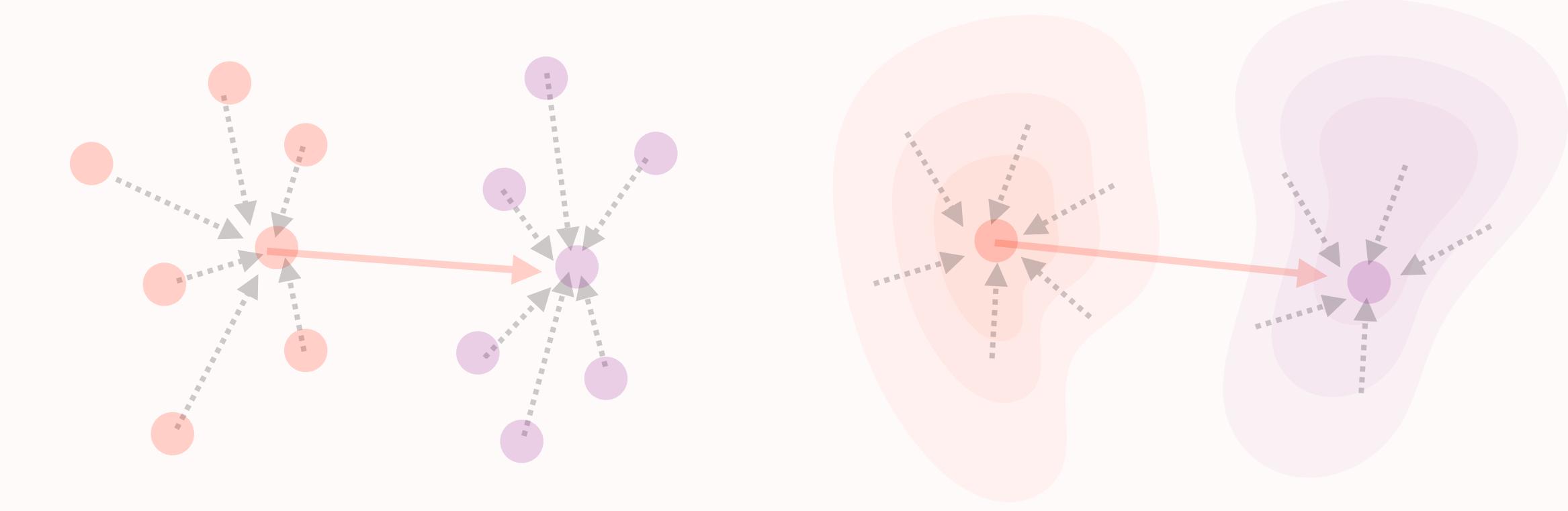
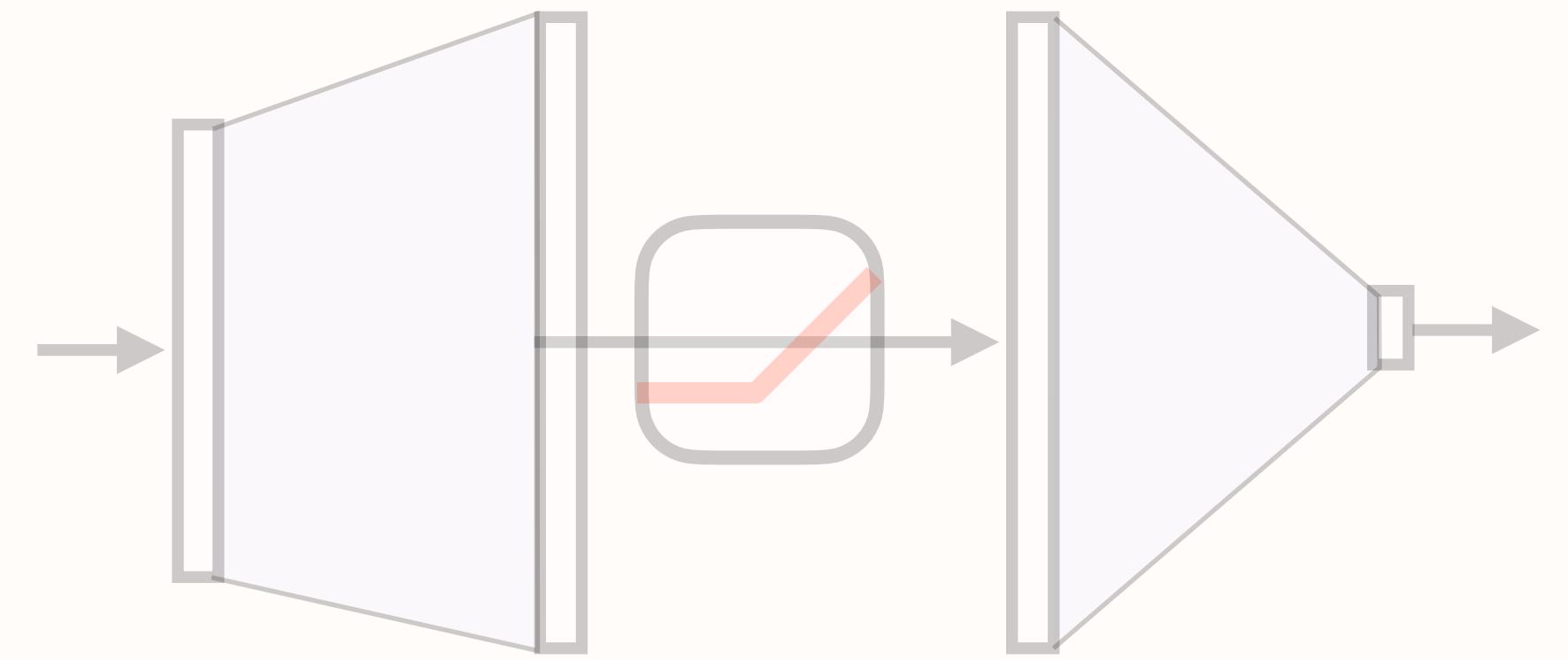
$$L_k(\mu)(b) := \int \frac{e^{by} y^k v}{\int e^{by'} d\mu(y')} d\mu(y)$$

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

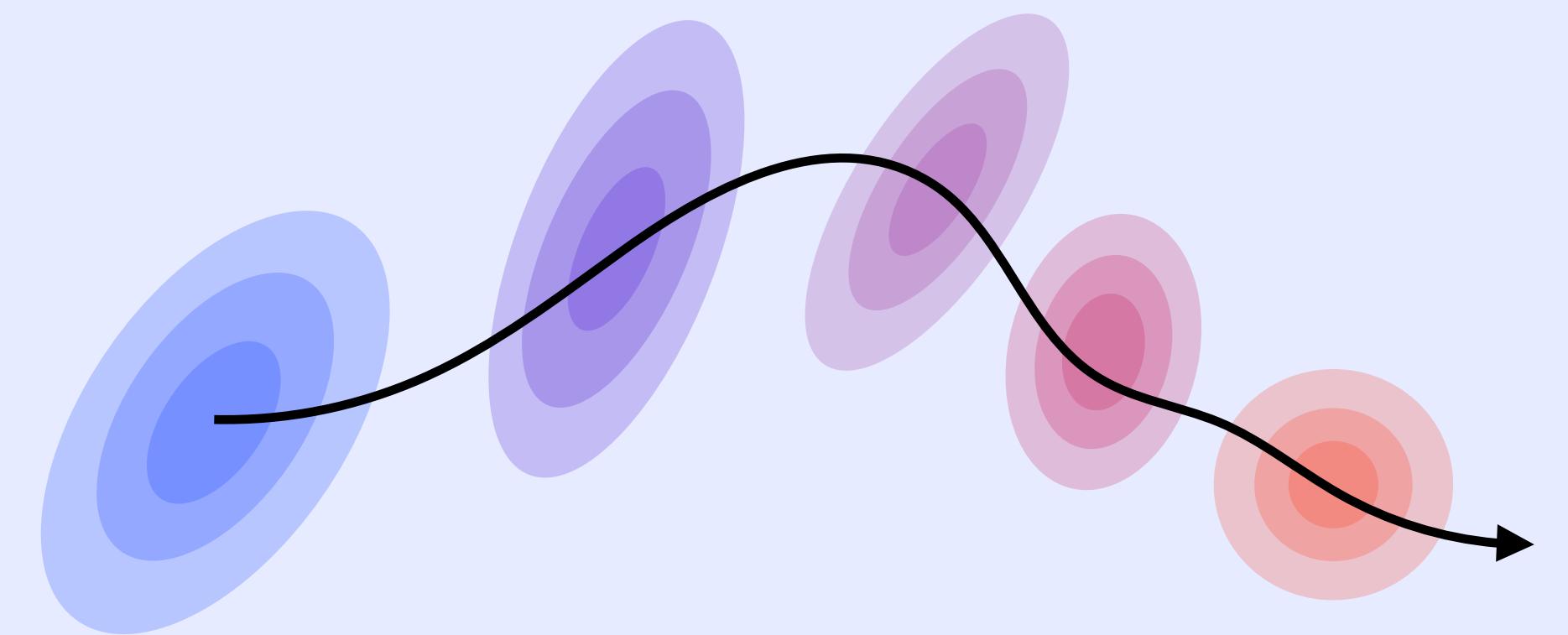
$$L_1(\mu) = L_1(\mu') \Rightarrow \forall k, L_k(\mu) = L_k(\mu') \Rightarrow \forall k, \int y^k d\mu(y) = \int y^k d\mu'(y)$$

*In higher dimensions:* use Radon transform.





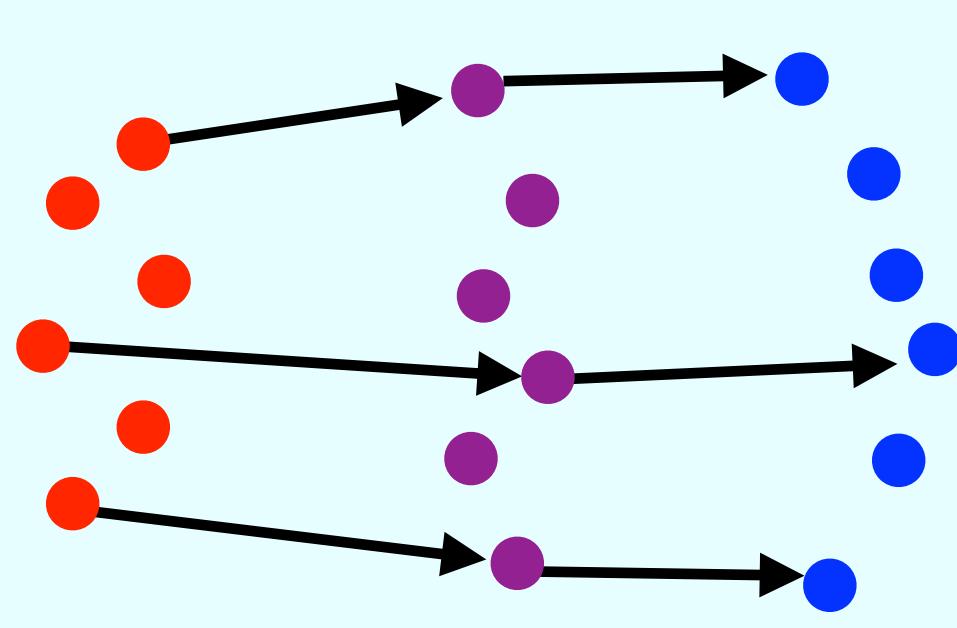
**Infinite depth  
and PDEs**



# Infinite Depth as a Neural PDE

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

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

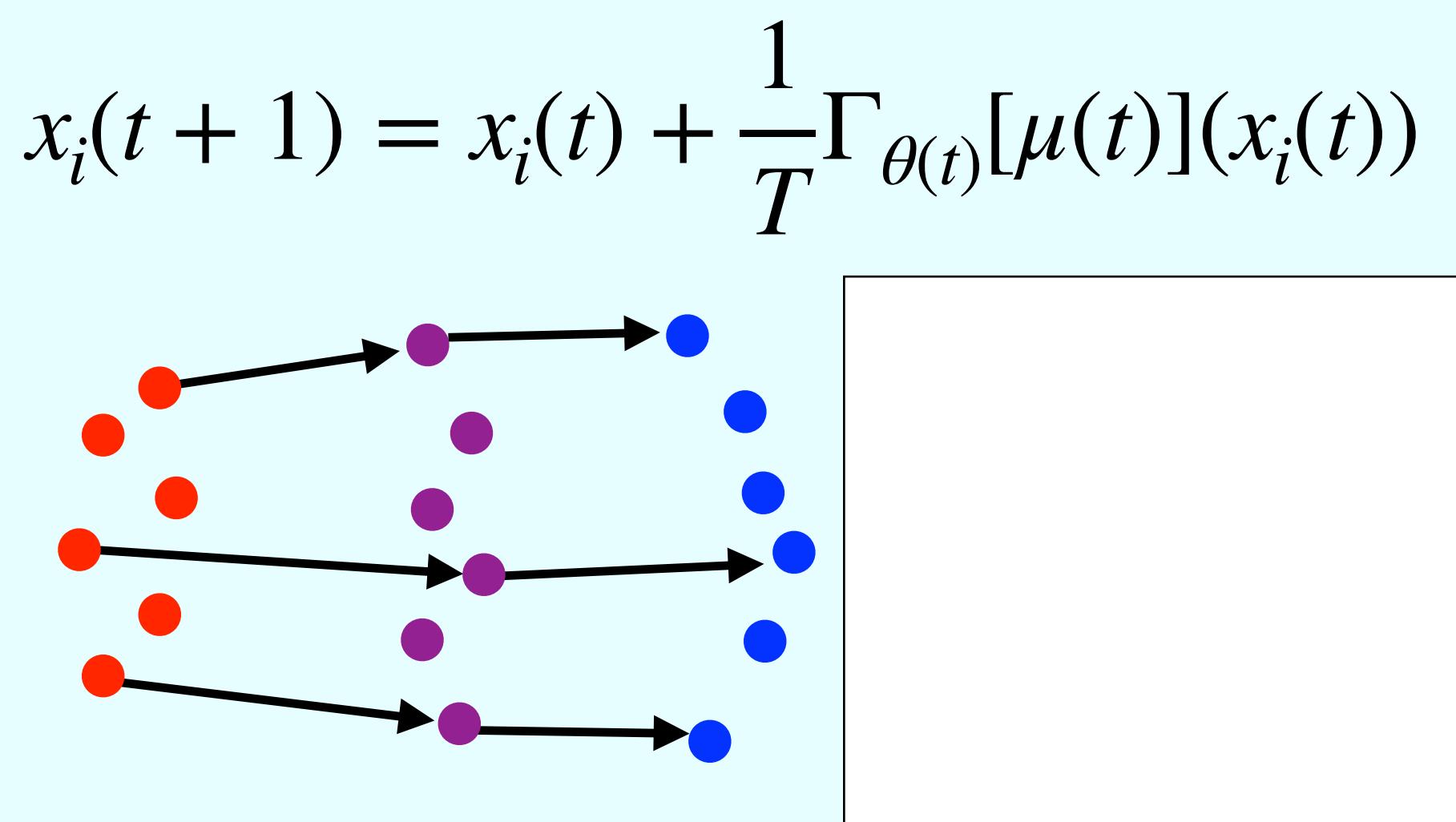


# Infinite Depth as a Neural PDE

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

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

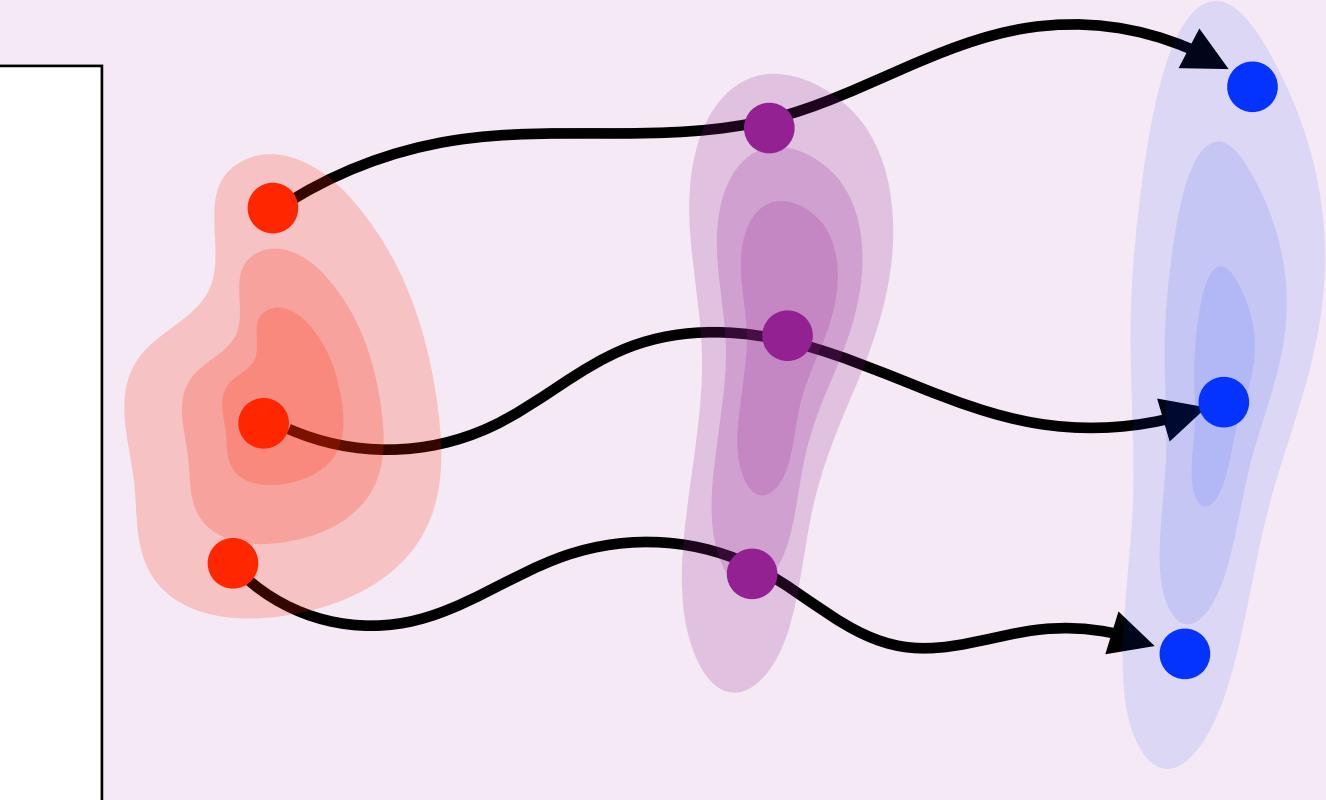
$$\mu(t) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$$



$T \rightarrow +\infty$   
Infinite depth

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

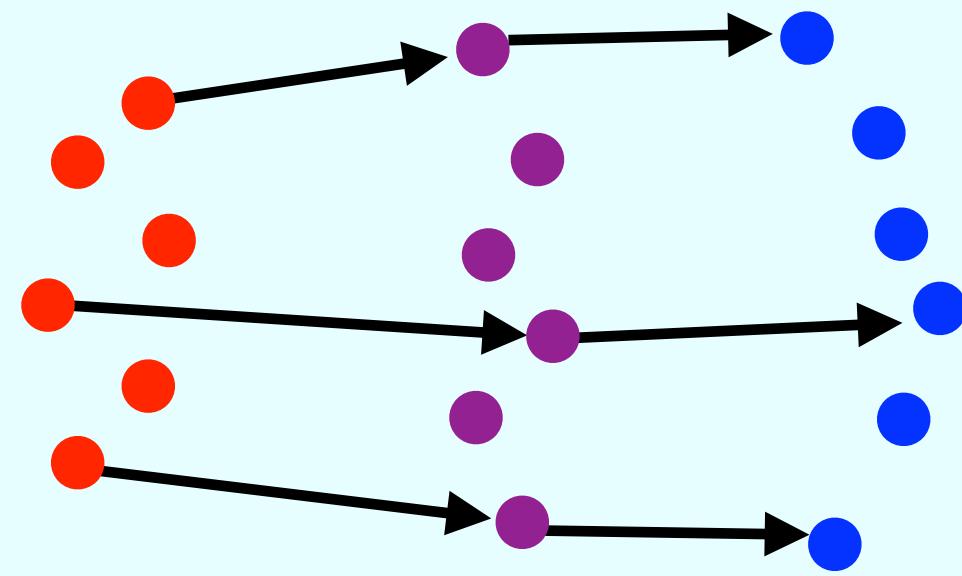
Coupled  
EDOs



# Infinite Depth as a Neural PDE

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

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



$T \rightarrow +\infty$   
Infinite depth

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

Coupled  
EDOs

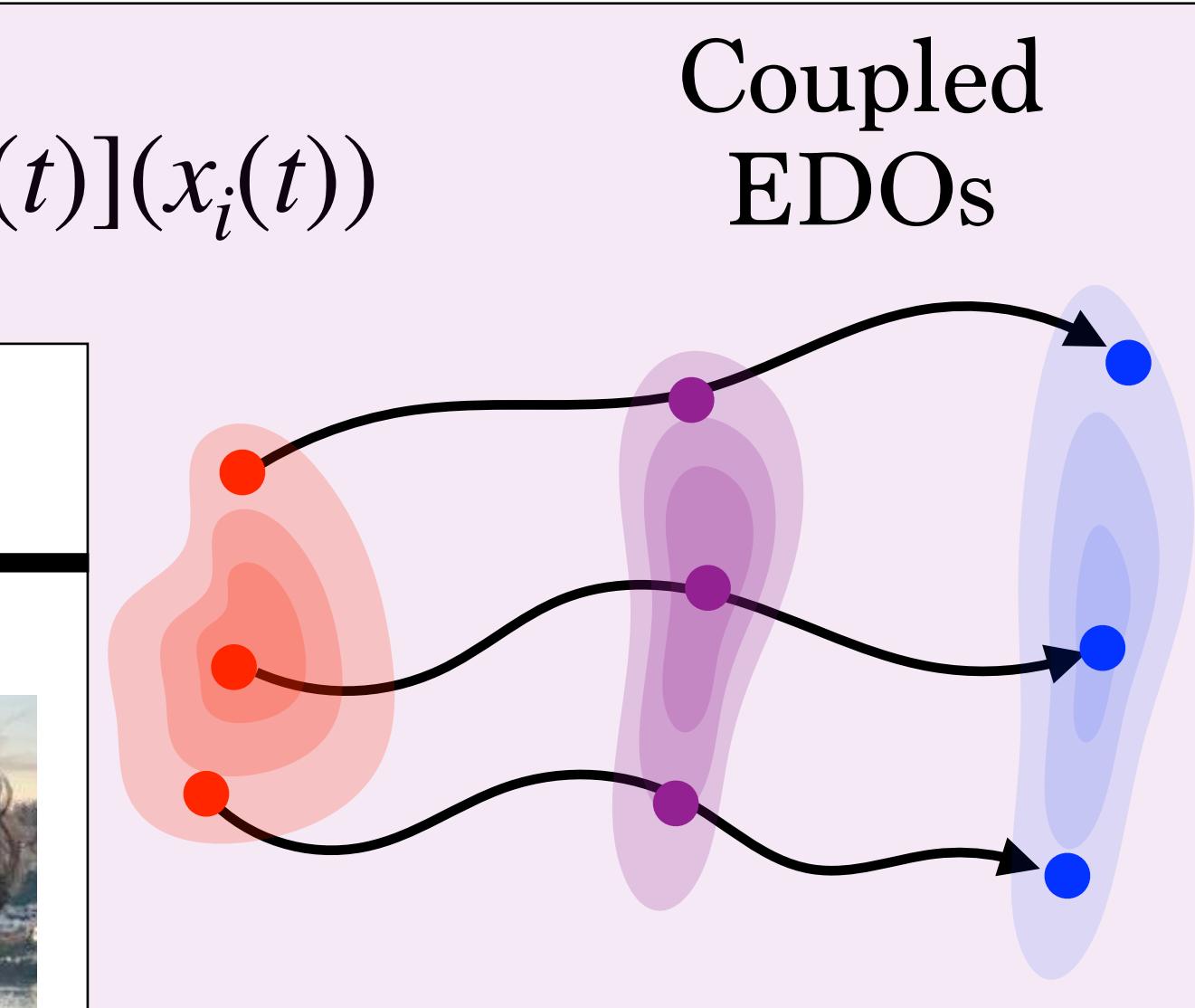
$$\frac{d\mu}{dt} + \text{div}(\mu \Gamma_\theta[\mu]) = 0$$

Transformer  
PDE

→ Not a Wasserstein flow.

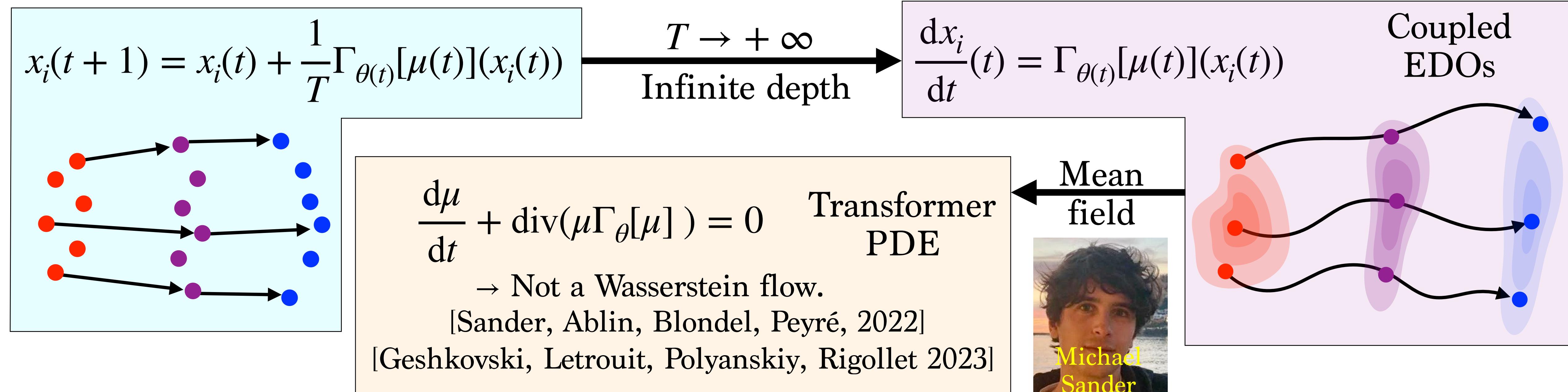
[Sander, Ablin, Blondel, Peyré, 2022]

[Geshkovski, Letrouit, Polyanskiy, Rigollet 2023]



# 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

# Gaussian Case and Clustering

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

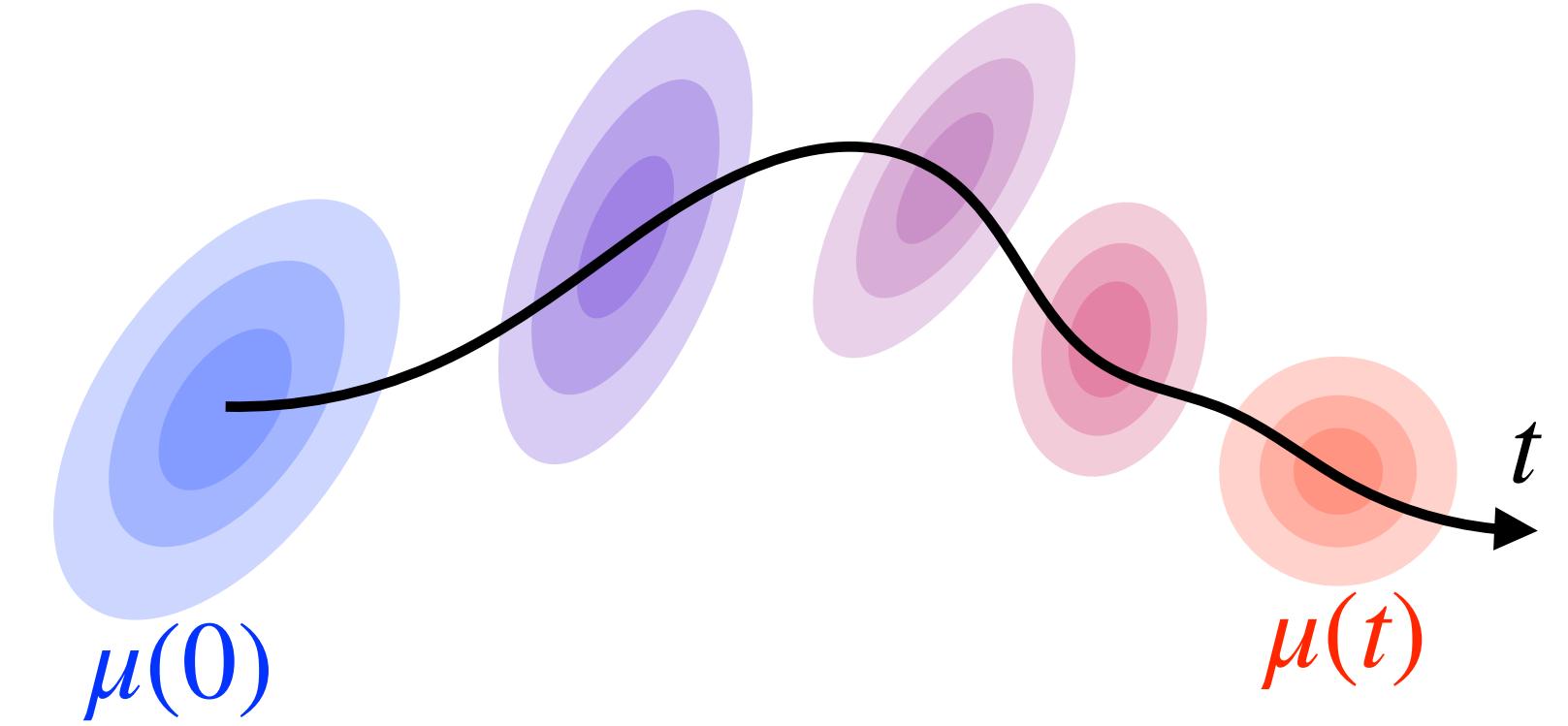
$$\frac{d\mu}{dt} + \operatorname{div}(\mu \Gamma_\theta[\mu]) = 0$$

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

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

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

$$\dot{\Sigma} = V\Sigma Q^\top K\Sigma + \Sigma K^\top Q\Sigma V^\top$$



# Gaussian Case and Clustering

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

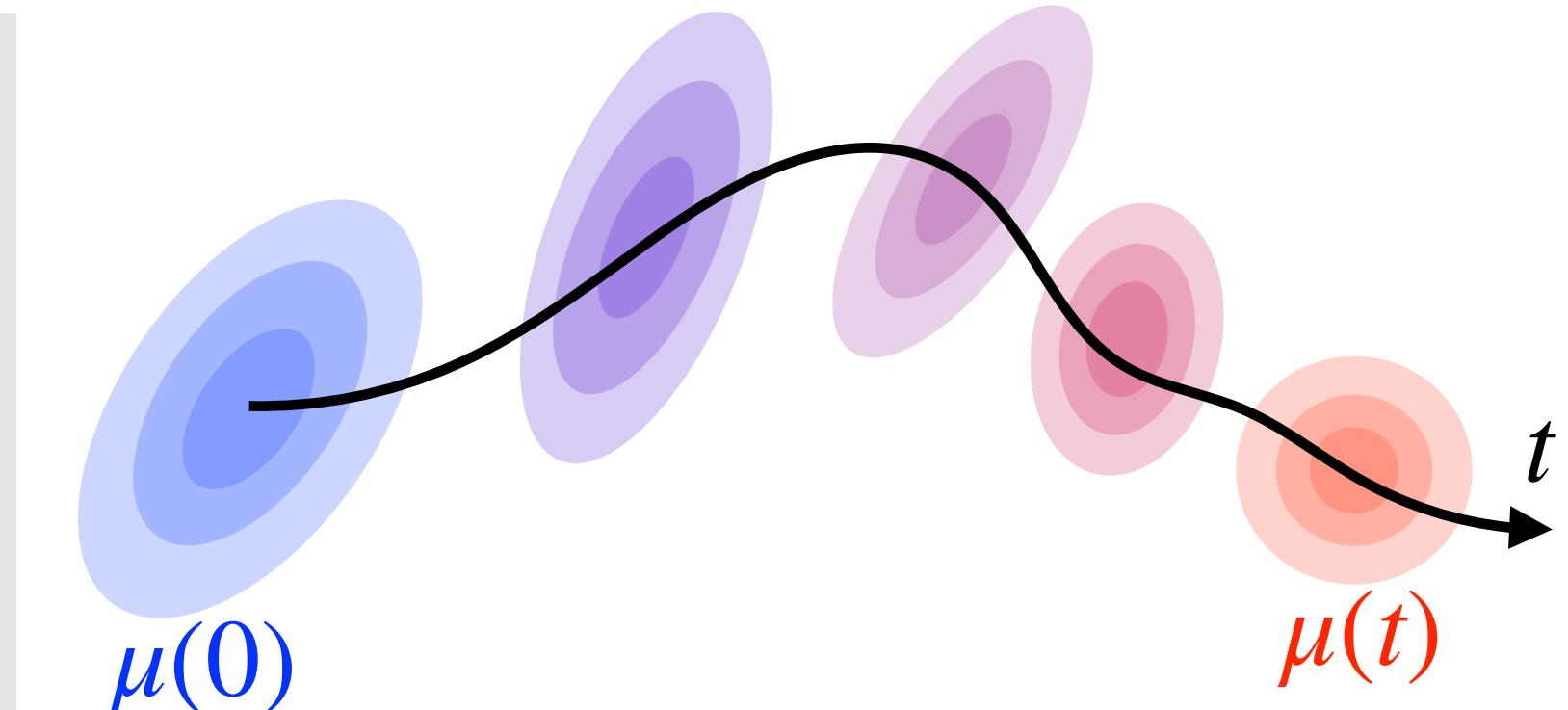
$$\theta(t) = (Q(t), K(t), V(t))$$

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

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

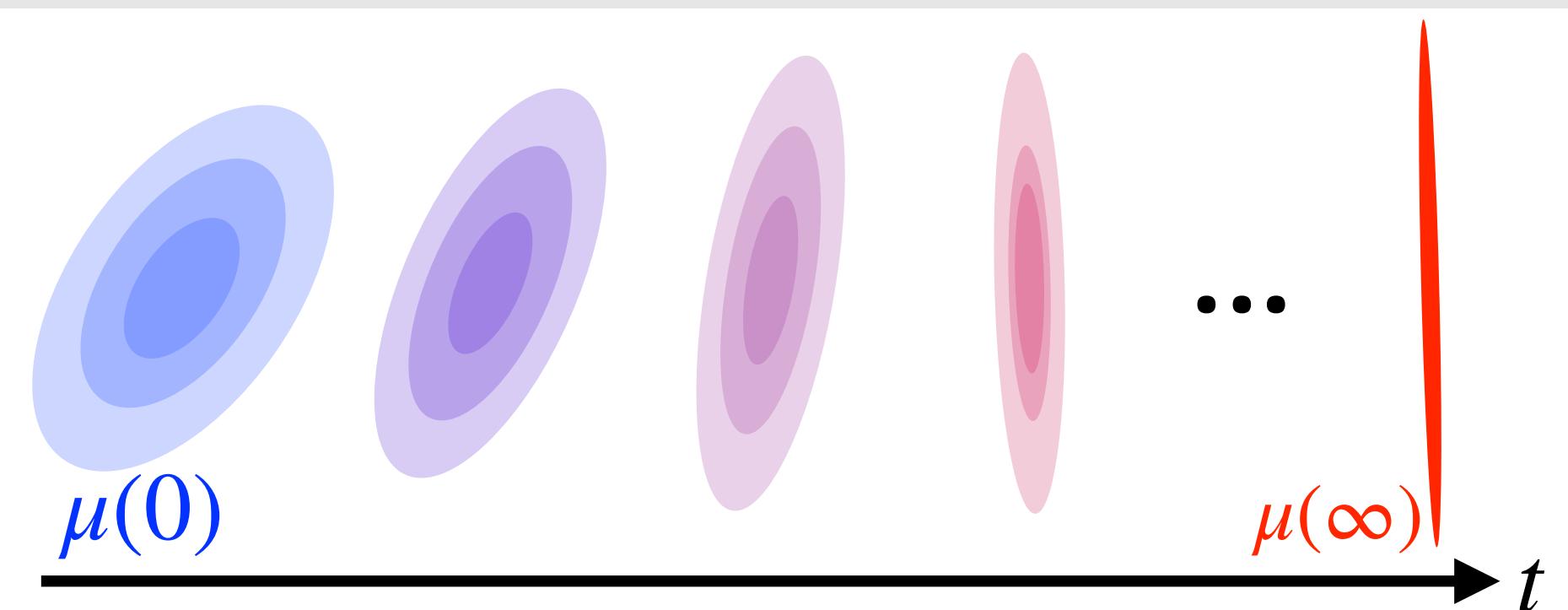
then  $\mu(s) = \mathcal{N}(\mathbf{m}(s), \Sigma(s))$

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



*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.

# Open Problems and their Solutions

*Universality:* Toward quantitative approximation bound, leverage smoothness.

*Optimisation:* Understand the structure of optimal  $(Q, K, V)$

Getting Albert interested in my problems

