

Attention with Markov: A Markovian Tale of Transformers

Ashok Vardhan Makkuva

Joint work with Marco Bondaschi, Nived Rajaraman, Adway Girish, Alliot Nagle, Martin Jaggi, Hyeji Kim, and Michael Gastpar



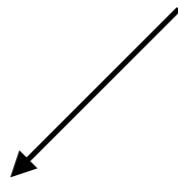
Attention with Markov

Attention with Markov

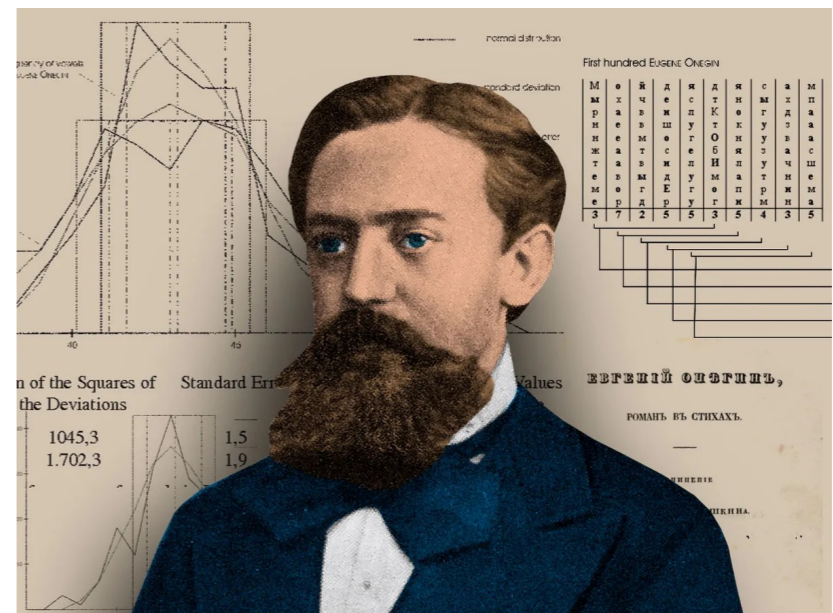
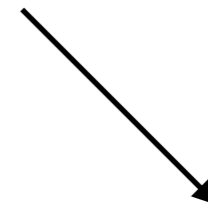
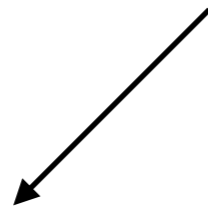
Attention with Markov

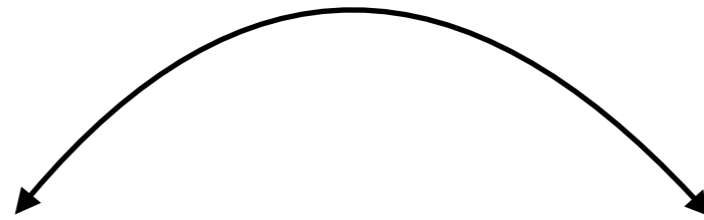
Attention

Attention

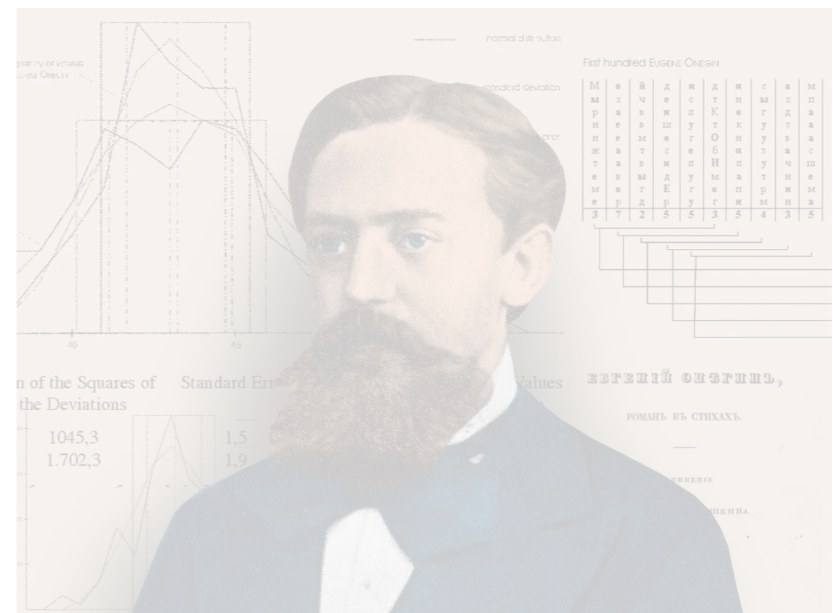
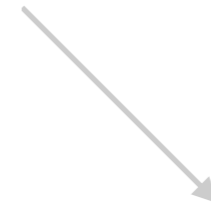
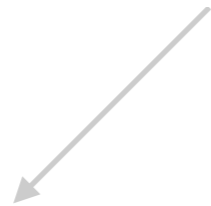


Attention with Markov

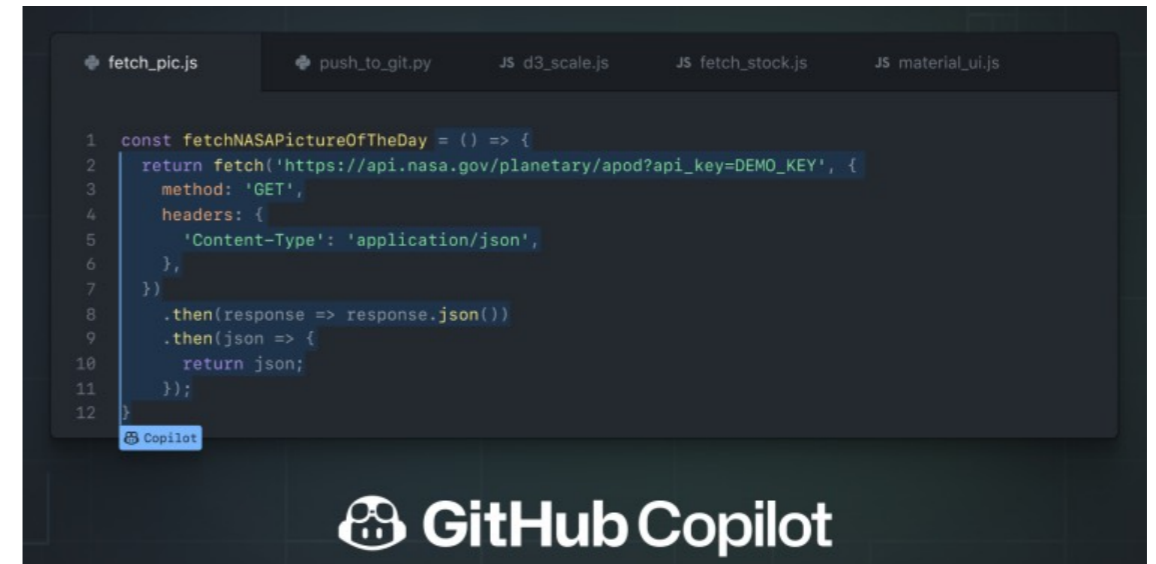




Attention with Markov



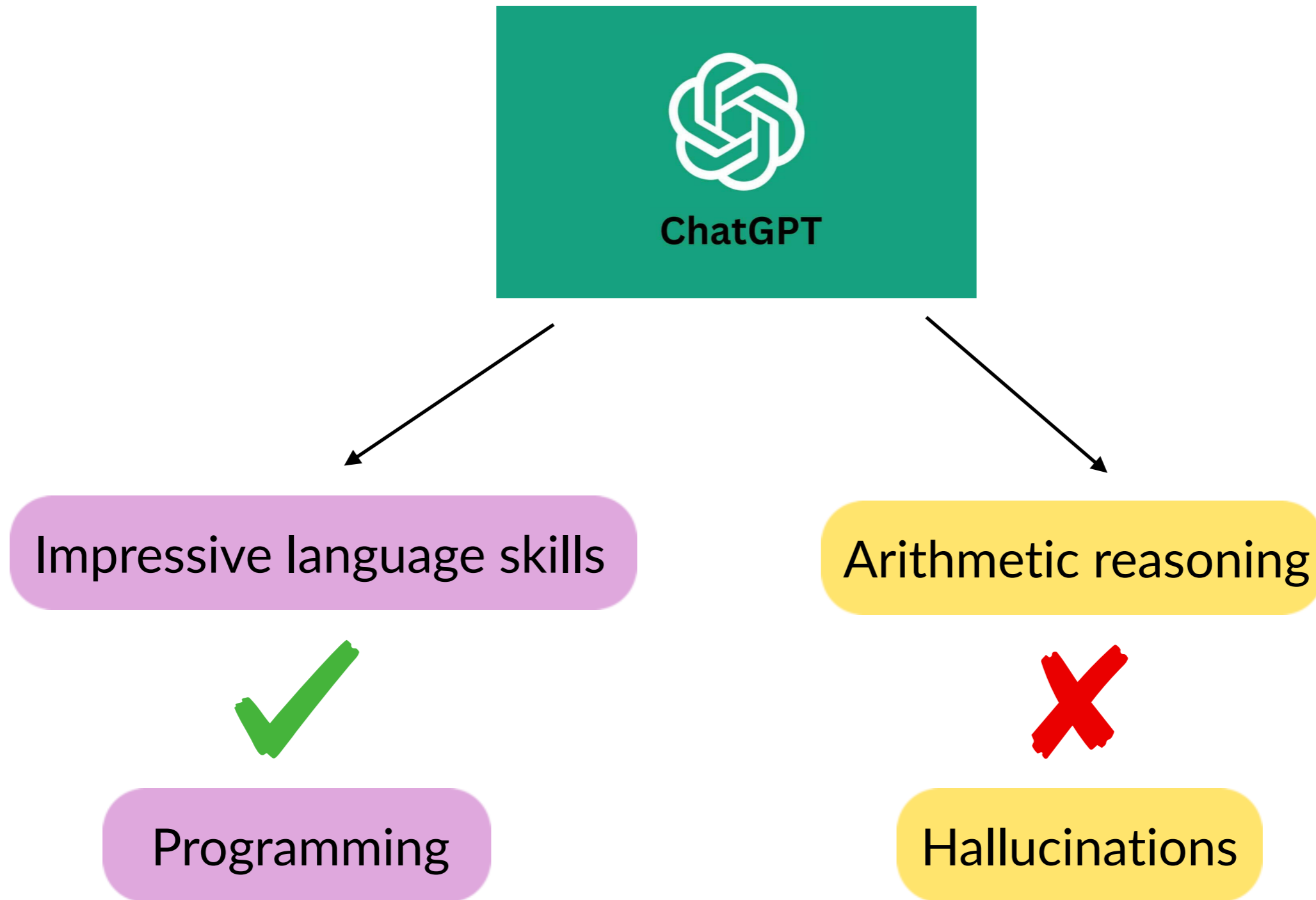
LLMs are part of daily lives



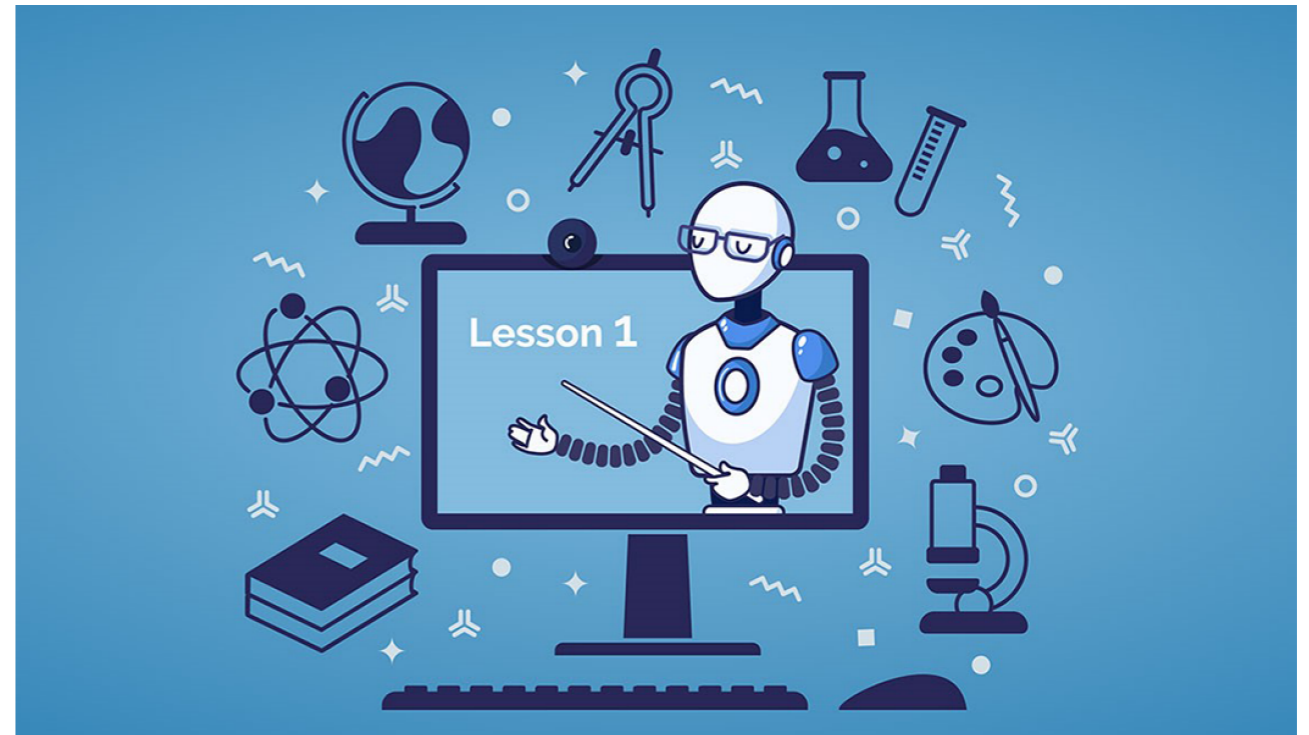
The good and the bad



The good and the bad



Critical domains



Need of the hour



Fundamental understanding



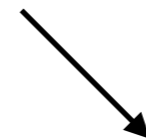
Fundamental understanding

What do they learn?

How do they learn?



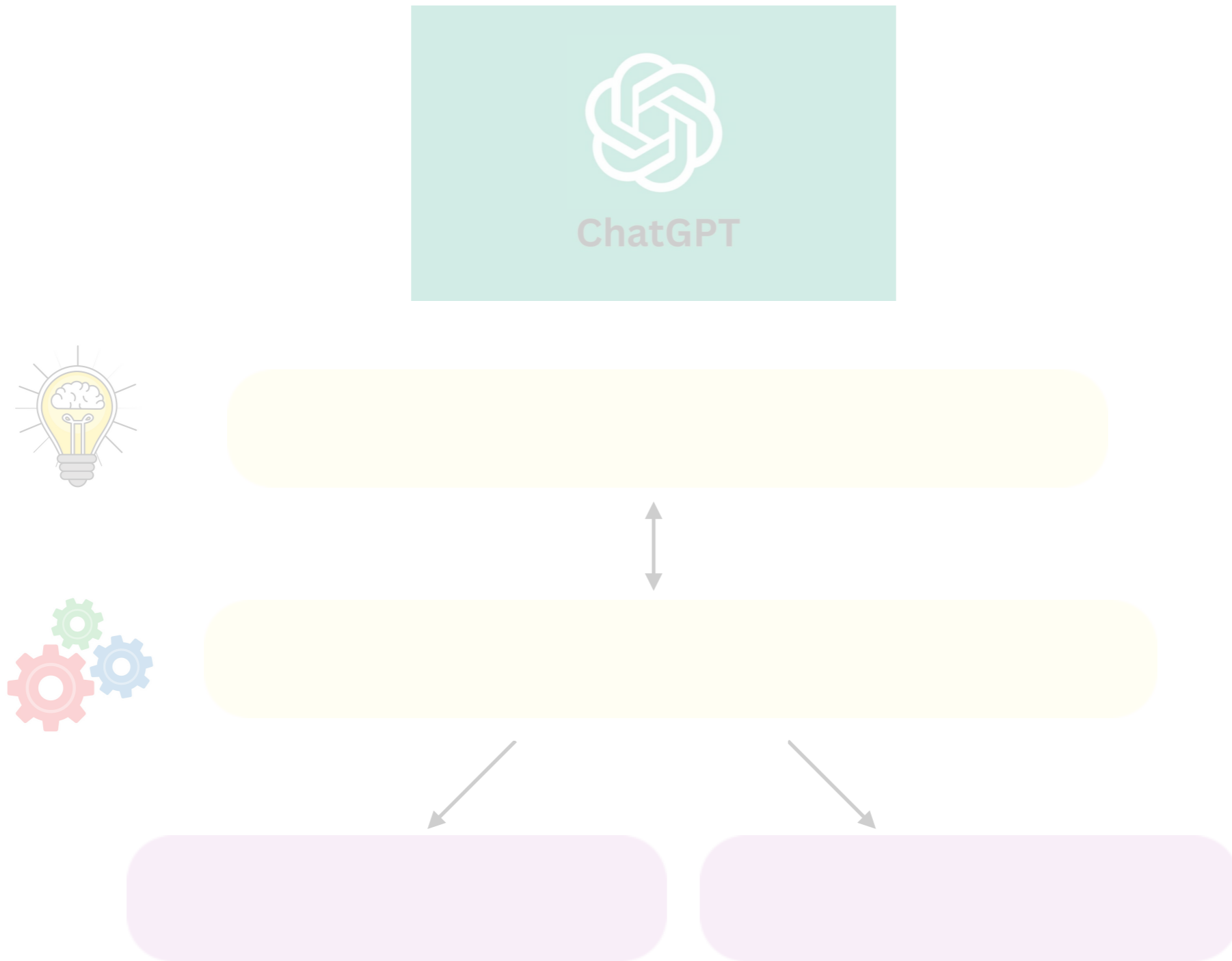
Principled frameworks and tools



What do they learn?

How do they learn?

Challenges

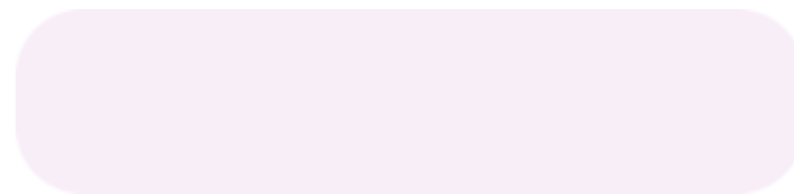
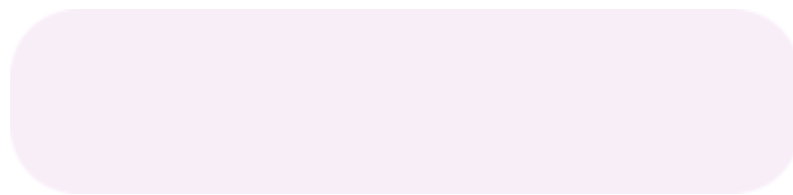
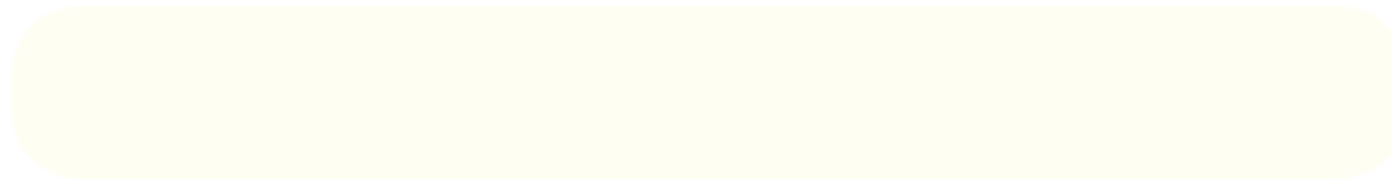
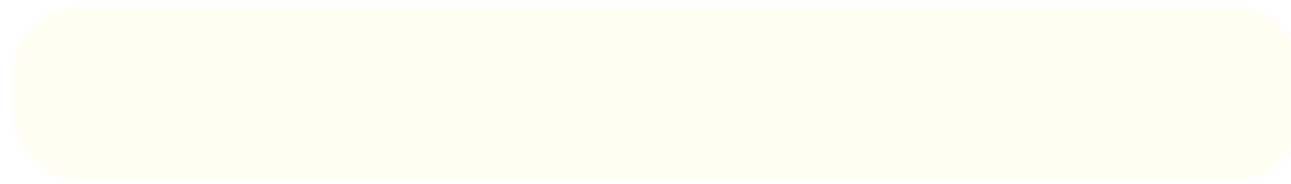


Challenges

Inherently complex



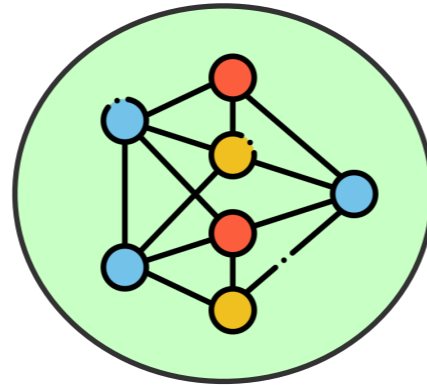
Mathematically intractable



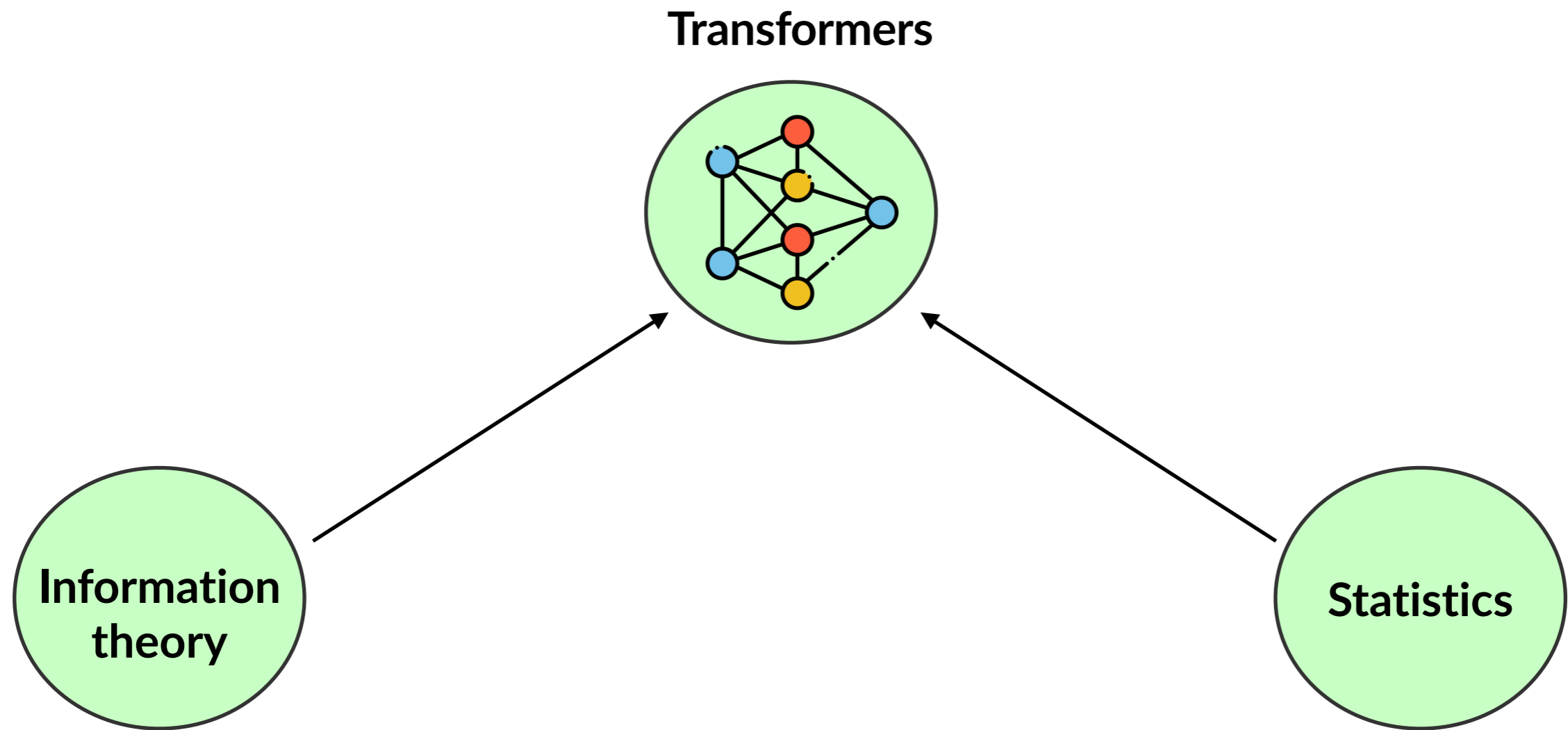
My research

My research

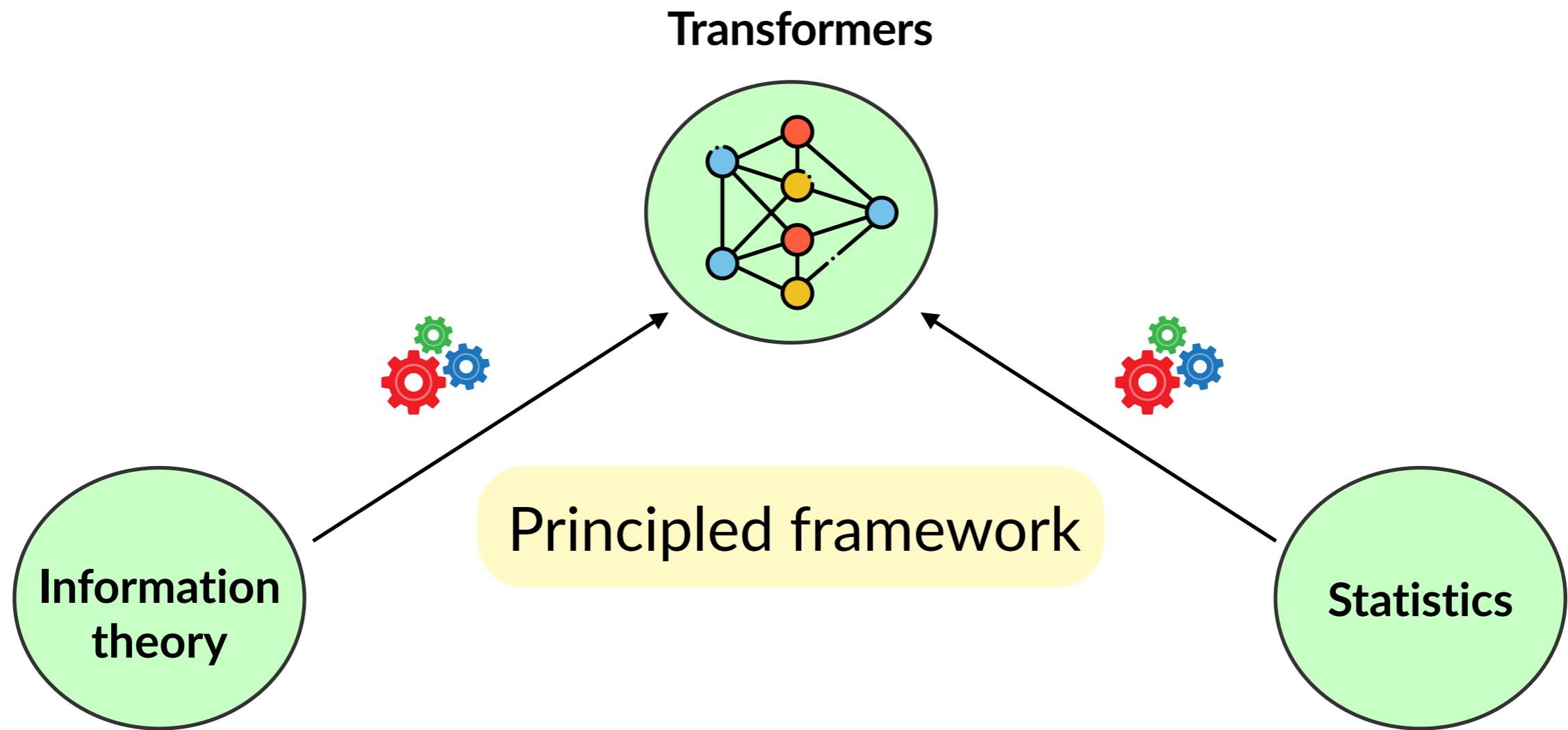
Transformers



My research



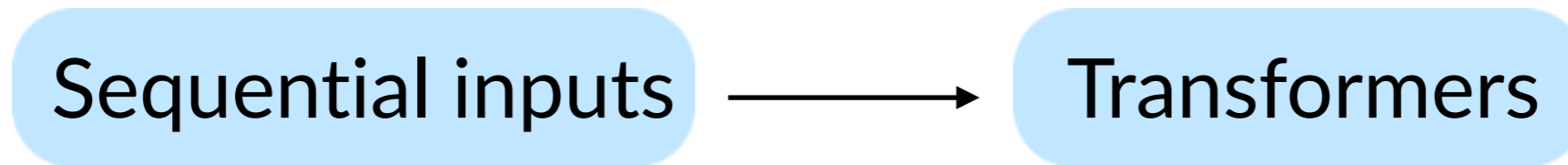
My research





Attention with Markov

Attention with Markov



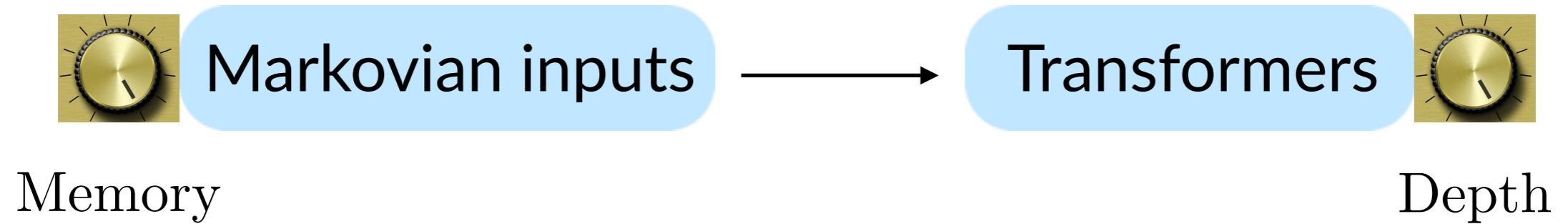
Attention with Markov

Markovian inputs

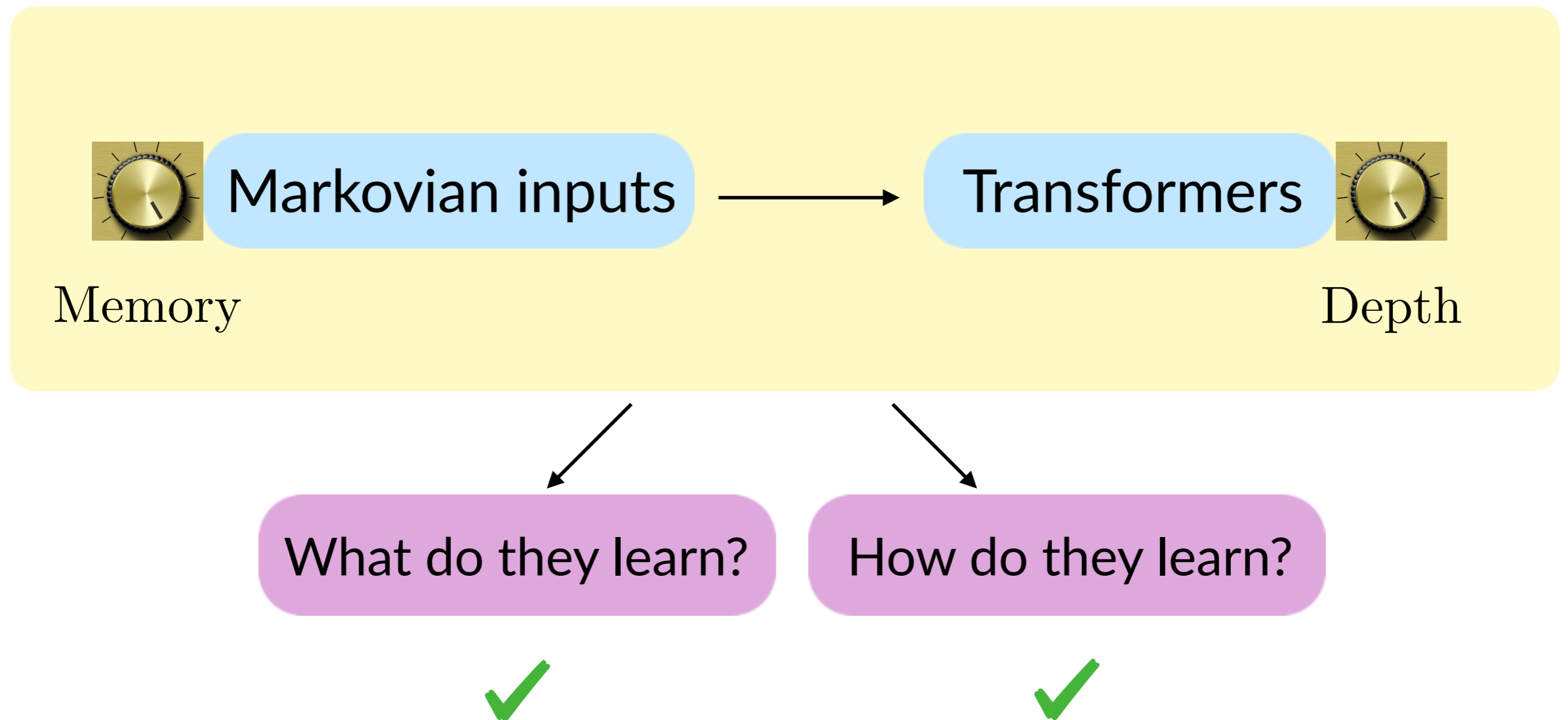


Transformers

Attention with Markov



Attention with Markov



Depth plays a crucial role in transformer functionality



Memory

1-layer transformer sometimes fails to learn even first-order Markov chains! *

Depth

3-layer transformer can learn Markov chains of all orders

This talk

Local to Global: Learning Dynamics and Effect of Initialization for Transformers

NeurIPS 2024

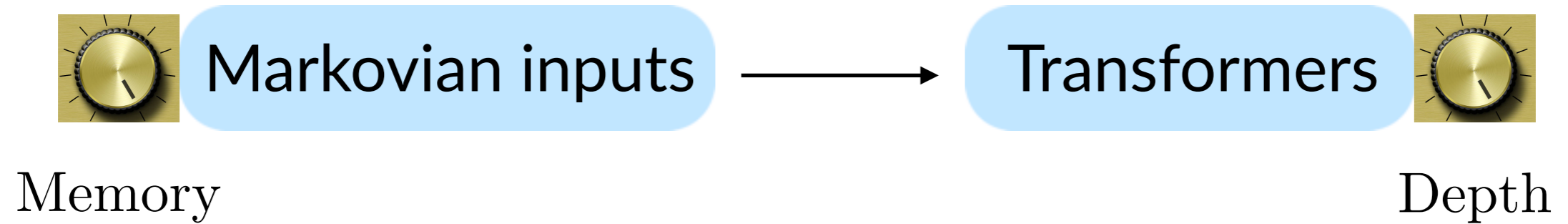
Transformers on Markov Data: Constant Depth Suffices

NeurIPS 2024

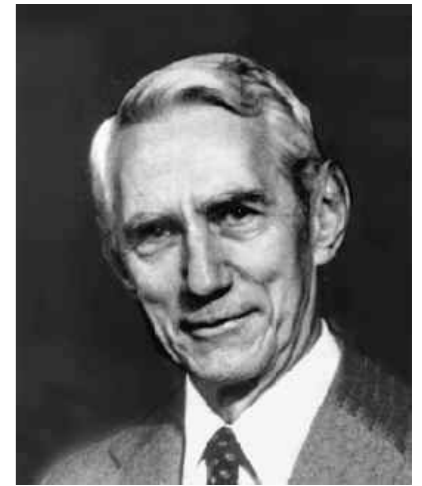
Attention with Markov: A Curious Case of Single-layer Transformers

*ICML 2024 Workshop,
Under review ICLR 2025
(All hail Reviewer #2)*

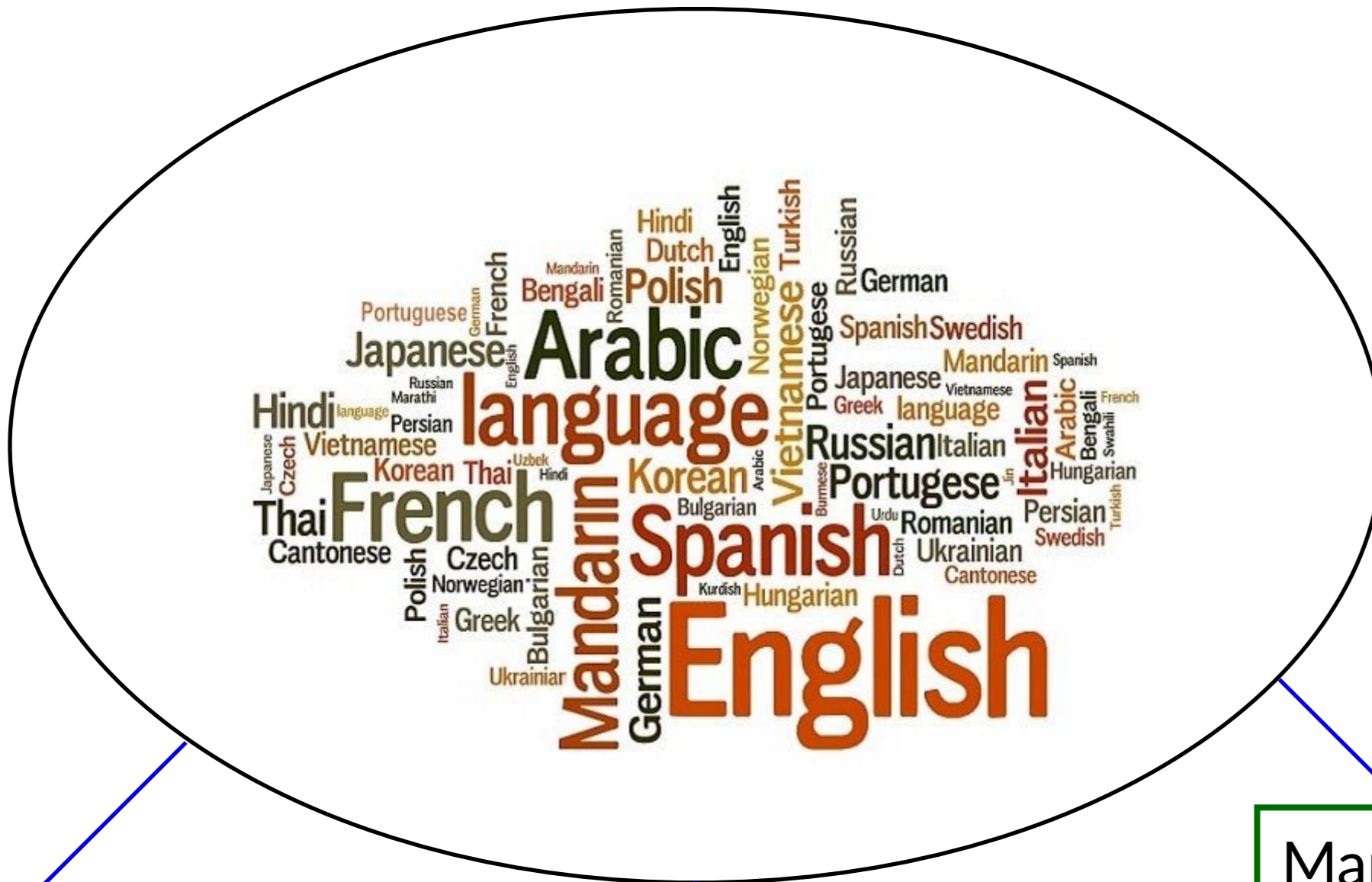
Attention with Markov



Why Markovian?



Shannon, 1948



Grammar

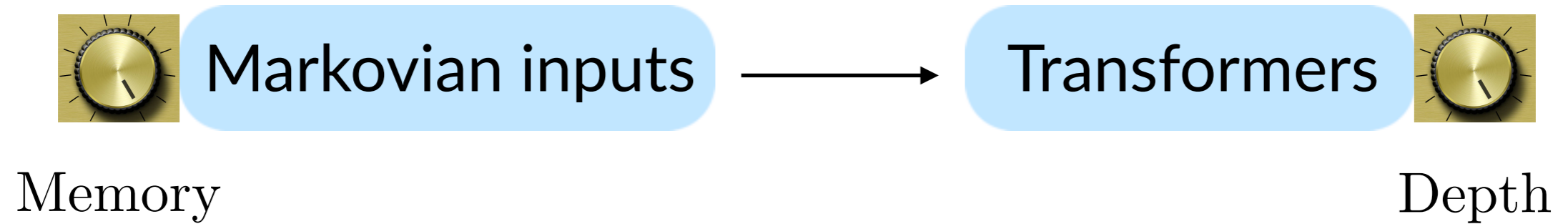
Syntax

Markovian

- Blue
- Black

Bl

Attention with Markov



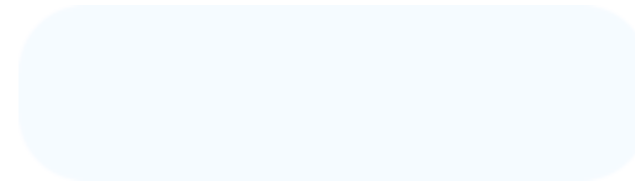
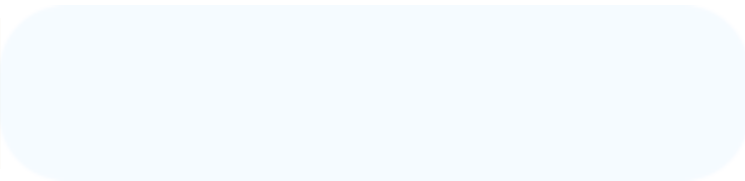
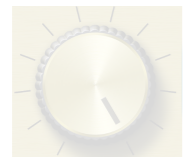
Part I



Memory = 1

Depth = 1

Part II



Memory

Depth > 1

Part I



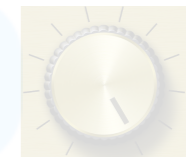
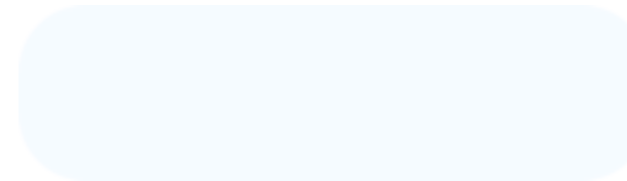
Memory = 1

Depth = 1

Part I



Markovian inputs

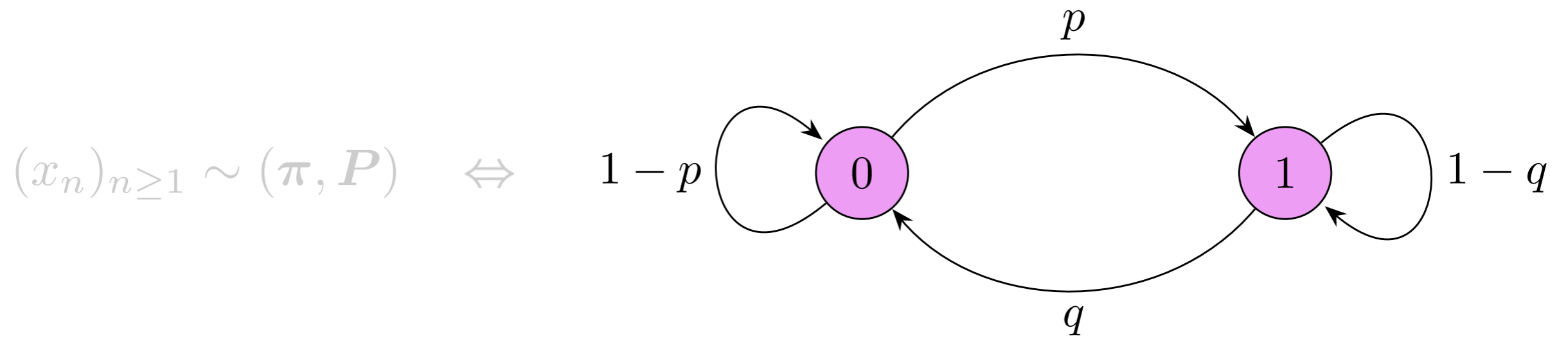


Memory = 1

Depth = 1

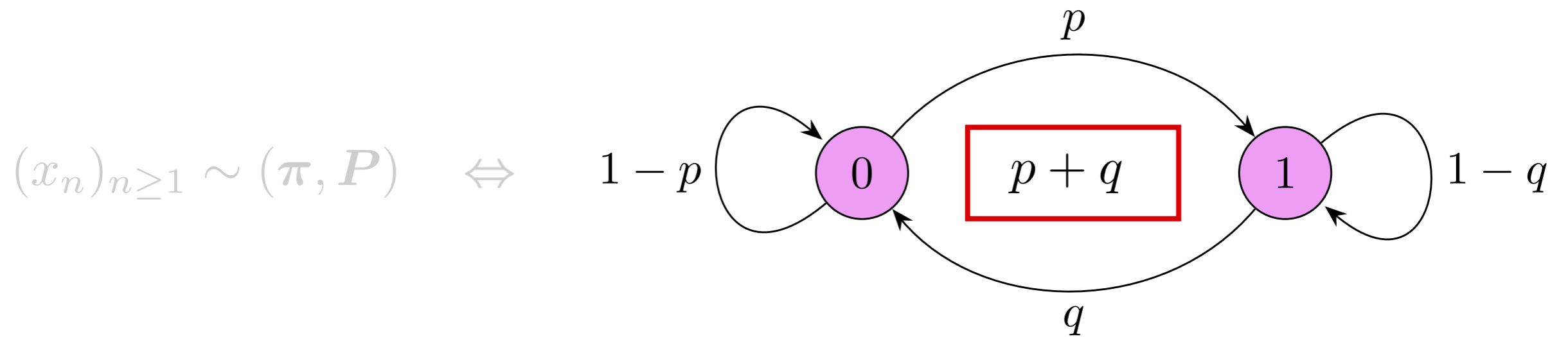
Input data: First-order Markov chain

Input data: First-order Markov chain



$$\pi = (\pi_0, \pi_1) = \left(\frac{q}{p+q}, \frac{p}{p+q} \right), \quad P = (P_{ij}) = \begin{bmatrix} 1-p & p \\ q & 1-q \end{bmatrix}.$$

Input data: First-order Markov chain

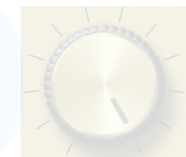
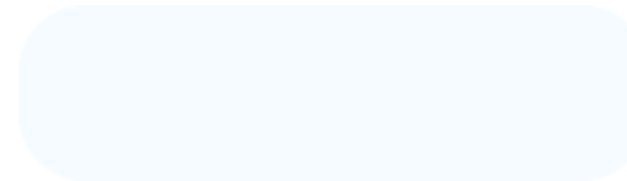


$$\pi = (\pi_0, \pi_1) = \left(\frac{q}{p+q}, \frac{p}{p+q} \right), \quad P = (P_{ij}) = \begin{bmatrix} 1-p & p \\ q & 1-q \end{bmatrix}.$$

Part I



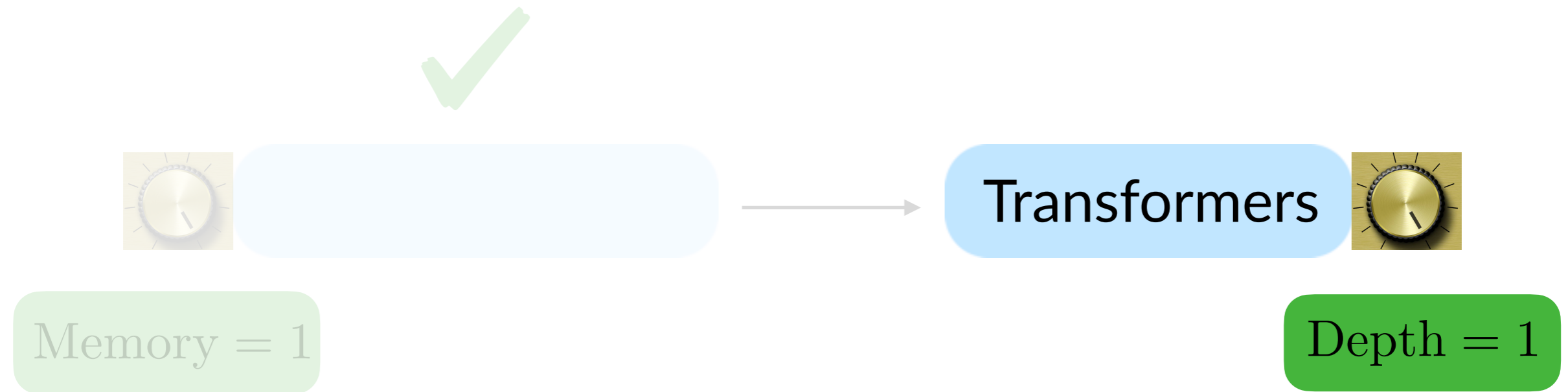
Markovian inputs



Memory = 1

Depth = 1

Part I



Transformers

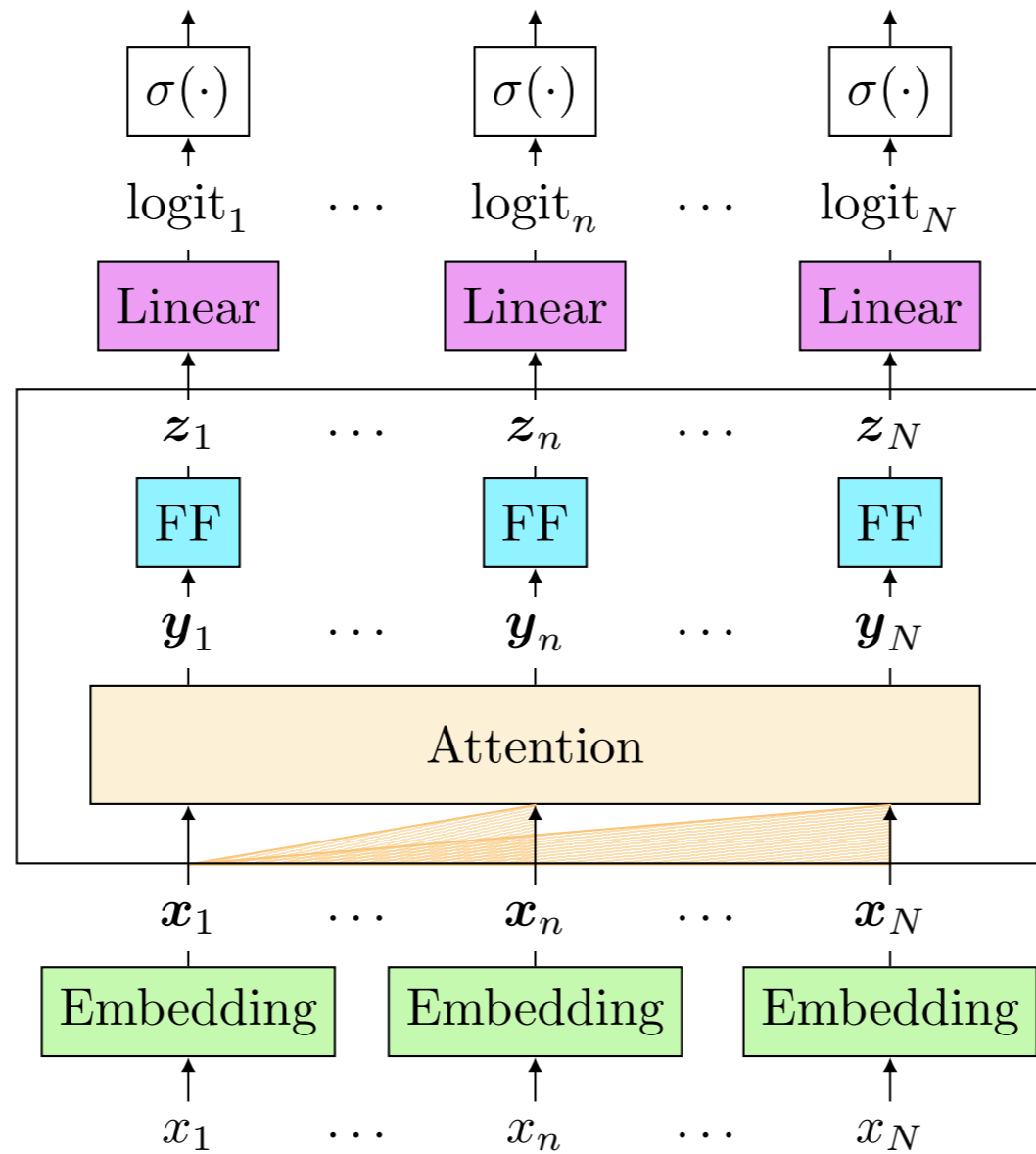
Transformers



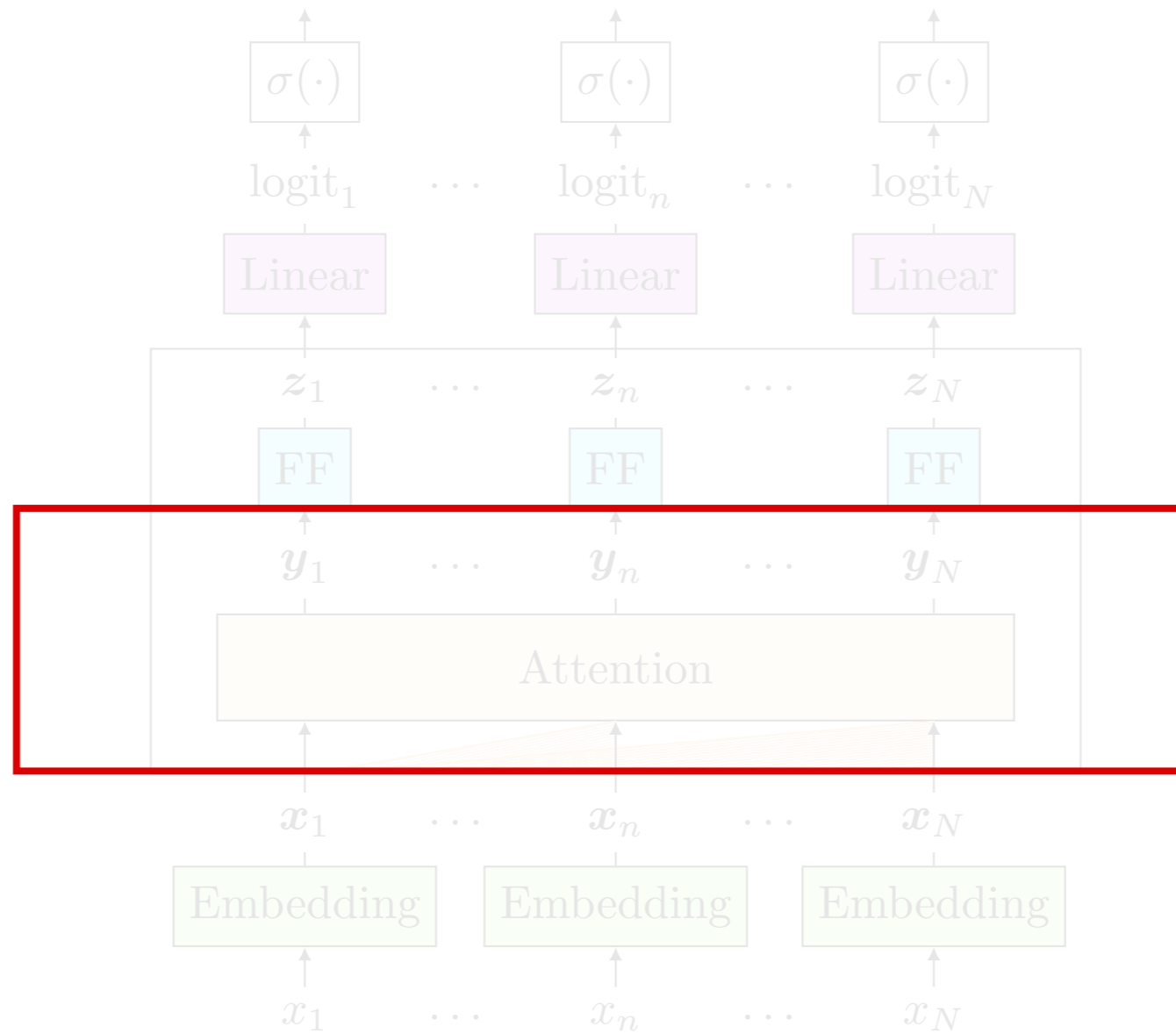
Transformers



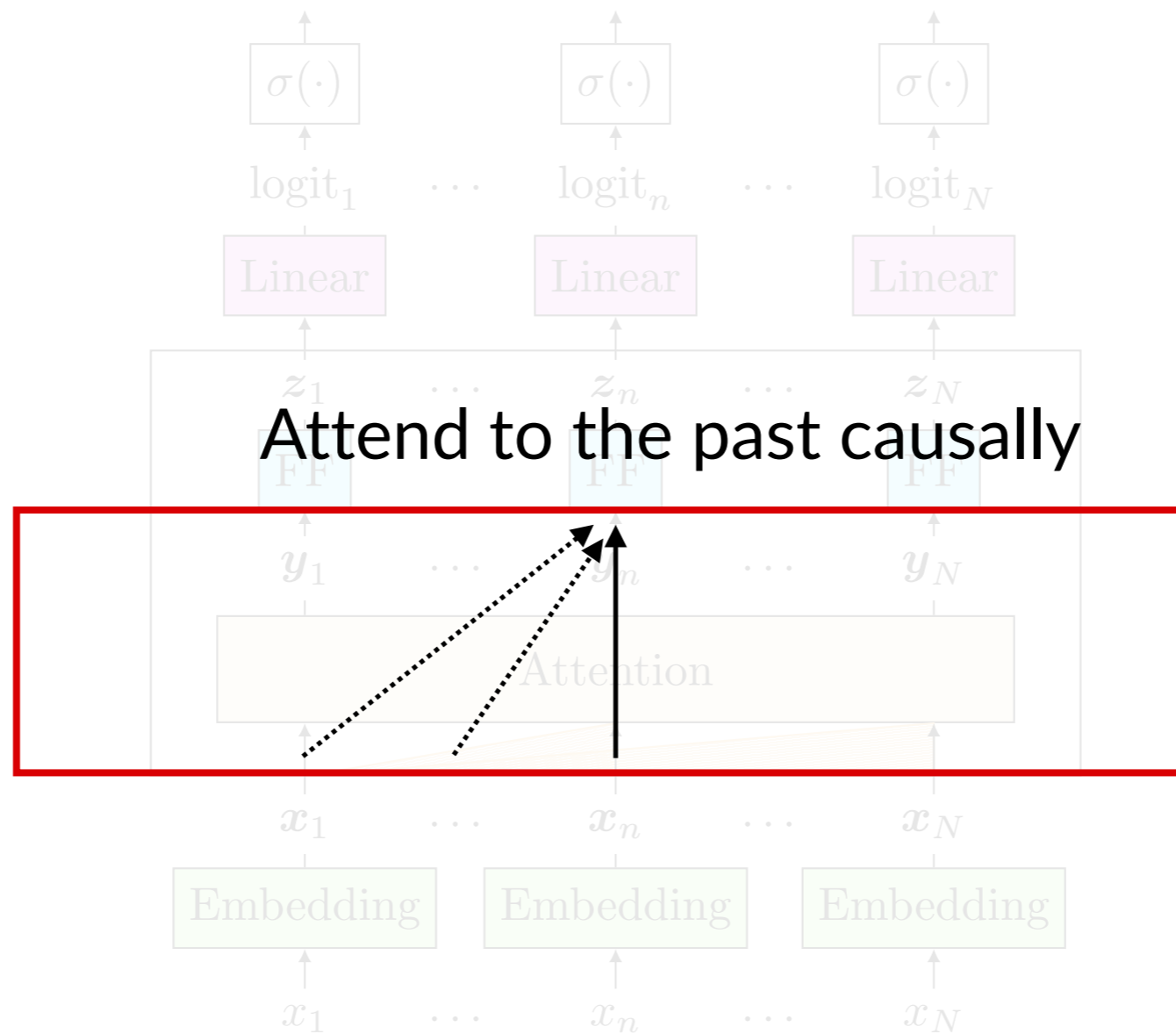
Transformers



Transformers



Transformers



Single-layer transformer

$$f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n)$$

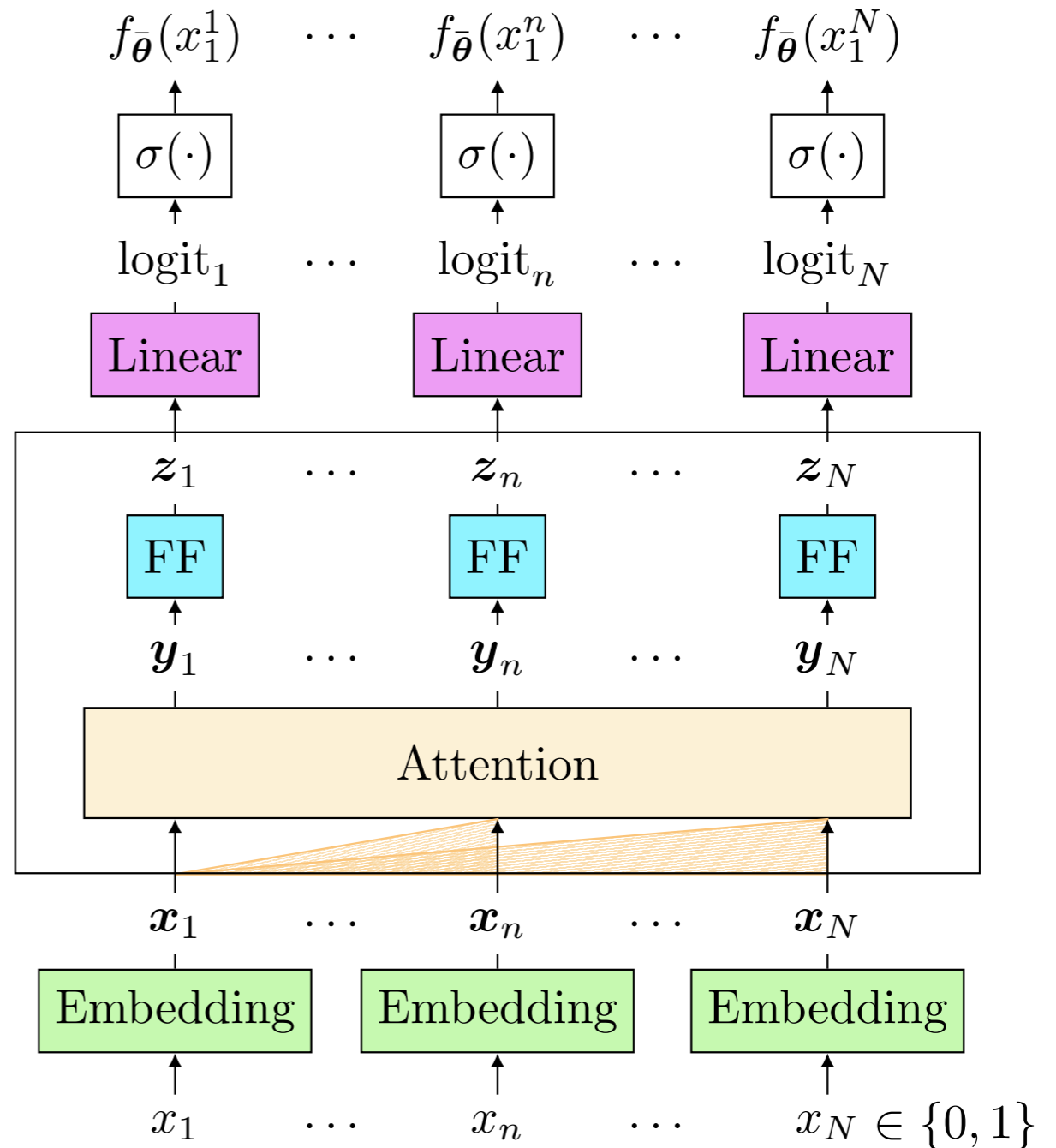
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$

$$\mathbf{z}_n = \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n)$$

$$\mathbf{y}_n = \mathbf{x}_n + \mathbf{W}_O \sum_{i=1}^n \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i$$

$$\mathbf{x}_n = x_n \cdot \mathbf{e} + \mathbf{p}_n$$

θ



Single-layer transformer

$$f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n)$$

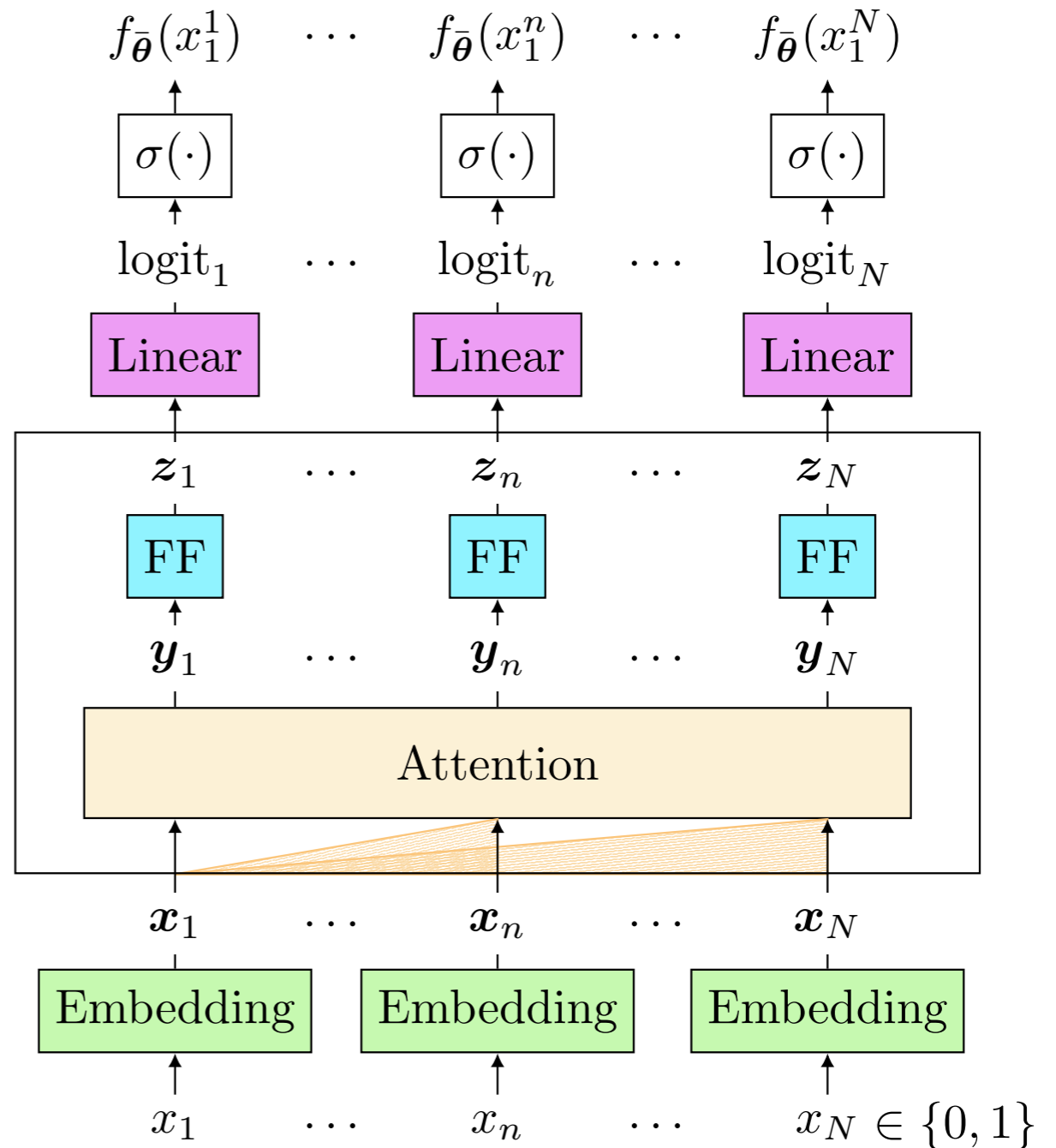
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$

$$\mathbf{z}_n = \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n)$$

$$\mathbf{y}_n = \mathbf{x}_n + \mathbf{W}_O \sum_{i=1}^n \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i$$

$$\mathbf{x}_n = x_n \cdot \mathbf{e} + \mathbf{p}_n$$

θ



Single-layer transformer

$$f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n)$$

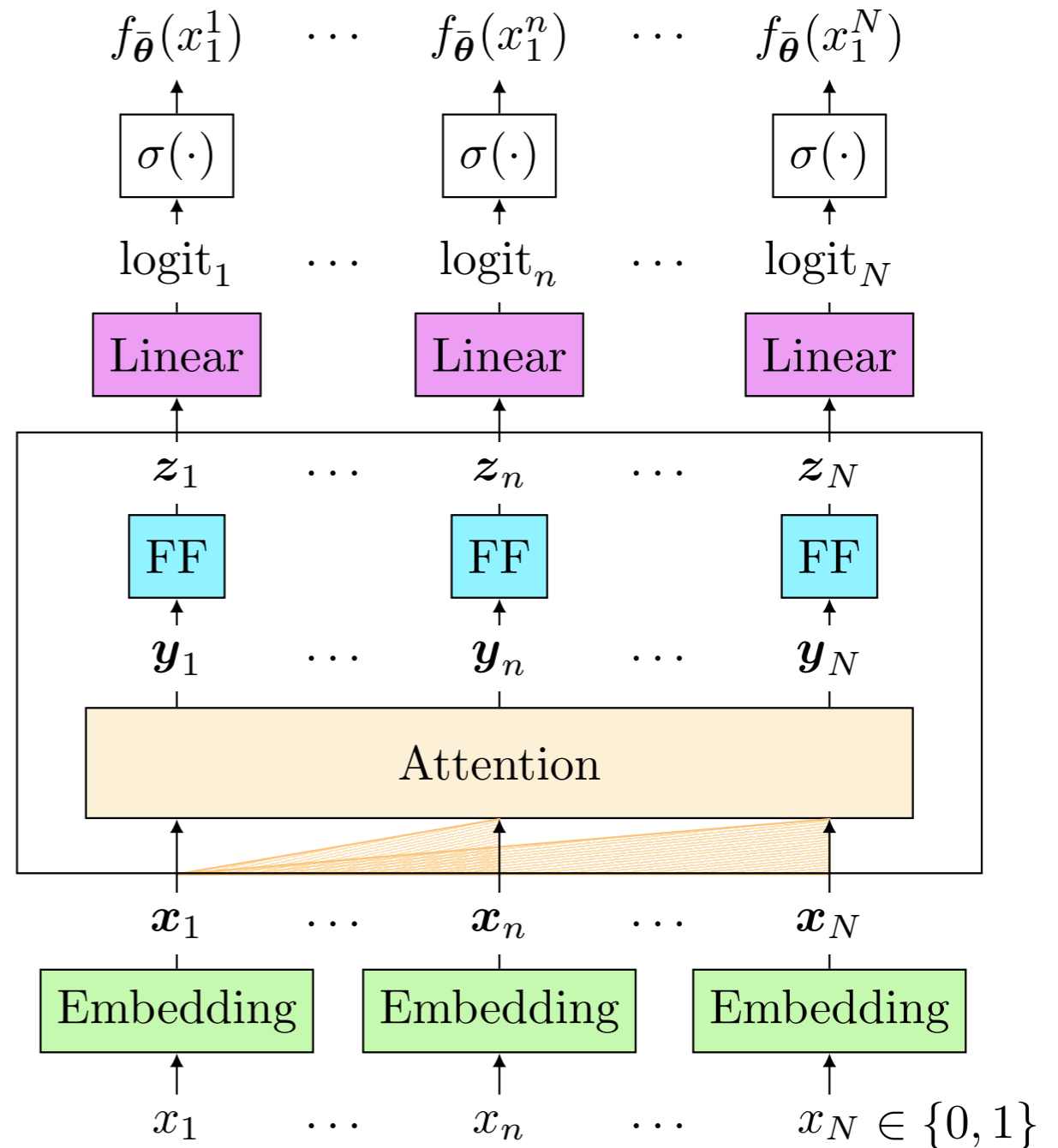
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$

$$\mathbf{z}_n = \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n)$$

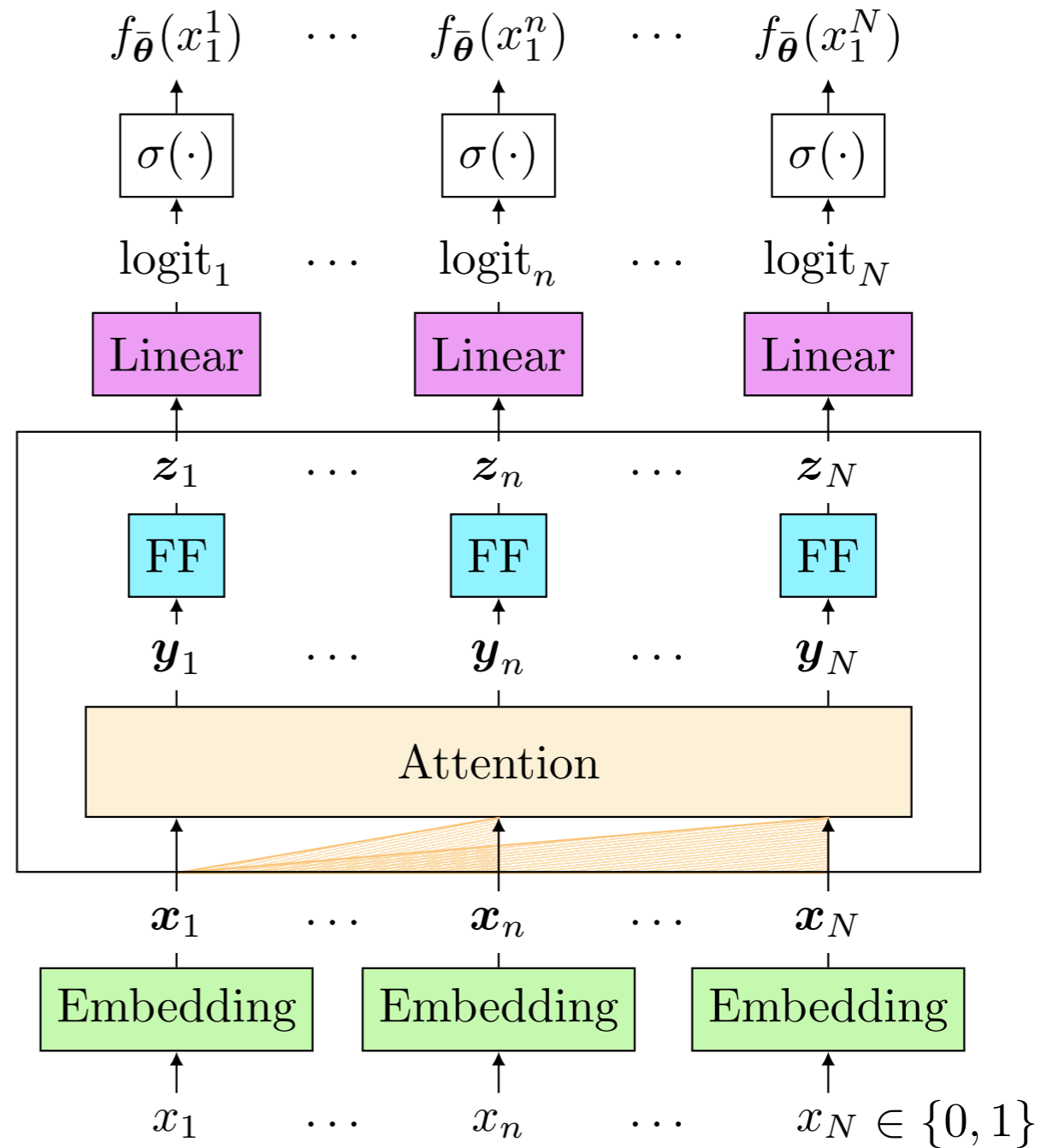
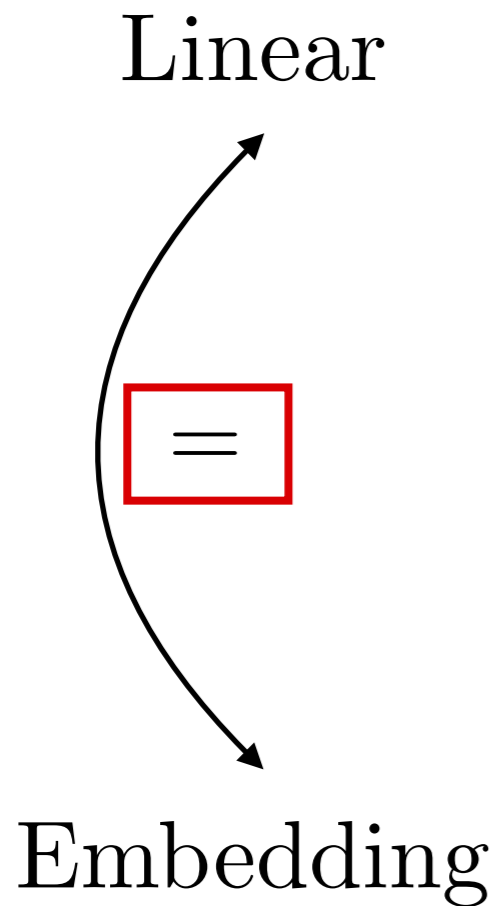
$$\mathbf{y}_n = \mathbf{x}_n + \mathbf{W}_O \sum_{i=1}^n \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i$$

$$\mathbf{x}_n = x_n \cdot \mathbf{e} + \mathbf{p}_n$$

θ



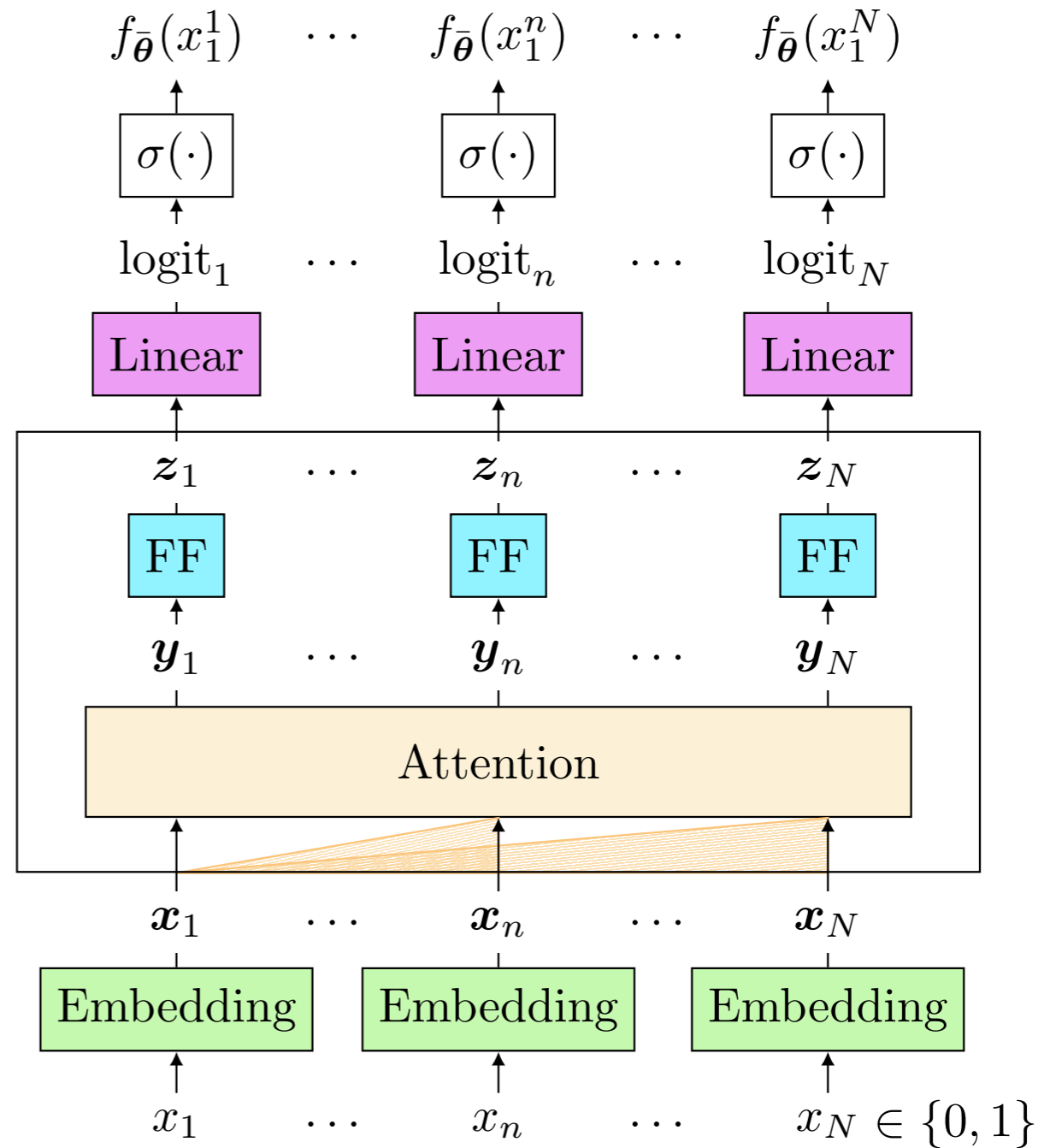
Weight tying



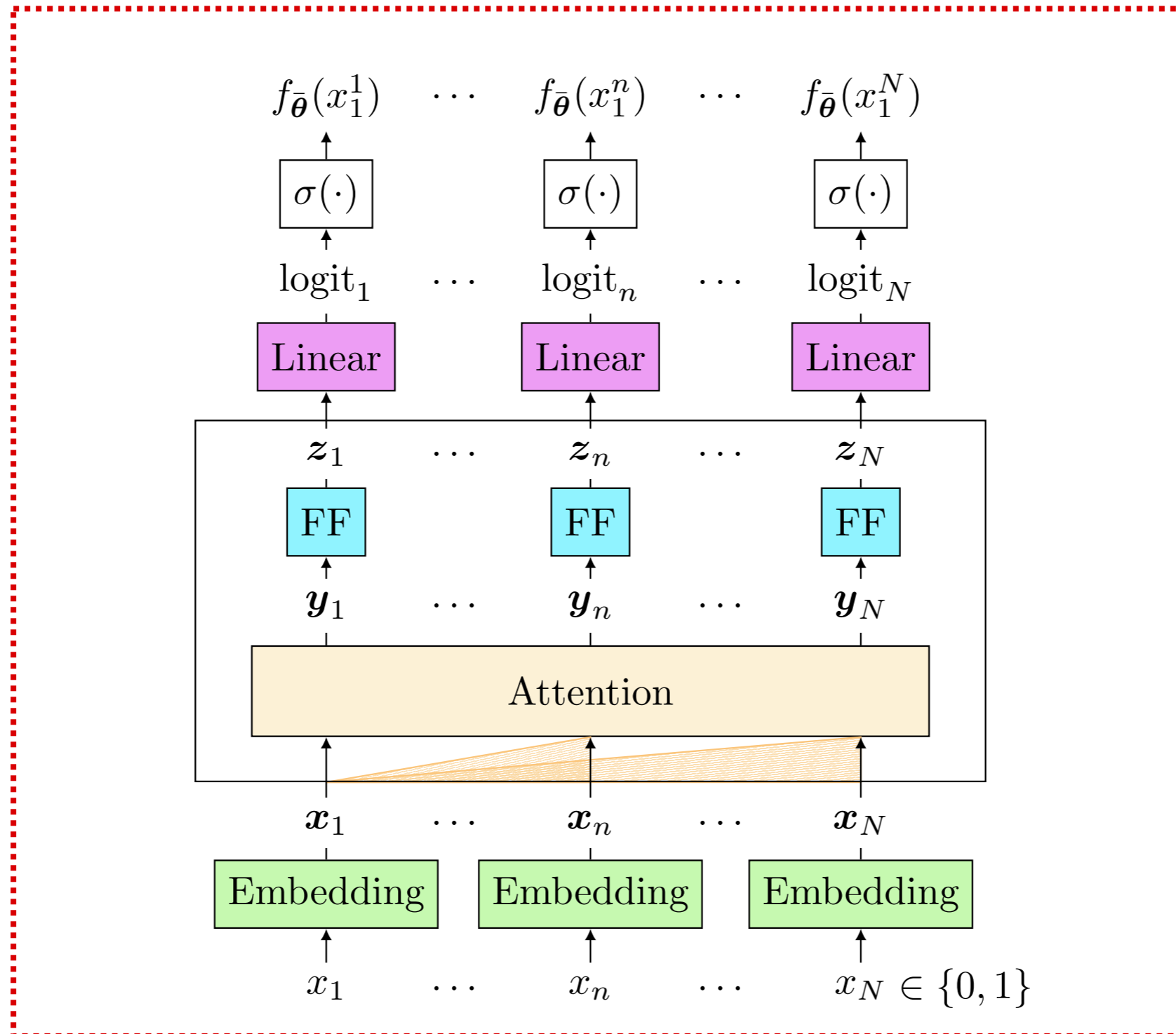
No weight tying

Linear

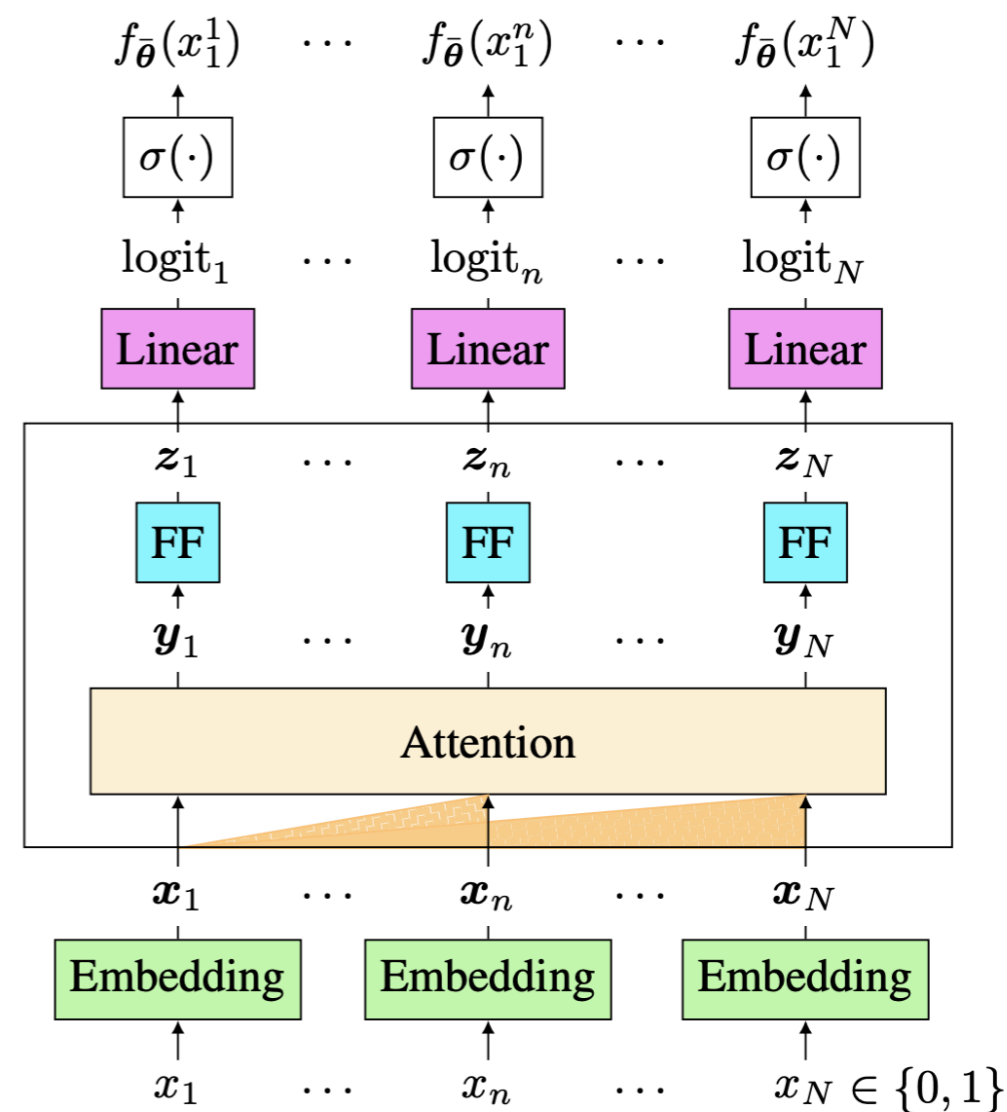
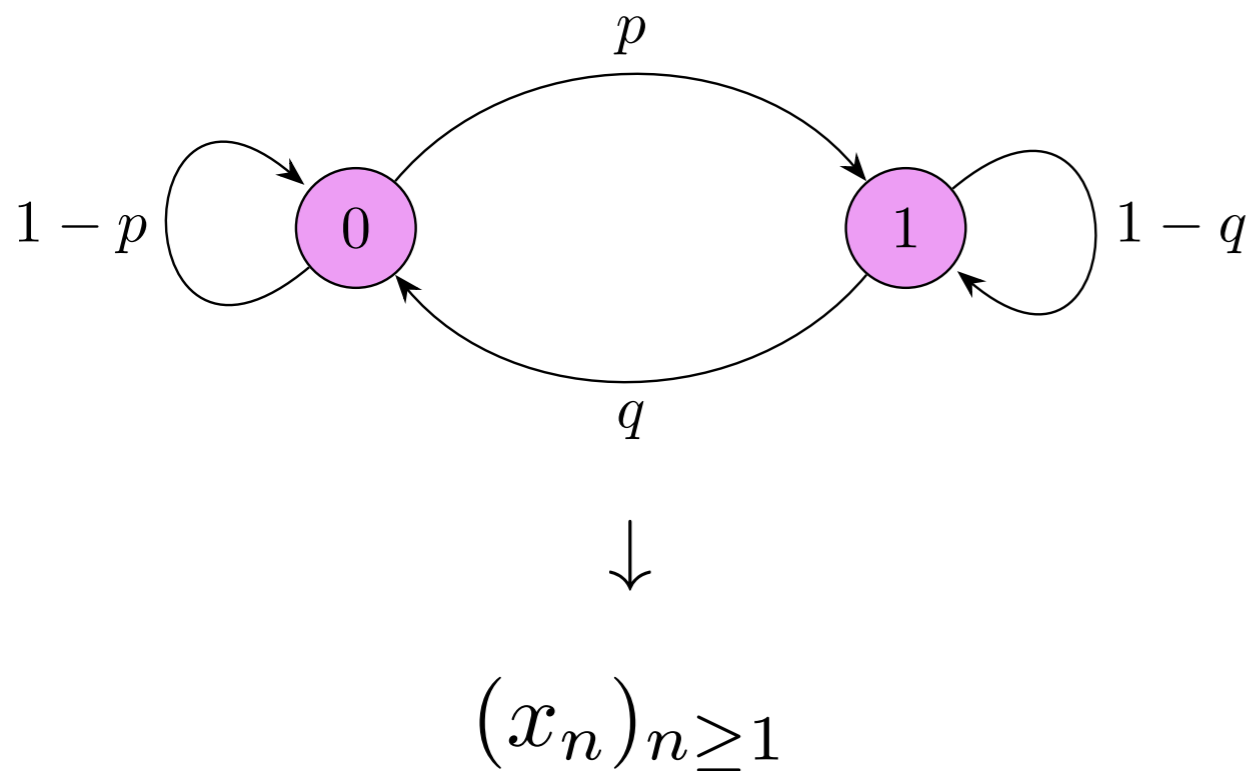
Embedding



Single-layer transformer



First-order Markov chain + Single-layer transformer



Input data: First-order Markov chain

Model: Depth = 1

Next-token prediction loss

Next-token prediction loss

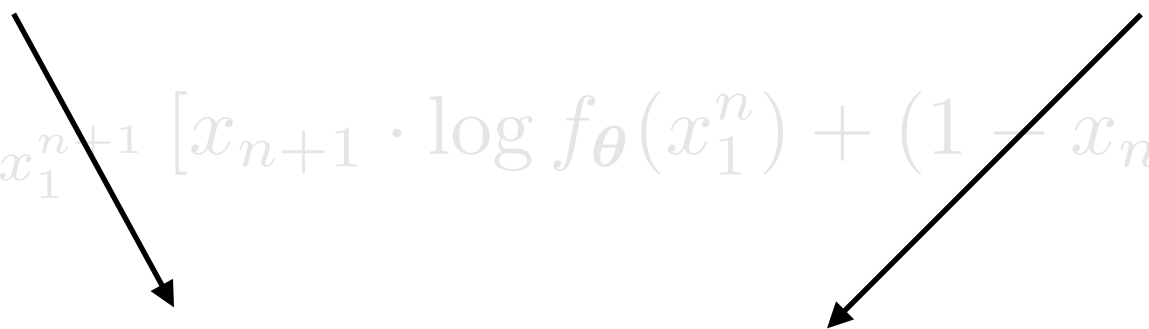
$$L(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{n=1}^N \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log f_{\boldsymbol{\theta}}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{\boldsymbol{\theta}}(x_1^n))]$$

Cross-entropy loss between x_{n+1} and prediction probability $f_{\boldsymbol{\theta}}(x_1^n)$

Ideally...

Prediction probability

Markov kernel

$$L(\theta) = -\frac{1}{N} \sum_{n=1}^N \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log f_{\theta}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{\theta}(x_1^n))]$$


$$f_{\theta}(x_1^n) \approx \mathbb{P}(x_{n+1} = 1 \mid x_n)$$

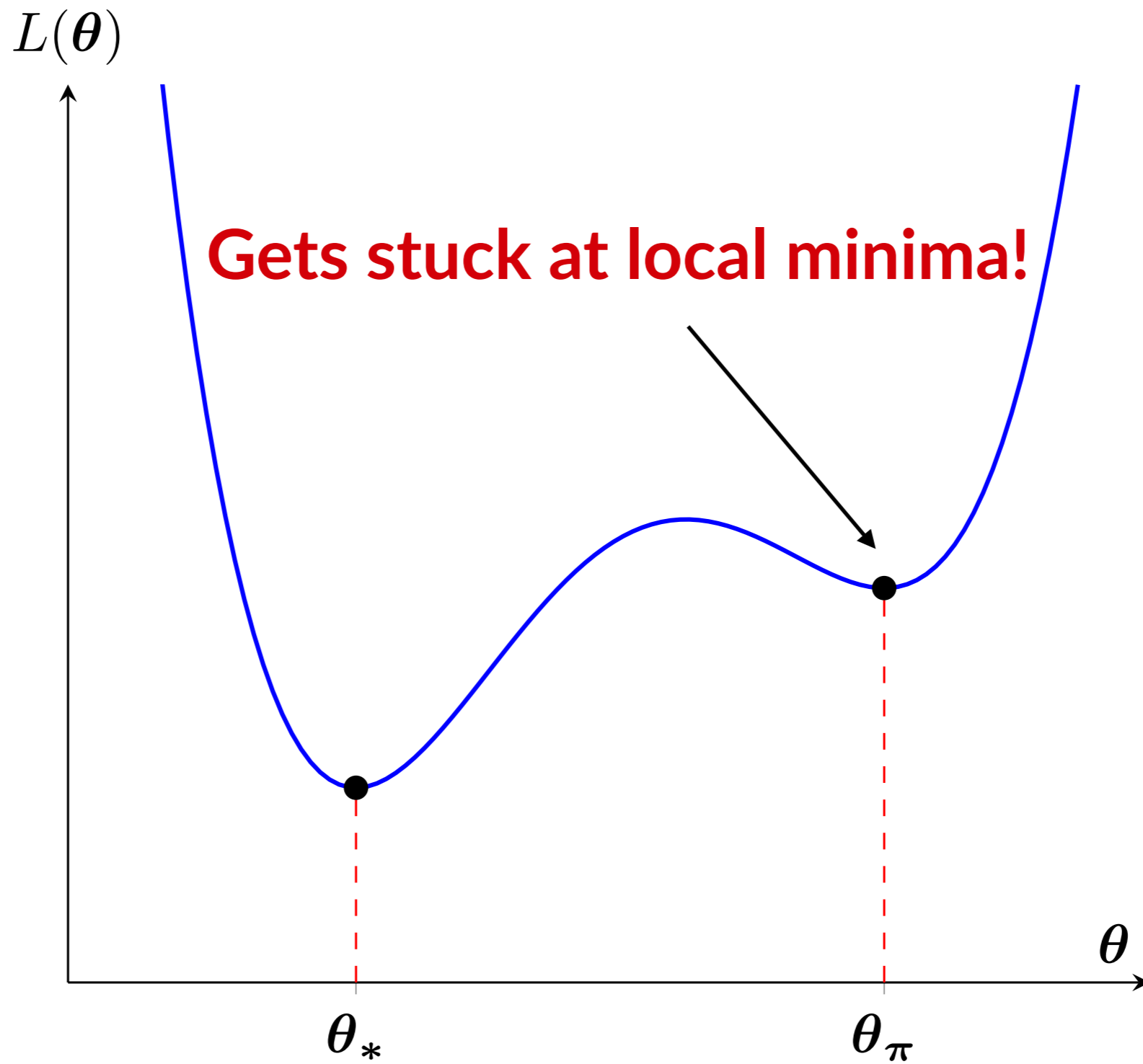
x_{n+1}

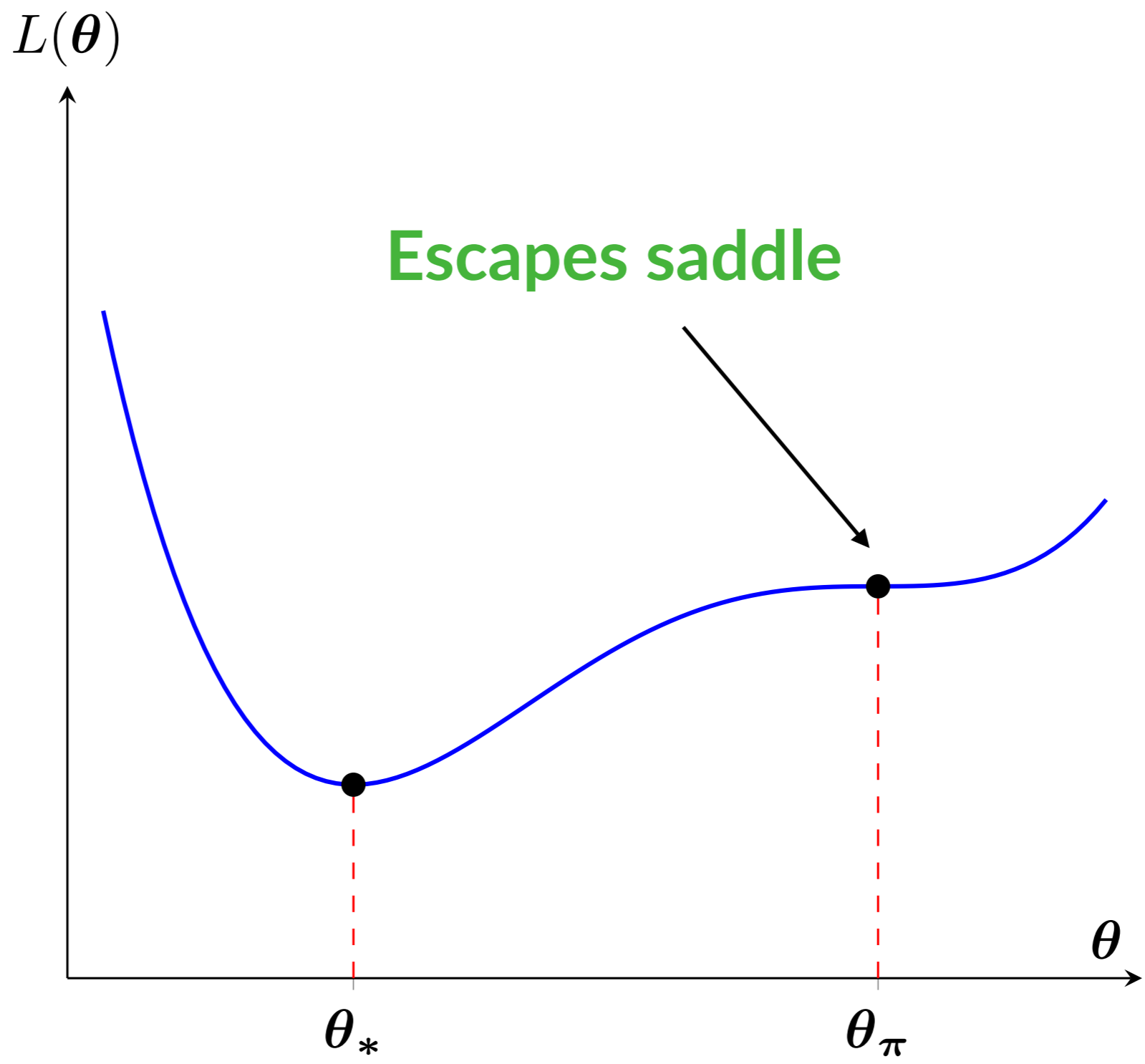
$f_{\theta}(x_1^n)$

But...

But...

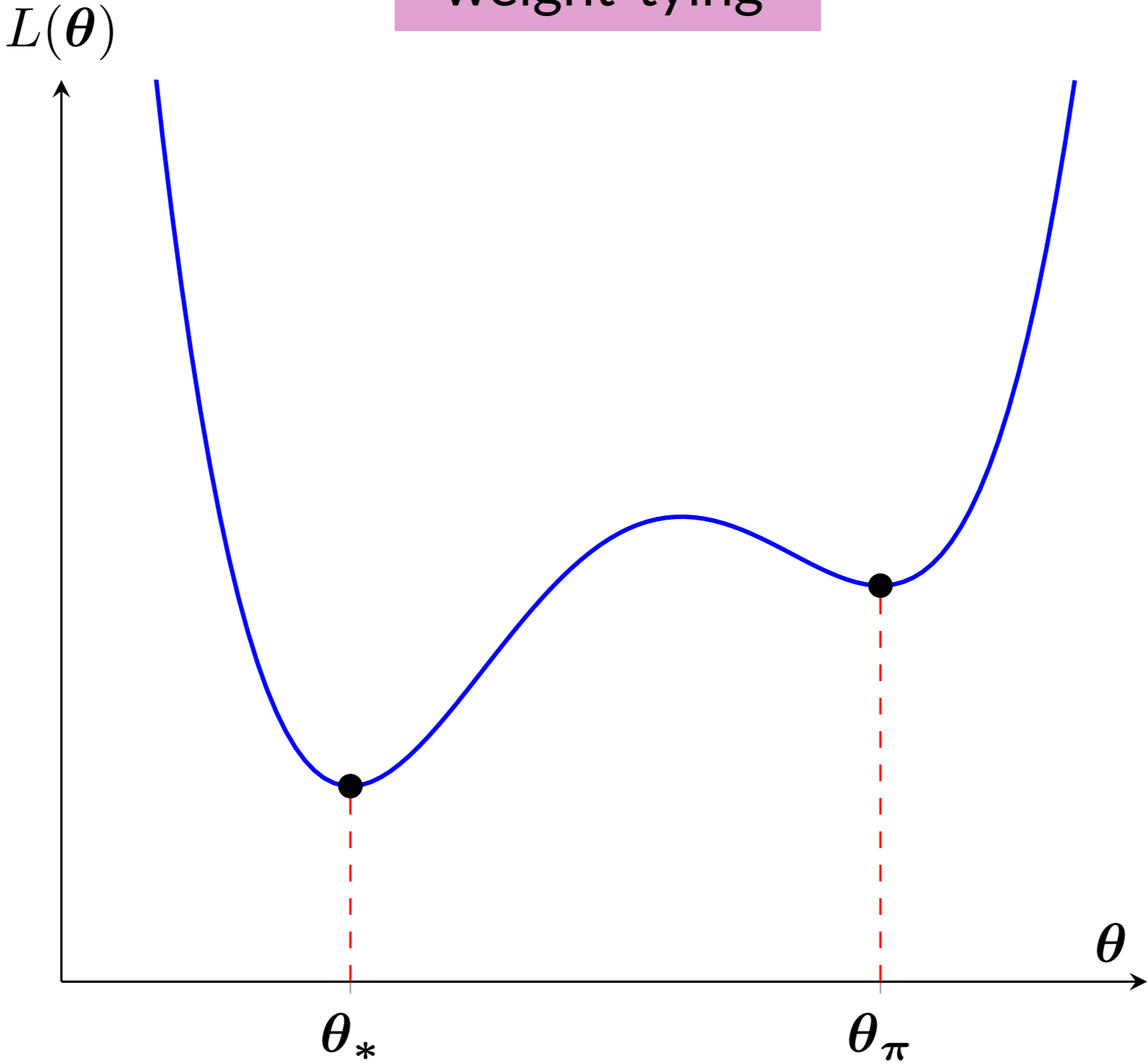
**Single-layer transformers sometimes fail
to learn even first-order Markov chains! ***



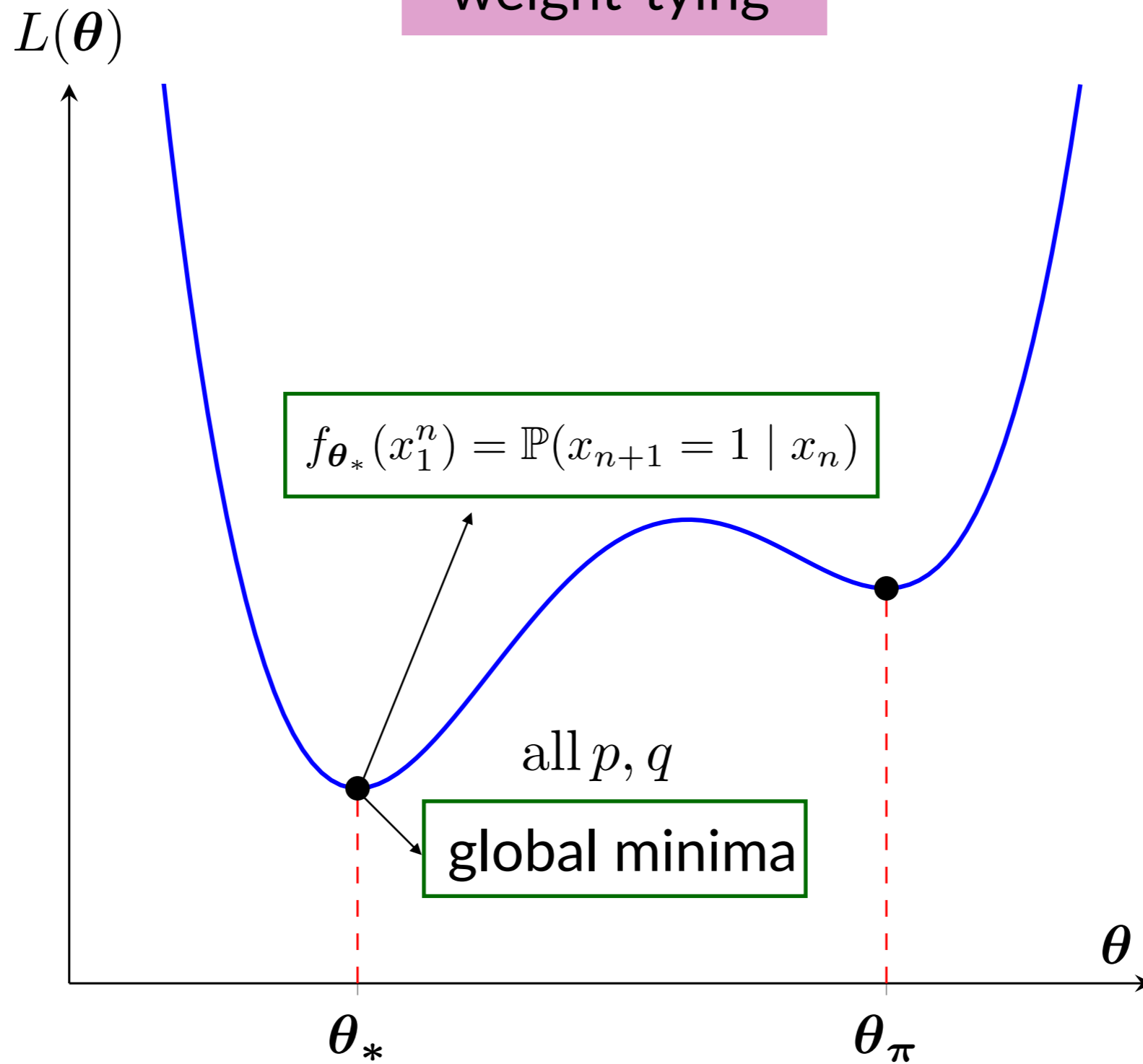


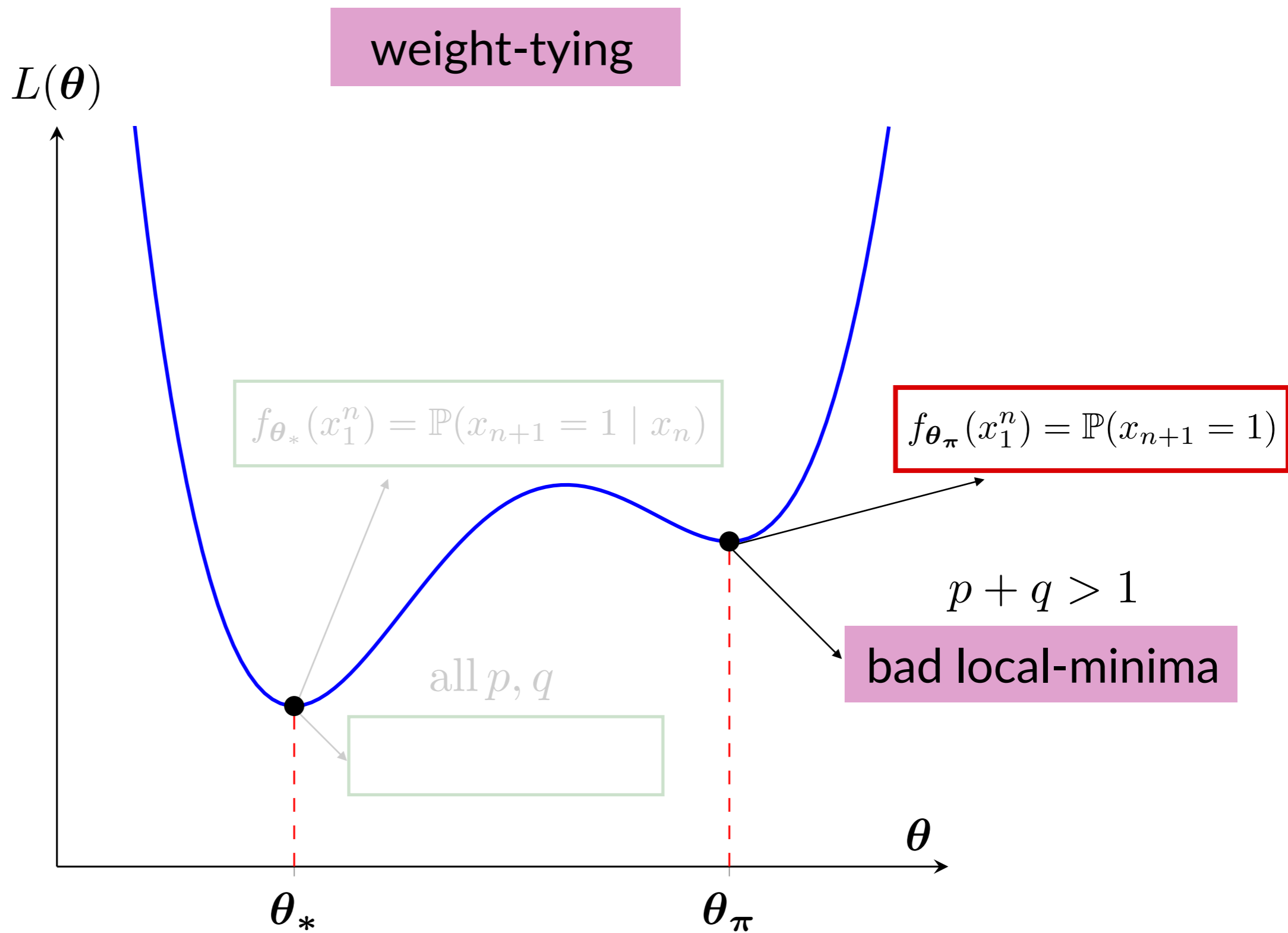
weight-tying

weight-tying

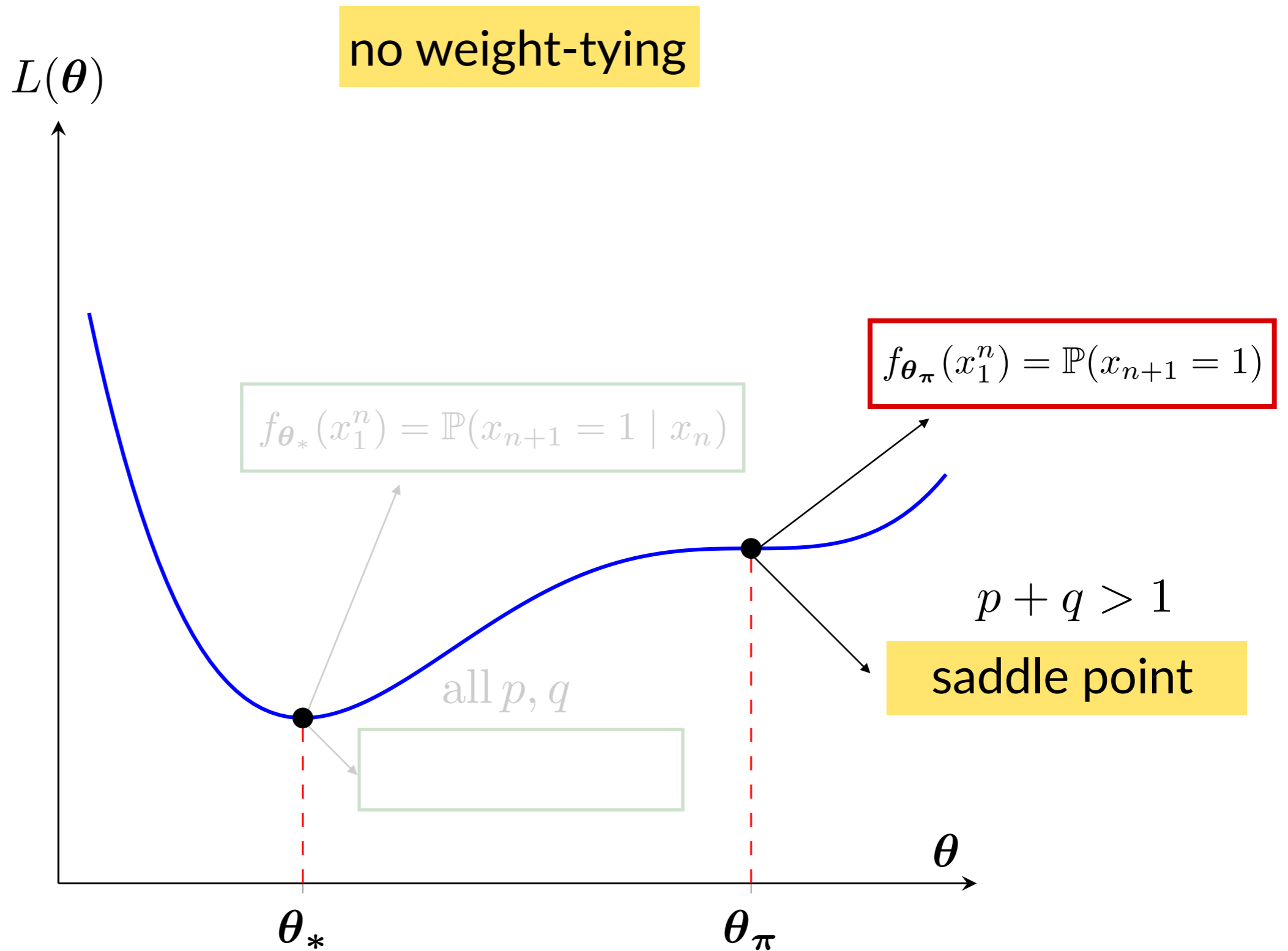


weight-tying





Interestingly...



Theoretical results

Theoretical results

Bad local minima (weight tying)

If $p + q > 1$ and the weights are tied, there exists a $\boldsymbol{\theta}_\pi$ with explicit construction such that:

- (i) $\boldsymbol{\theta}_\pi$ is a bad local minima for $L(\cdot)$ with $L(\boldsymbol{\theta}_\pi) > L(\boldsymbol{\theta}_*)$
- (ii) $f_{\boldsymbol{\theta}_\pi}(x_1^n) = \mathbb{P}(x_{n+1} = 1)$, the marginal distribution
- (iii) $L(\boldsymbol{\theta}_\pi) = H(\boldsymbol{\pi})$, the entropy of the stationary distribution
- (iv) $\nabla L(\boldsymbol{\theta}_\pi) = 0$, i.e. $\boldsymbol{\theta}_\pi$ is a stationary point

Theoretical results

If $p + q > 1$ and the weights are tied, there exists a θ_π with explicit construction such that:

- (i) θ_π is a bad local minima for $L(\cdot)$ with $L(\theta_\pi) > L(\theta_*)$
- (ii) $f_{\theta_\pi}(x_1^n) = \mathbb{P}(x_{n+1} = 1)$, the marginal distribution
- (iii) $L(\theta_\pi) = H(\pi)$, the entropy of the stationary distribution
- (iv) $\nabla L(\theta_\pi) = 0$, i.e. θ_π is a stationary point

Saddle point (no weight tying)

Under the same setting as above with the weights not tied, θ_π becomes a saddle point. It satisfies the same properties.

Theoretical results

Bad local minima (weight tying)

If $p + q > 1$ and the weights are tied, there exists a $\boldsymbol{\theta}_\pi$ with explicit construction such that:

- (i) $\boldsymbol{\theta}_\pi$ is a bad local minima for $L(\cdot)$ with $L(\boldsymbol{\theta}_\pi) > L(\boldsymbol{\theta}_*)$
- (ii) $f_{\boldsymbol{\theta}_\pi}(x_1^n) = \mathbb{P}(x_{n+1} = 1)$, the marginal distribution
- (iii) $L(\boldsymbol{\theta}_\pi) = H(\boldsymbol{\pi})$, the entropy of the stationary distribution
- (iv) $\nabla L(\boldsymbol{\theta}_\pi) = 0$, i.e. $\boldsymbol{\theta}_\pi$ is a stationary point

Saddle point (no weight tying)

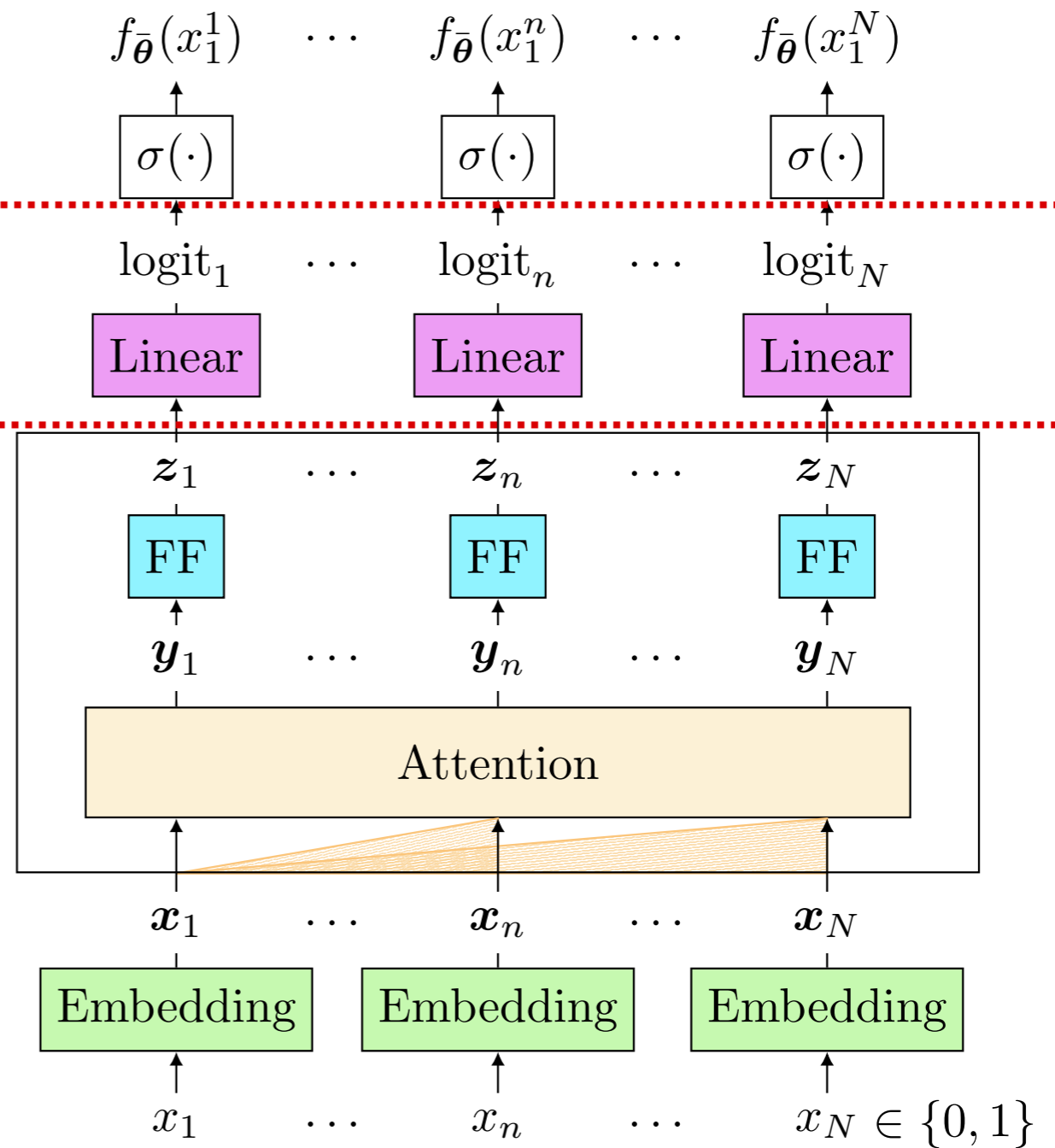
Under the same setting as above with the weights not tied, $\boldsymbol{\theta}_\pi$ becomes a saddle point. It satisfies the same properties.

Intuition

$$f_{\theta}(x_1^n) = \sigma(b) = \pi_1$$

$$\uparrow a = 0$$

$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$



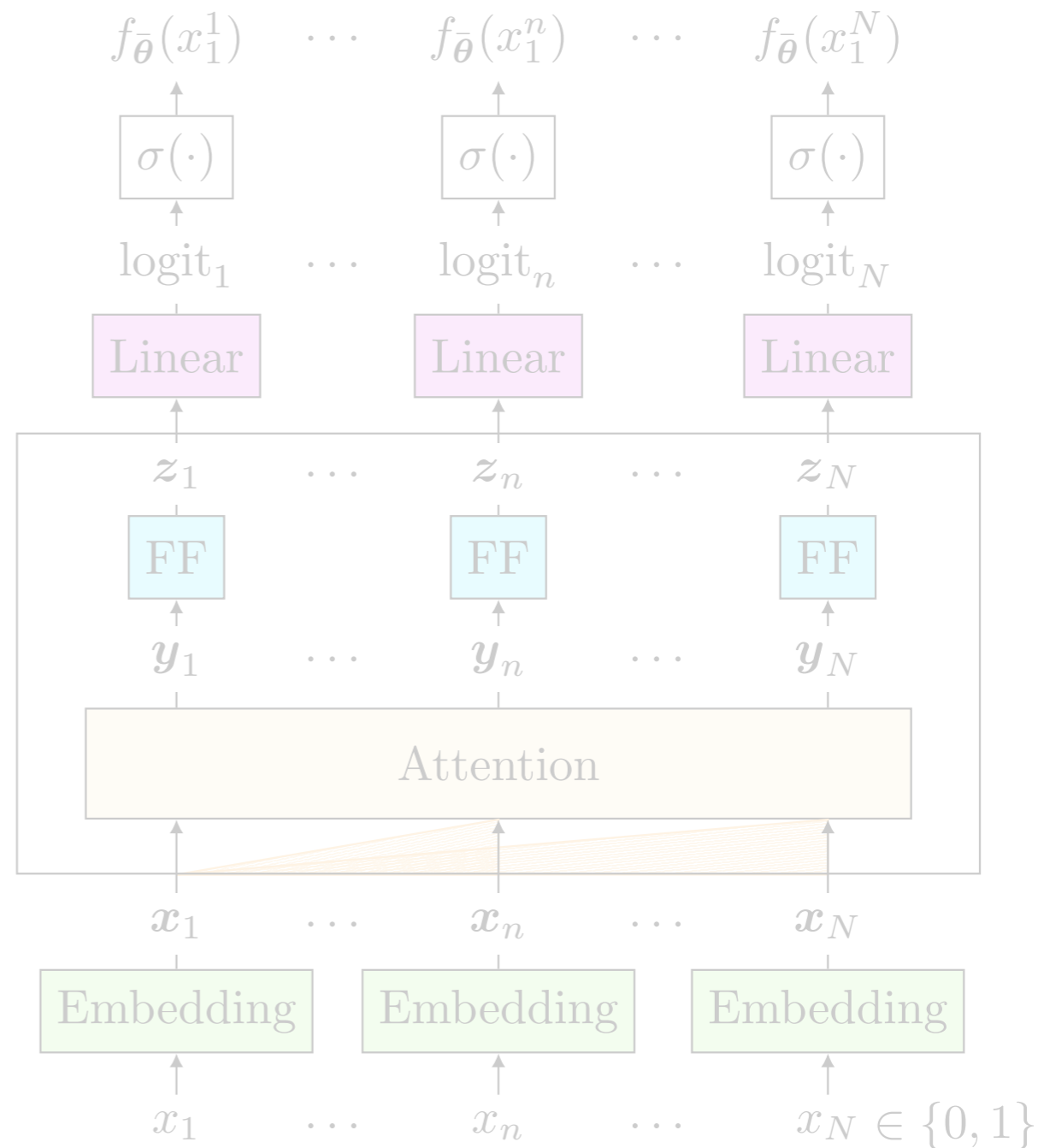
Intuition

weight-tying

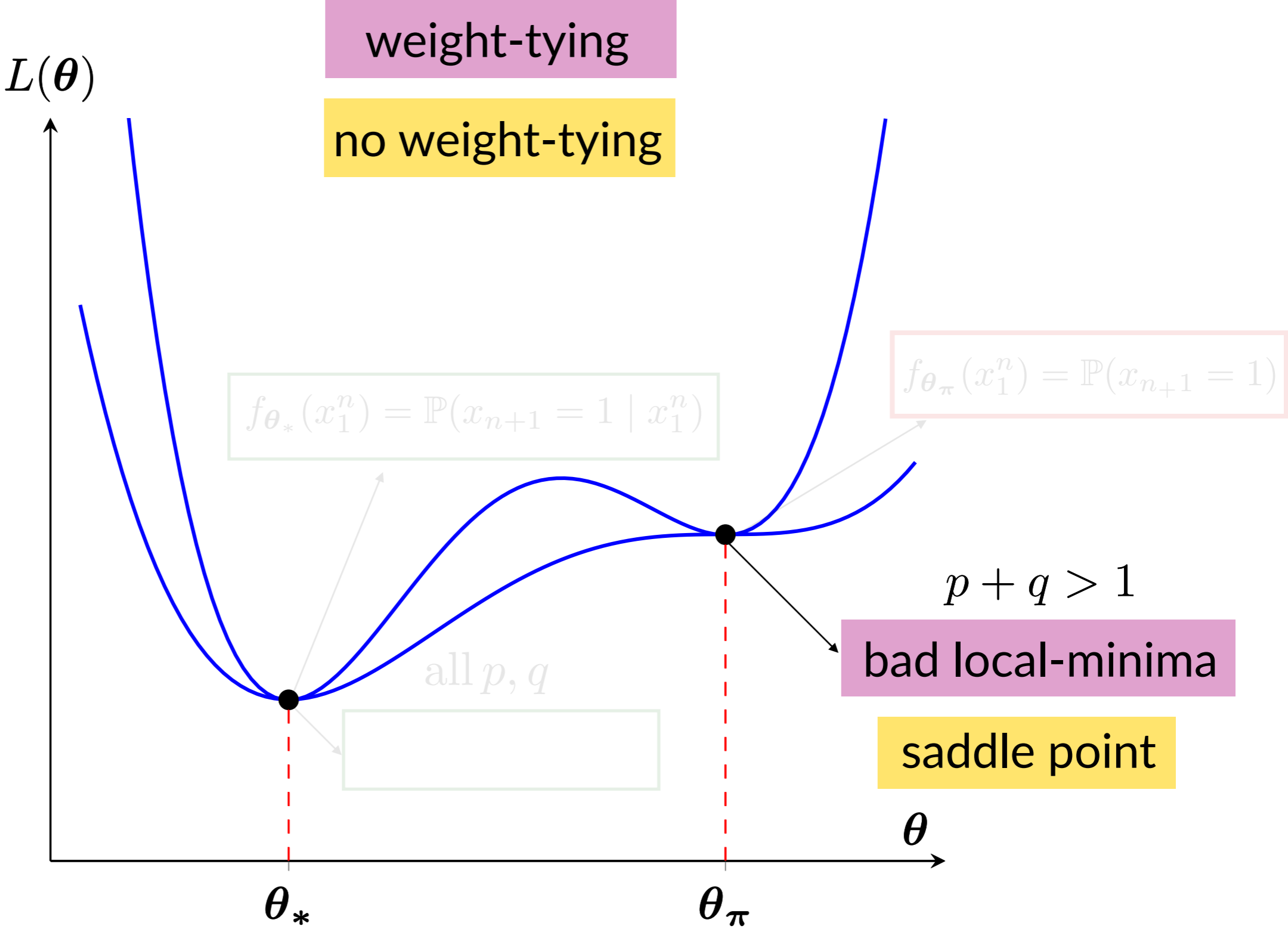
Hessian is (almost) positive-definite

no weight-tying

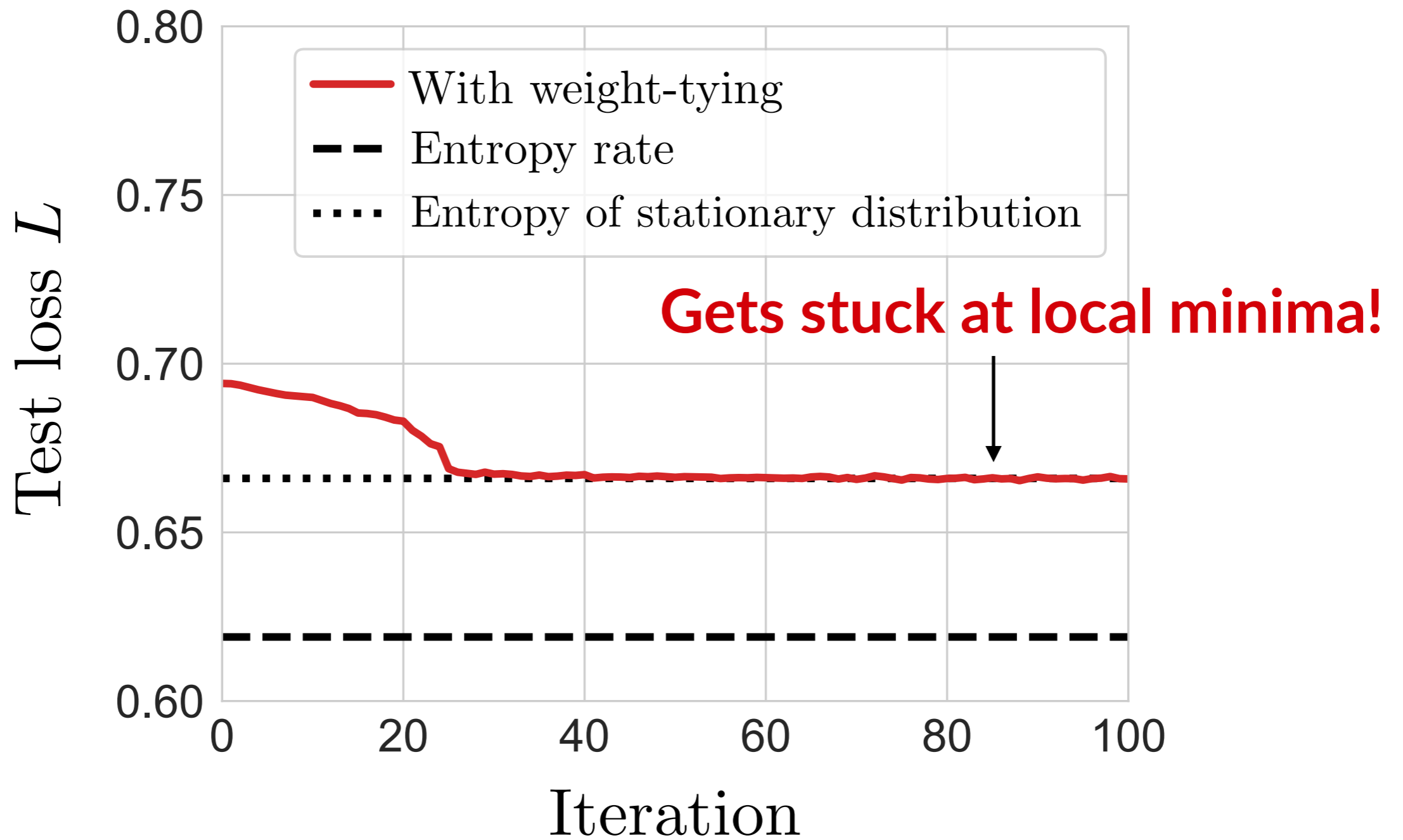
Hessian is indefinite



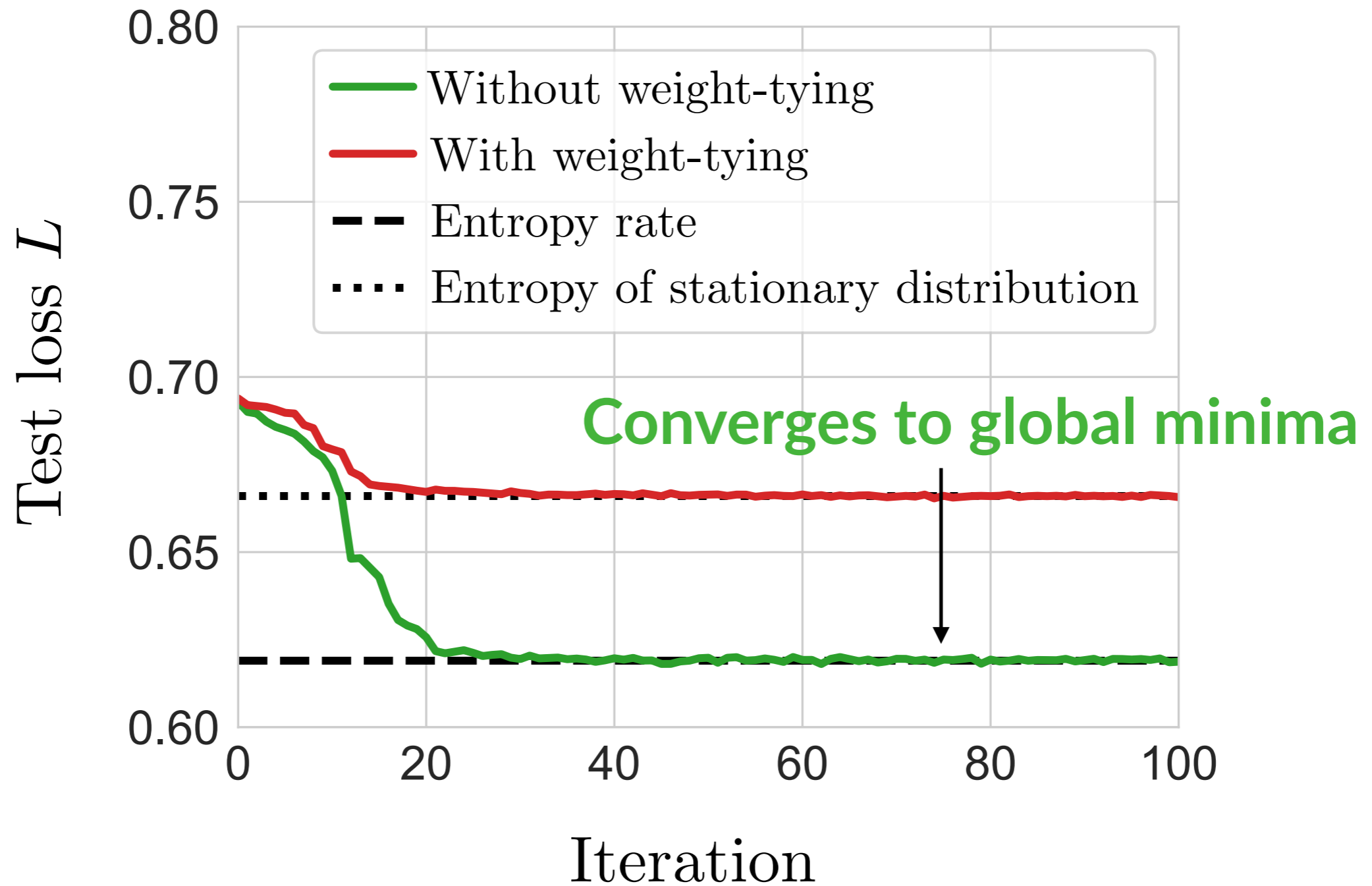
Main results



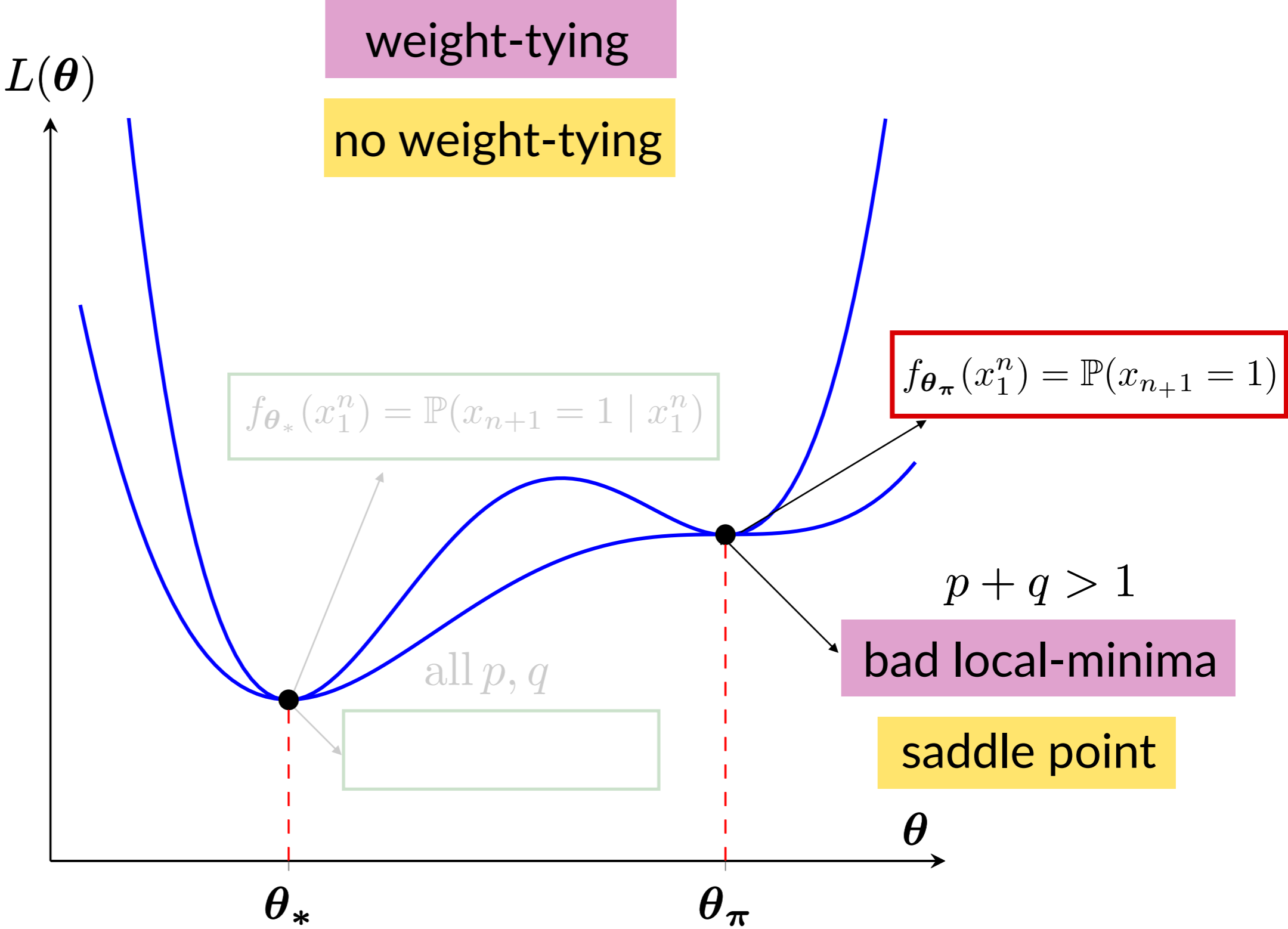
Weight tying



Without weight tying



Main results





Markovian inputs



Transformers



Memory = 1

Depth = 1



What do they learn?





Markovian inputs

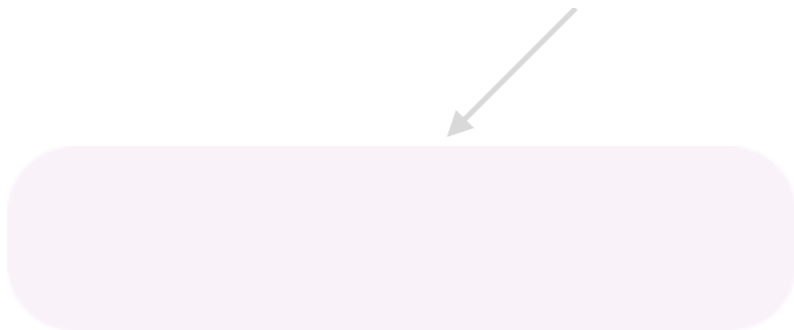


Transformers



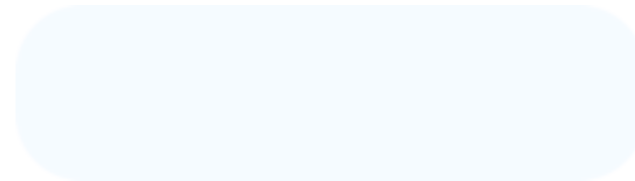
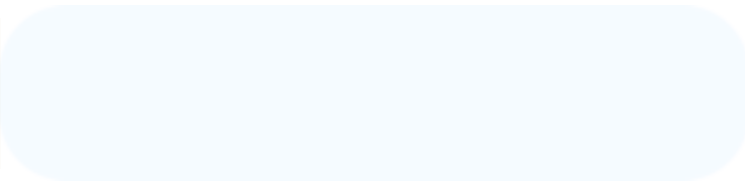
Memory = 1

Depth = 1



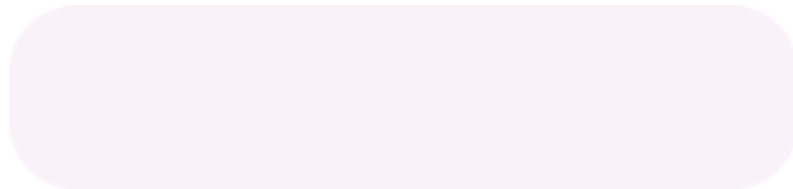
How do they learn?





Memory = 1

Depth = 1



How do they learn?



Learning dynamics

└ Gradient-flow

Gradient flow

Gradient flow

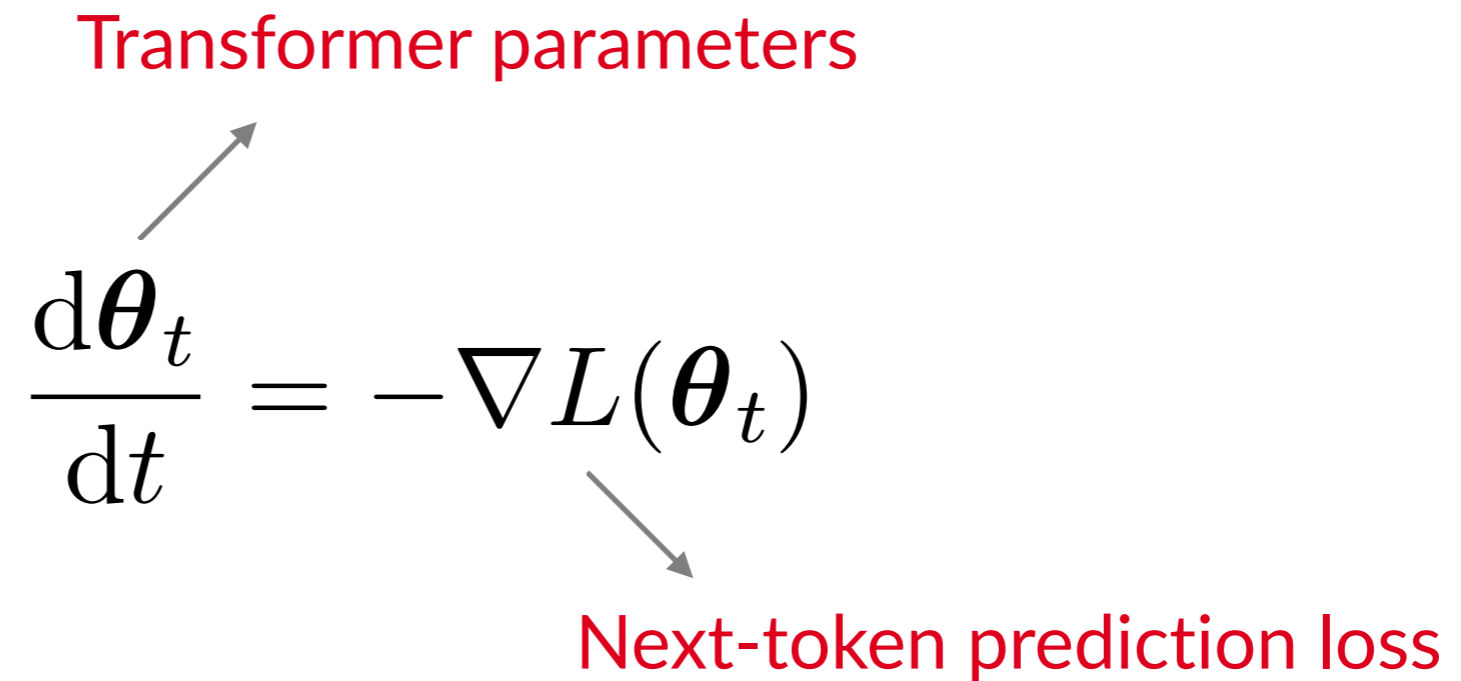
$$\frac{d\boldsymbol{\theta}_t}{dt} = -\nabla L(\boldsymbol{\theta}_t)$$

Gradient flow

Transformer parameters

$$\frac{d\boldsymbol{\theta}_t}{dt} = -\nabla L(\boldsymbol{\theta}_t)$$

Next-token prediction loss

The diagram features the mathematical equation $\frac{d\boldsymbol{\theta}_t}{dt} = -\nabla L(\boldsymbol{\theta}_t)$ centered on the page. A red arrow points from the text 'Transformer parameters' above to the $\boldsymbol{\theta}_t$ term in the denominator of the derivative. Another red arrow points from the text 'Next-token prediction loss' below to the $L(\boldsymbol{\theta}_t)$ term in the gradient.

Recall

$$f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n)$$

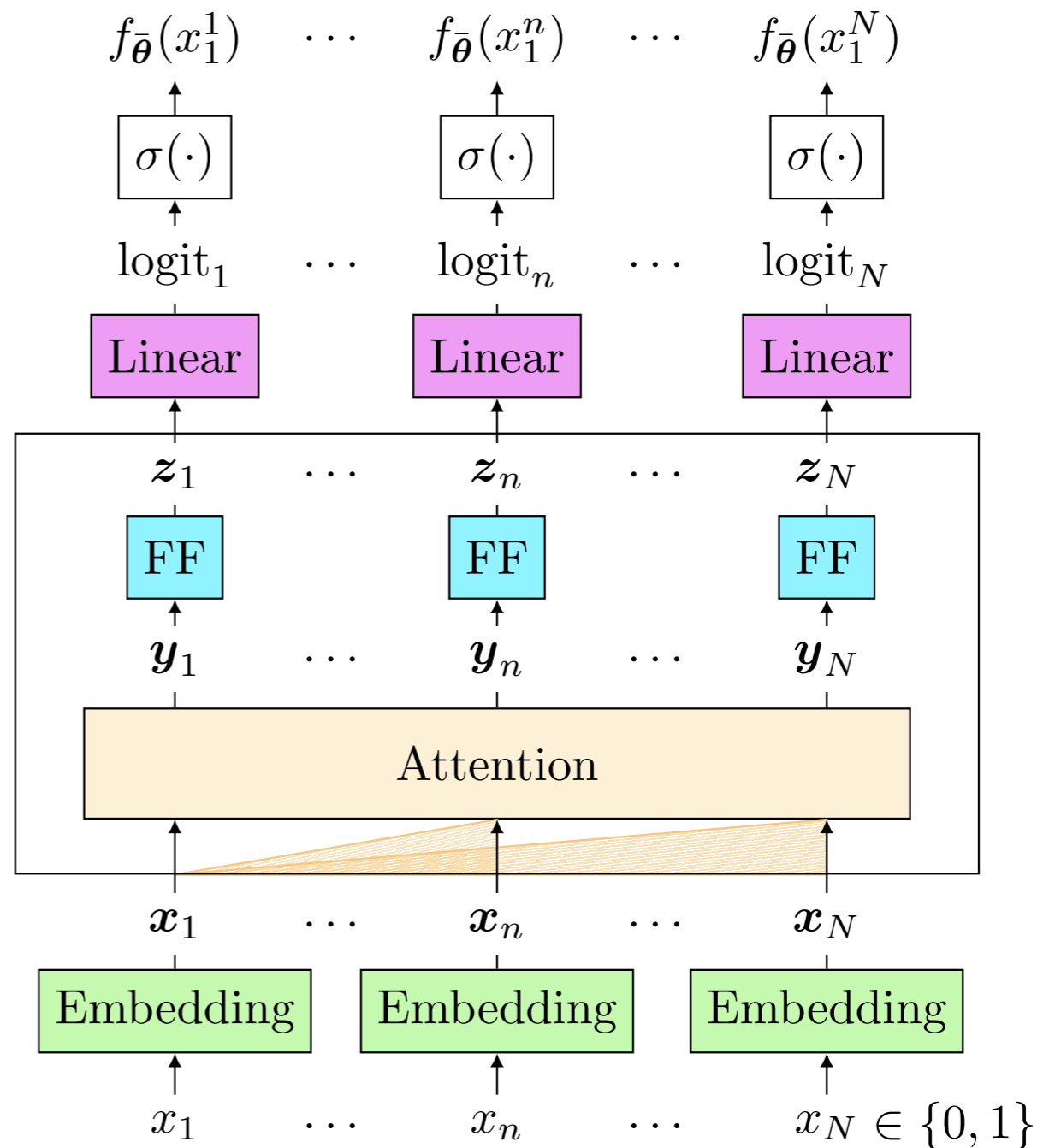
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$

$$\mathbf{z}_n = \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n)$$

$$\mathbf{y}_n = \mathbf{x}_n + \mathbf{W}_O \sum_{i=1}^n \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i$$

$$\mathbf{x}_n = x_n \cdot \mathbf{e} + \mathbf{p}_n$$

θ



Low-rank structure

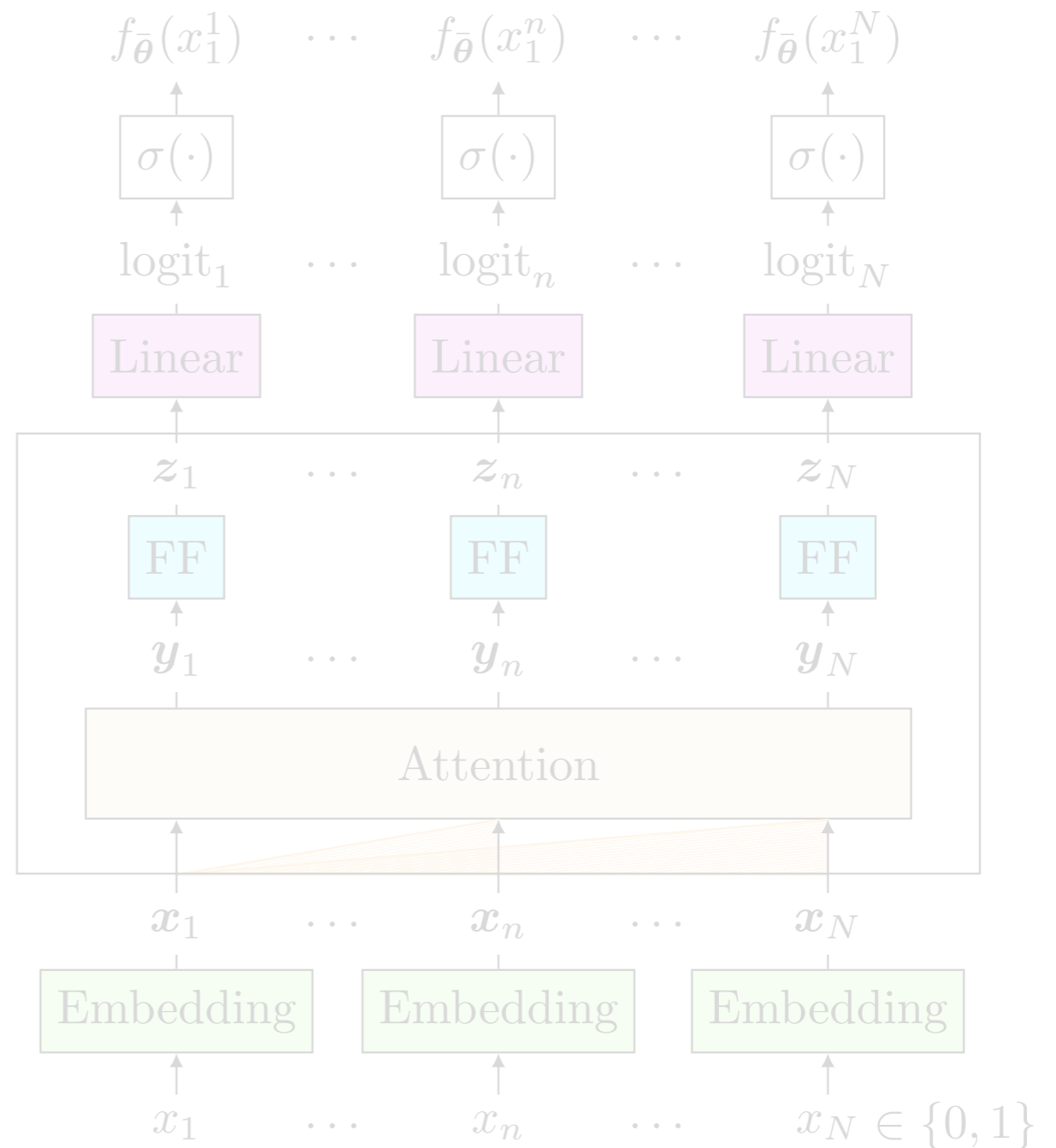
$$f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n)$$

$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$

$$\mathbf{z}_n = \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n)$$

$$\mathbf{y}_n = \mathbf{x}_n + \mathbf{W}_O \sum_{i=1}^n \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i$$

$$\mathbf{x}_n = x_n \cdot \mathbf{e} + \mathbf{p}_n$$



Reparametrization

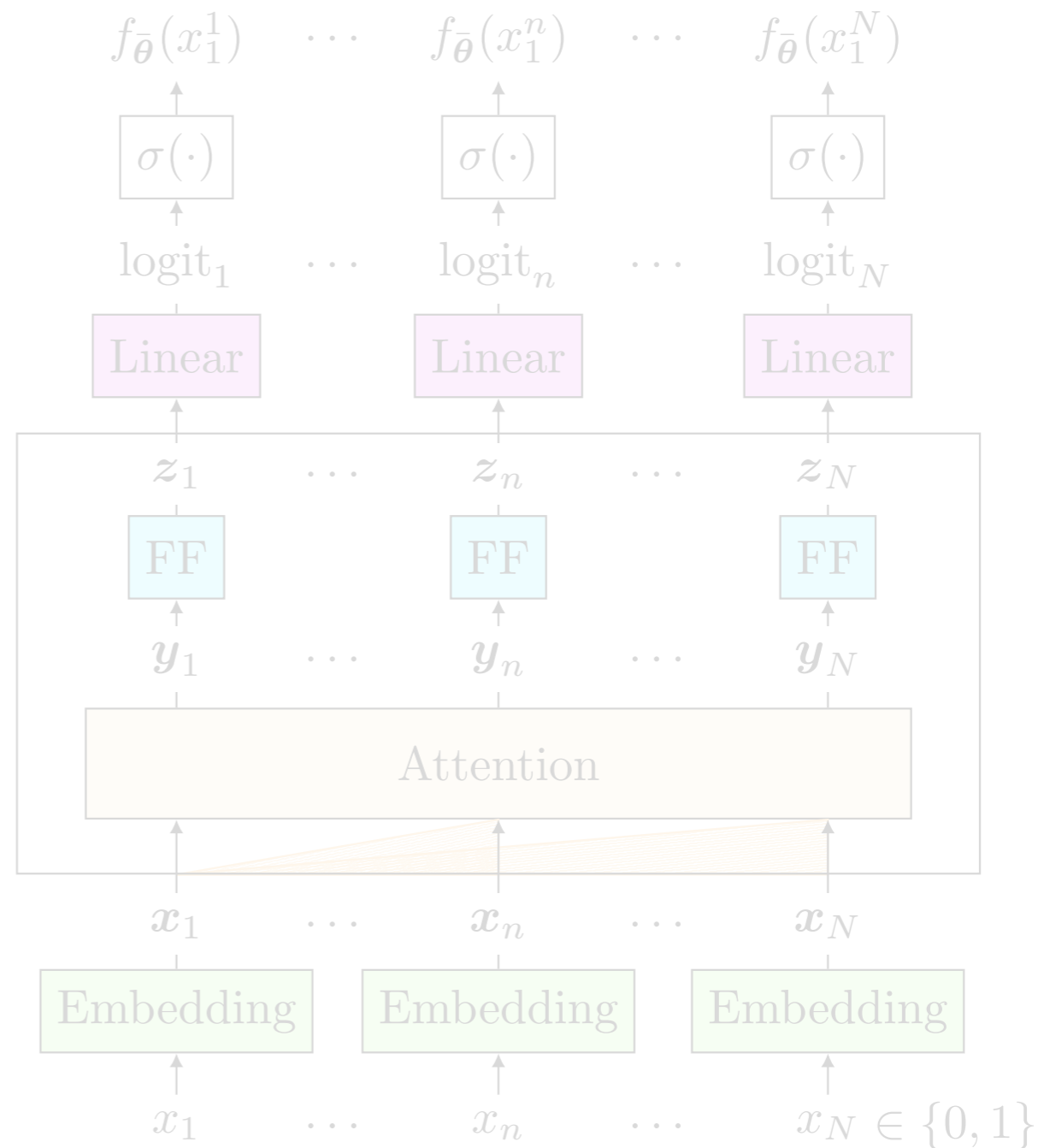
$$f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n)$$

$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R} \quad \mathbf{w}$$

$$\mathbf{z}_n = \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n)$$

$$\mathbf{y}_n = \mathbf{x}_n + \mathbf{W}_O \sum_{i=1}^n \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i$$

$$\mathbf{x}_n = x_n \cdot \mathbf{e} + \mathbf{p}_n \quad \mathbf{a} \quad \mathbf{e}$$



Reparametrization

$$f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n)$$

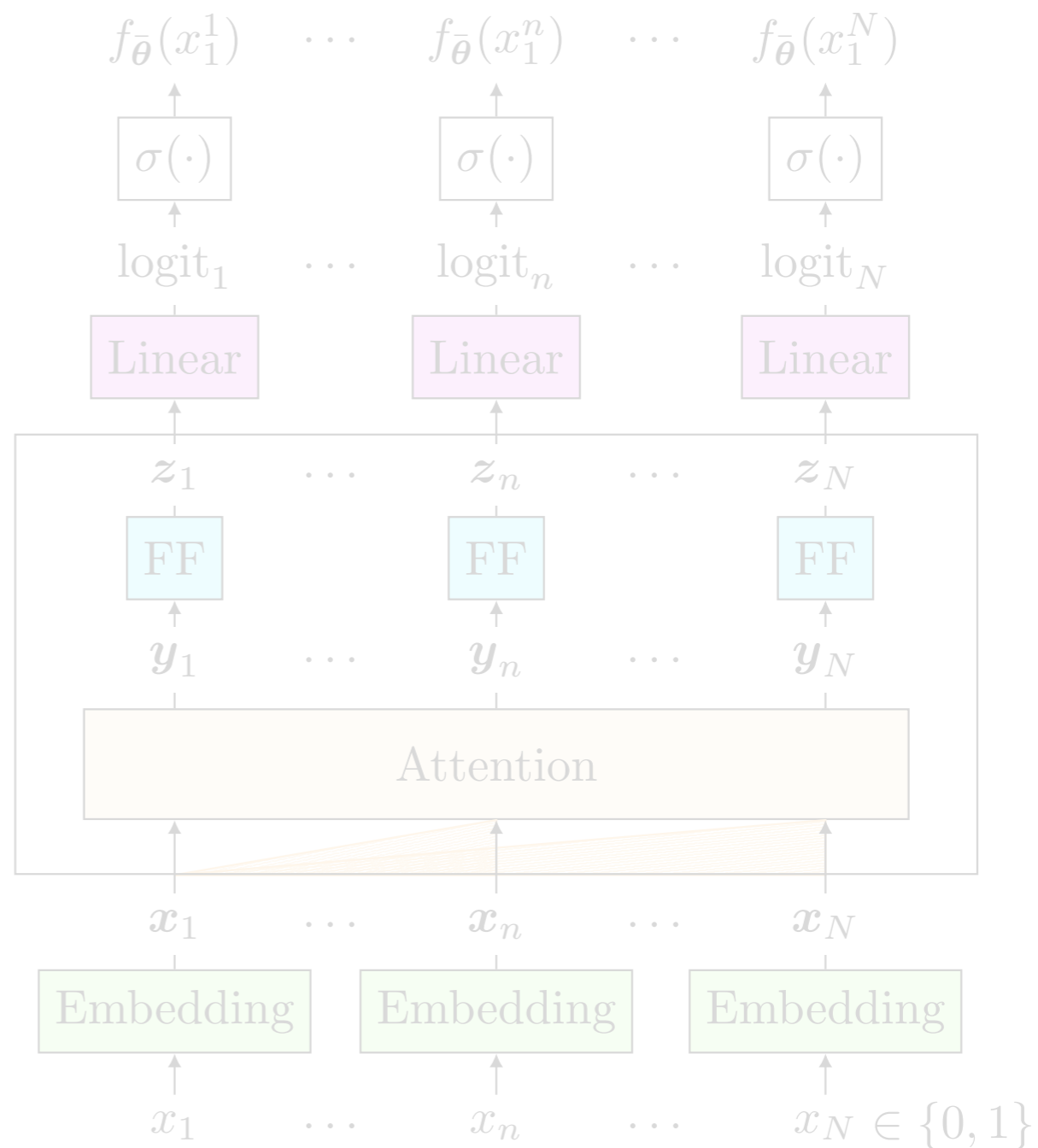
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R} \quad \mathbf{w}$$

$$\mathbf{z}_n = \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n)$$

$$\mathbf{y}_n = \mathbf{x}_n + \mathbf{W}_O \sum_{i=1}^n \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i$$

$$\mathbf{x}_n = x_n \cdot \mathbf{e} + \mathbf{p}_n \quad \mathbf{a}$$

$$\theta = (e, w, a) \in \mathbb{R}^3$$



Gradient flow

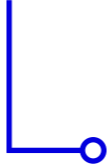
$$\frac{d\boldsymbol{\theta}_t}{dt} = -\nabla L(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_t = (e_t, w_t, a_t) \in \mathbb{R}^3$$

$$\frac{d\boldsymbol{\theta}_t}{dt} = -\nabla L(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_t = (e_t, w_t, a_t) \in \mathbb{R}^3$$

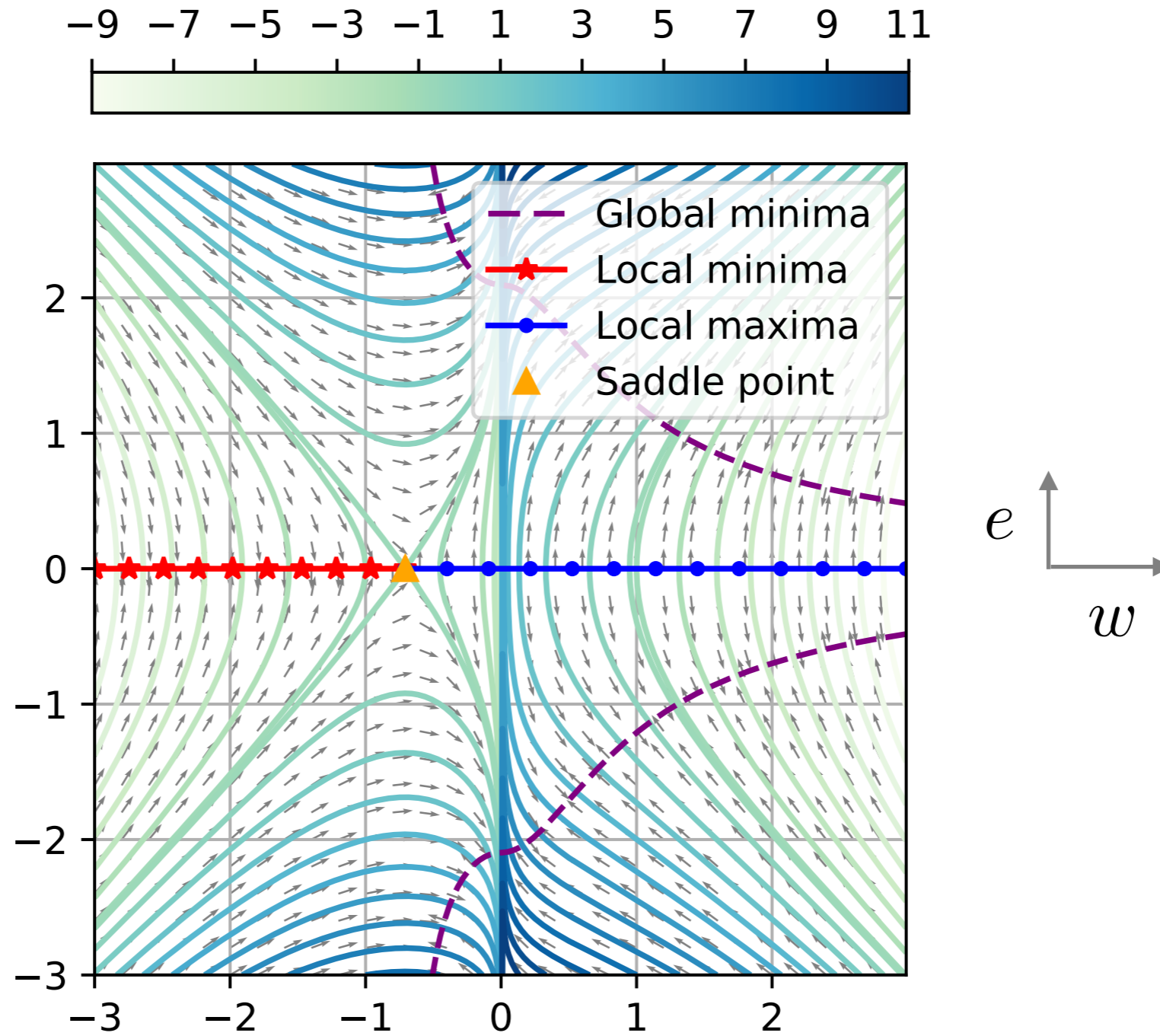
How does the flow look like?

$$\frac{d\boldsymbol{\theta}_t}{dt} = -\nabla L(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_t = (e_t, w_t, a_t) \in \mathbb{R}^3$$

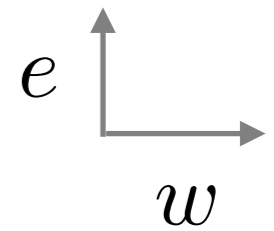
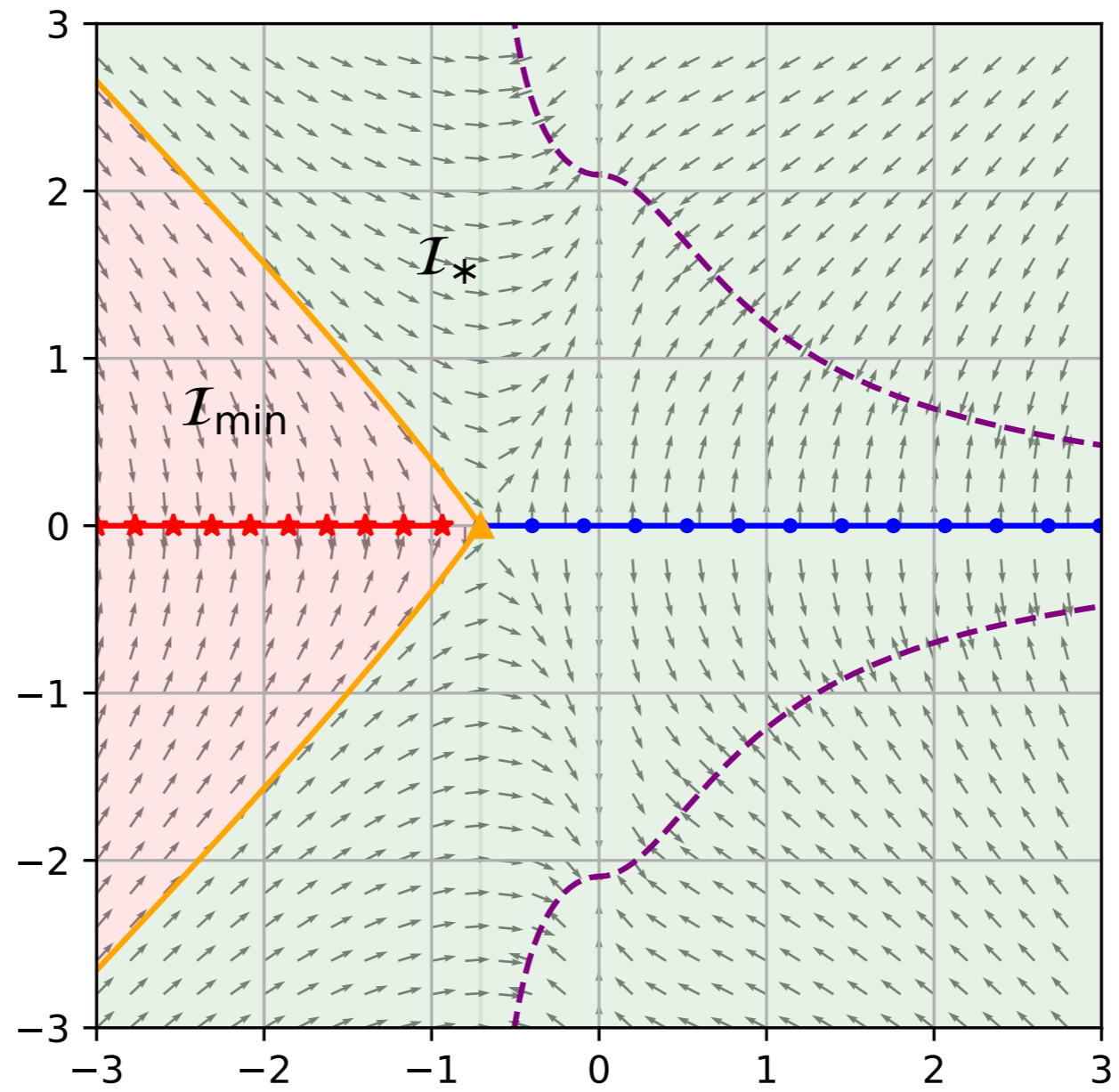
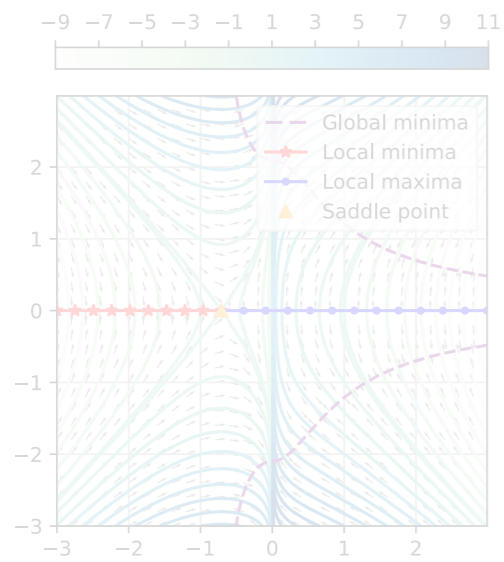
How does the flow look like?


$$a = 0 \rightarrow \boldsymbol{\theta} = (e, w)$$

$p+q < 1$

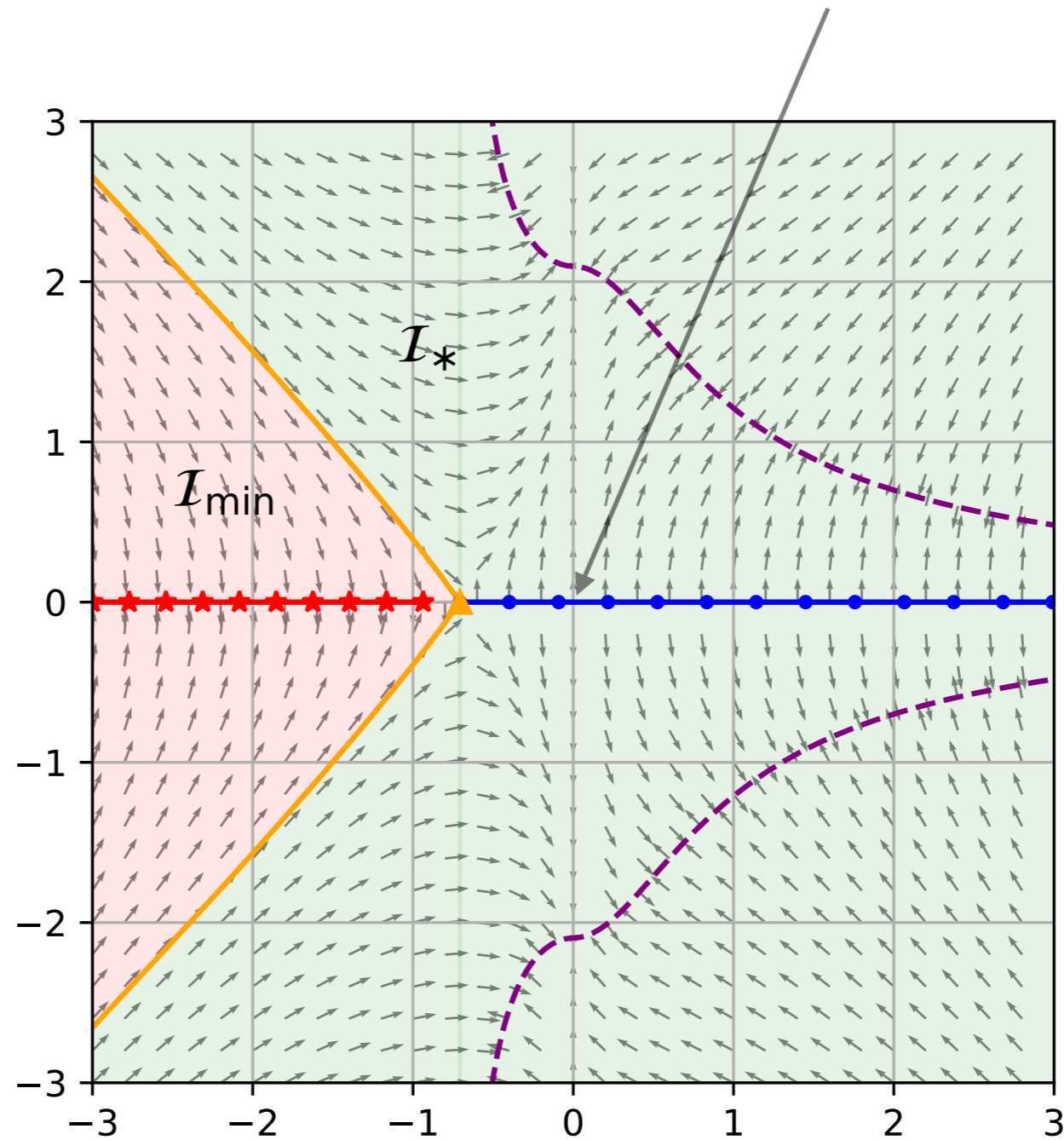
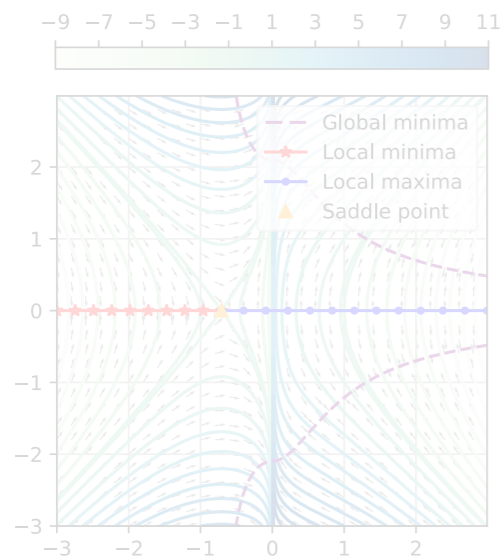


$$p+q < 1$$

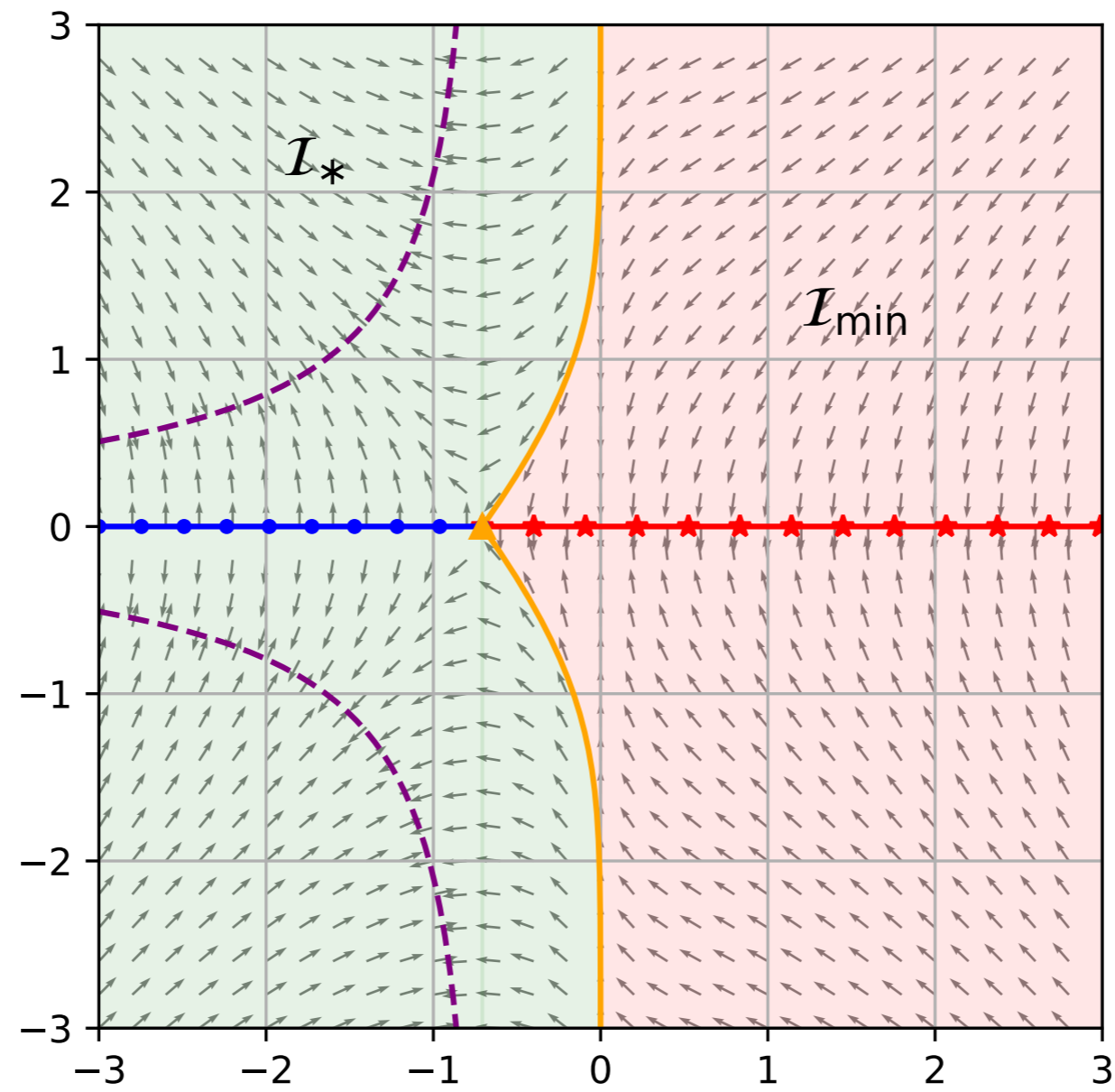


$p+q < 1$

Gaussian init. converges to global minima

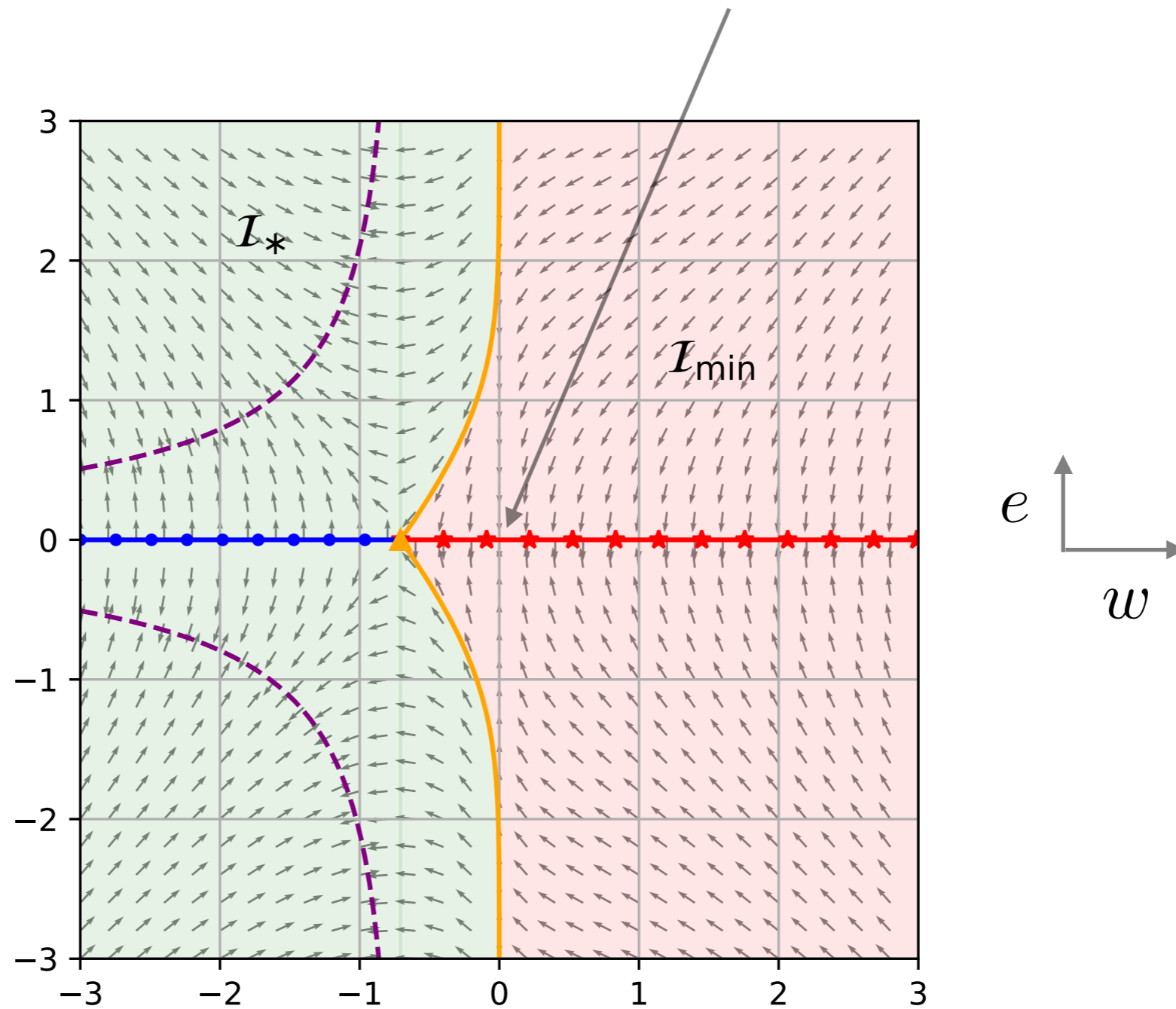


$p+q > 1$

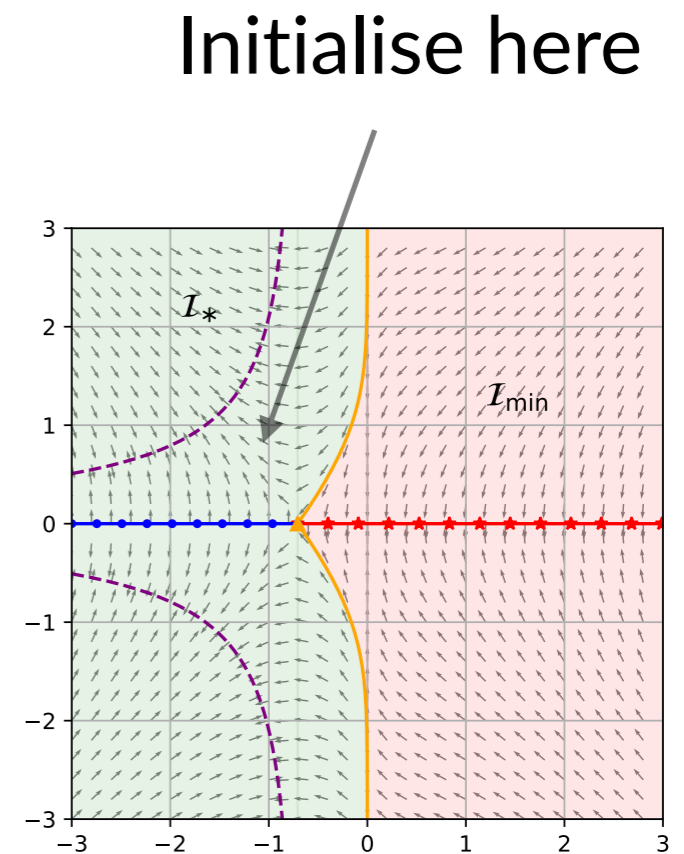
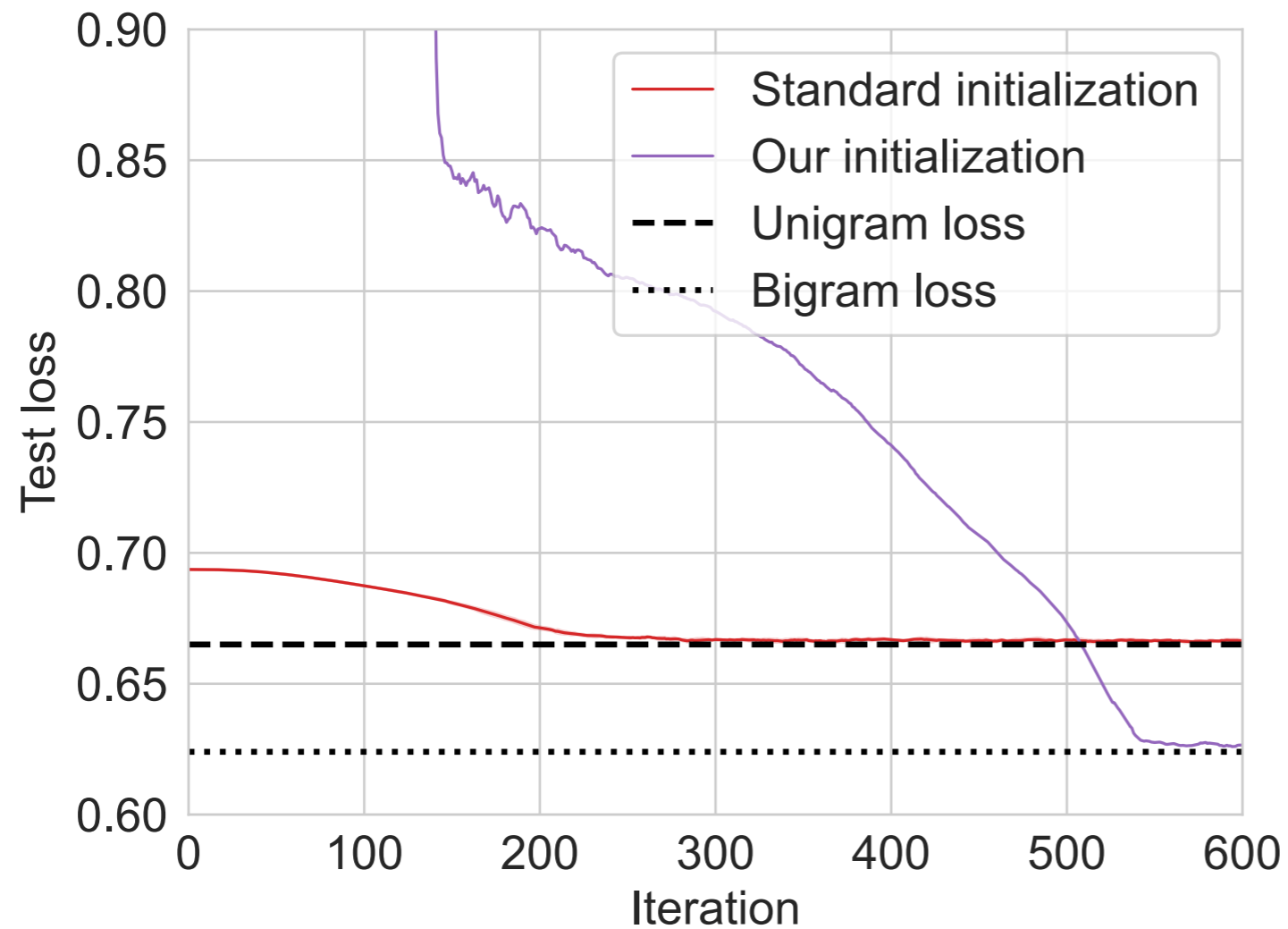


$p+q > 1$

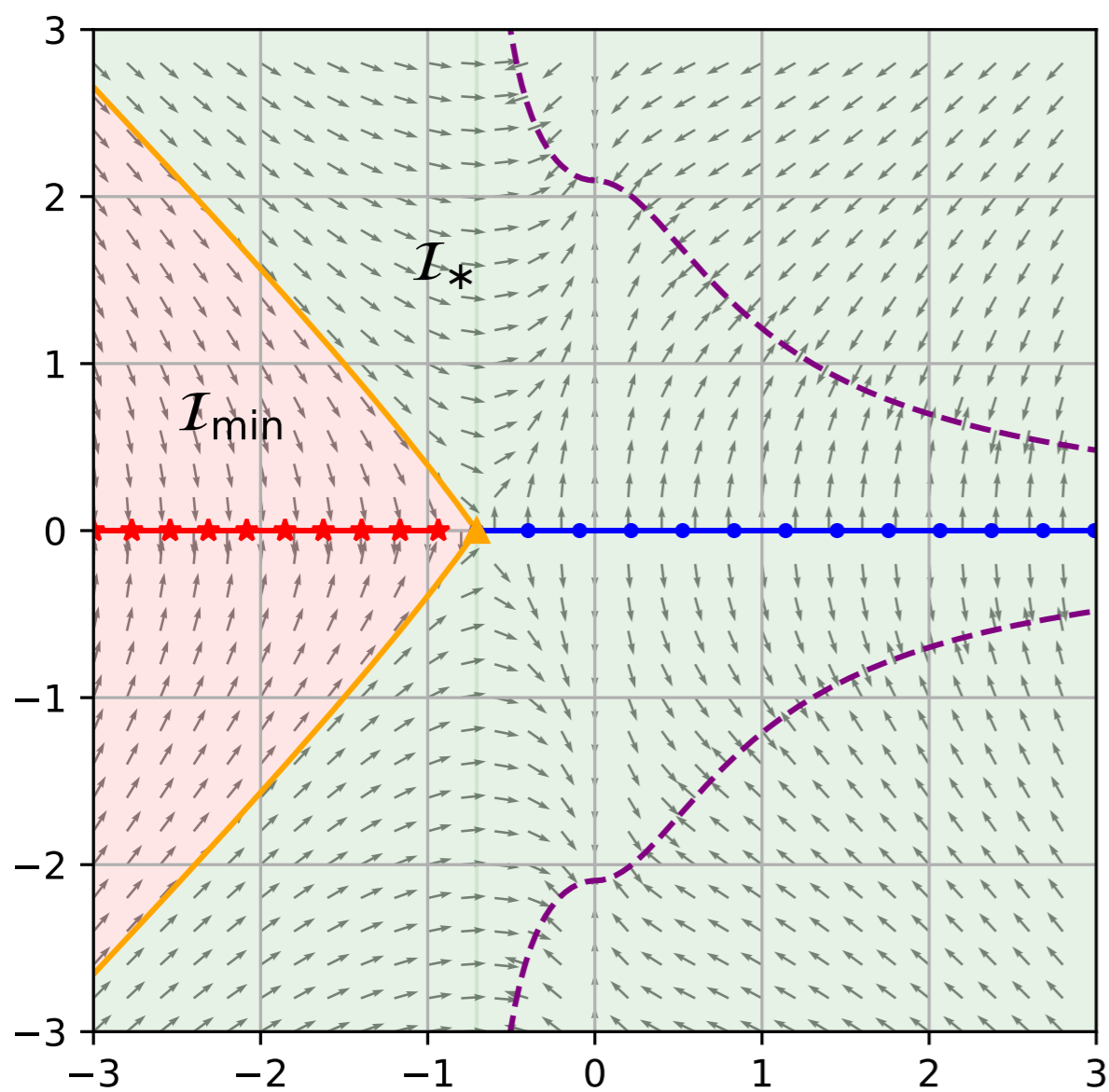
Gets stuck at local minima!



Can we escape it?

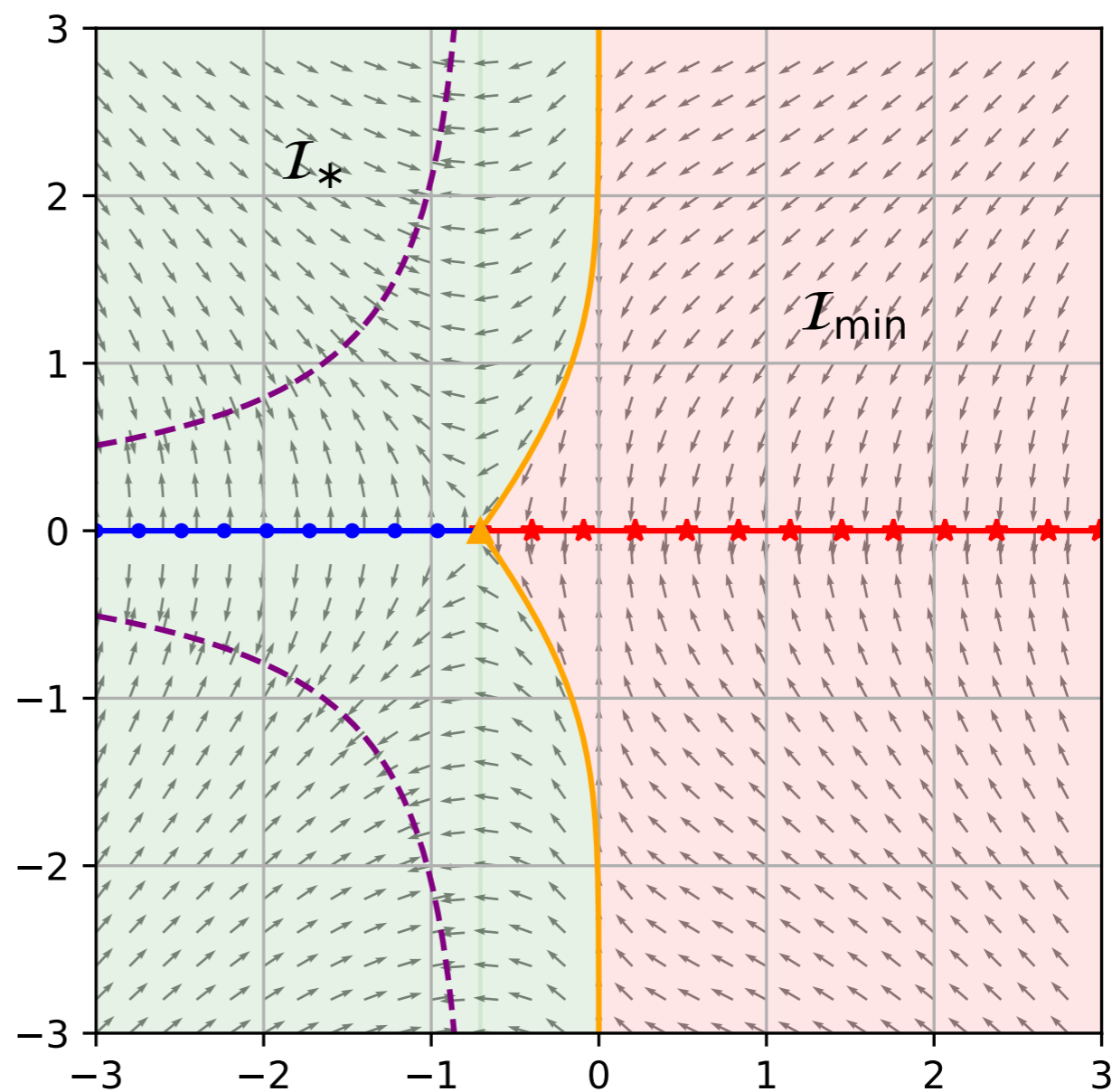


Converges to global minima

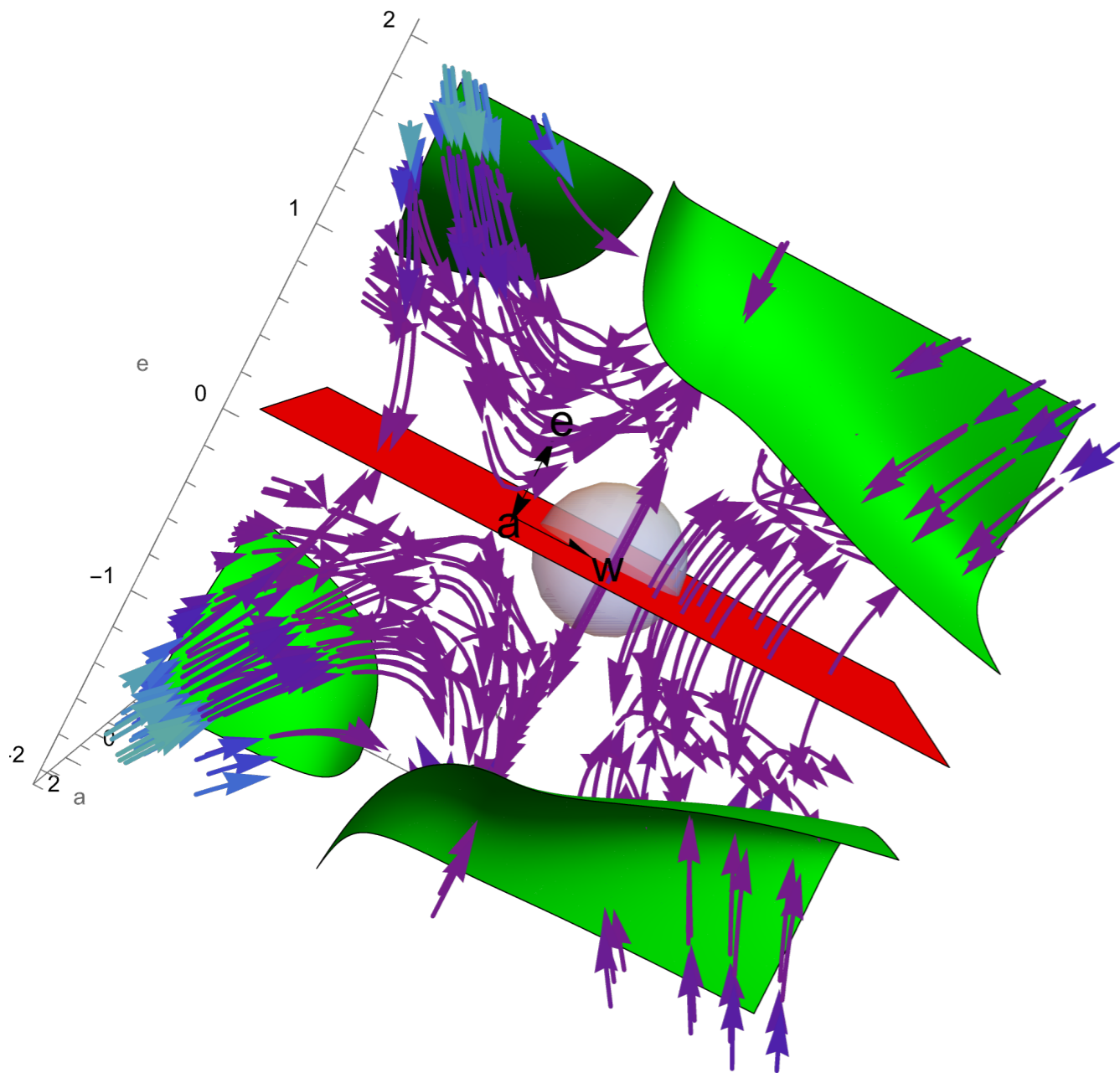


$$p + q < 1$$

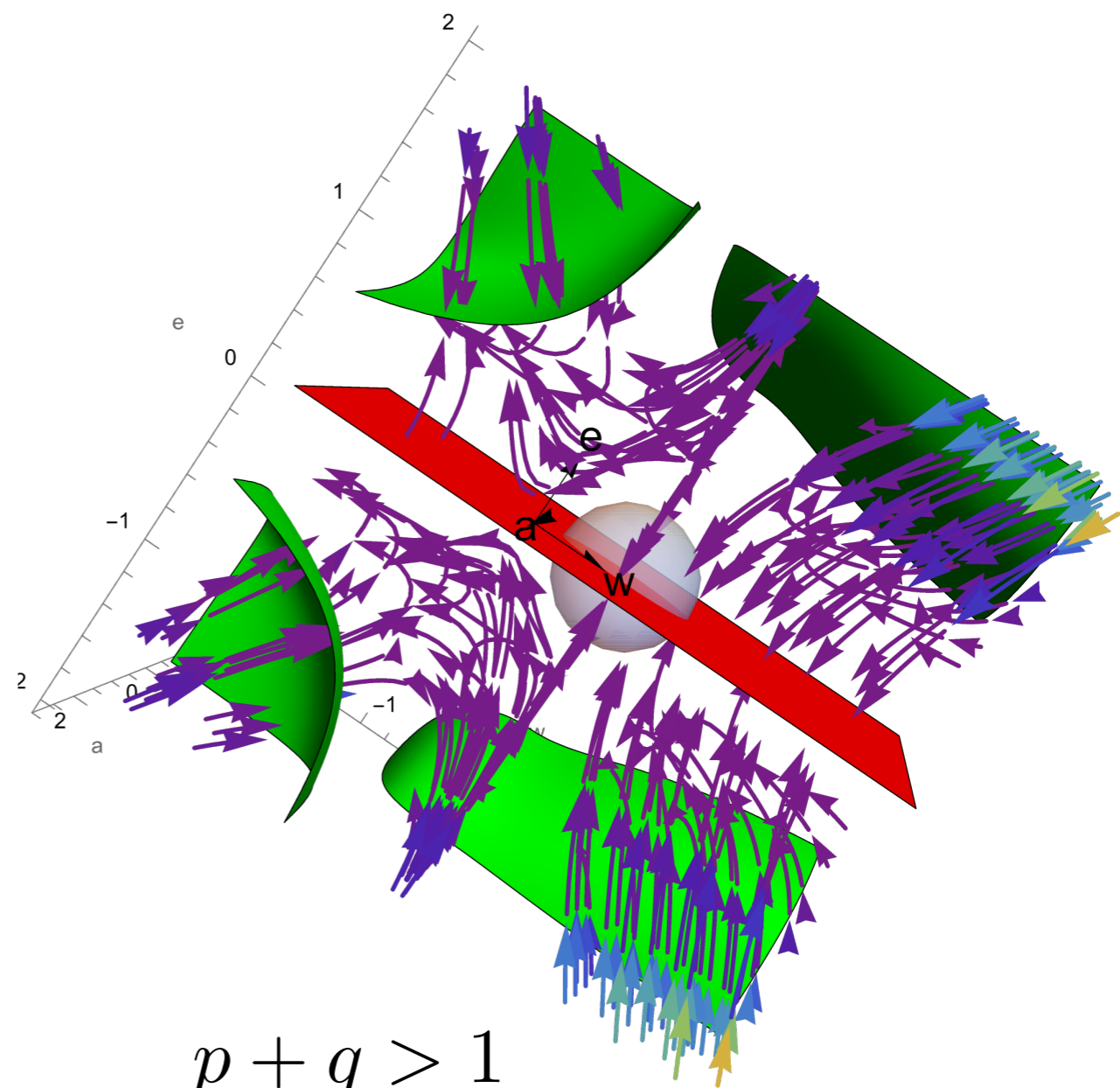
Gets stuck at local minima!



$$p + q > 1$$



$$p + q < 1$$



$$p + q > 1$$

■ Global minima
 ■ Local minima
 ■ Ball around origin



Markovian inputs



Transformers



Memory = 1

Depth = 1



What do they learn?



How do they learn?



Single-layer transformers sometimes fail to learn even first-order Markov chains!



Memory = 1

Depth = 1

Markovian switching and initialization play a key role in the learning dynamics



Part II



Markovian inputs

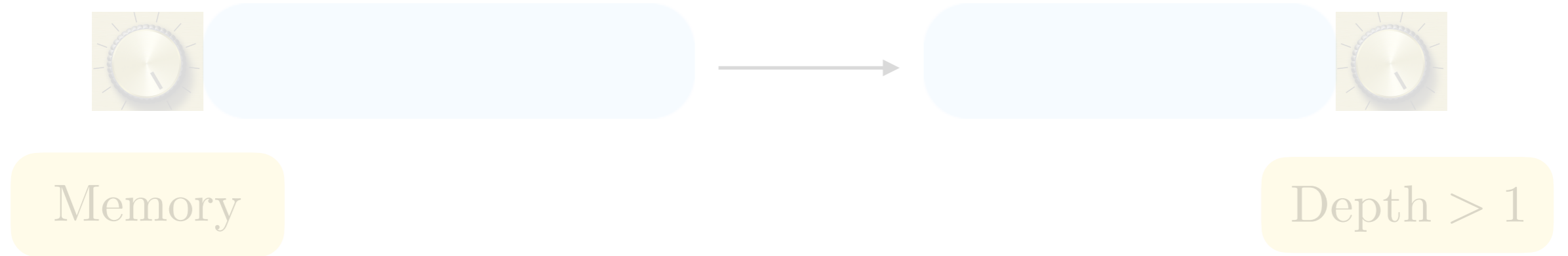


Transformers



Memory

Depth > 1

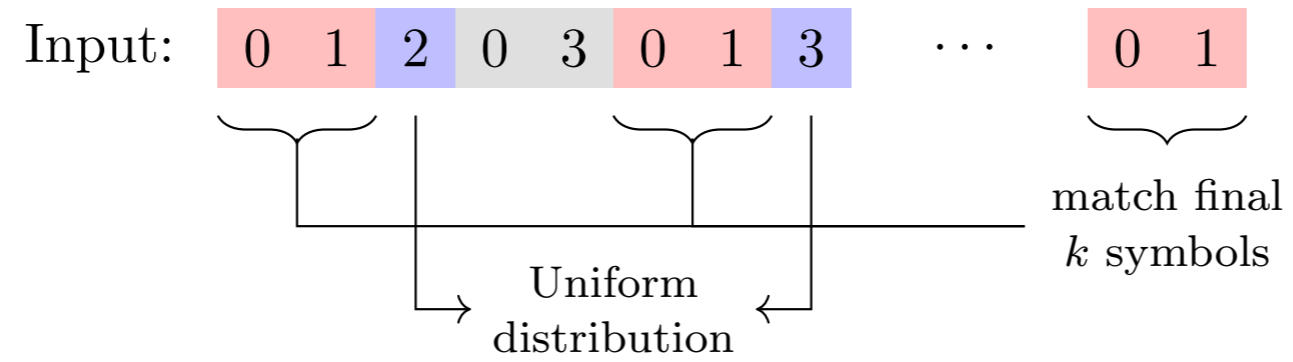


In-context learning (ICL) emerges!

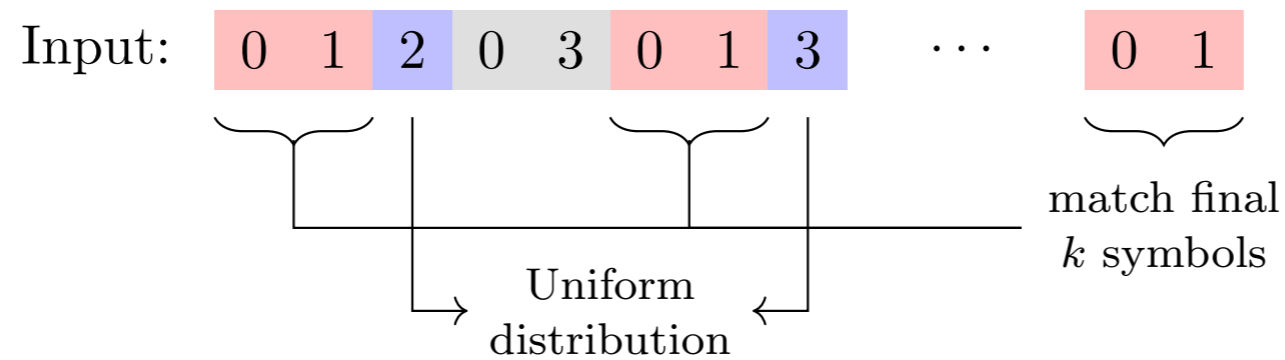
ICL

Mr and Mrs **Dursley**, of number four, Privet Drive, were proud to say that they were perfectly normal, thank you very much. They were the last people you'd expect to be involved in anything strange or mysterious, because they just didn't hold with such nonsense. Mr **Dursley** was the director of a firm called Grunnings, which made drills. He was a big, beefy man with hardly any neck, although he did have a very large moustache. Mrs **Durs**_____

ICL for Markov inputs

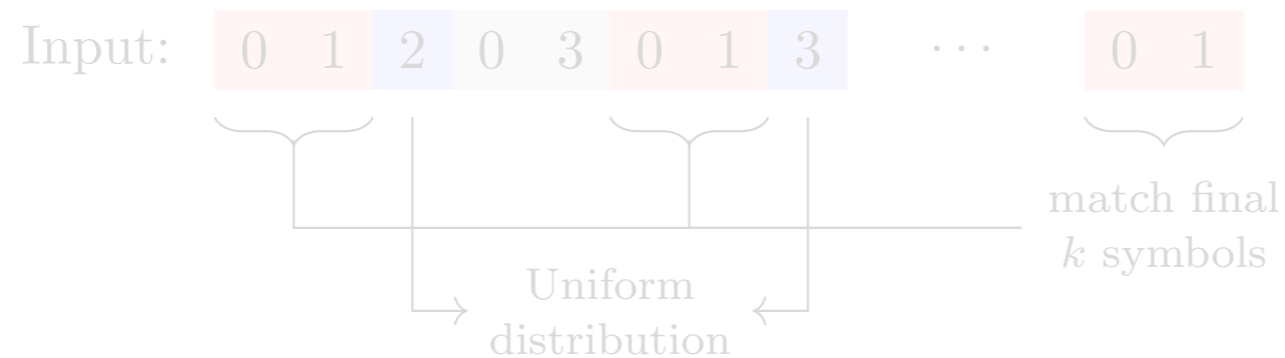


ICL for kth-order Markov



Main idea: Use historical tokens of same context as x_t to predict x_{t+1}

ICL for kth-order Markov



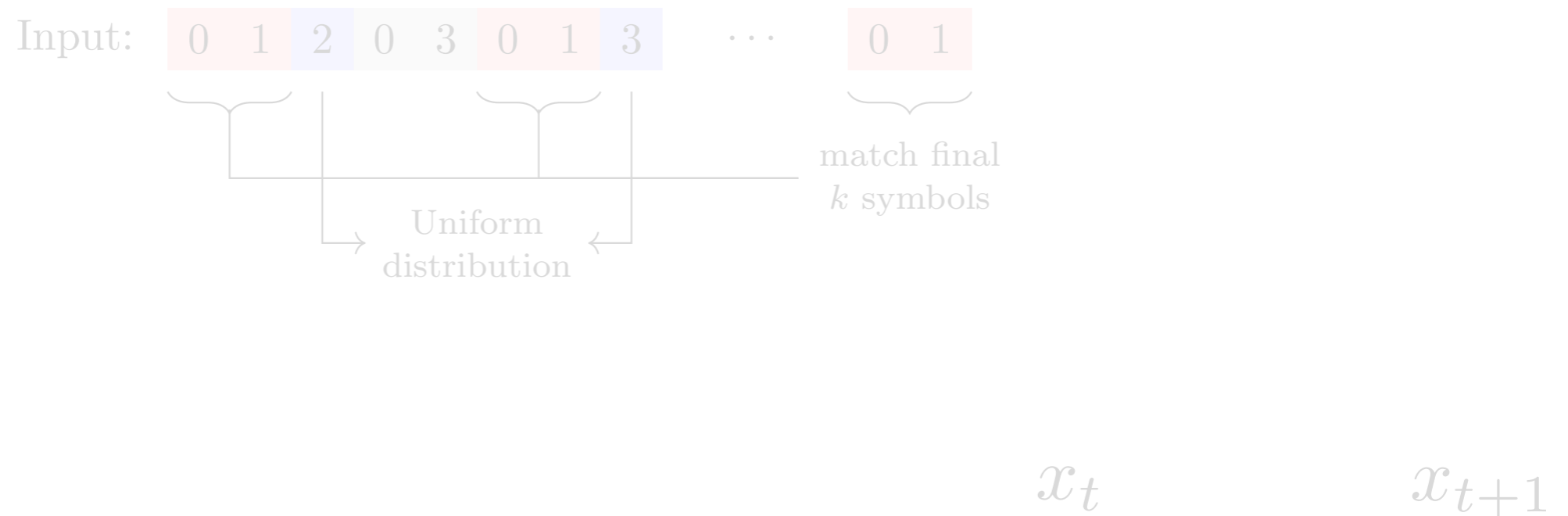
x_t

x_{t+1}

In-context estimator:

$$\widehat{\text{Pr}}_k(x \mid x_1, \dots, x_t) \triangleq \frac{\sum_{i=k+1}^t \mathbb{I}(x_i = x, x_{i-1} = x_t, \dots, x_{i-k} = x_{t-k+1})}{\sum_{i=k+1}^t \mathbb{I}(x_{i-1} = x_t, \dots, x_{i-k} = x_{t-k+1})}$$

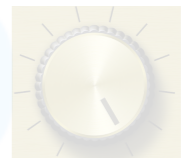
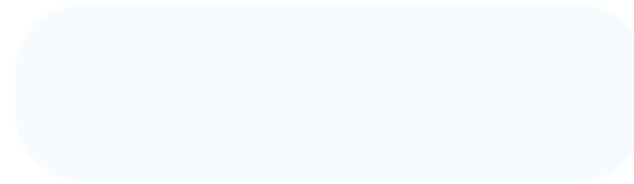
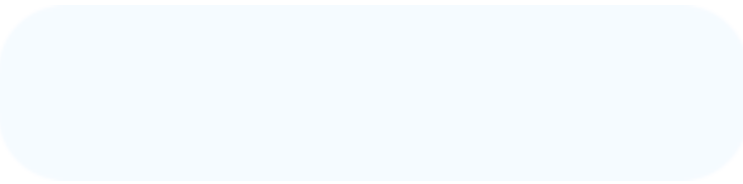
ICL for kth-order Markov



In-context estimator:

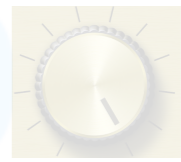
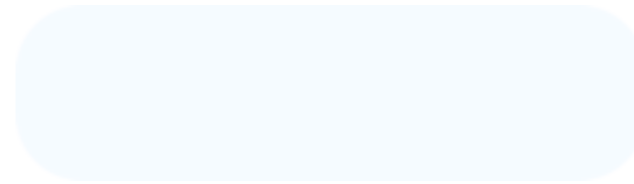
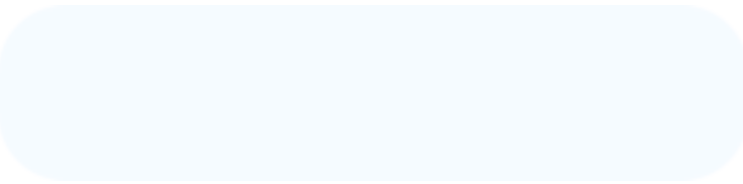
Context-matching

$$\widehat{\text{Pr}}_k(x \mid x_1, \dots, x_t) \triangleq \frac{\sum_{i=k+1}^t \mathbb{I}(x_i = x, x_{i-1} = x_t, \dots, x_{i-k} = x_{t-k+1})}{\sum_{i=k+1}^t \mathbb{I}(x_{i-1} = x_n, \dots, x_{i-k} = x_{t-k+1})}$$



Memory = k

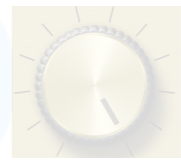
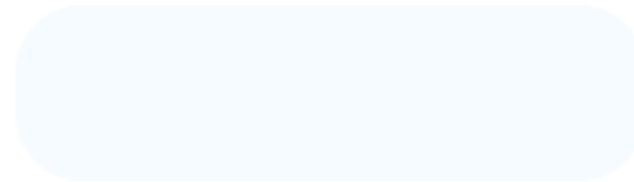
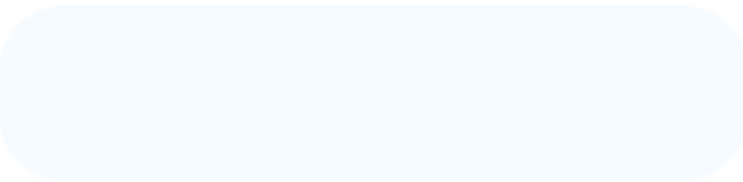
Depth



Memory = k

Depth?

What we know so far

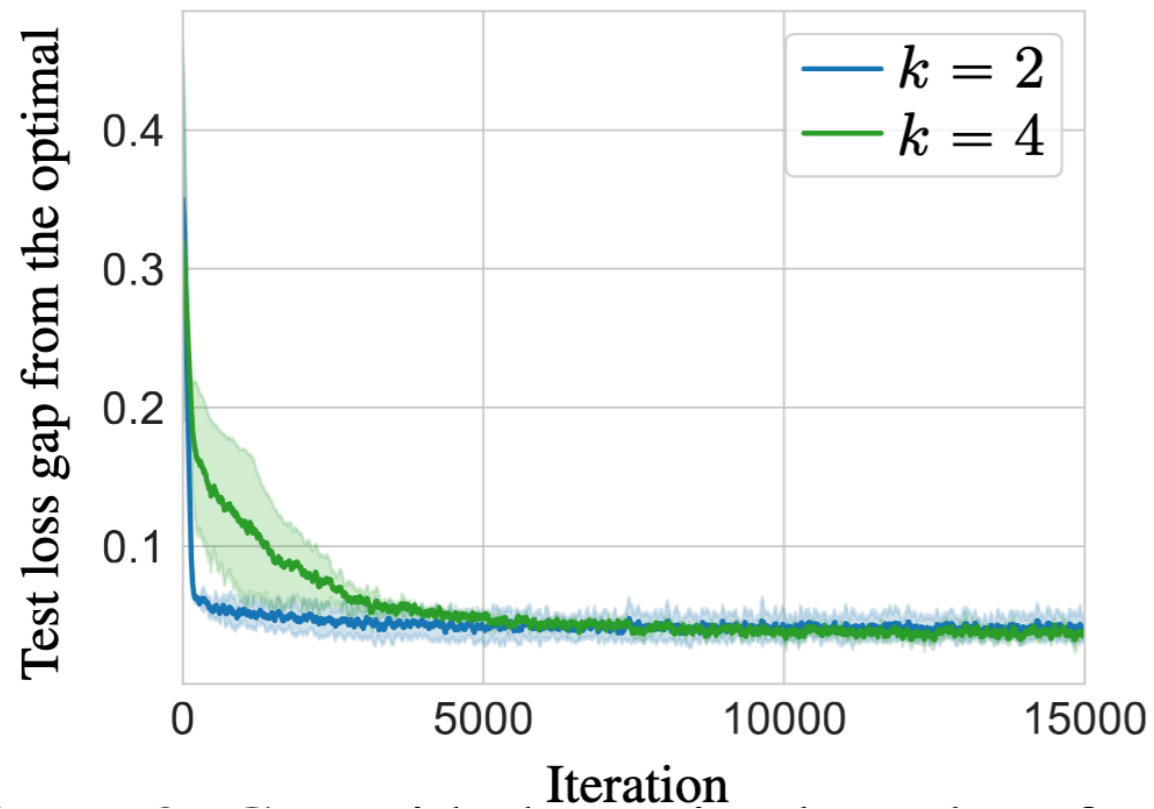


Memory = k

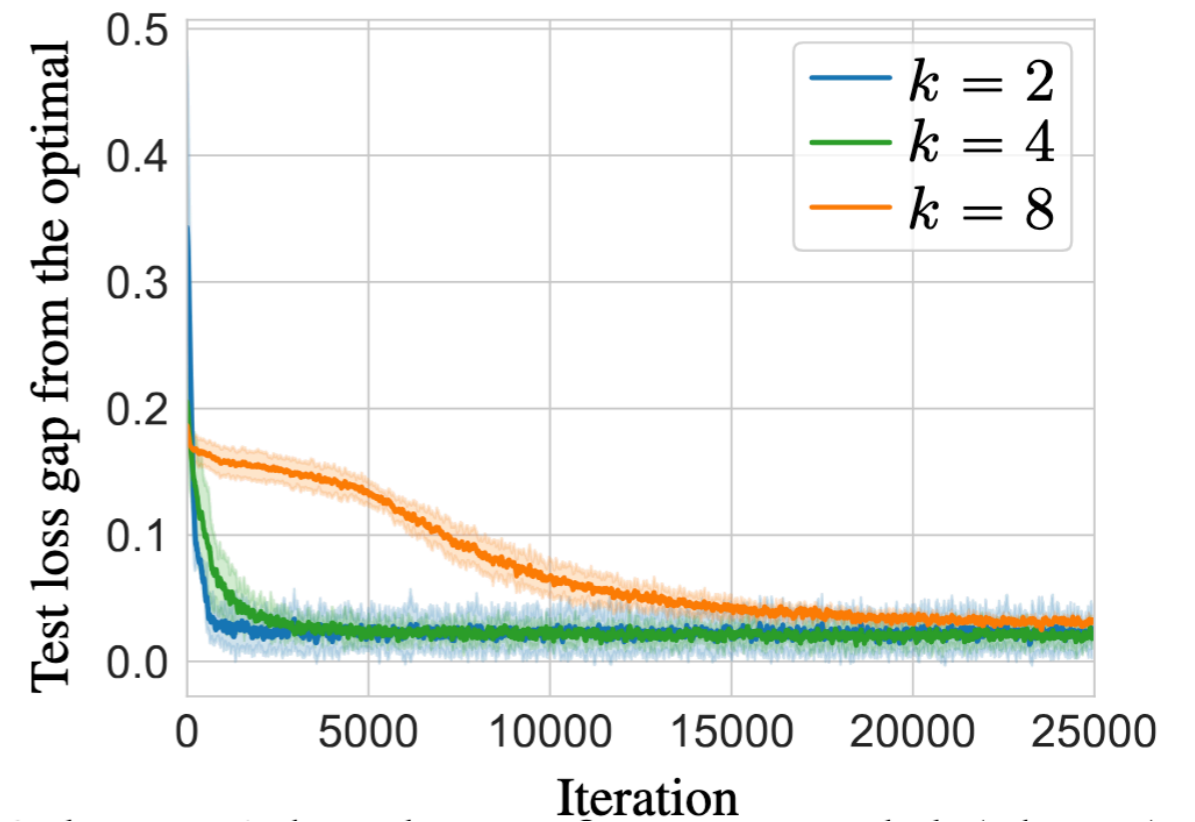
Depth?

- ▶ (Edelman et al., 2024 & Niching et al., 2024): Number of heads/layers should scale with k

But...



2-layer transformer



3-layer transformer

What's happening

What's happening

Main result (Constant depth suffices)

Any order k in-context estimator can be represented by a transformer with 3 layers, 1 head per layer, relative positional encodings and layer norm

What's happening

Main result (Constant depth suffices)

Any order k in-context estimator can be represented by a transformer with 3 layers, 1 head per layer, relative positional encodings and layer norm

Without non-linearities, you need logarithmic depth

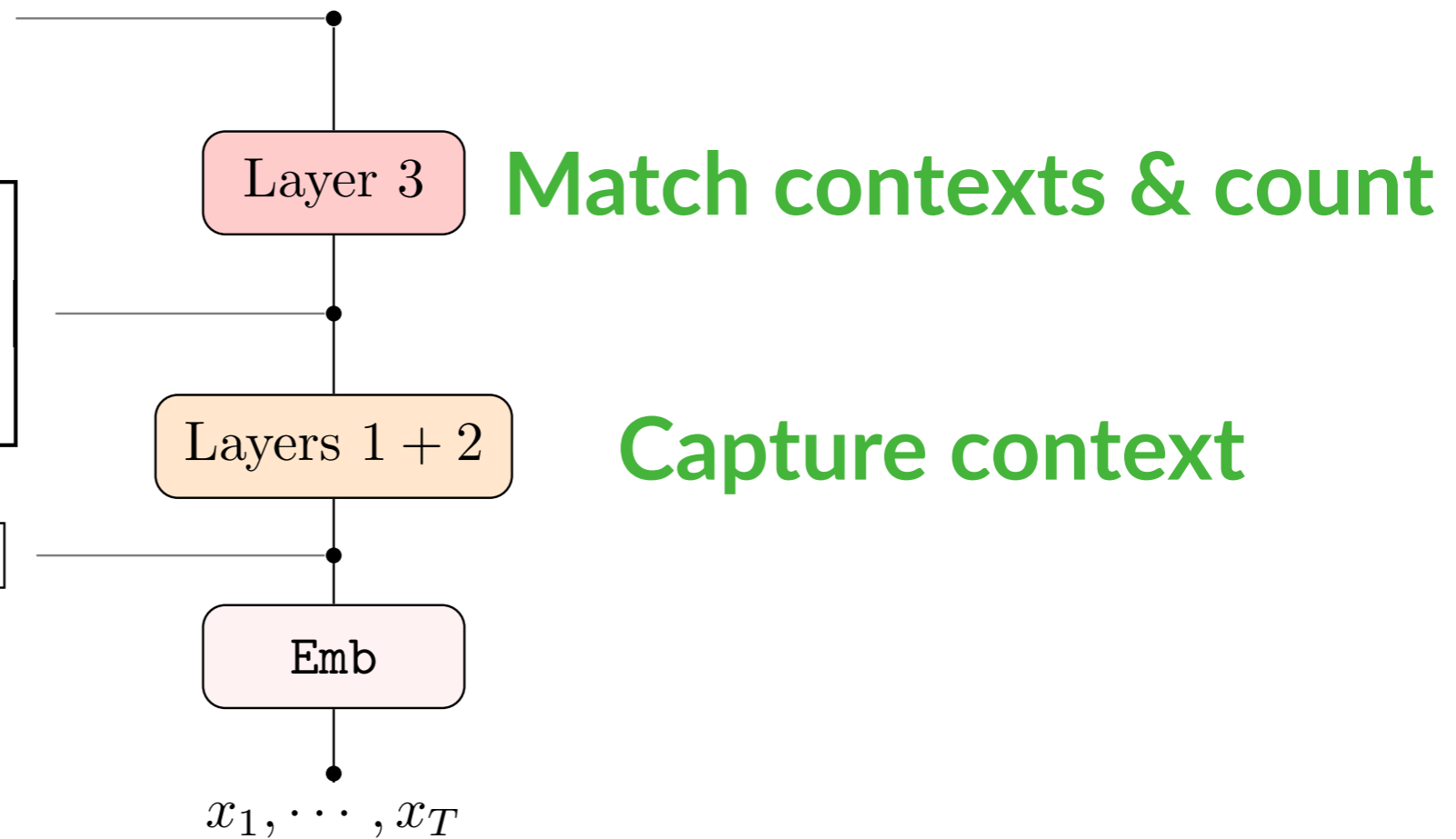
Intuition

$$\text{att}_{T,n} \propto \exp\left(\kappa \frac{\langle \mathbf{v}_n, \mathbf{u}_T \rangle}{\|\mathbf{v}_n\|_2 \|\mathbf{u}_T\|_2}\right)$$

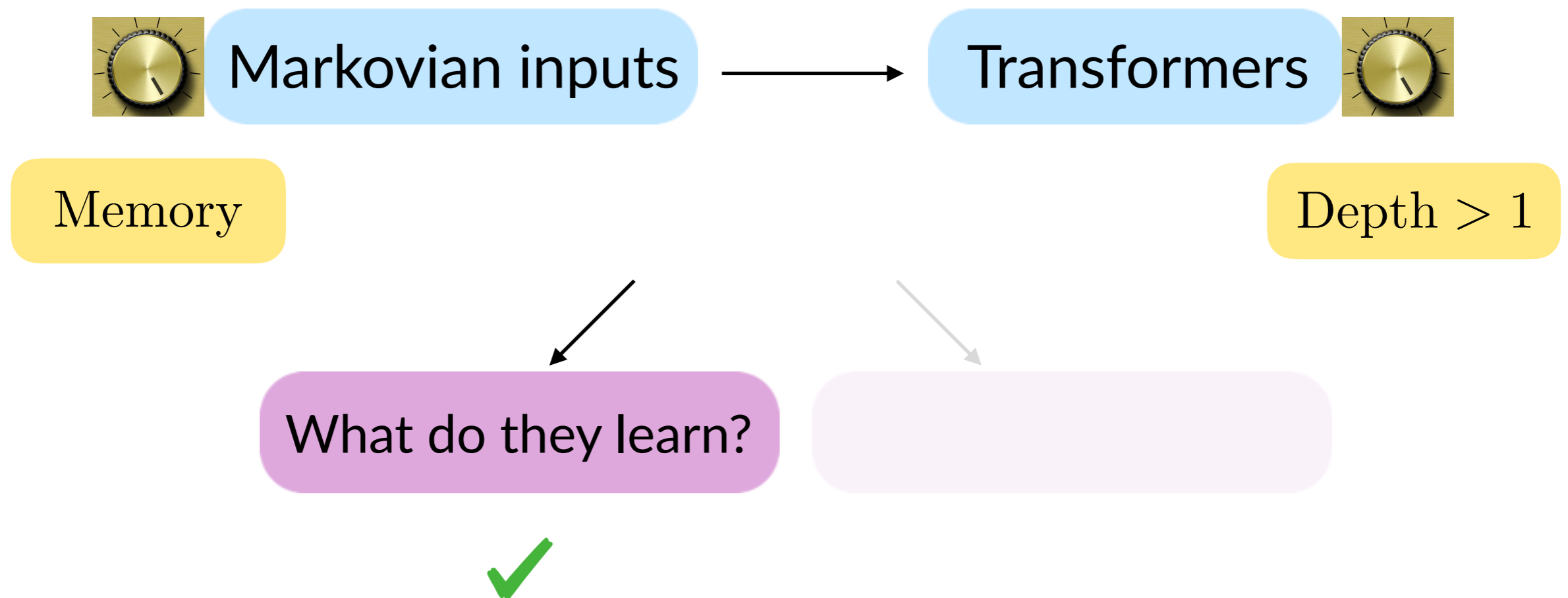
(realizes a k^{th} -order induction head)

$$\left[\begin{array}{c|c} \text{Emb}(x_n) & \text{Emb}(x_T) \\ \hline \dots & \dots \\ \frac{\mathbf{u}_n}{\|\mathbf{u}_n\|_2} & \frac{\mathbf{u}_T}{\|\mathbf{u}_T\|_2} \\ \frac{\mathbf{v}_n}{\|\mathbf{v}_n\|_2} & \frac{\mathbf{v}_T}{\|\mathbf{v}_T\|_2} \end{array} \right]$$

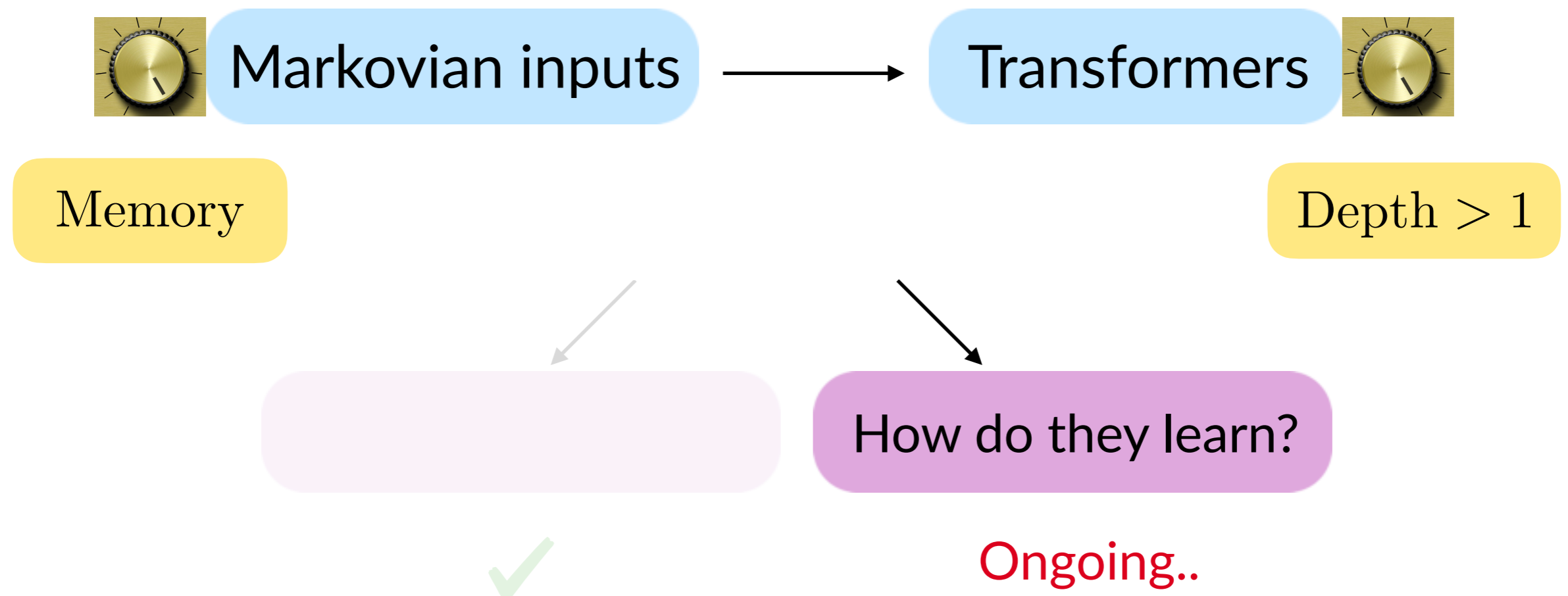
$$\left[\dots \mid \text{Emb}(x_n) \mid \dots \mid \text{Emb}(x_T) \right]$$

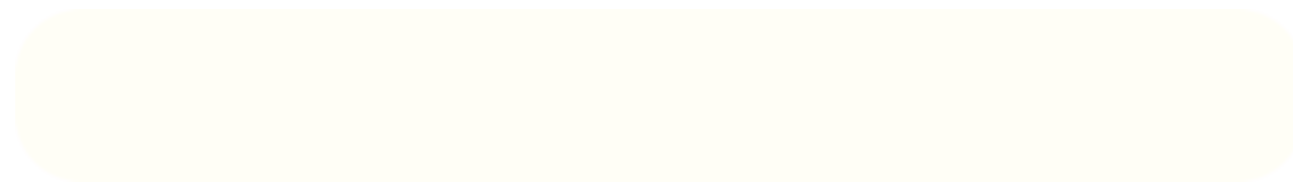


Summary

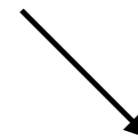


Summary





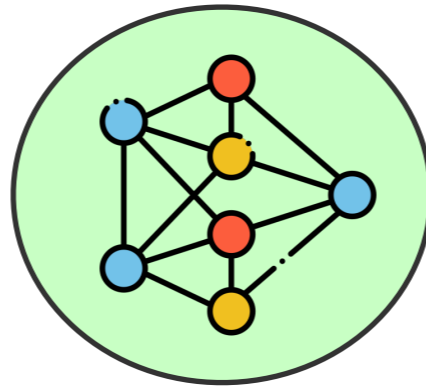
Principled frameworks and tools



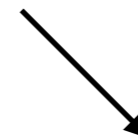
What do they learn?

How do they learn?

Transformers



Attention with Markov



What do they learn?

How do they learn?



Future directions

- ▶ More realistic data models: Hidden-Markov Models
- ▶ State-space models
- ▶ More practical guidelines for training LLMs

