

From Markov to Laplace:

A Markovian Tale of Large Language Models (LLMs)

Ashok Vardhan Makkuva

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



Berkeley
NEVER EXPOSE THE WRONGDOING OF THE
UNIVERSITY OF CALIFORNIA



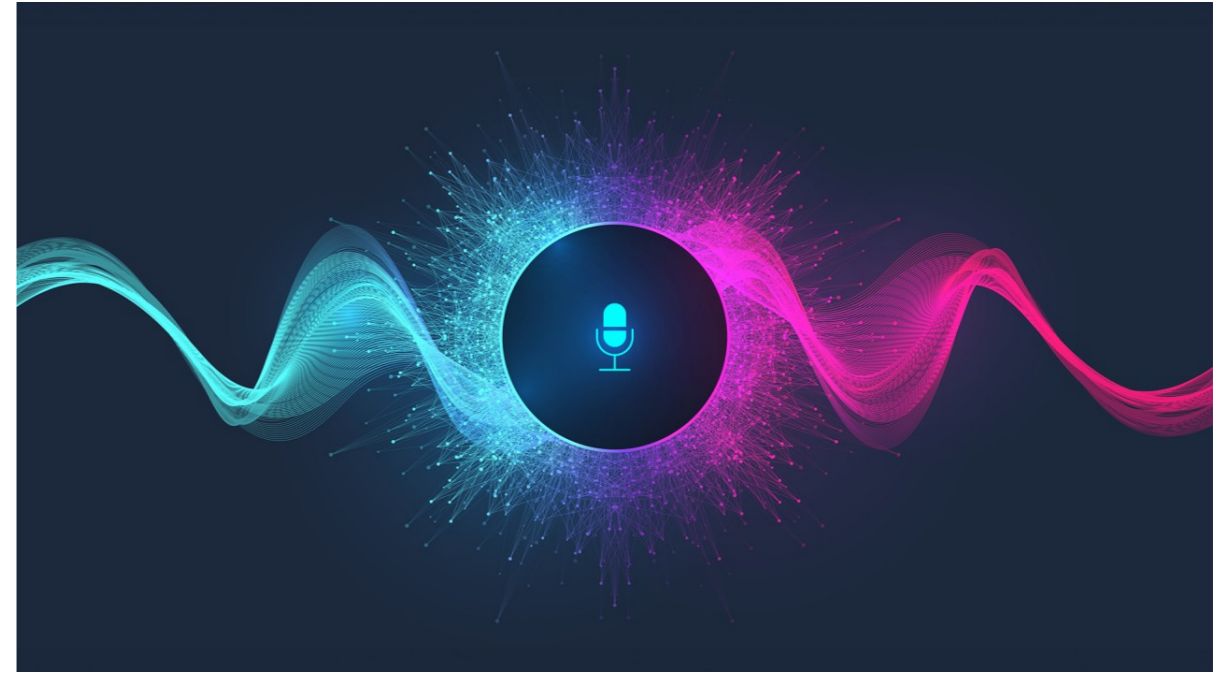
Massachusetts
Institute of
Technology

EPFL

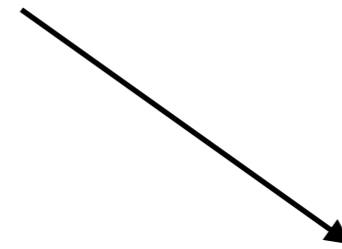
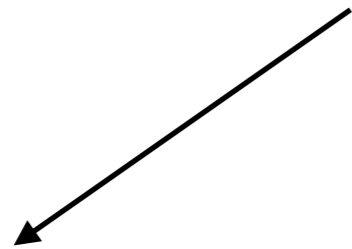


TEXAS
The University of Texas at Austin

LLMs are part of daily lives



Not fully Trustworthy



Impressive language skills

Arithmetic reasoning

Programming

Hallucinations



Need of the hour



LLM Interpretability



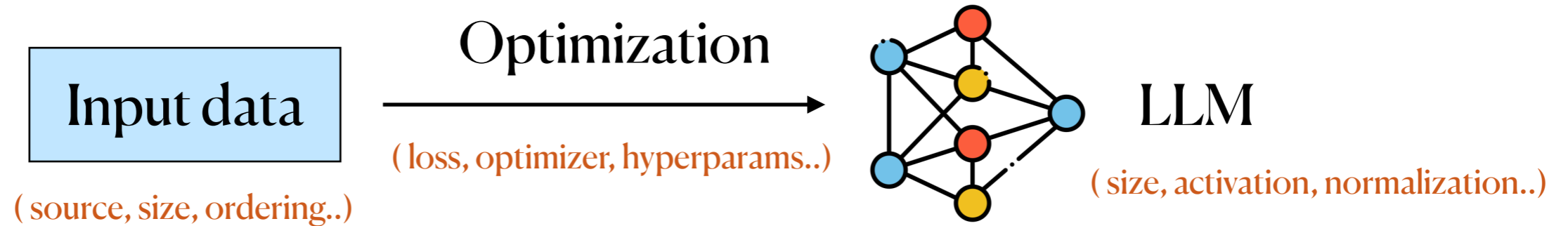
What do they learn?

How do they learn?

Challenges



Challenges

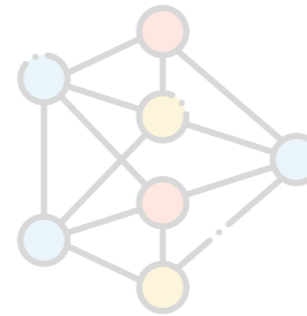
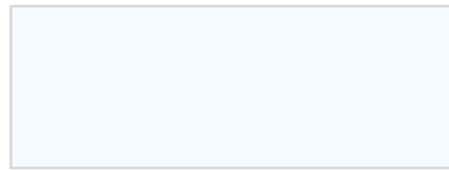


- Complex and opaque
- How/why do they work?
- Expensive



Too many tuning knobs \$63M

Need



- Complex and opaque

Principles



Clarity & scientific understanding

- How/why do they work?

Transparency



Core learning mechanisms

- Expensive

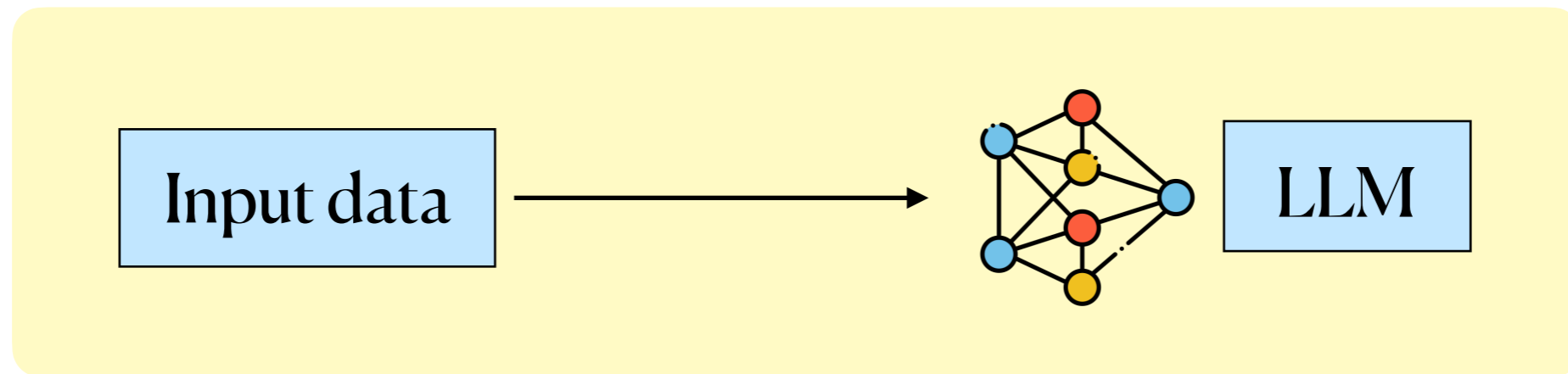
Efficiency



Lighter & efficient models, algos.

LLM Interpretability

State of the art



What data can they represent?

How do they learn them?

- **Key limitation**

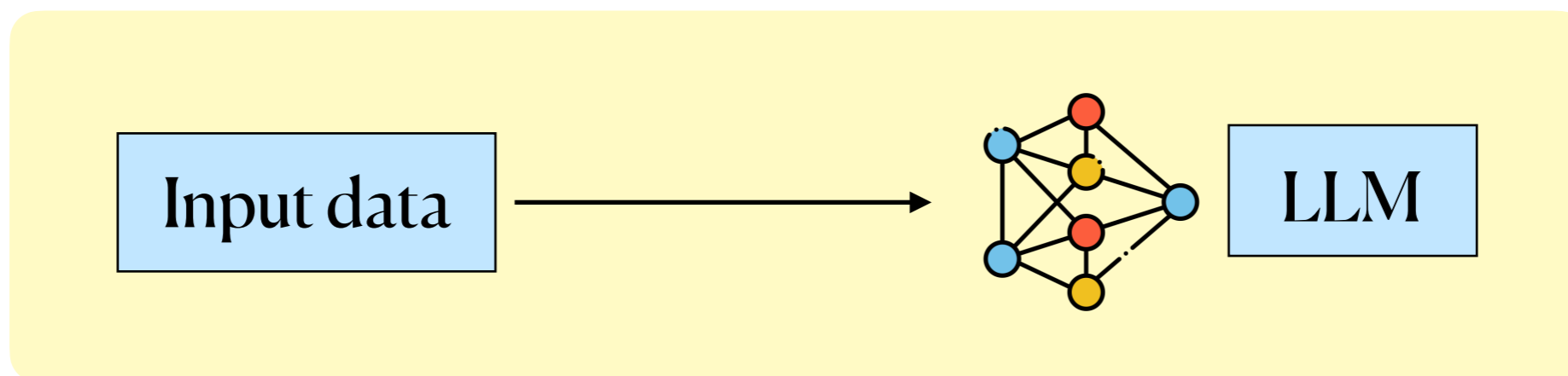
- ▶ **Lack generality:** specific data and models (e.g. transformers), restrictive assumptions

Our approach



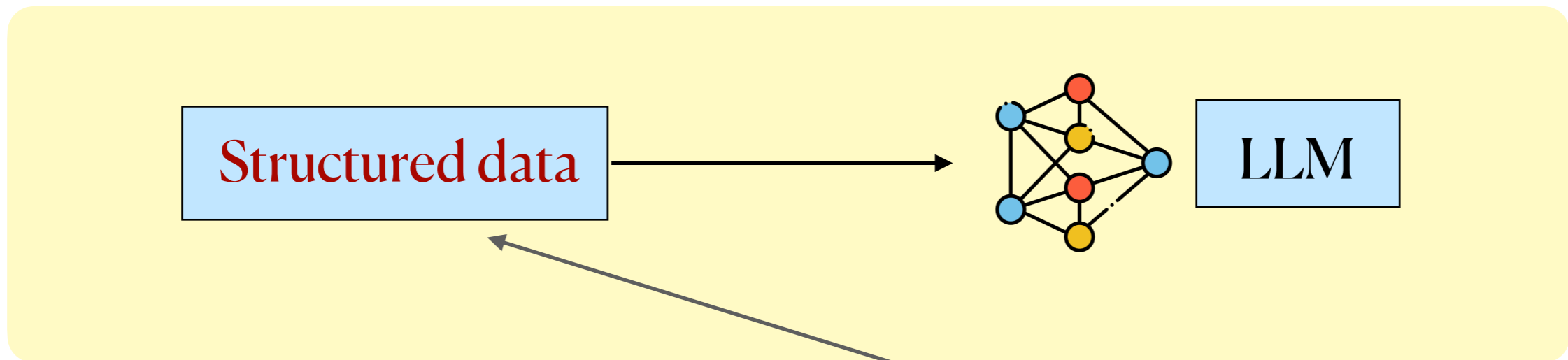
LLM Interpretability via Structured Data

Our approach



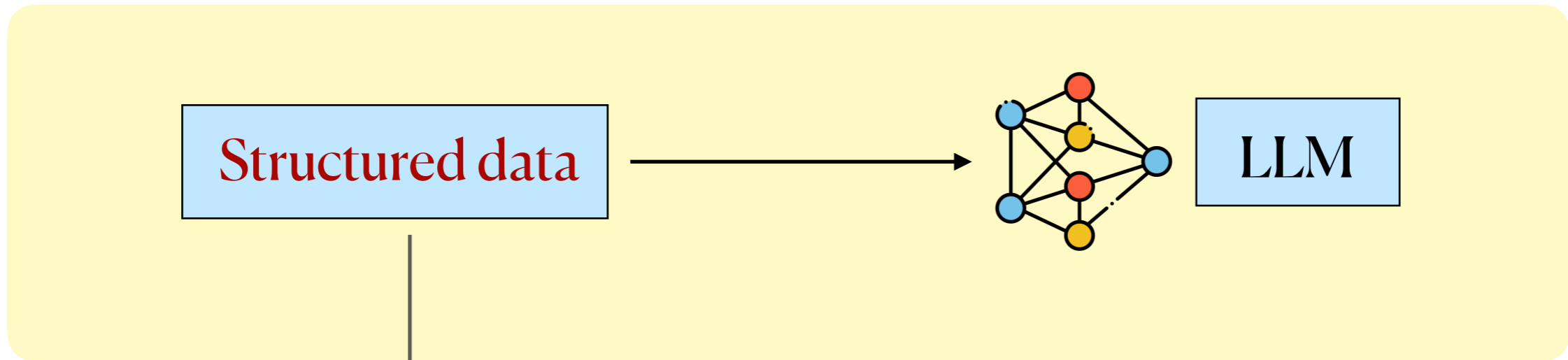
LLM Interpretability via Structured Data

Our approach

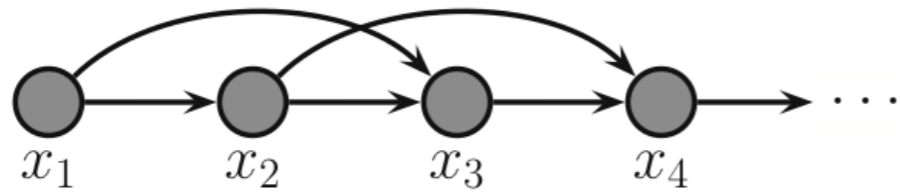


LLM Interpretability via Structured Data

Our approach



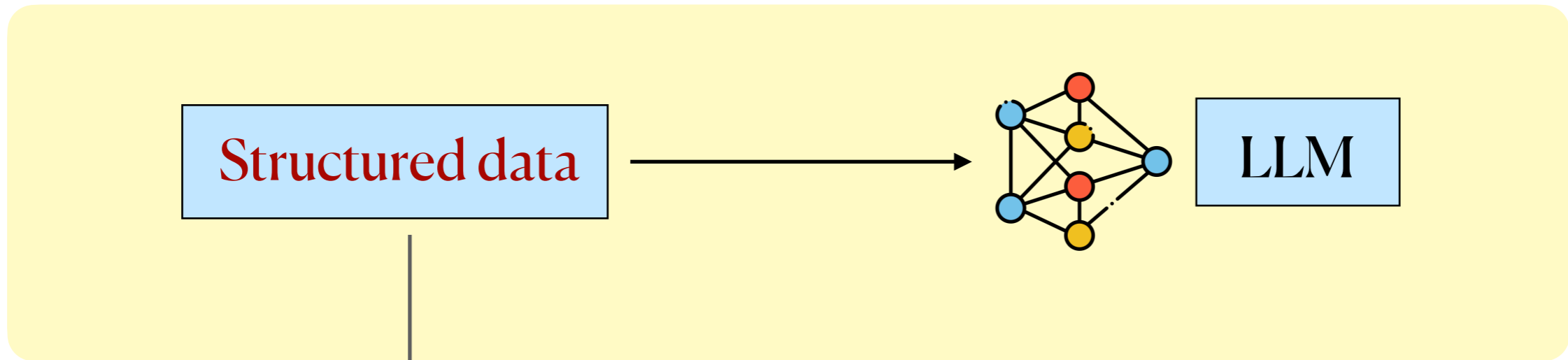
Markov:



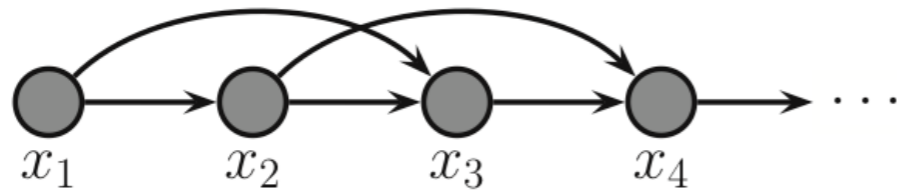
Mathematically controlled

Insightful \rightarrow complex real-world data

Our approach

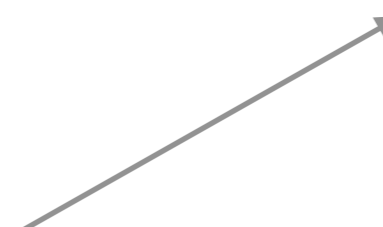


Markov:

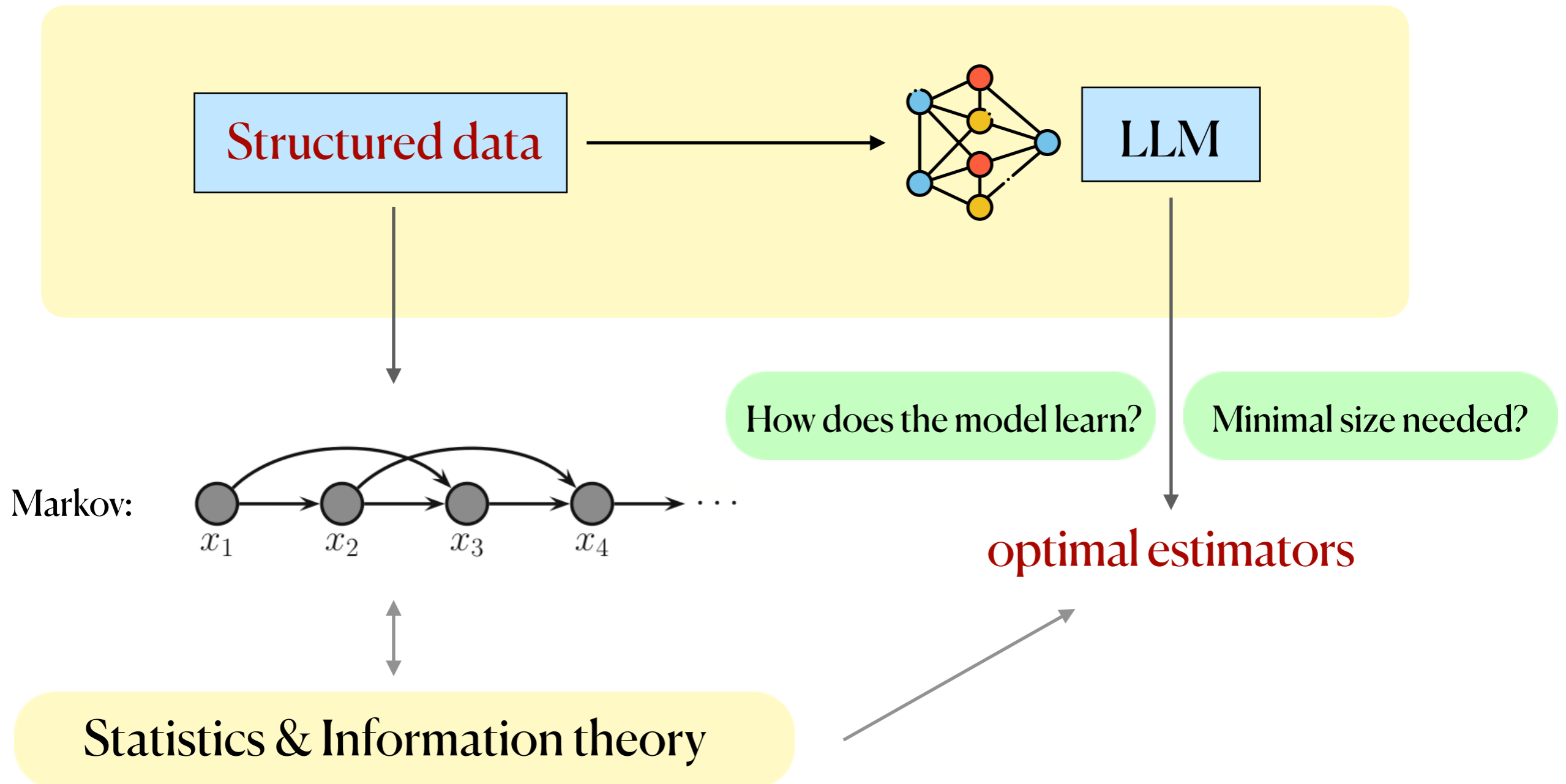


well studied
optimal estimators

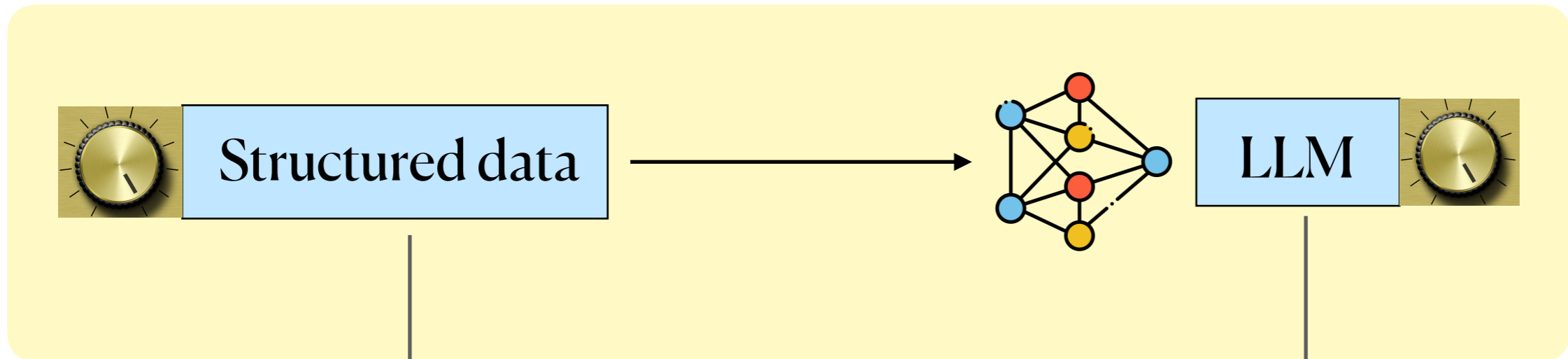
Statistics & Information theory



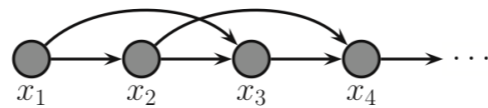
Our approach



This talk



Markov processes



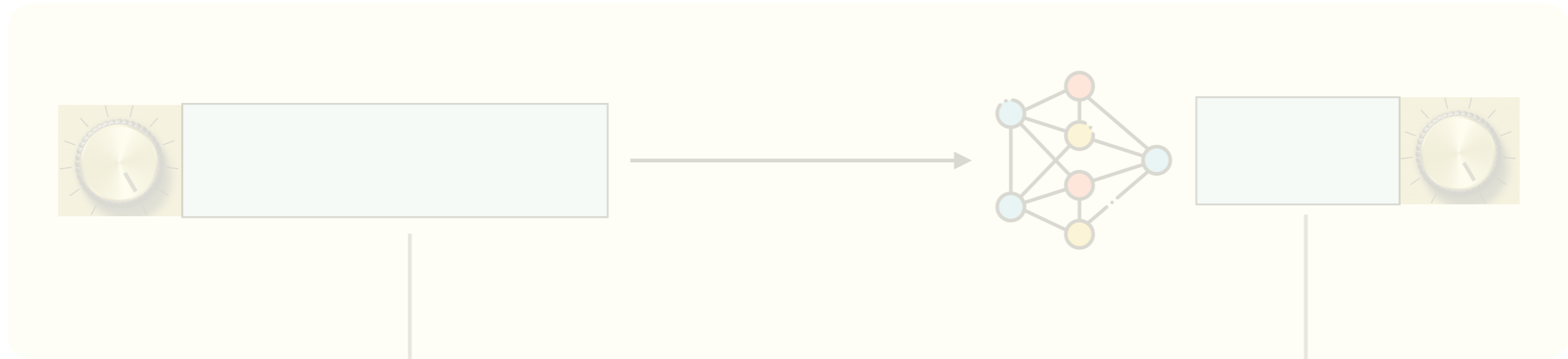
Transformers



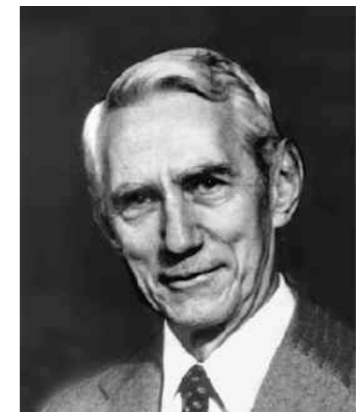
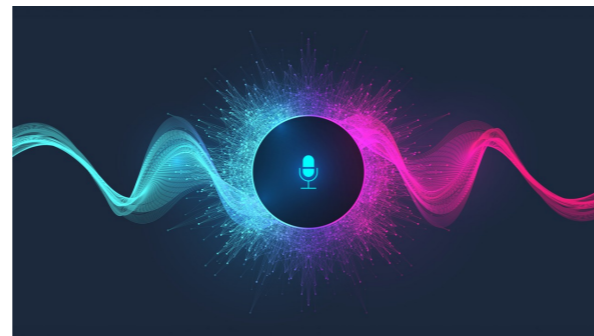
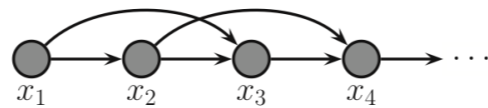
State-Space Models



Why Markovian?

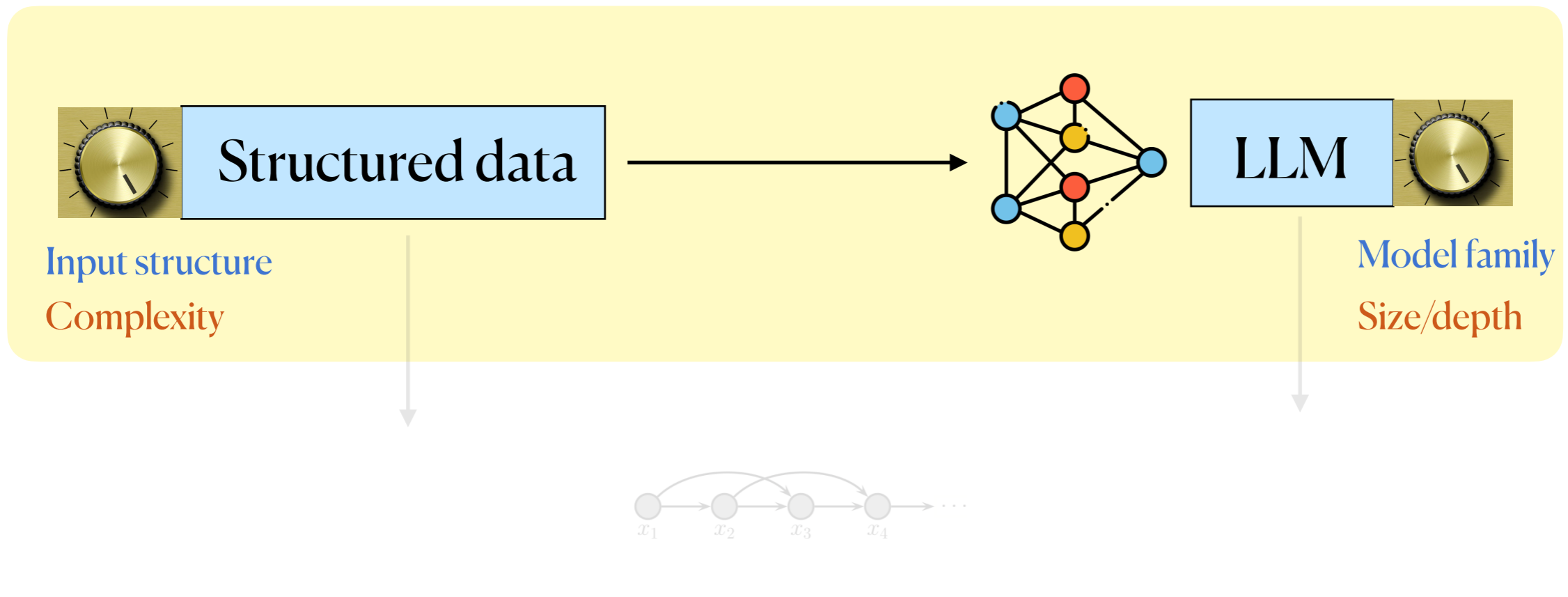


Markov processes

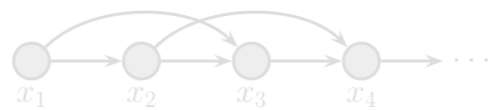
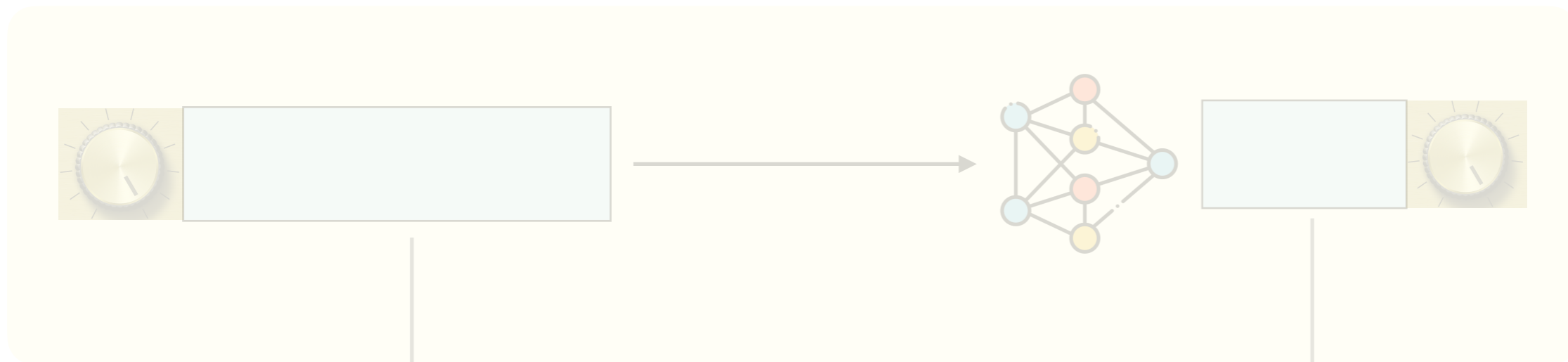


Shannon, 1948

Key features



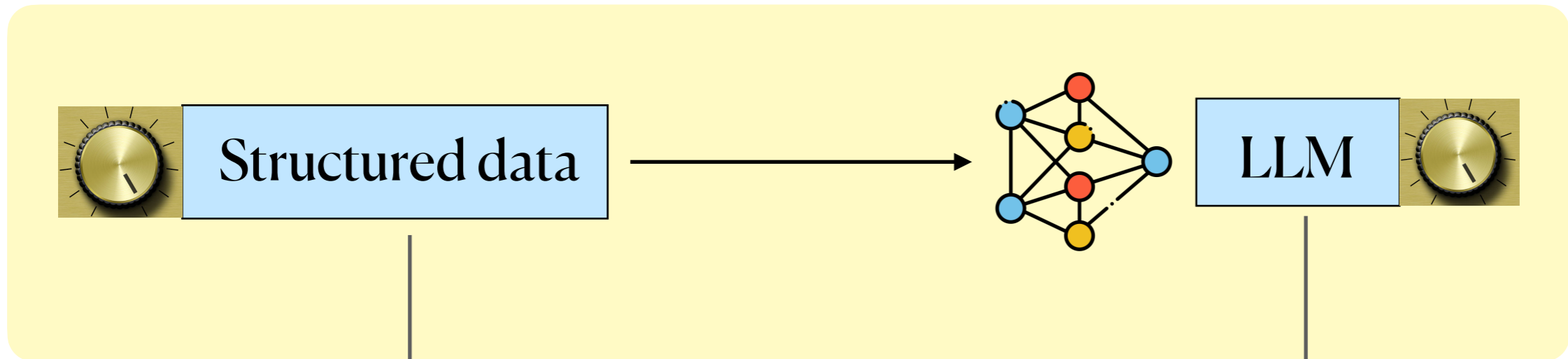
- ▶ **Comprehensive:** broad family of (data, models) → **compare** various architectures
- ▶ **Fine-grained analysis**



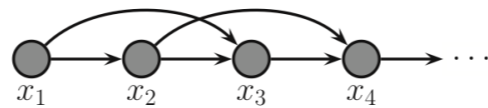
What do they learn?

How do they learn?

This talk



Markov processes



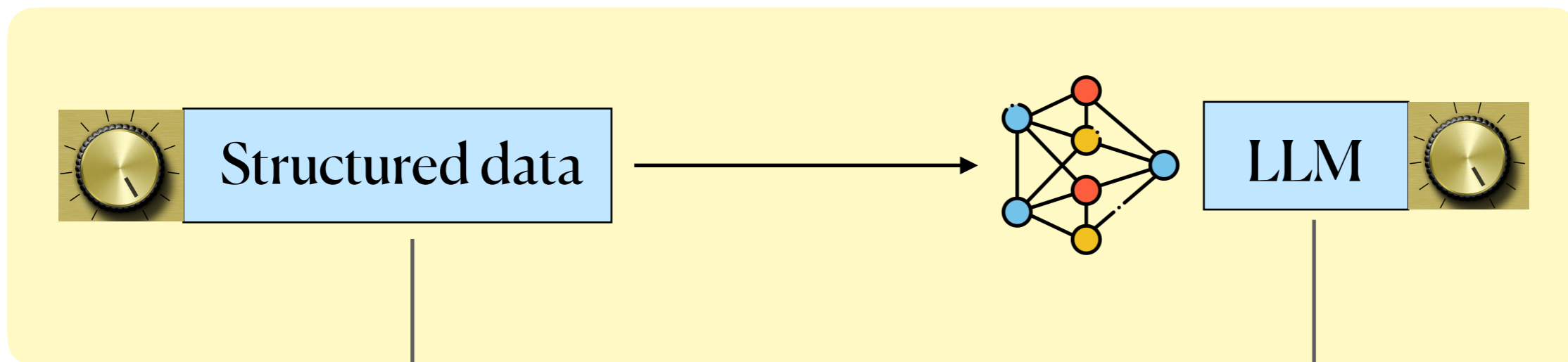
Transformers



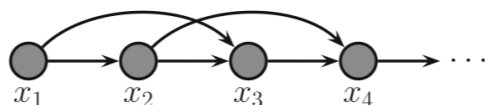
State-Space Models



Part I



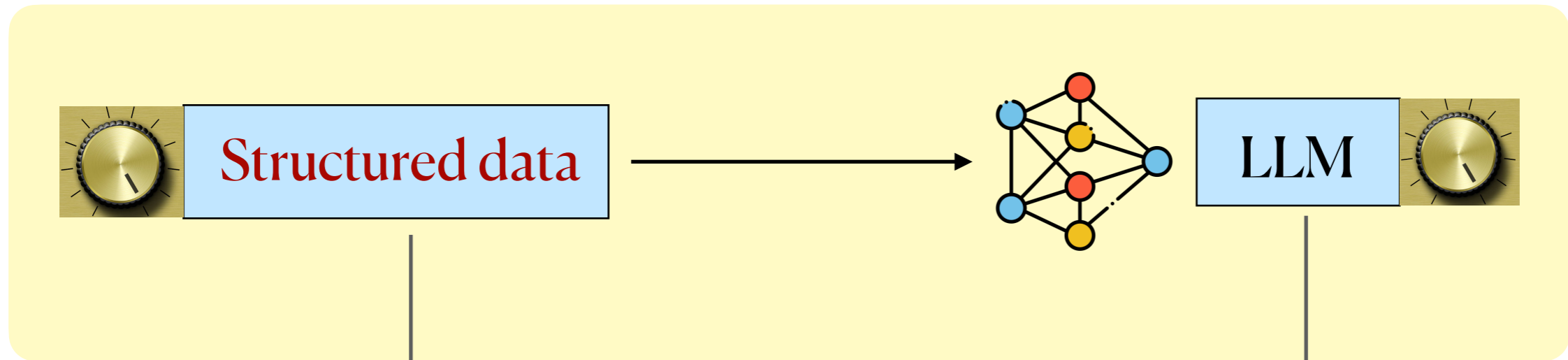
Markov processes



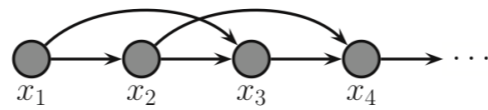
Transformers



Part II



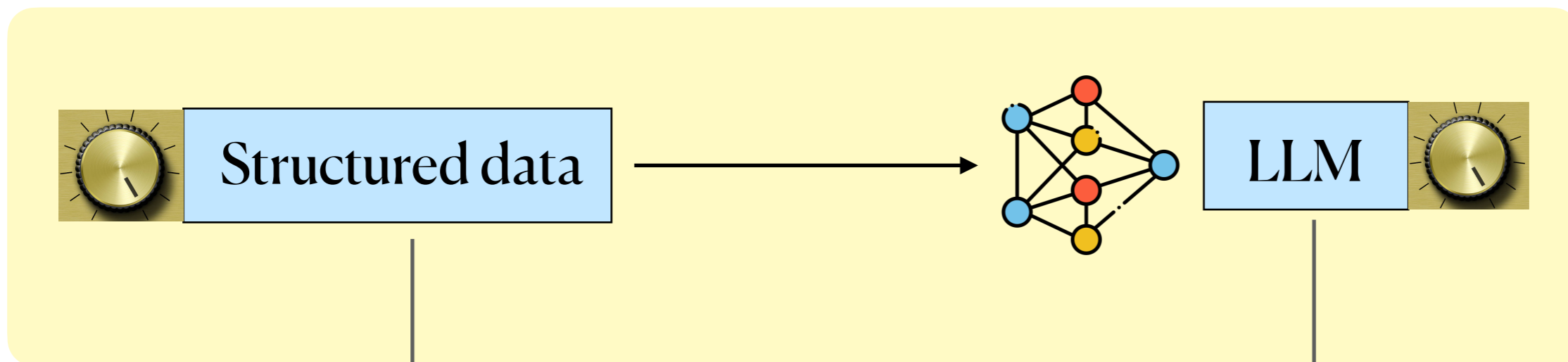
Markov processes



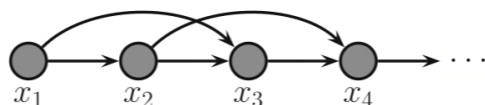
State-Space Models



Part I



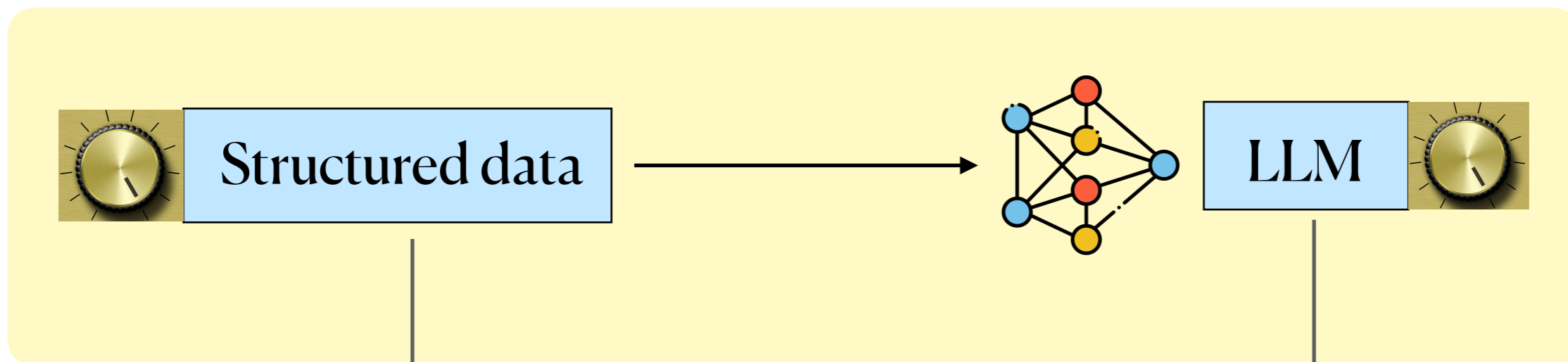
Markov processes



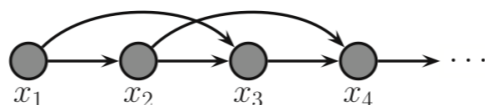
Transformers



Part I



Markov processes

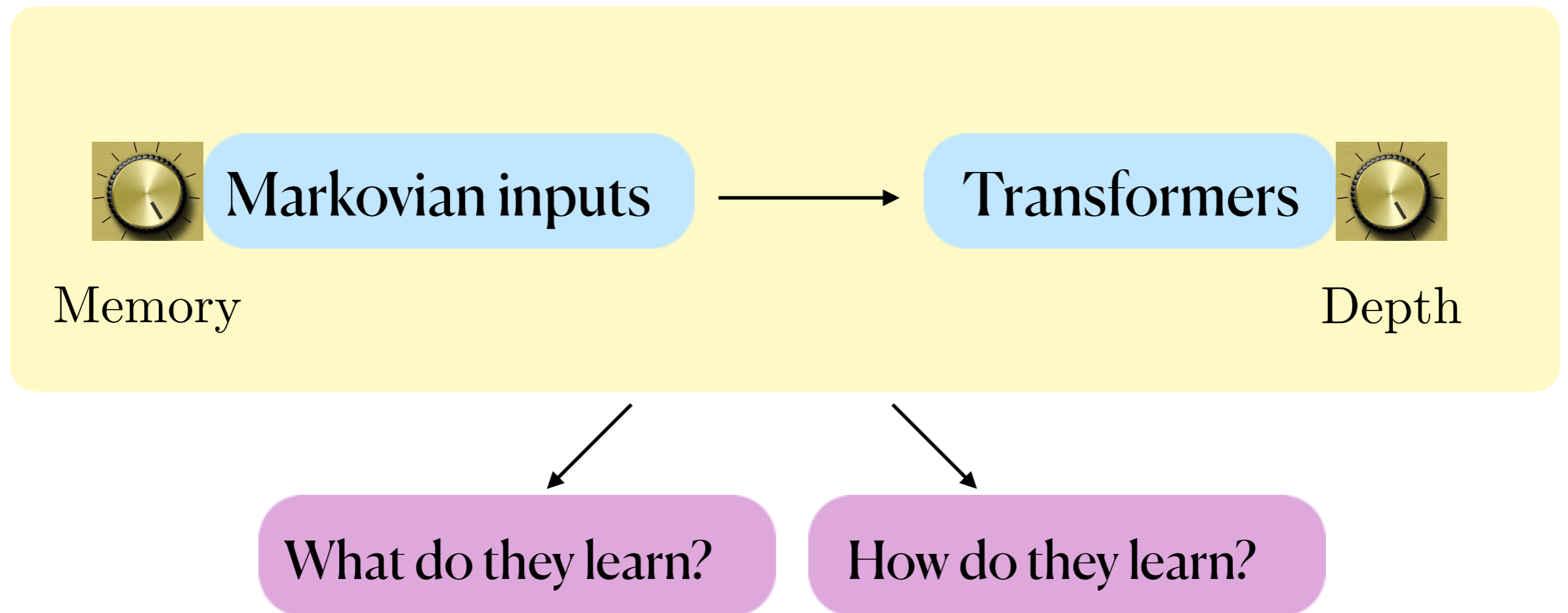


Transformers



Attention with Markov

Attention with Markov



**Depth and non-linearity play a crucial role
in transformer functionality**

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

**2-layer transformer can learn Markov
chains of all orders!**

**What One Cannot Two Can: Two-Layer Transformers
Provably Represent Induction Heads on Any-Order
Markov Chains**

NeurIPS 2025 Spotlight

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

ICLR 2025 Spotlight

**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



Markovian inputs



Transformers



Memory

Depth

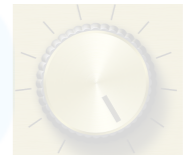
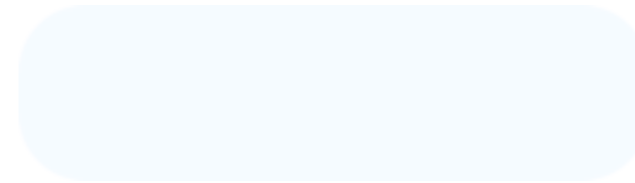
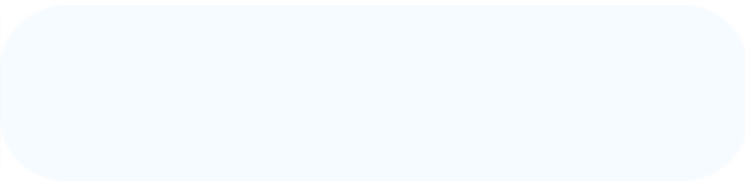
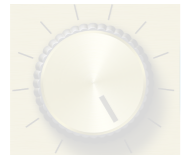
Part I.a



Memory = 1

Depth = 1

Part I.b



Memory

Depth > 1

Part I.a



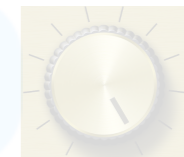
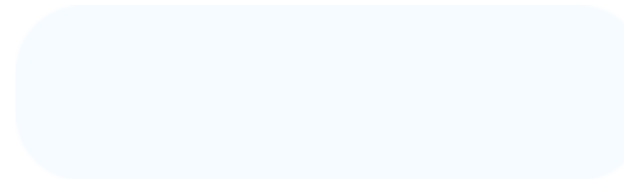
Memory = 1

Depth = 1

Part I.a



Markovian inputs

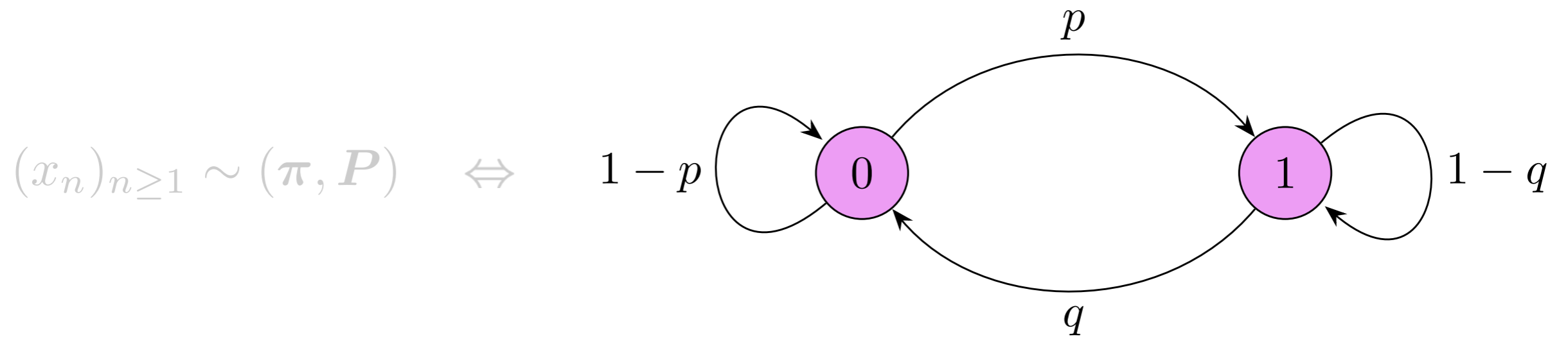


Memory = 1

Depth = 1

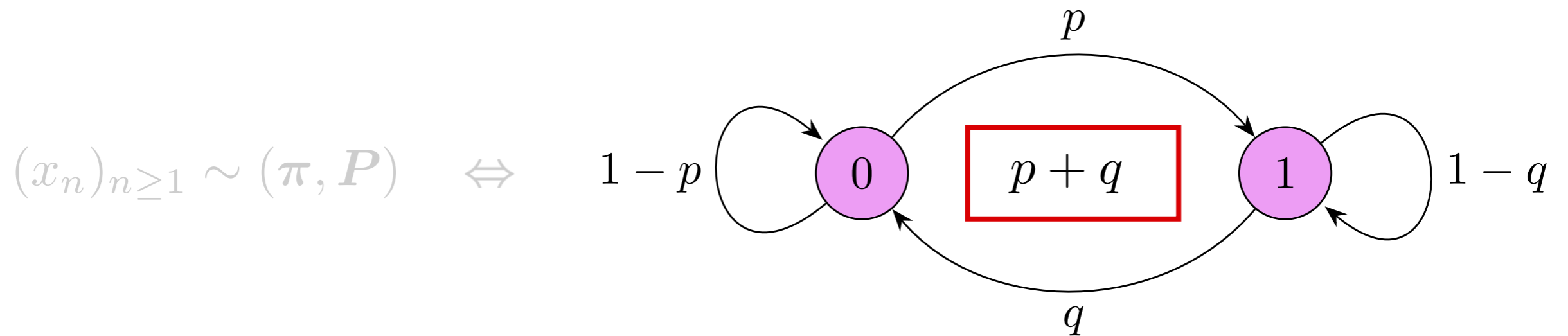
Input data: First-order Markov chain

Input data: First-order Markov chain



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

Input data: First-order Markov chain

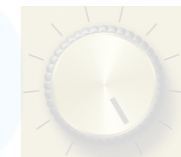
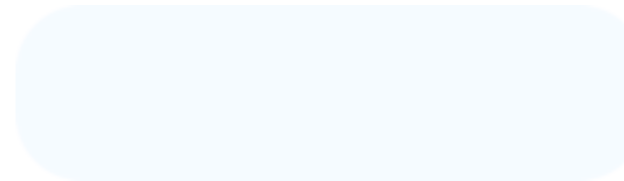


$$\pi = (\pi_0, \pi_1) = \left(\frac{q}{p+q}, \frac{p}{p+q} \right), \quad P = (P_{ij}) = \begin{bmatrix} 1-p & p \\ q & 1-q \end{bmatrix}.$$

Part I.a



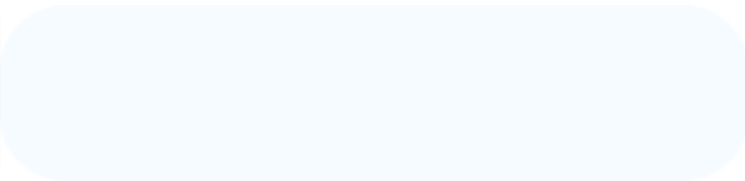
Markovian inputs



Memory = 1

Depth = 1

Part I.a



Memory = 1

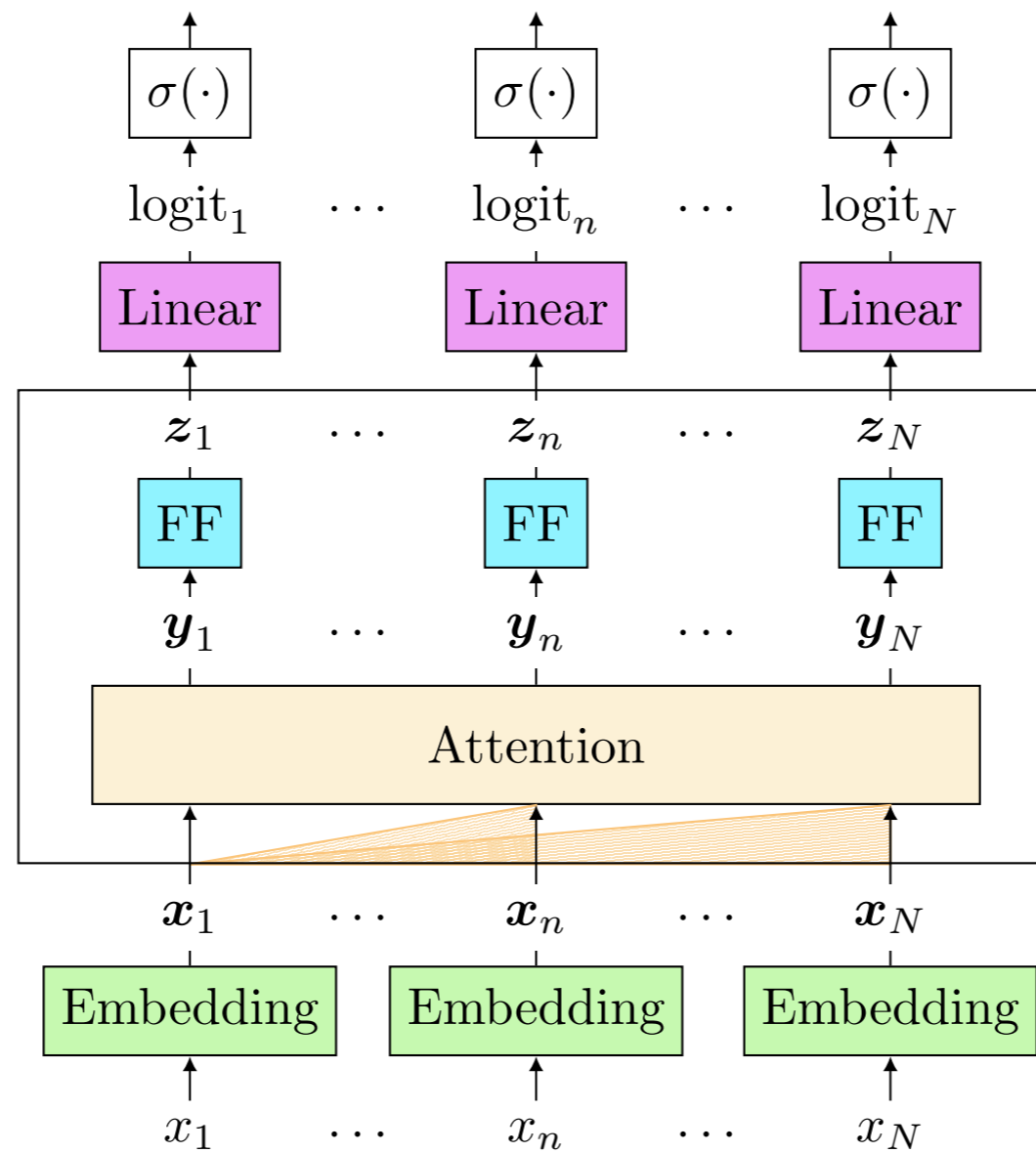
Depth = 1

Transformers

Transformers

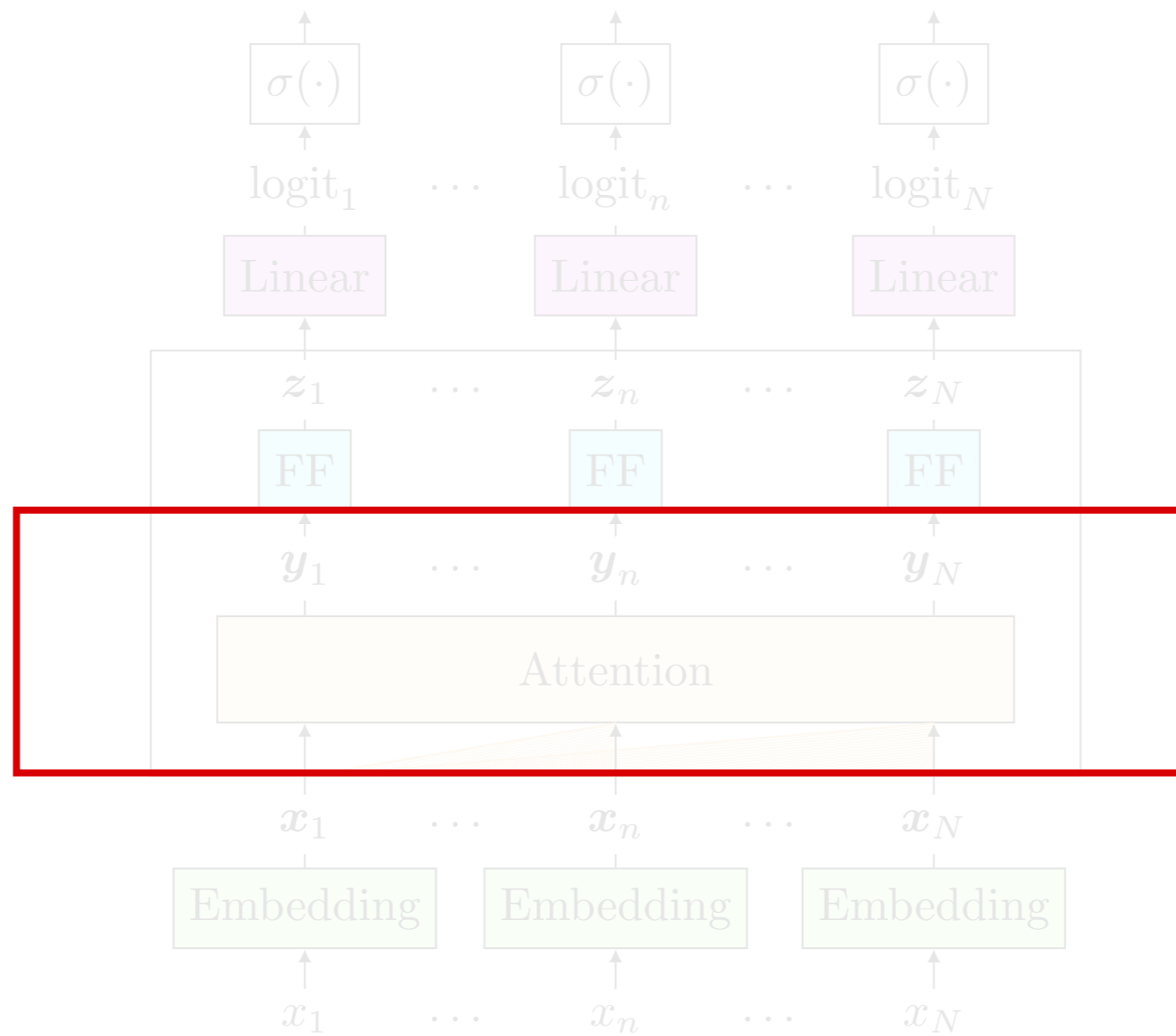


Transformers

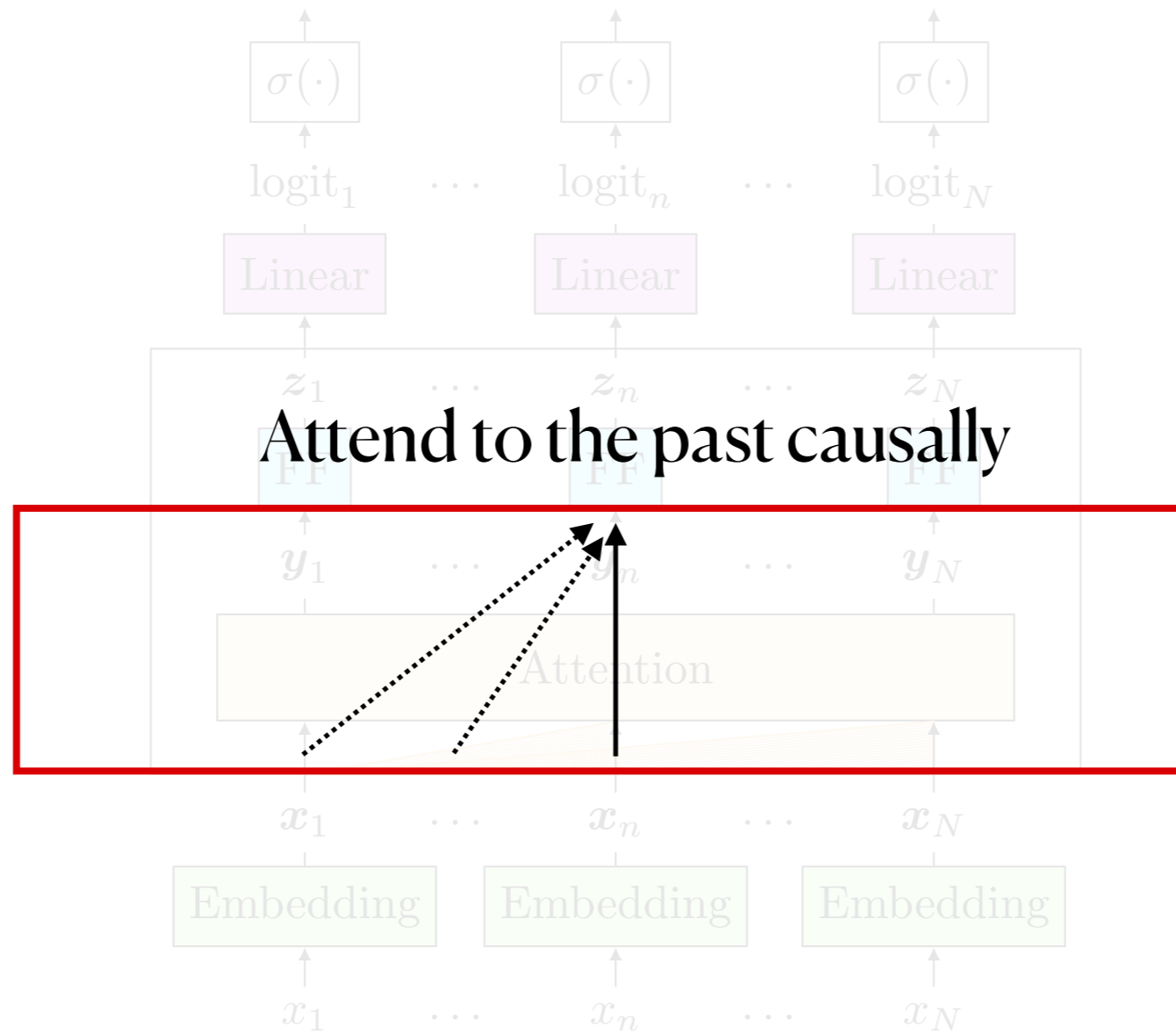


(Vaswani et. al., 2017)

Transformers



Transformers



Single-layer transformer

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

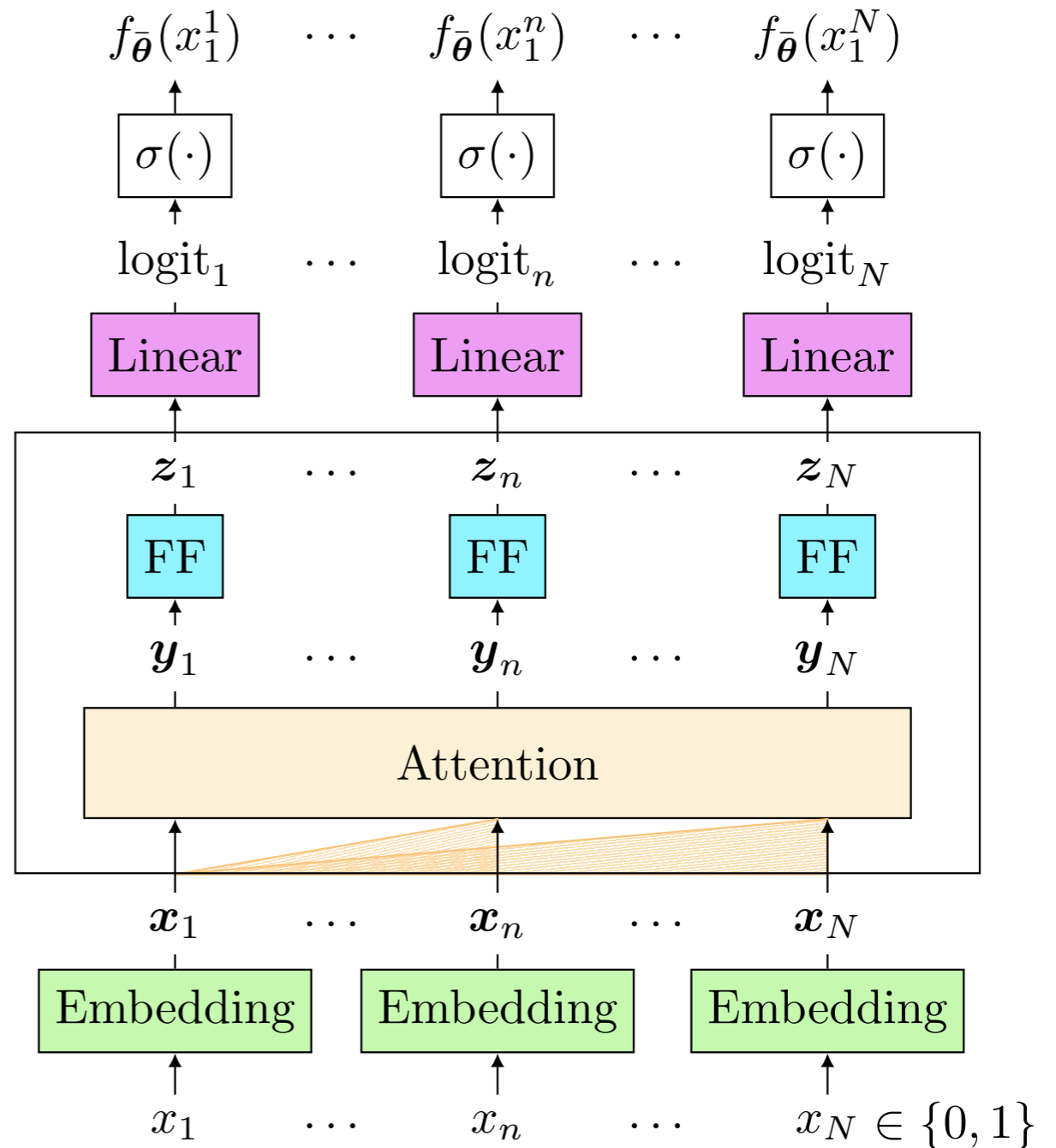
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$

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

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

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

θ



Single-layer transformer

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

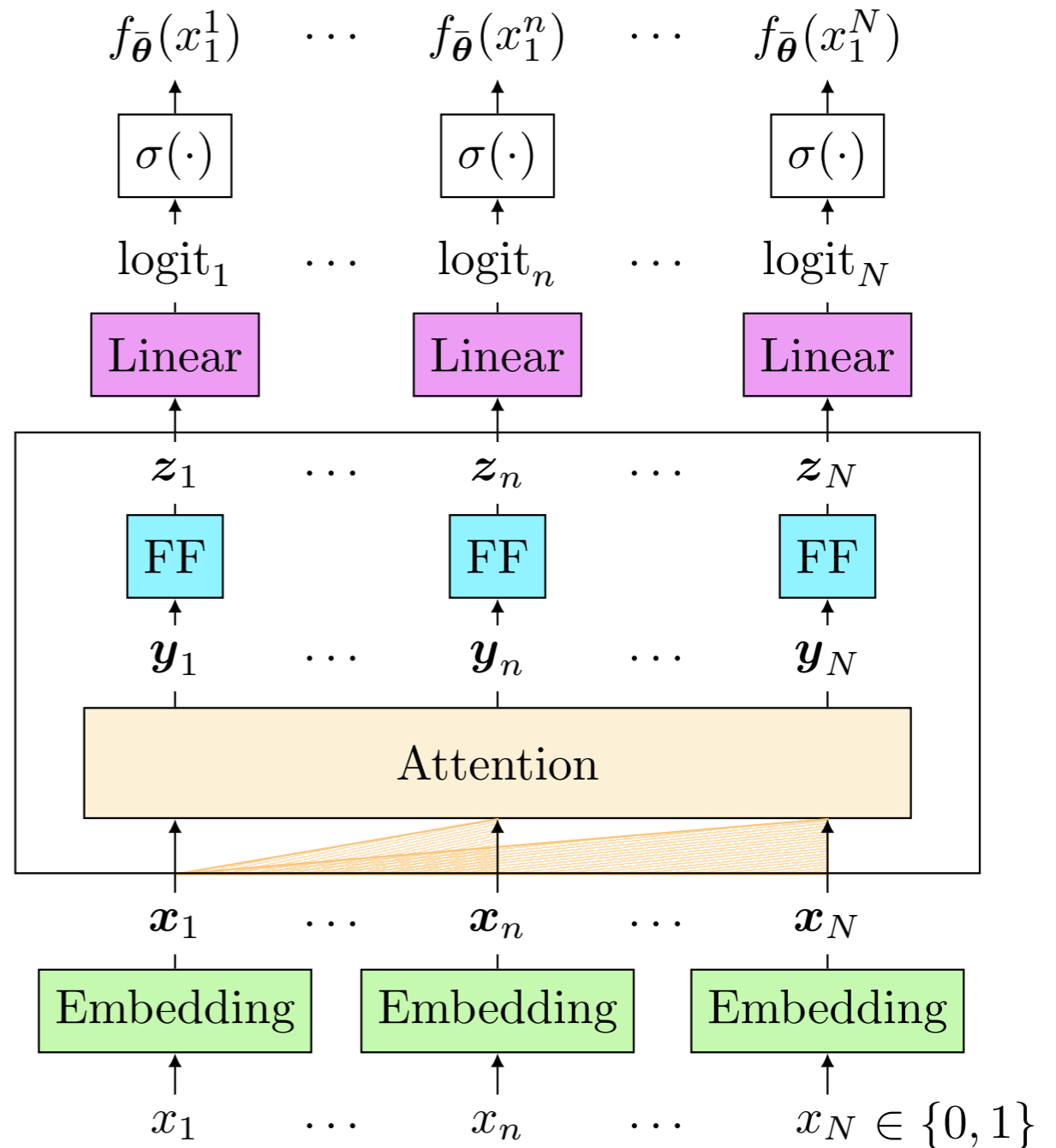
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$

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

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

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

θ



Single-layer transformer

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

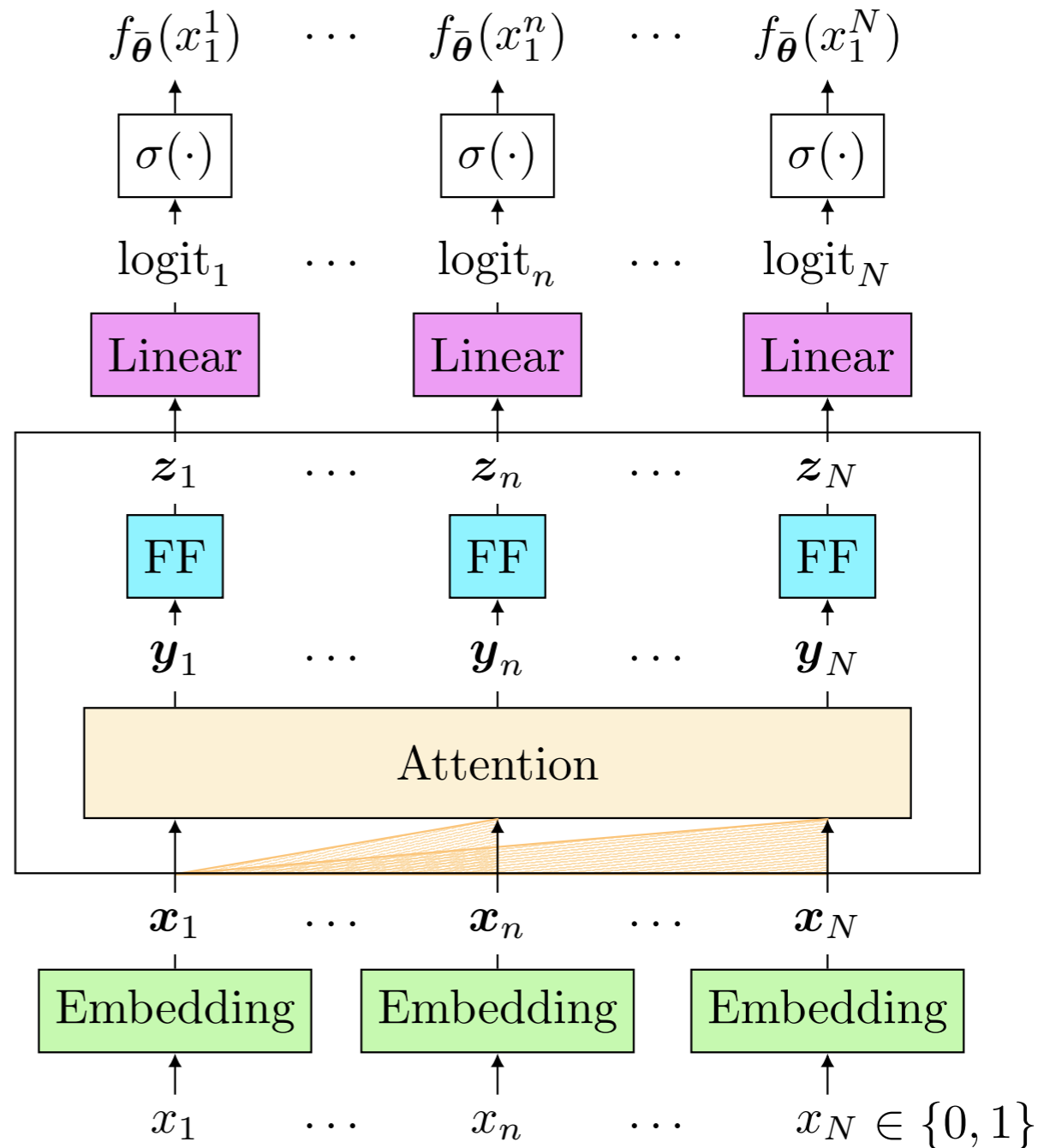
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$

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

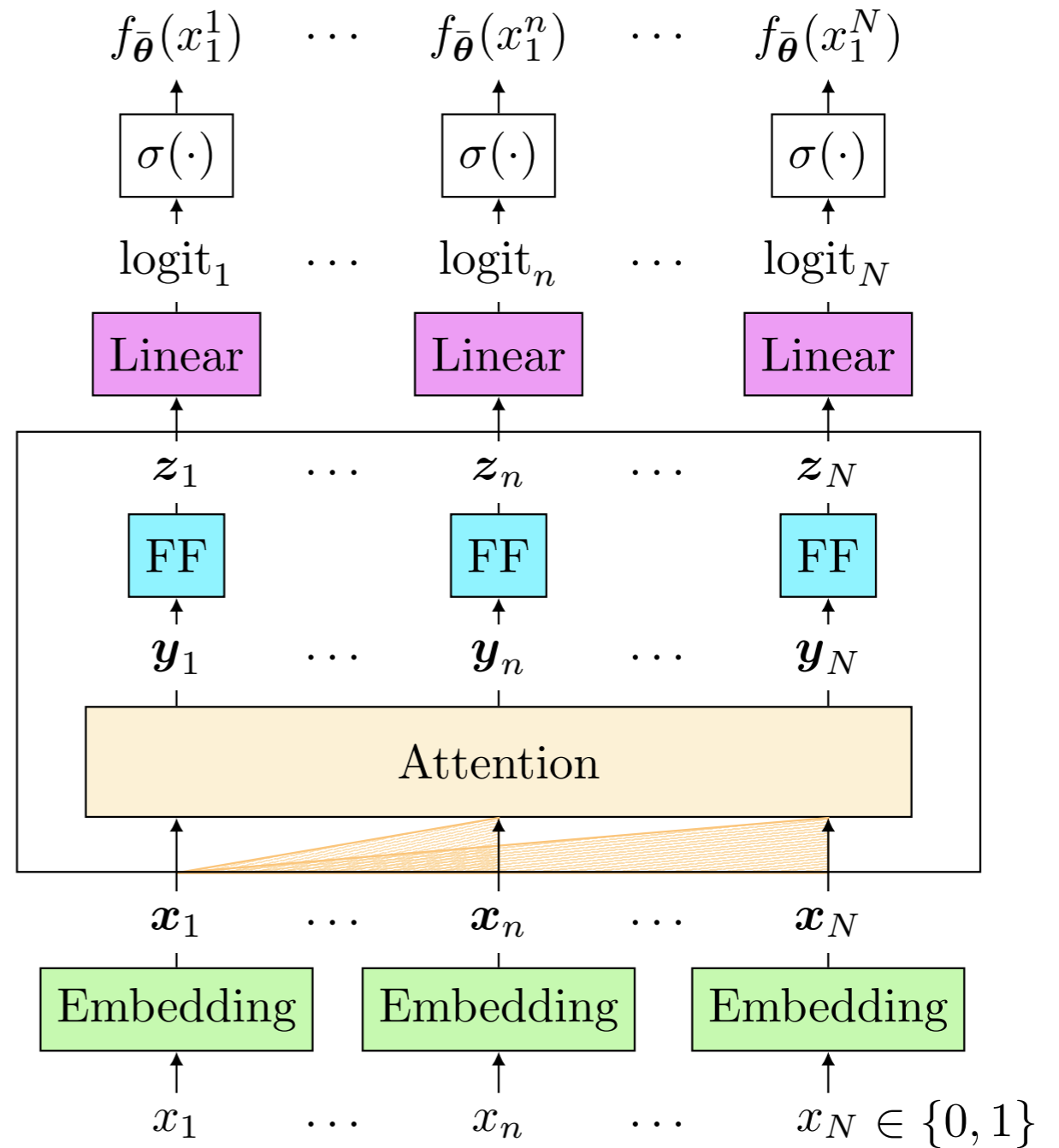
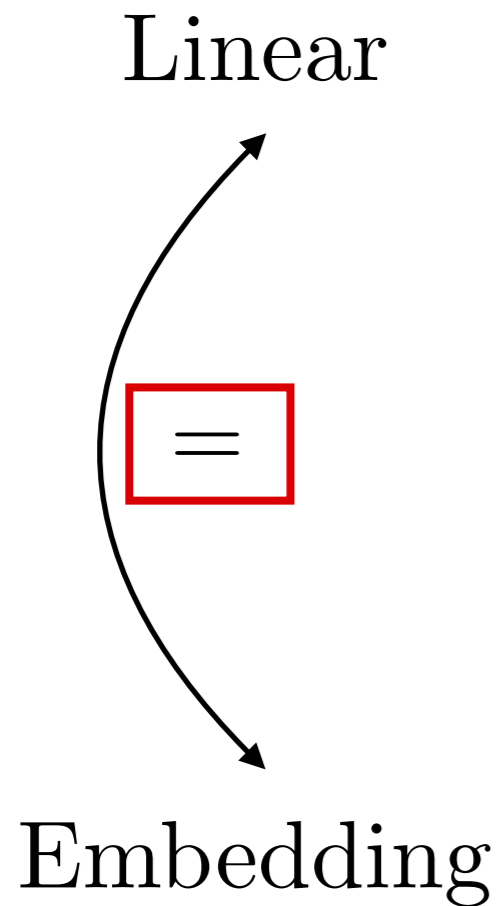
$$\mathbf{y}_n = \mathbf{x}_n + \mathbf{W}_O \sum_{i=1}^n \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i$$

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

θ

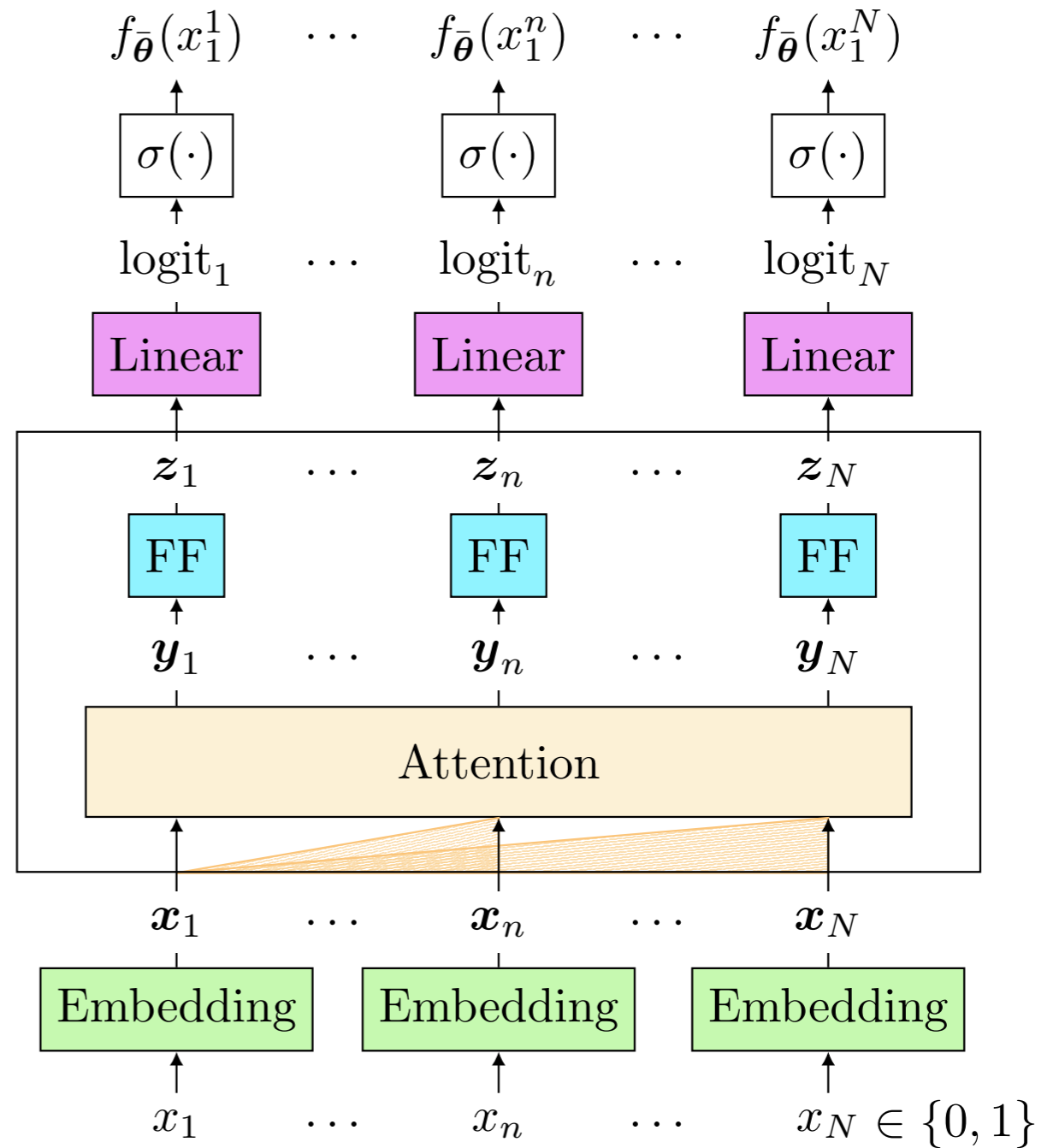


Weight tying



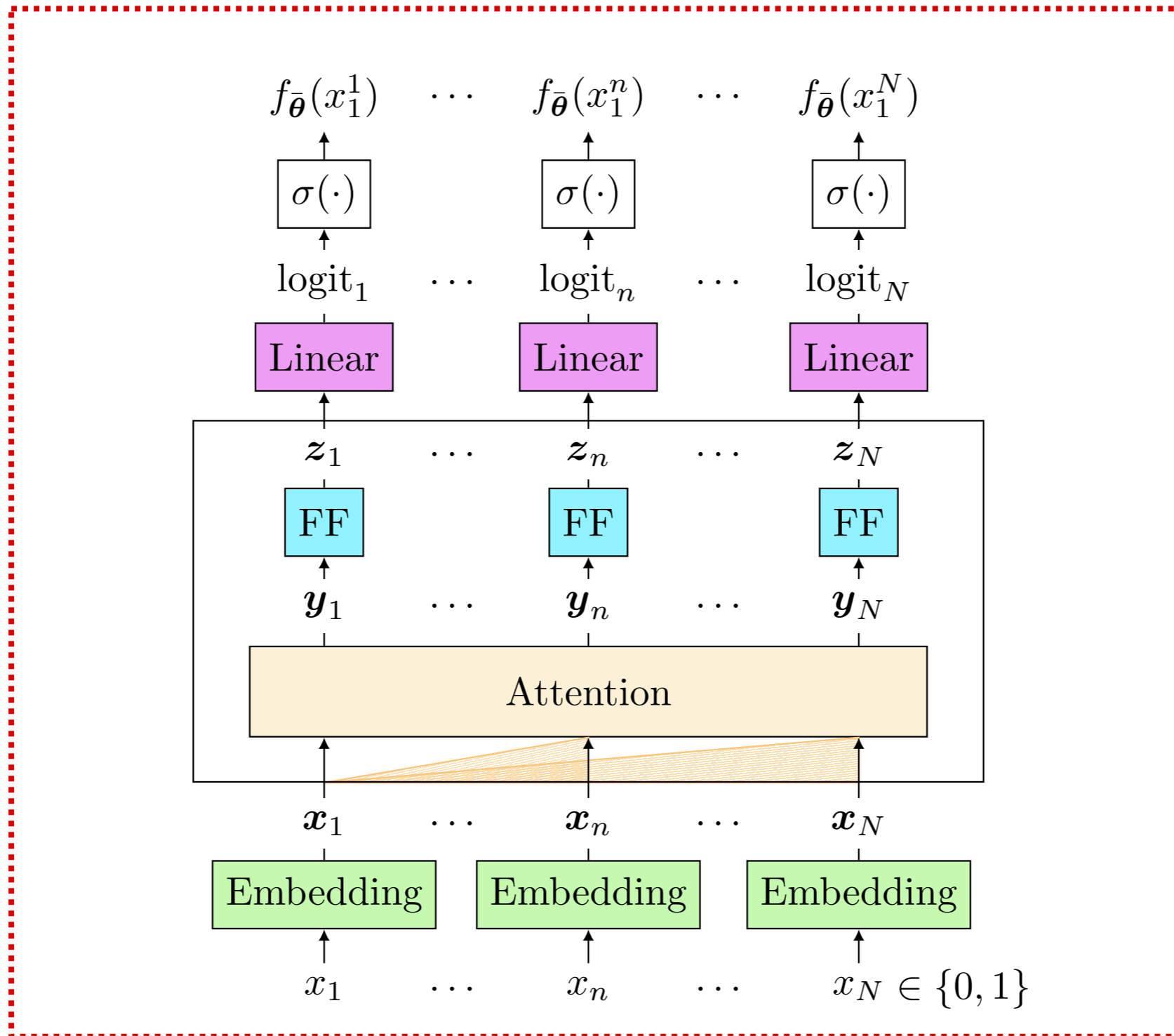
No weight tying

Linear



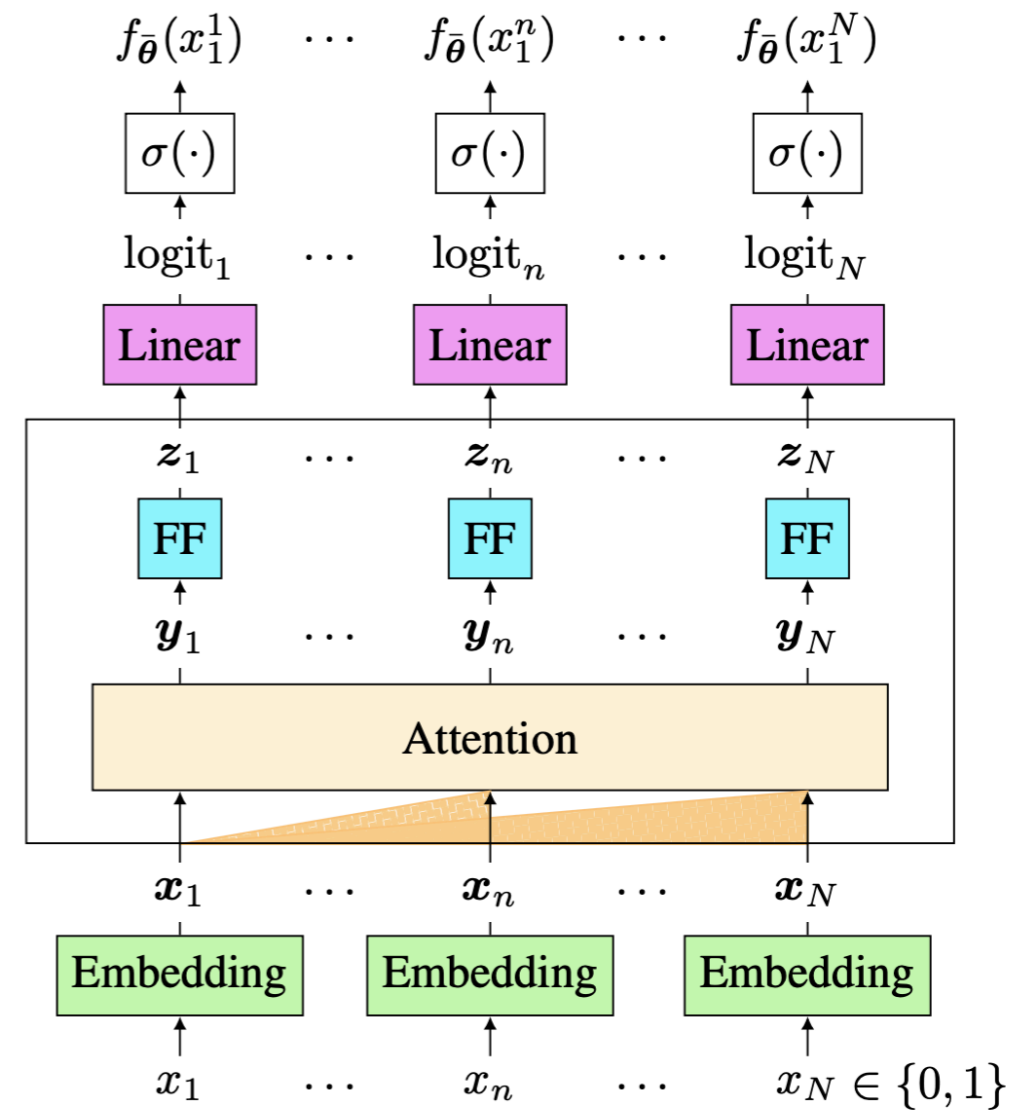
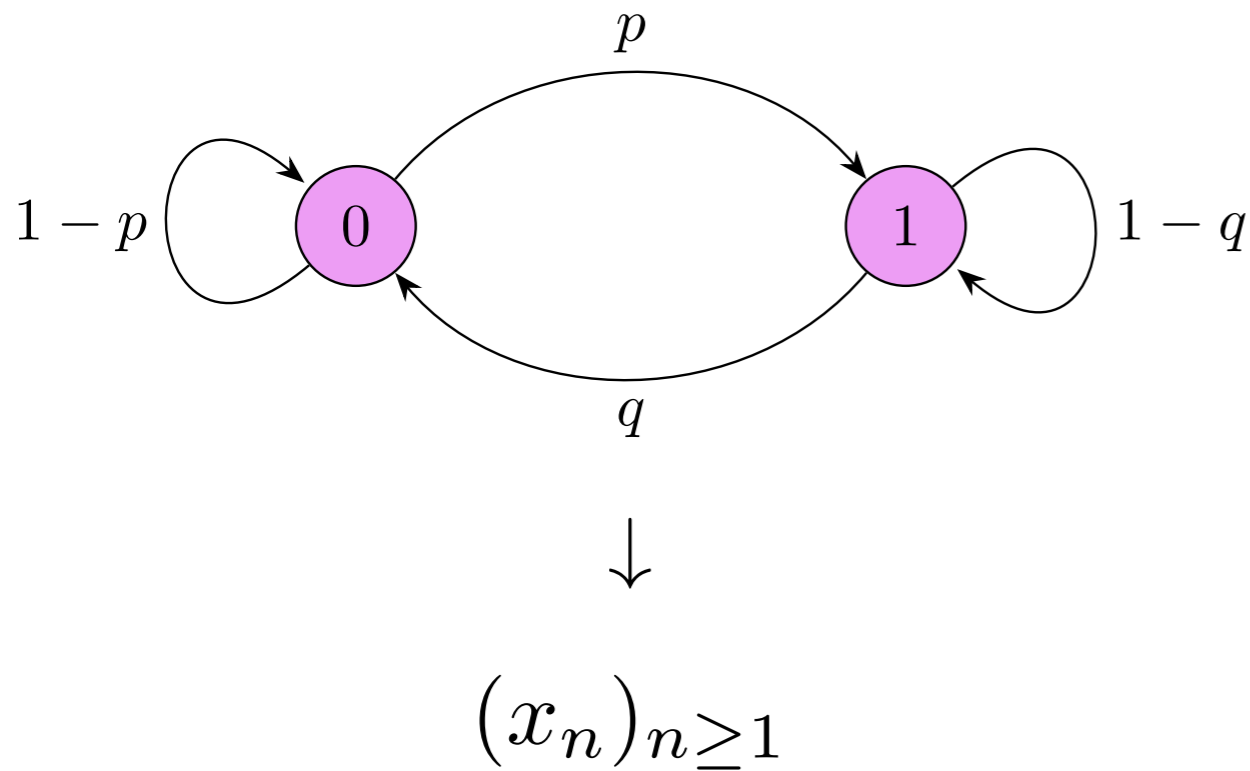
Embedding

Single-layer transformer



θ

First-order Markov chain + Single-layer transformer



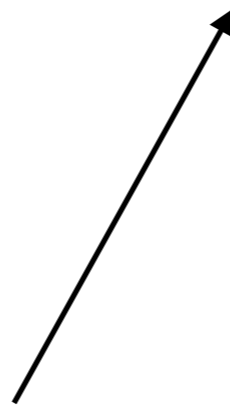
Input data: First-order Markov chain

Model: Depth = 1

Next-token prediction loss

Next-token prediction loss

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

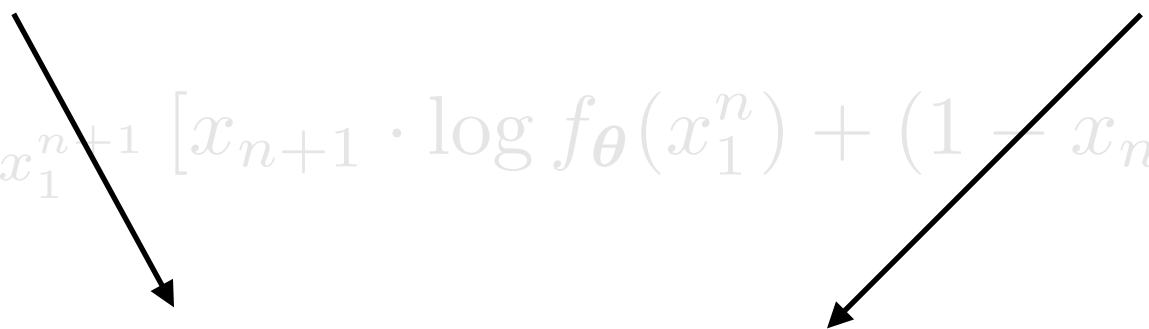


Cross-entropy loss between x_{n+1} and prediction probability $f_{\boldsymbol{\theta}}(x_1^n)$

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}

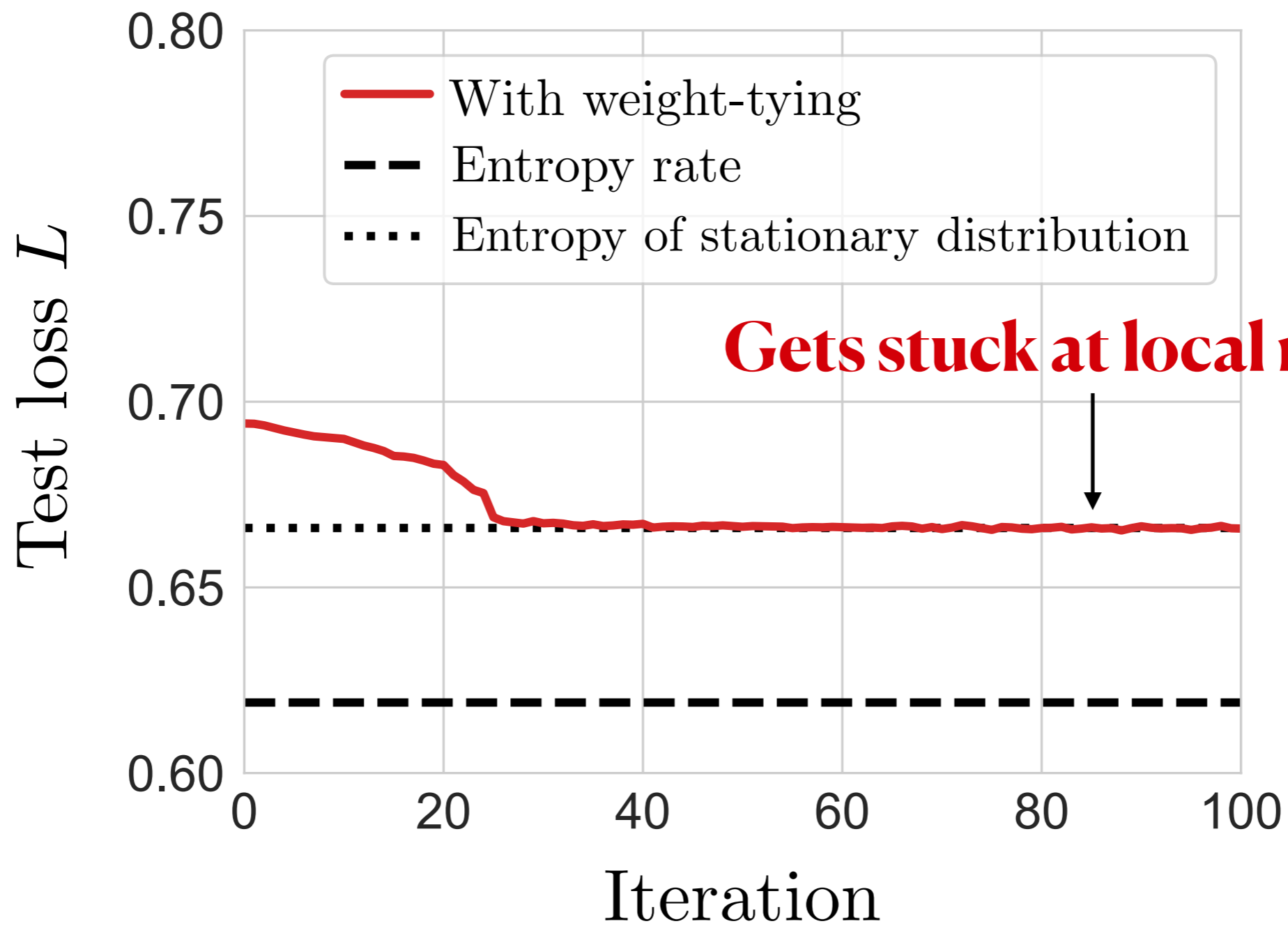
$f_{\theta}(x_1^n)$

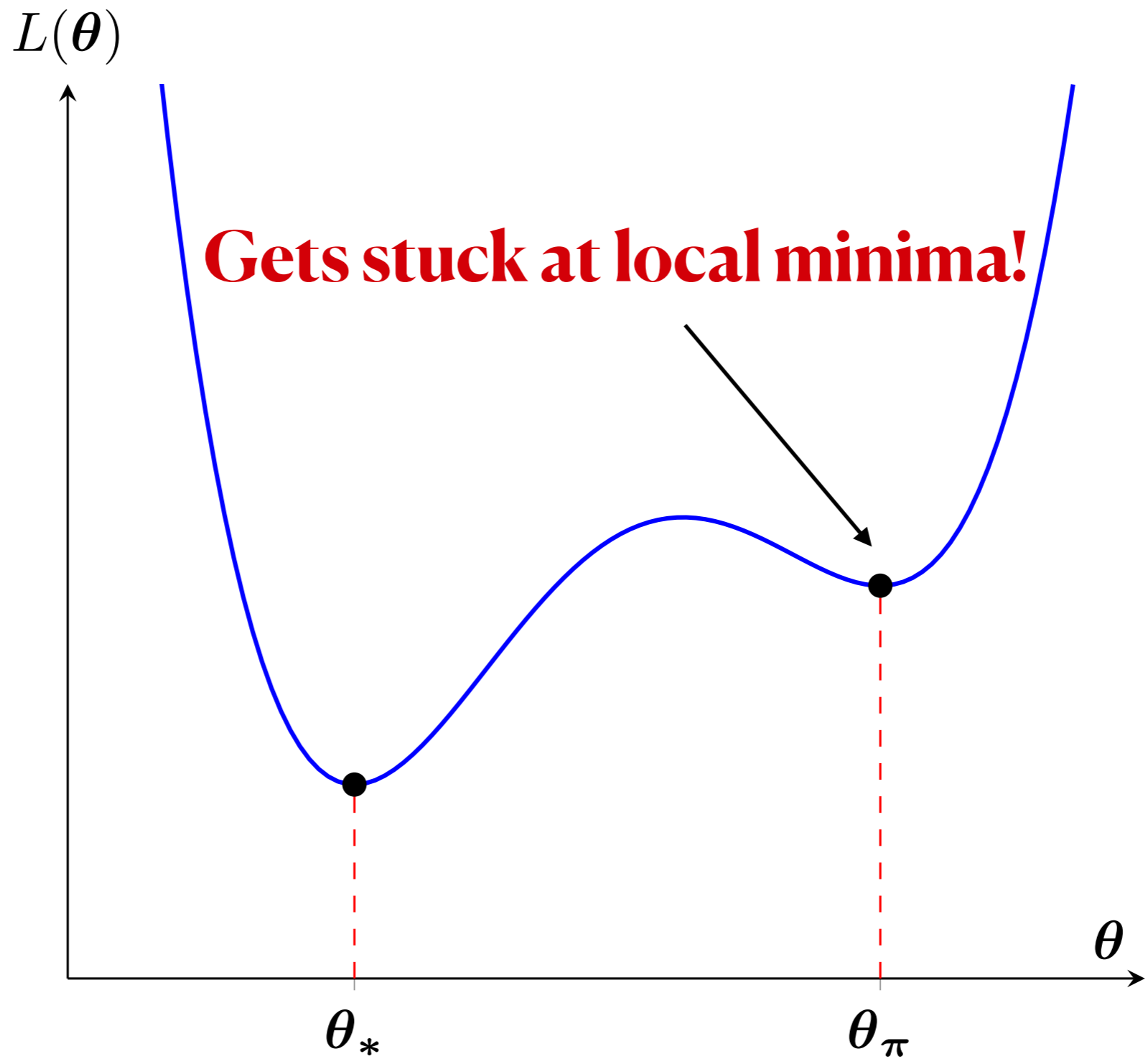
But...

But...

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

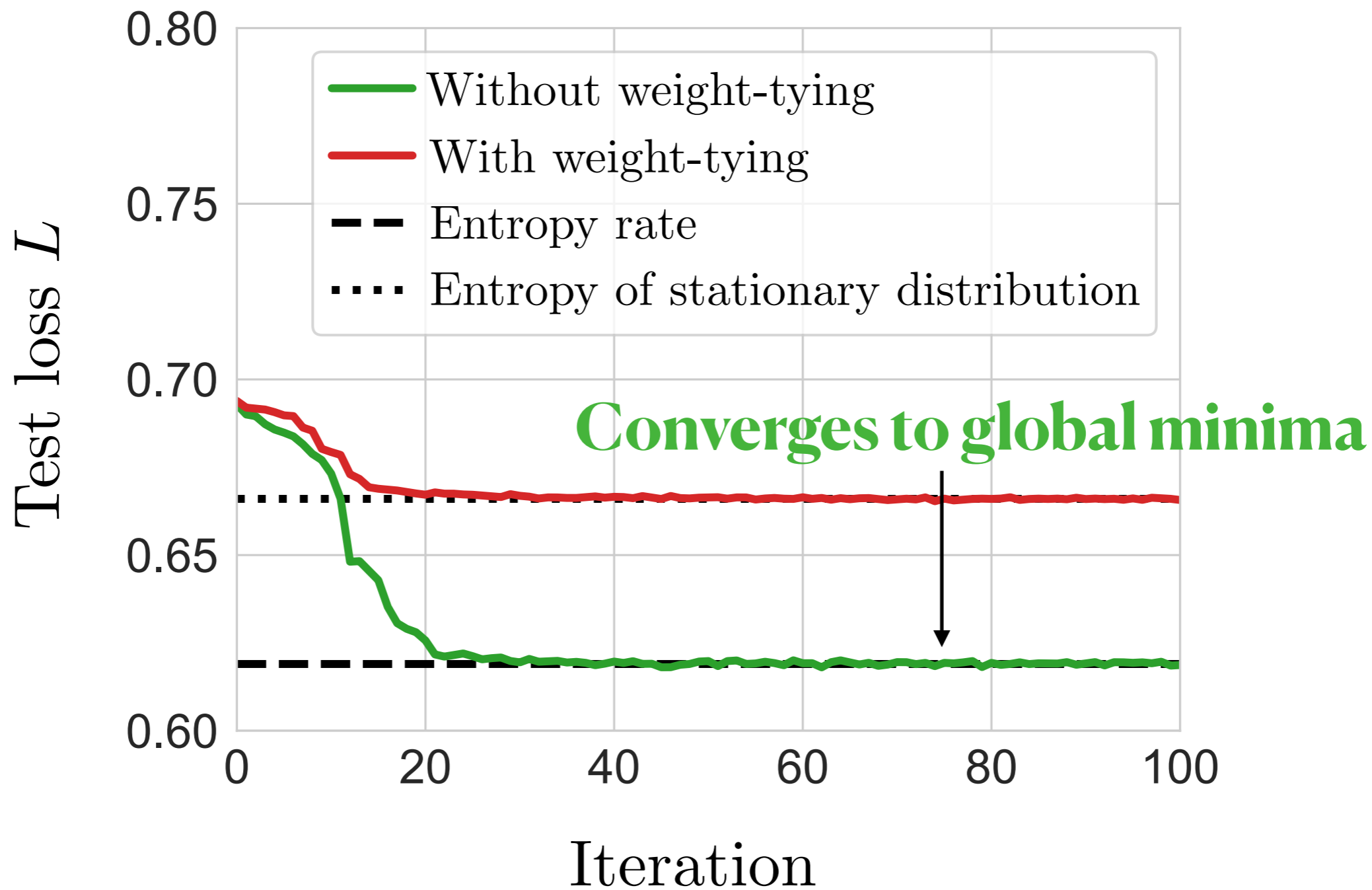
Weight tying

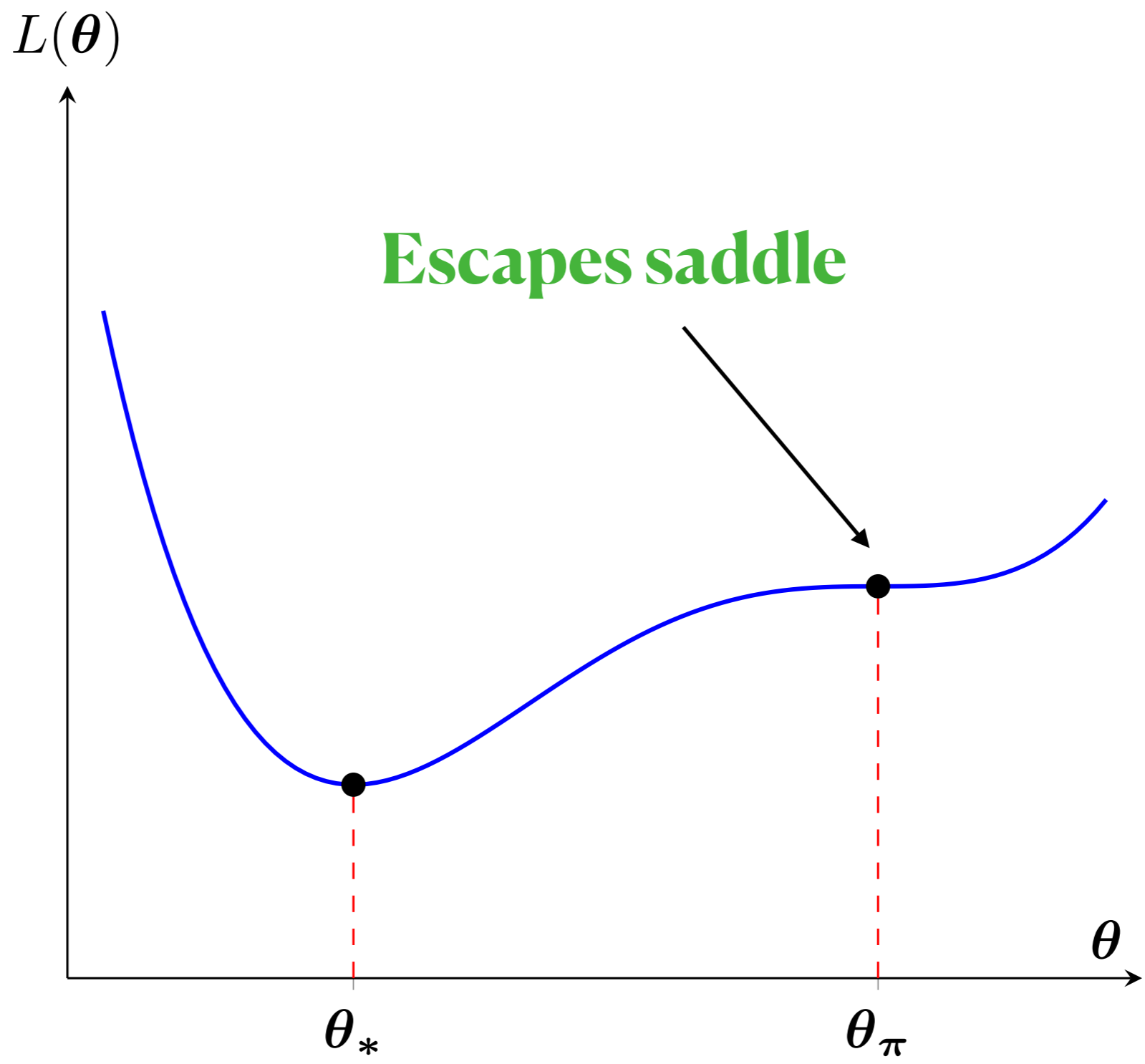




Gets stuck at local minima!

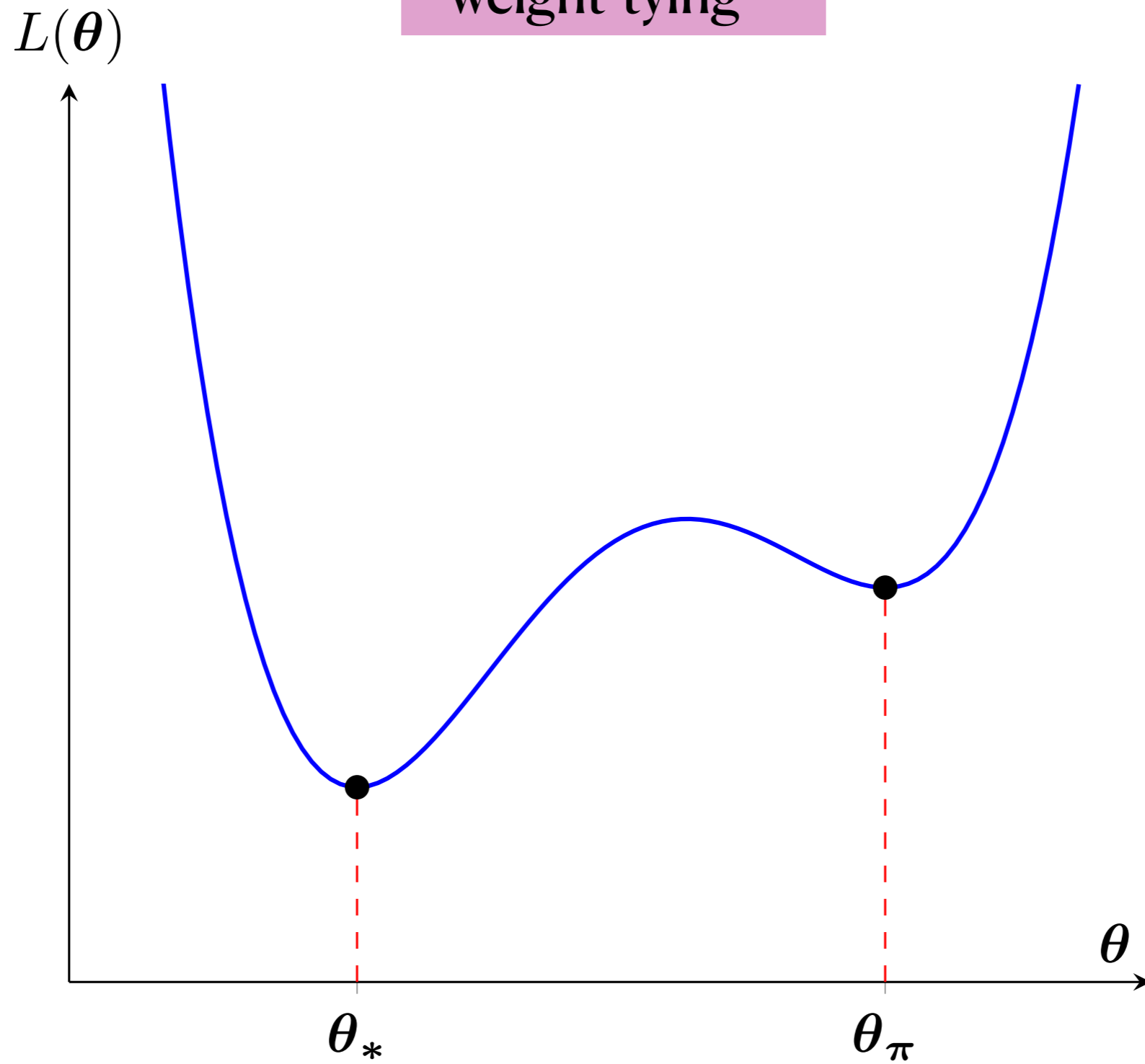
Without weight tying





What's happening

weight-tying



weight-tying

$L(\theta)$

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

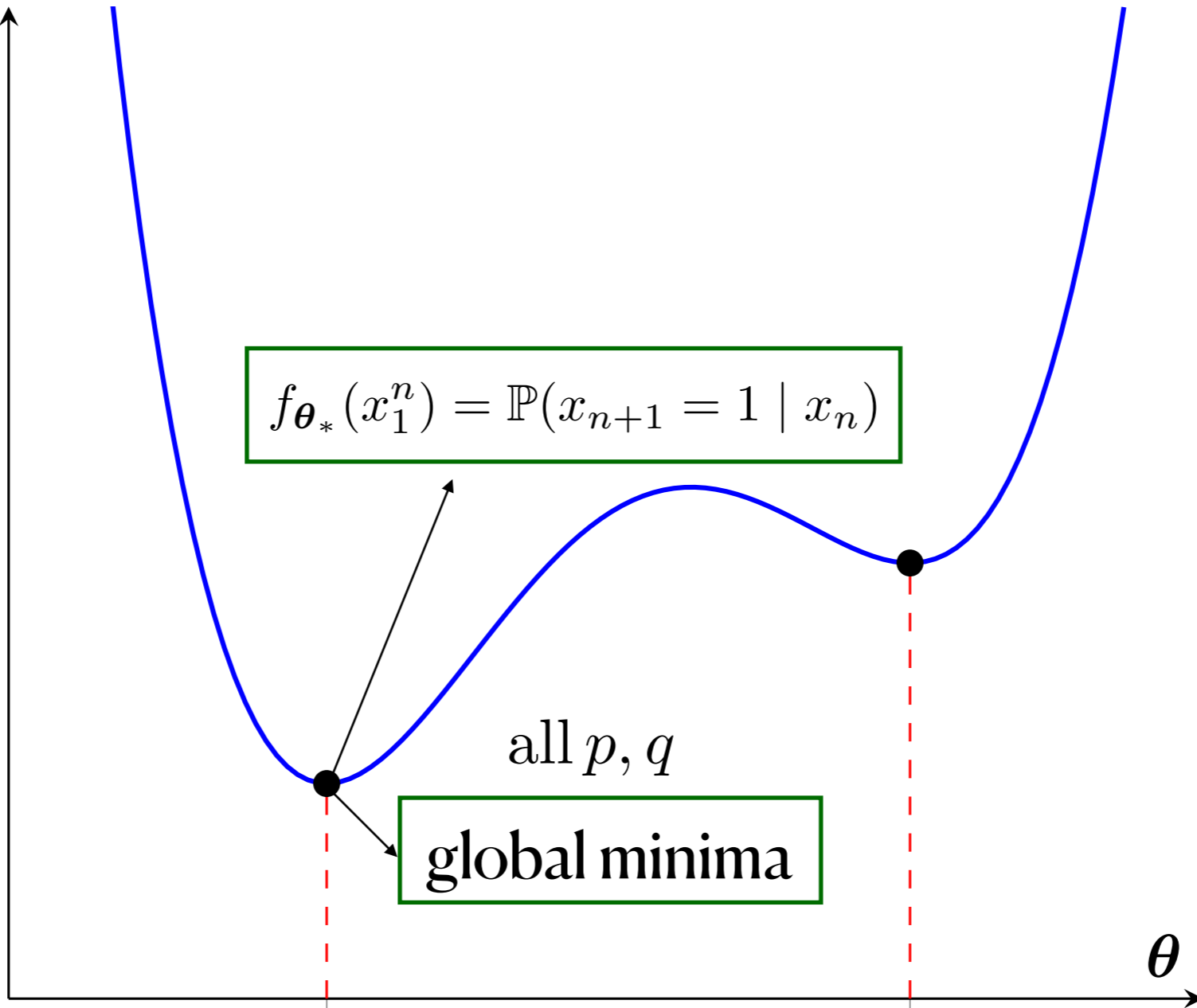
all p, q

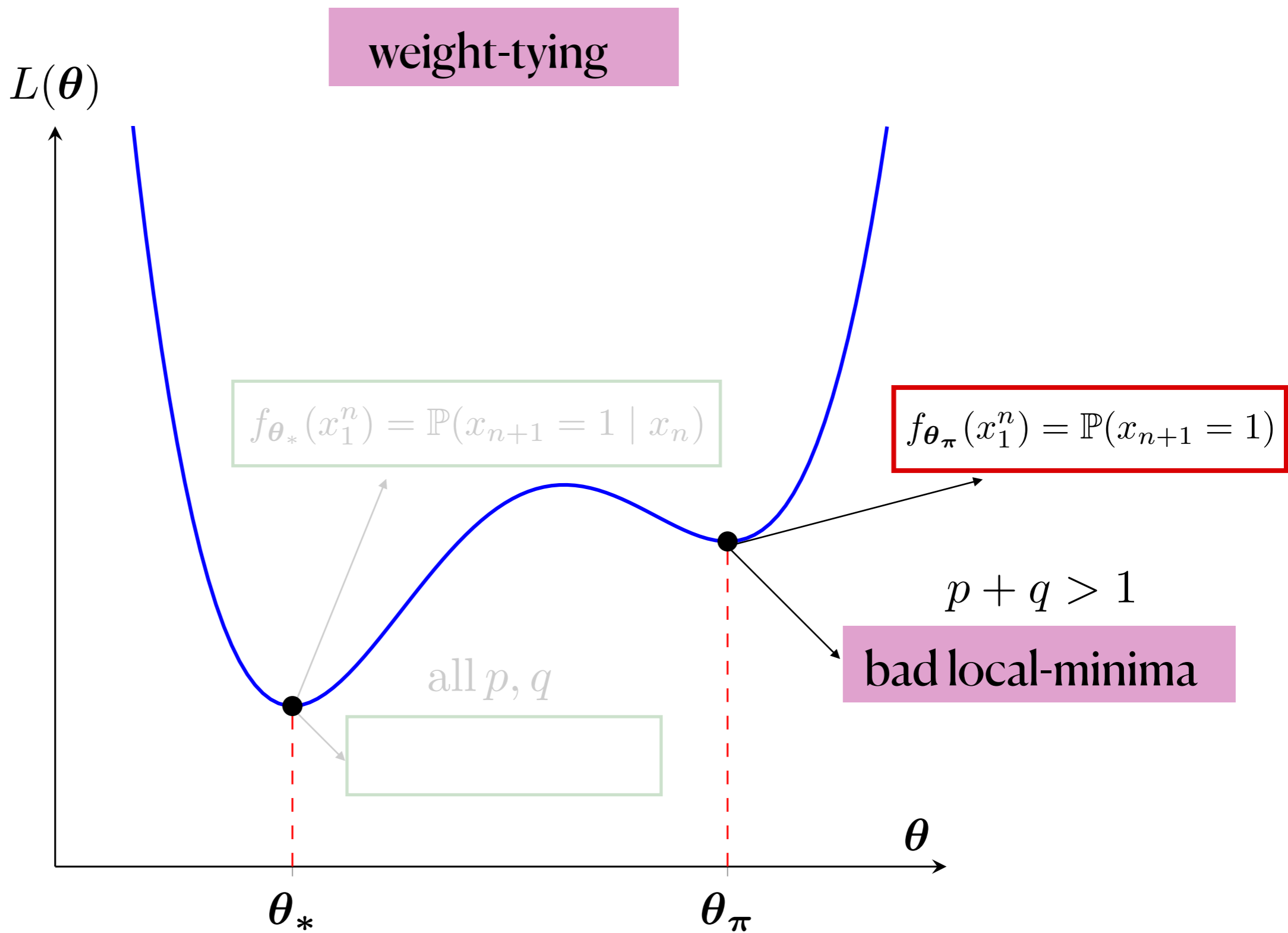
global minima

θ_*

θ_π

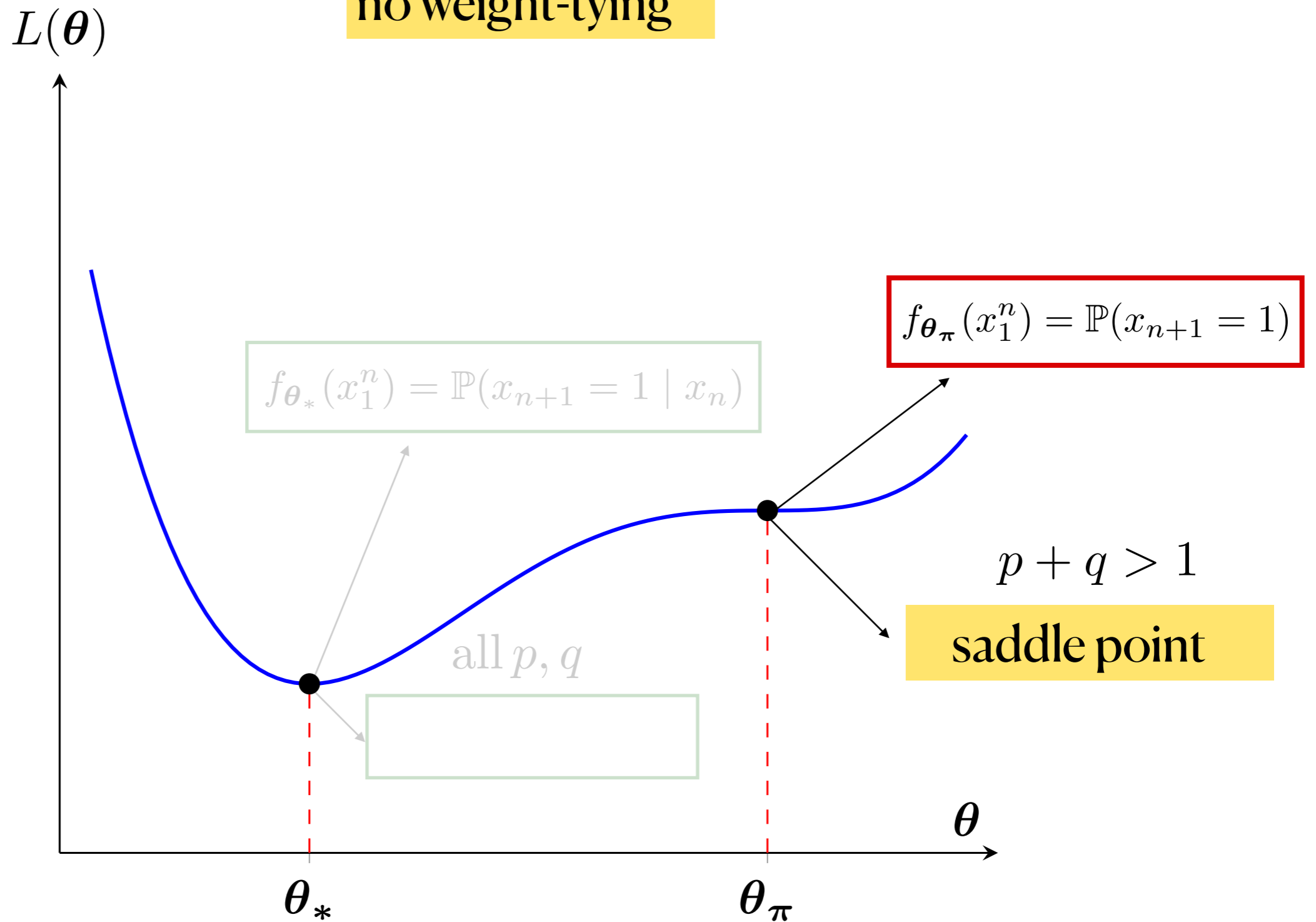
θ





Interestingly...

no weight-tying



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 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)

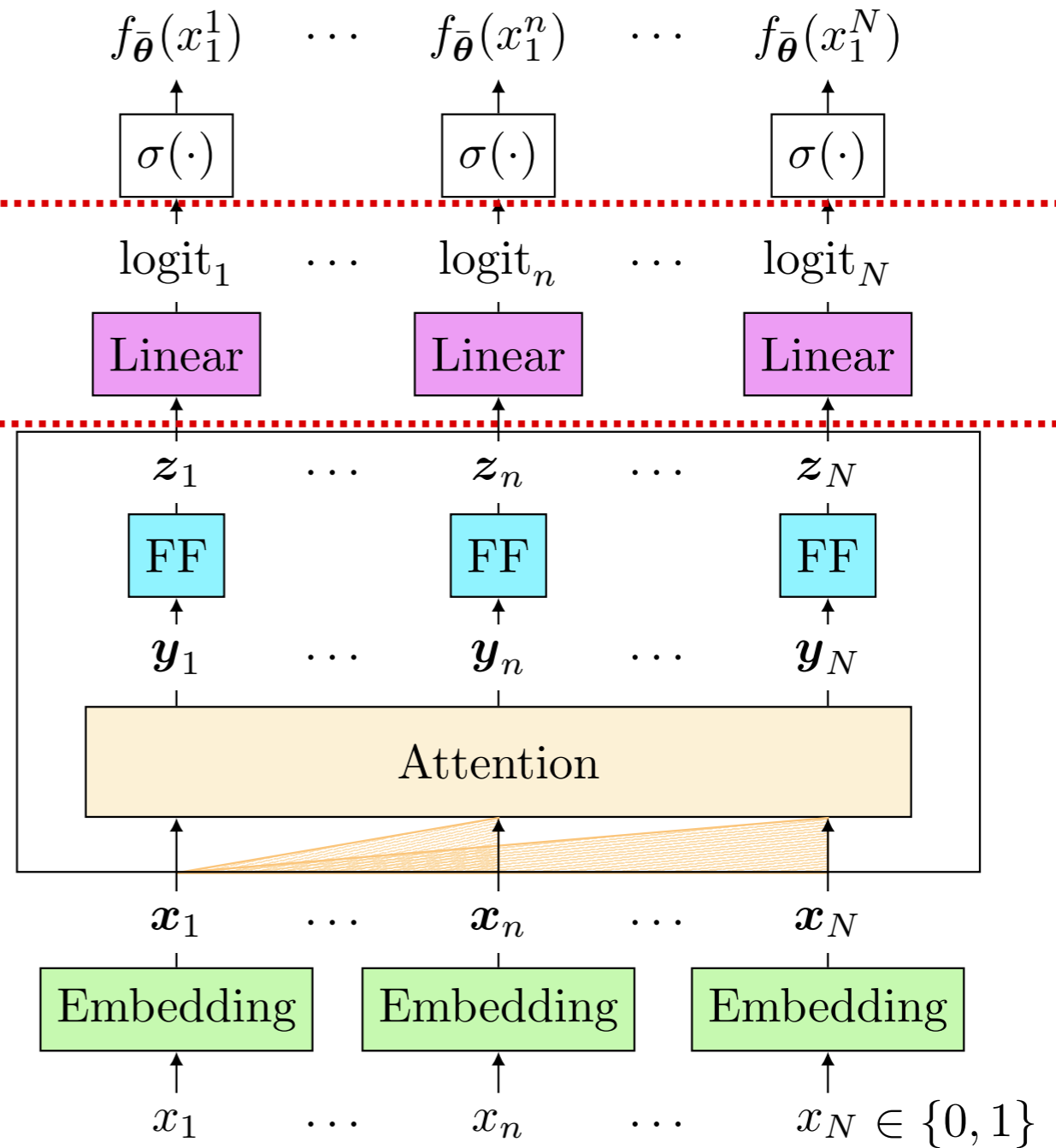
Under the same setting as above with the weights not tied, θ_π becomes a saddle point. It satisfies the same properties.

Intuition

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

$$\uparrow a = 0$$

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



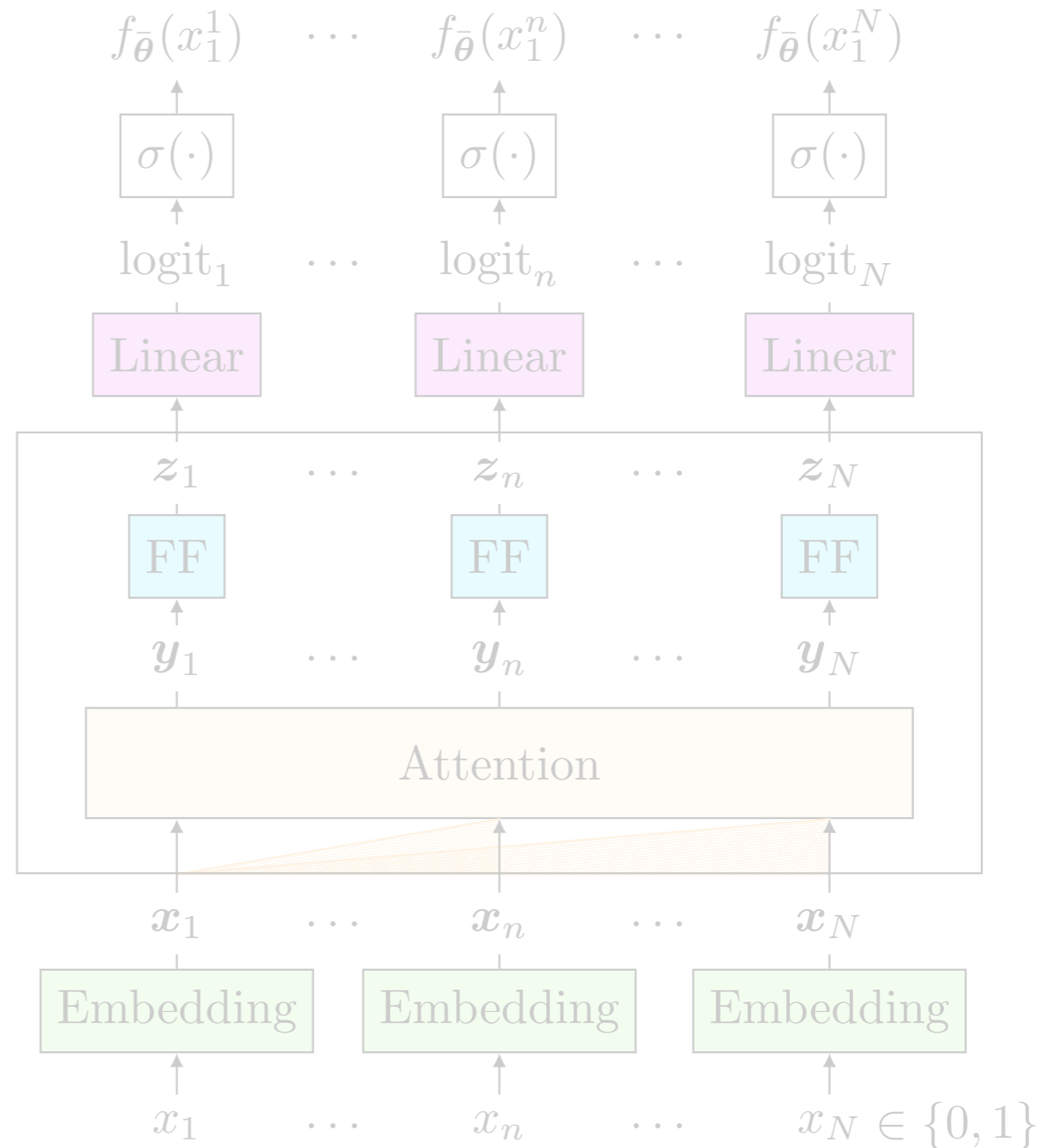
Intuition

weight-tying

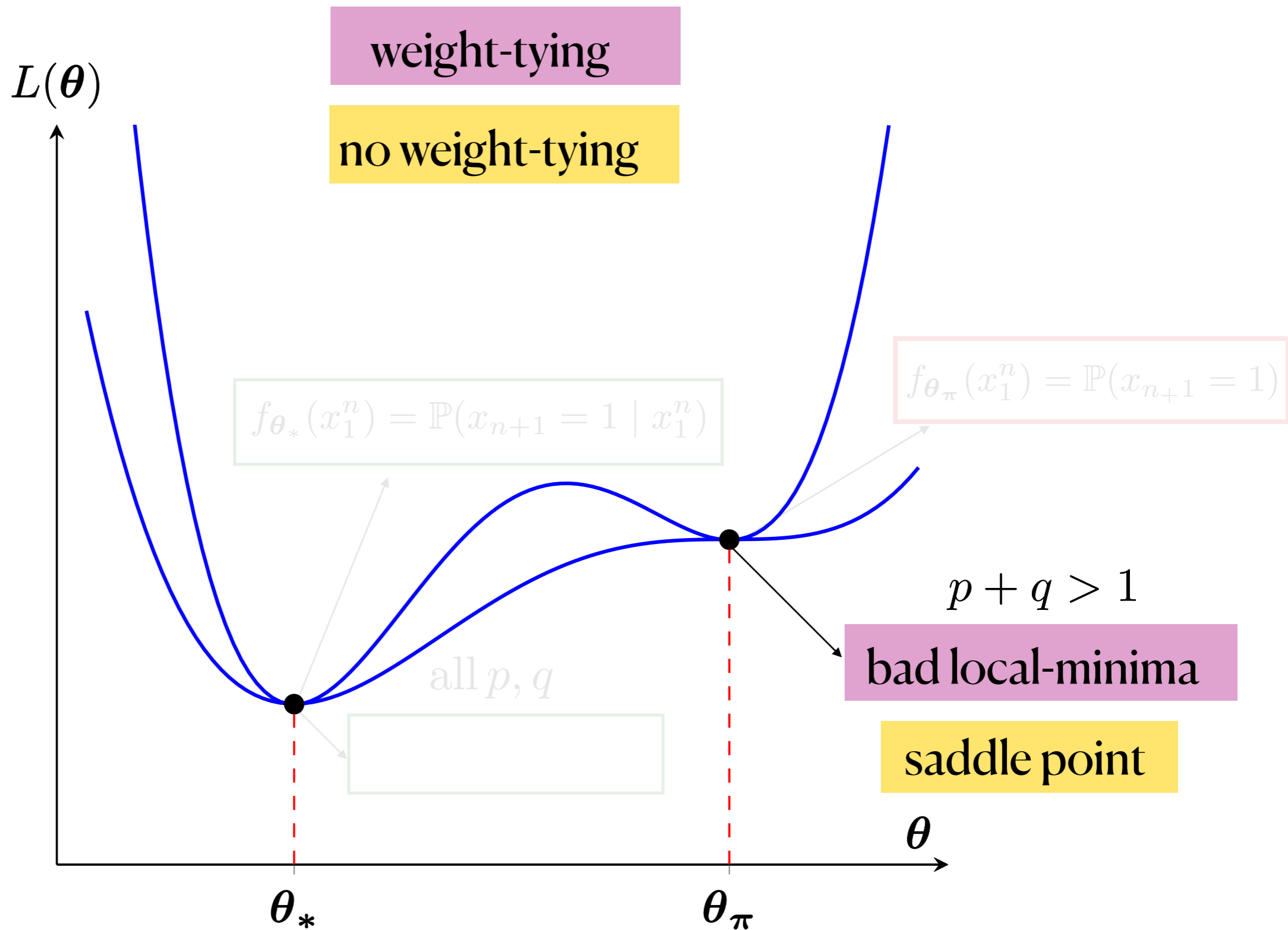
Hessian is (almost) positive-definite

no weight-tying

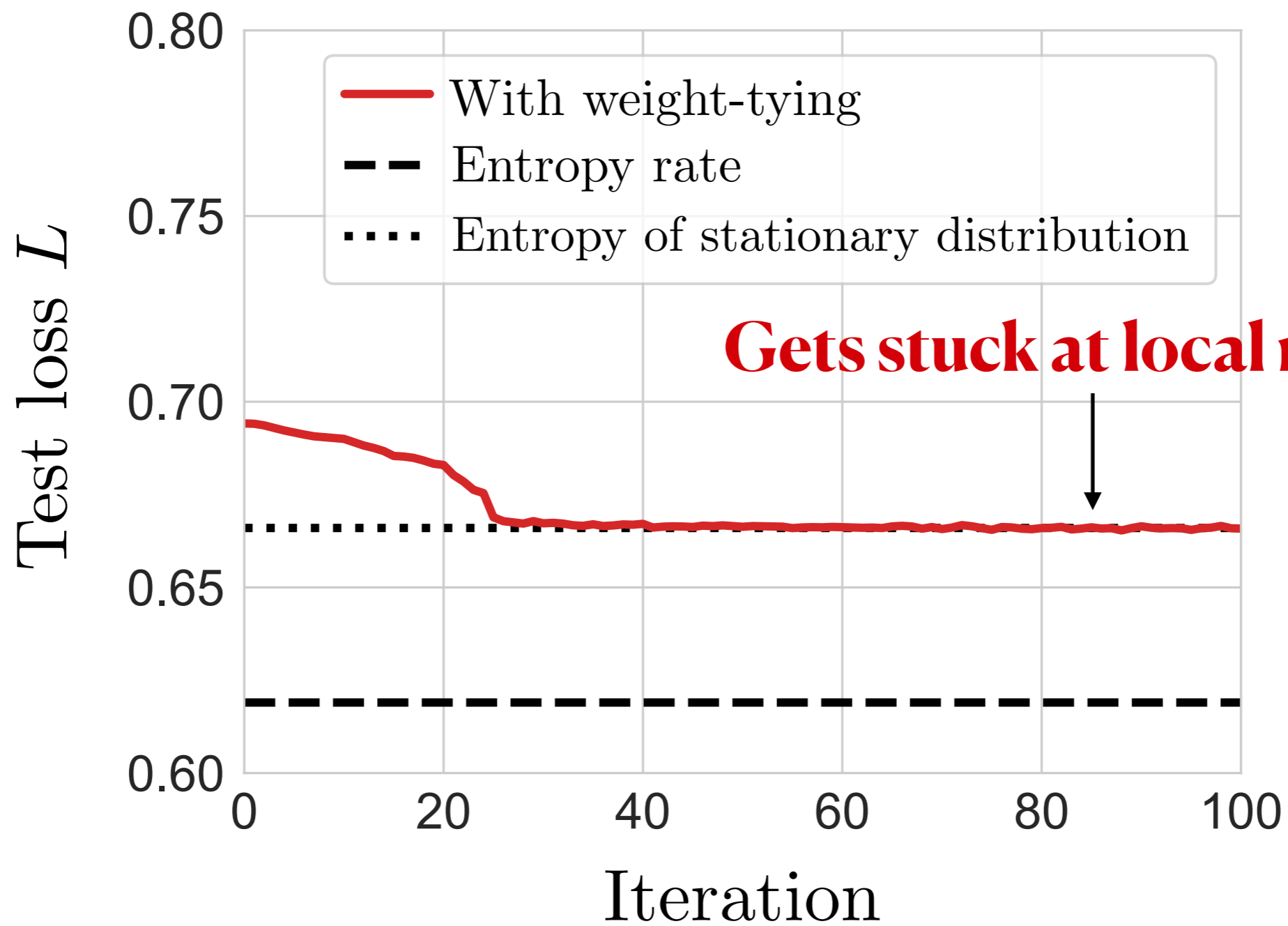
Hessian is indefinite



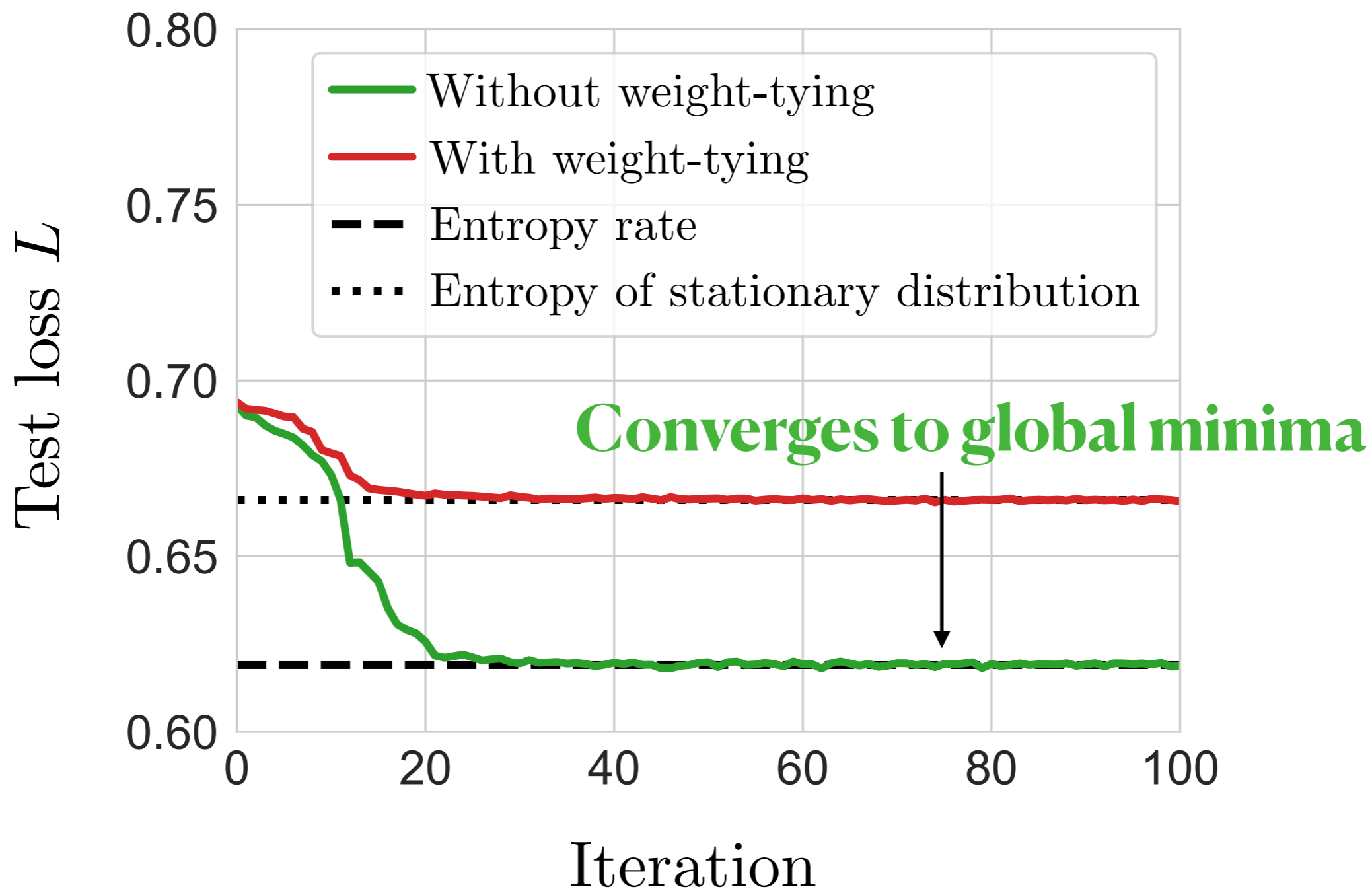
Main results



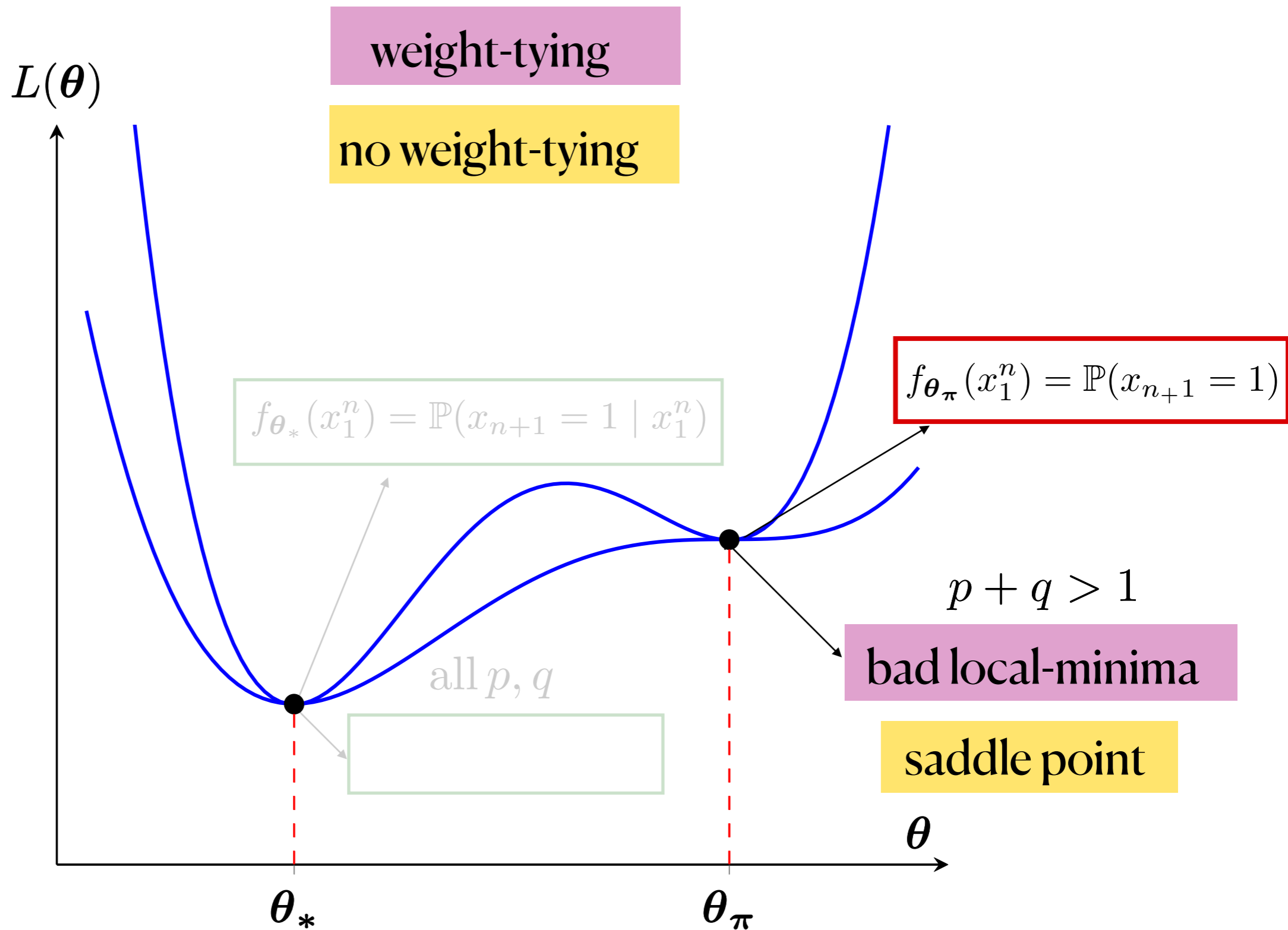
Weight tying



Without weight tying



Main results





Markovian inputs



Transformers



Memory = 1

Depth = 1



What do they learn?





Markovian inputs

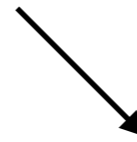


Transformers



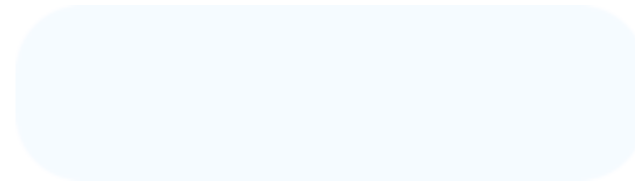
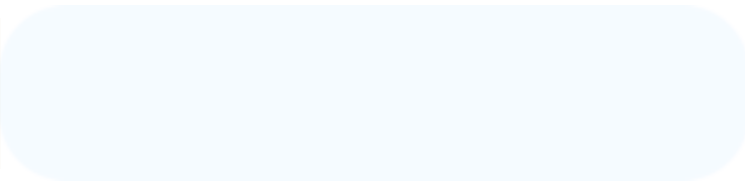
Memory = 1

Depth = 1



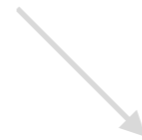
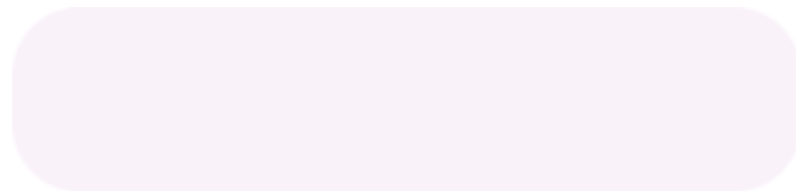
How do they learn?





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

Recall

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

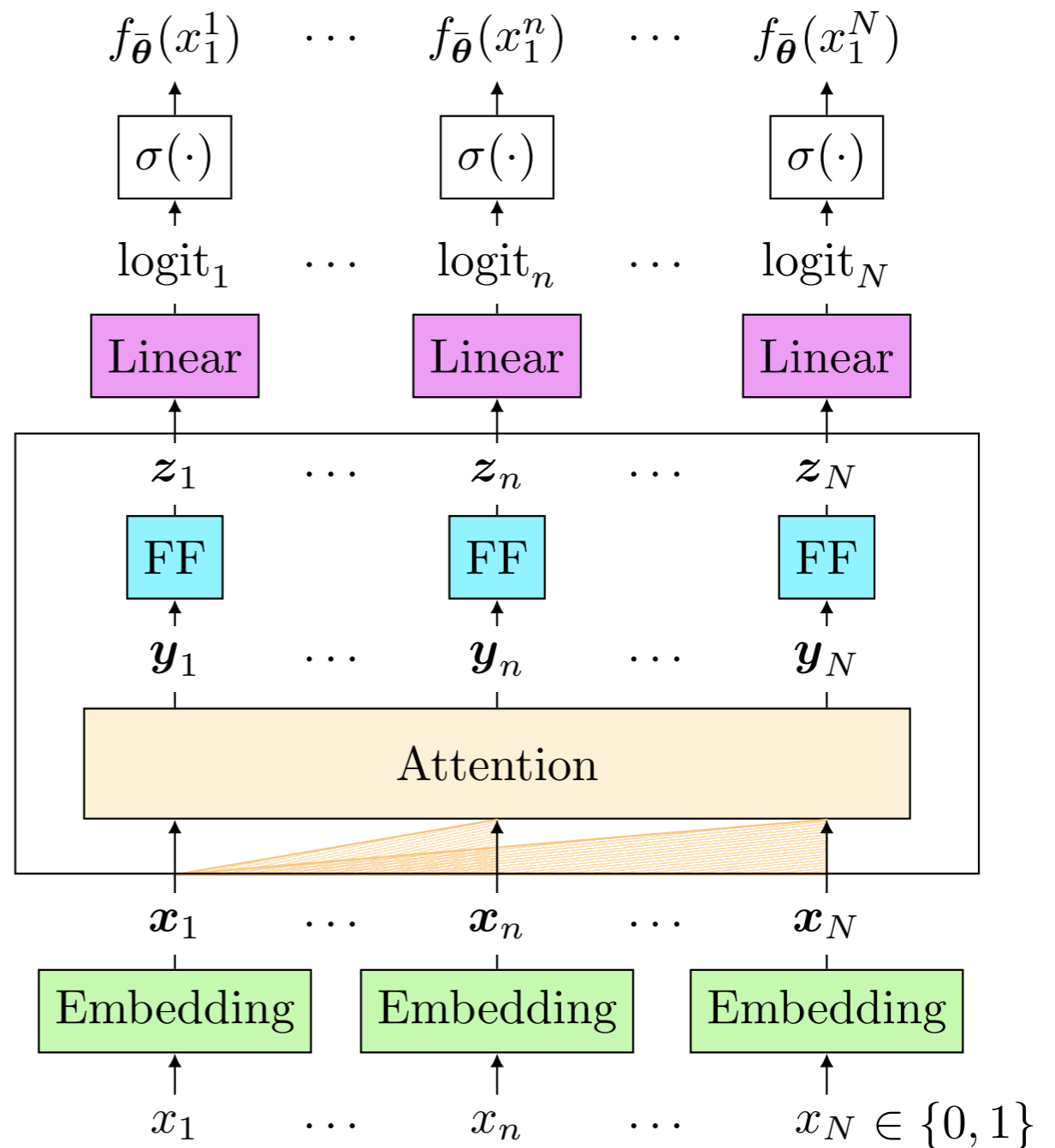
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}$$

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

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

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

θ



Low-rank structure

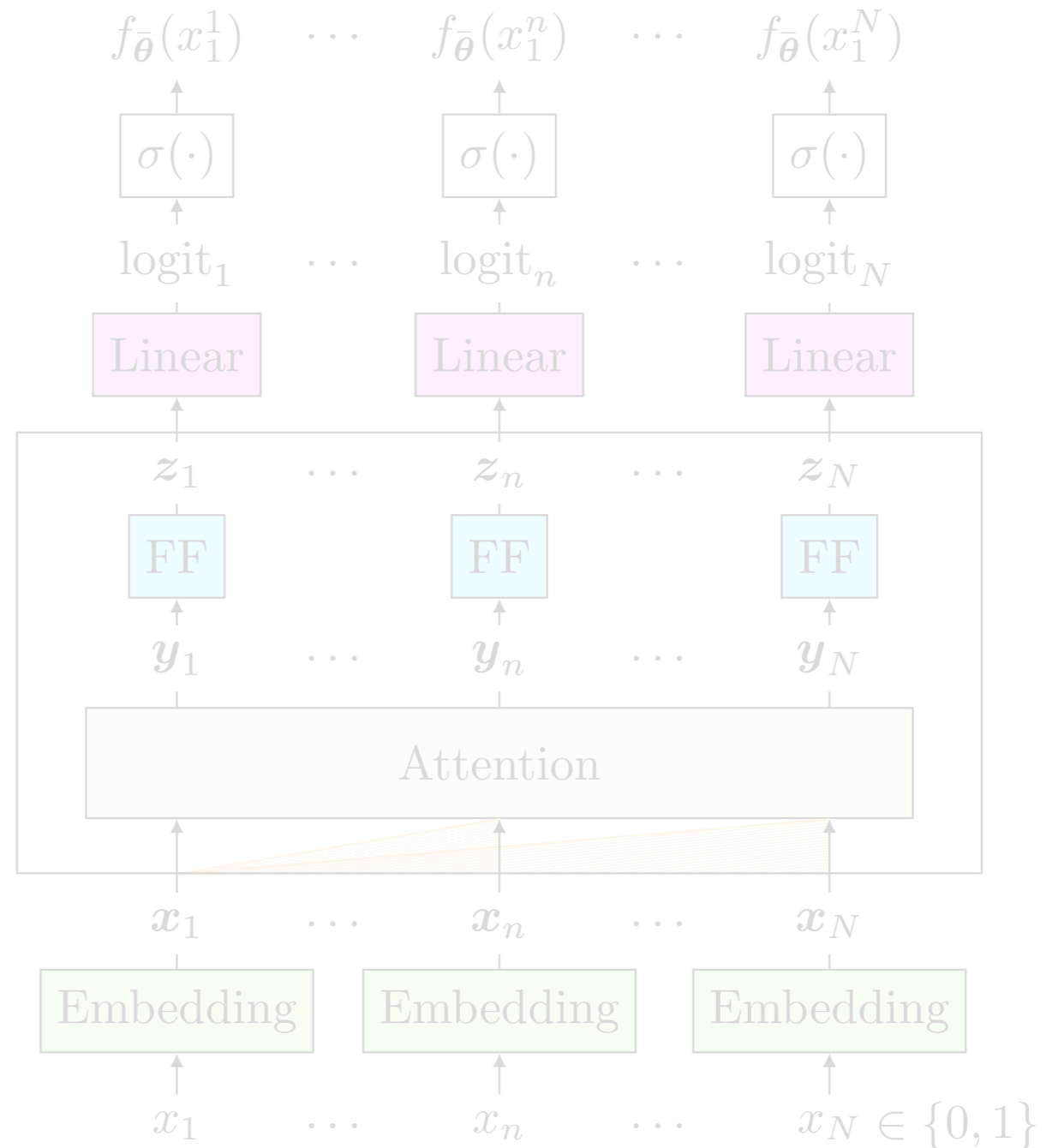
$$f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n)$$

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

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

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

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



Reparametrization

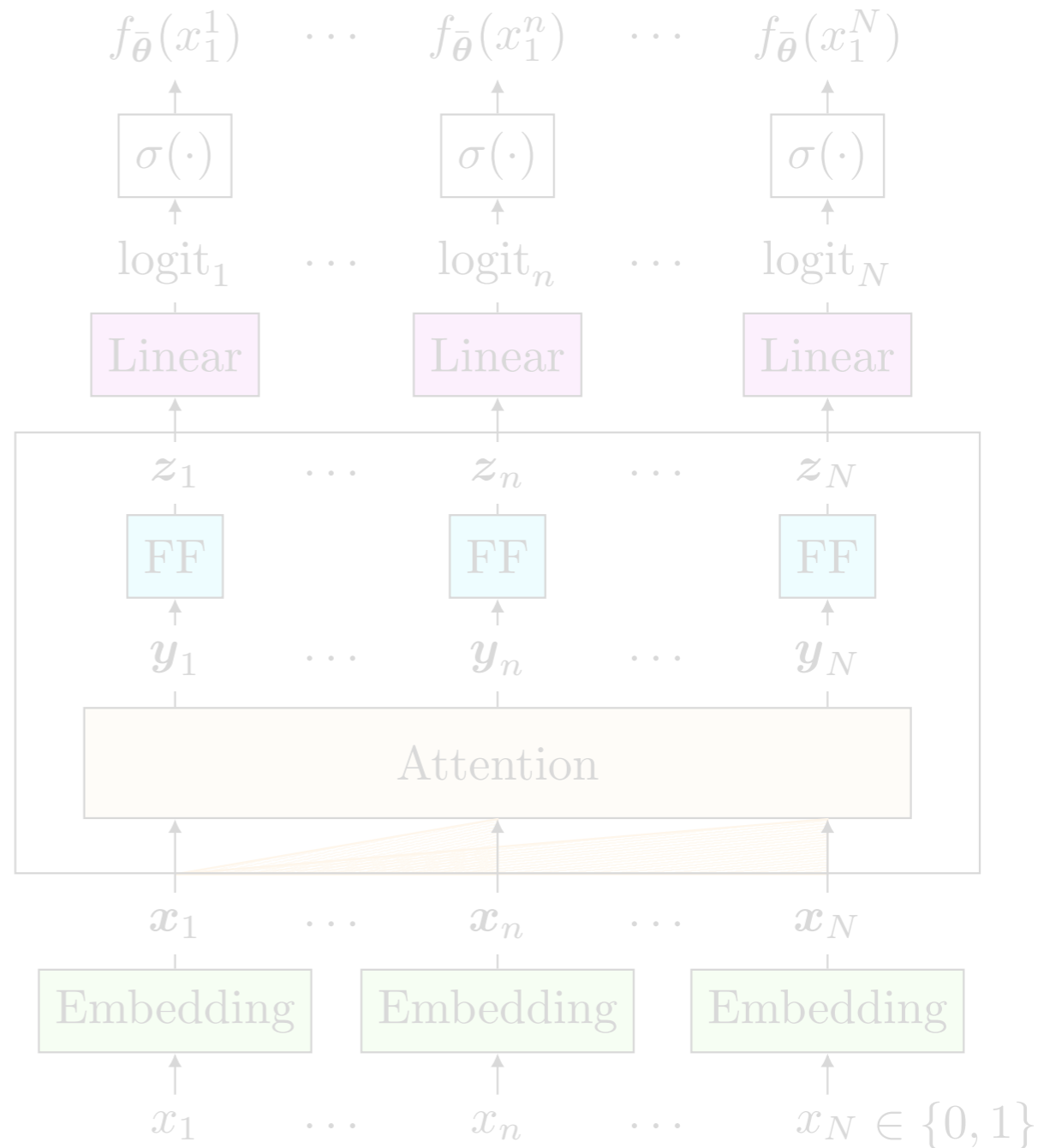
$$f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n)$$

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

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

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

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



Reparametrization

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

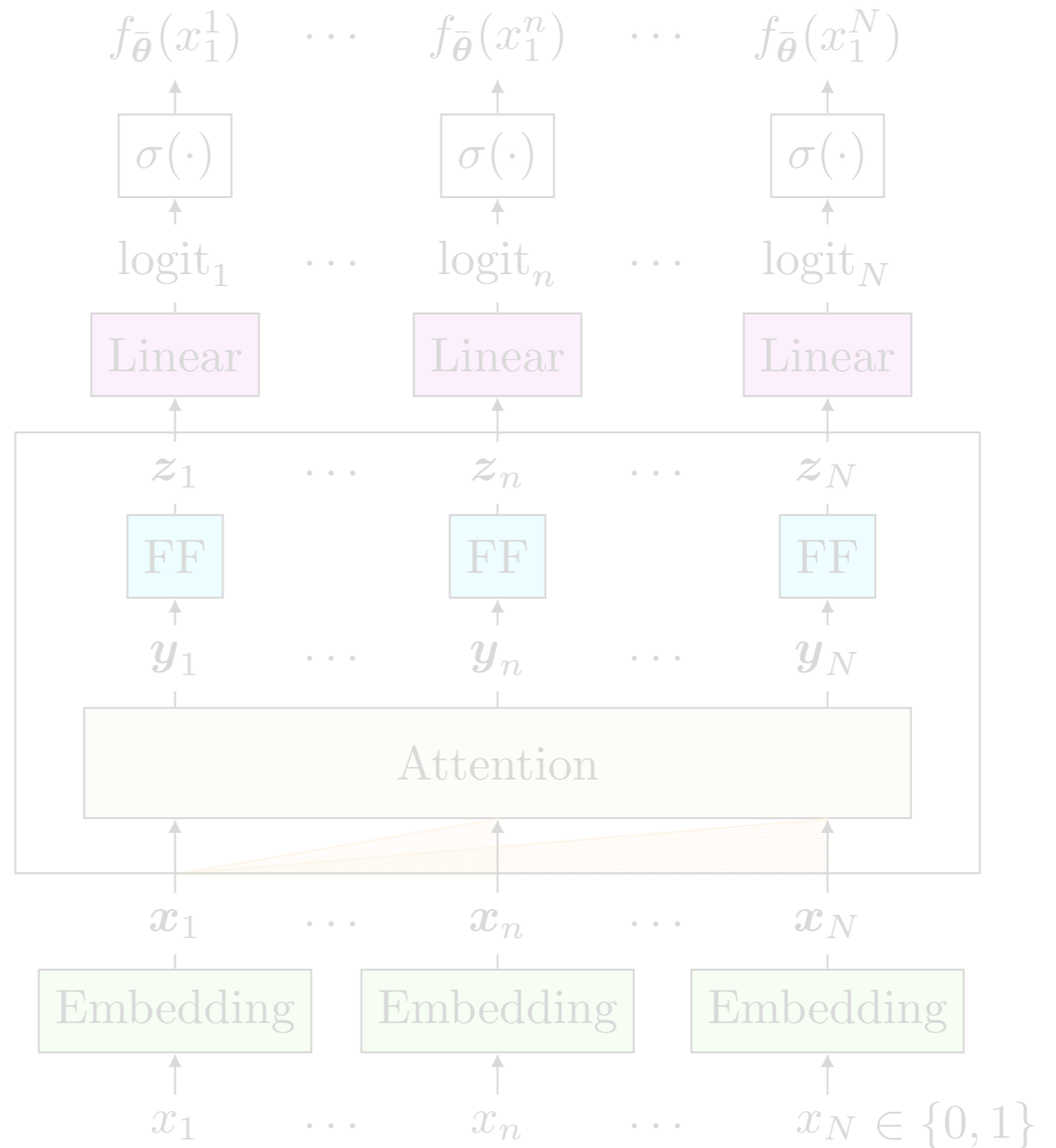
$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R} \quad \mathbf{w}$$

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

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

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

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



Gradient flow

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

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

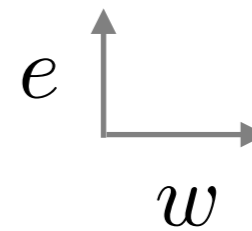
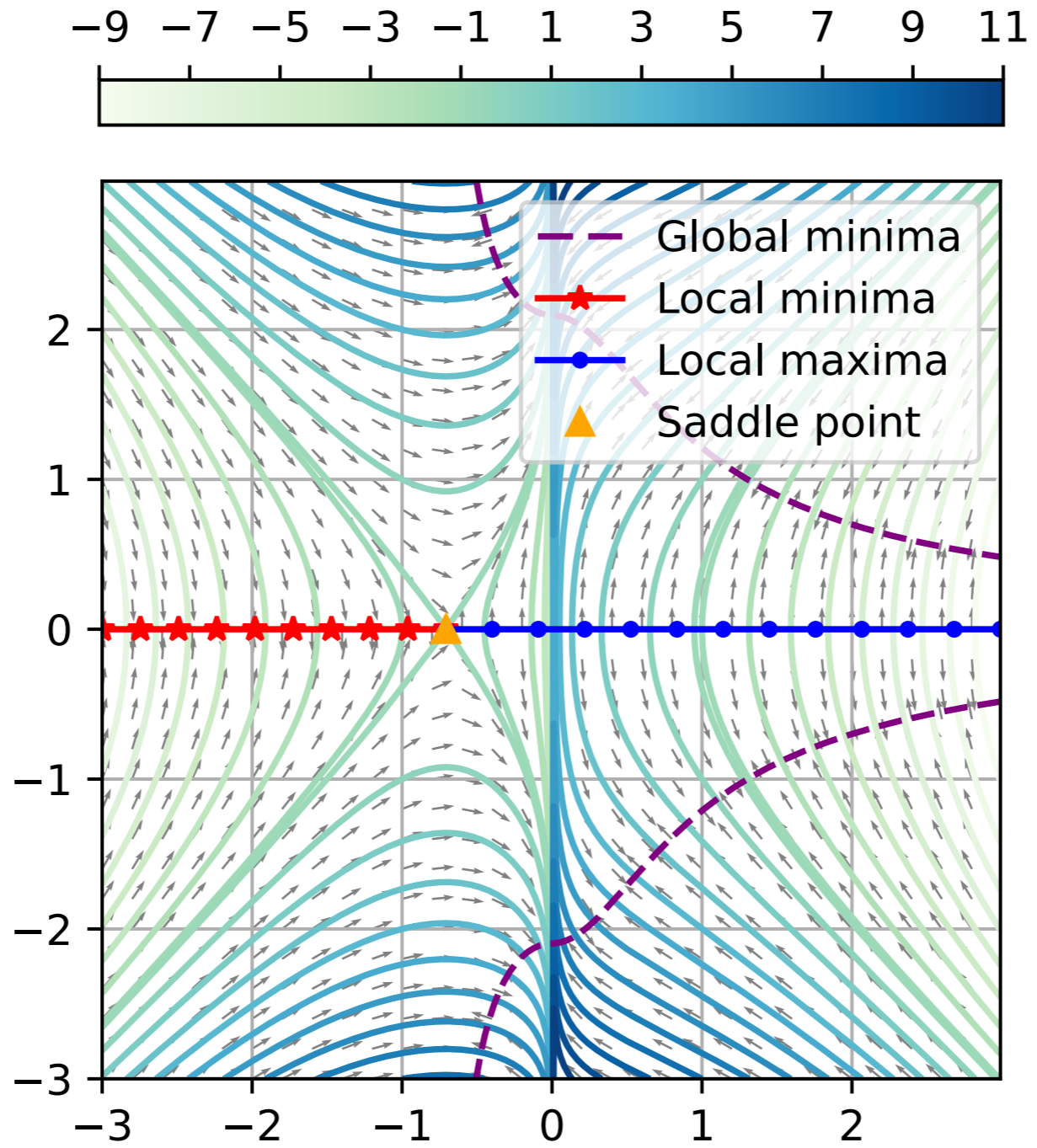
How does the flow look like?

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

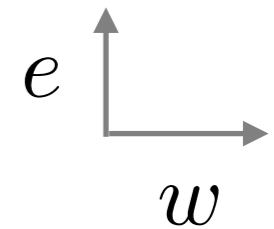
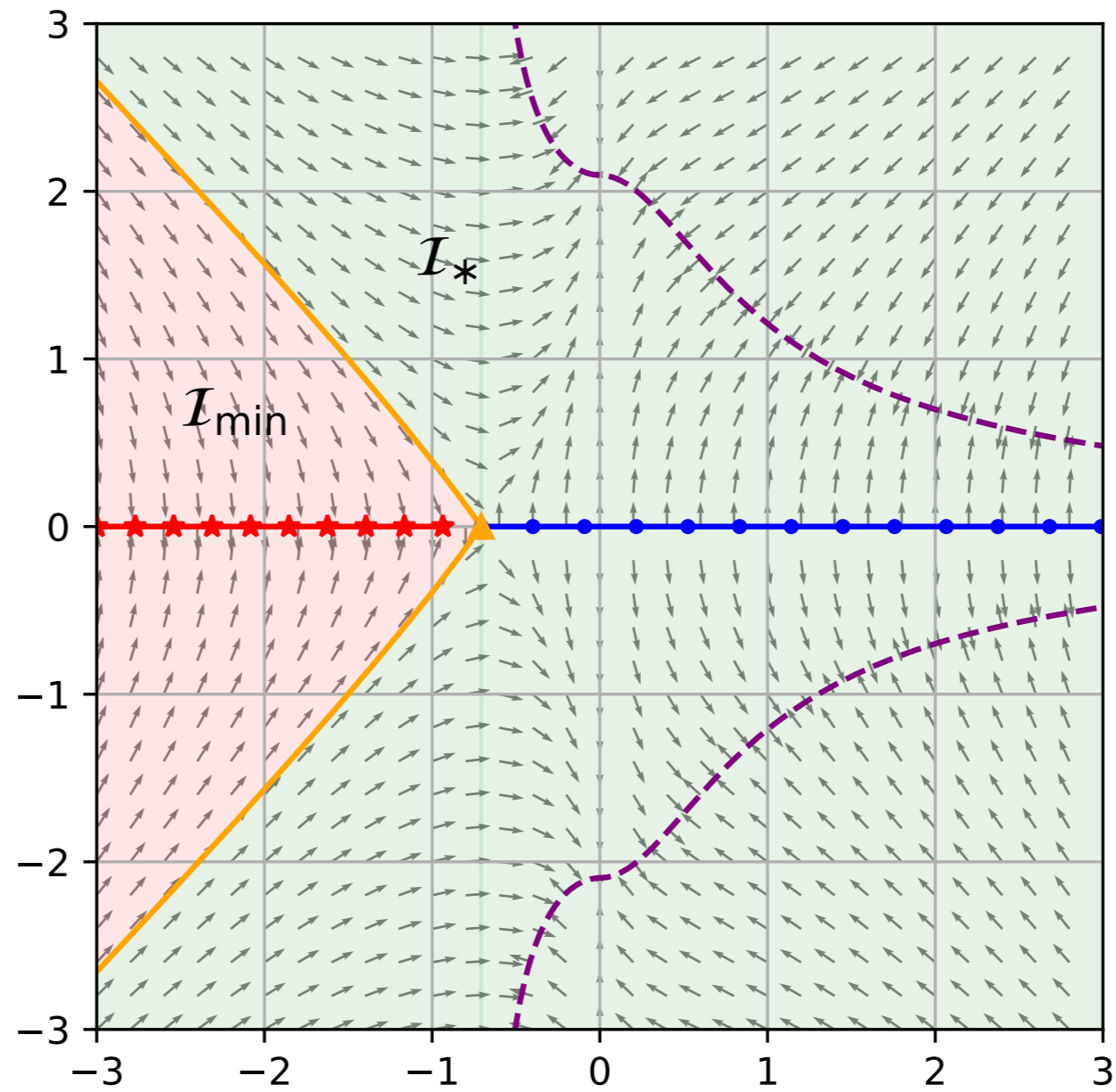
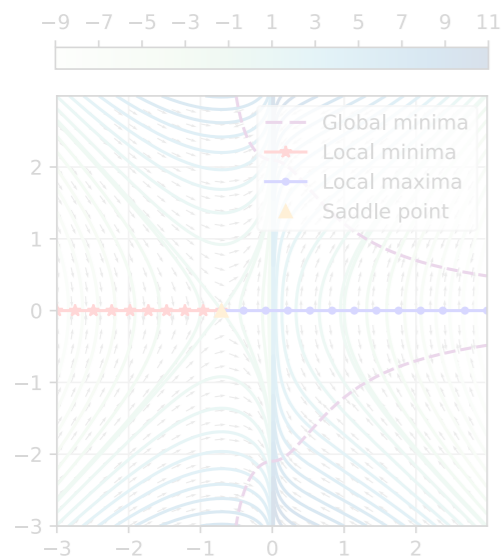
How does the flow look like?

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

$$p+q < 1$$

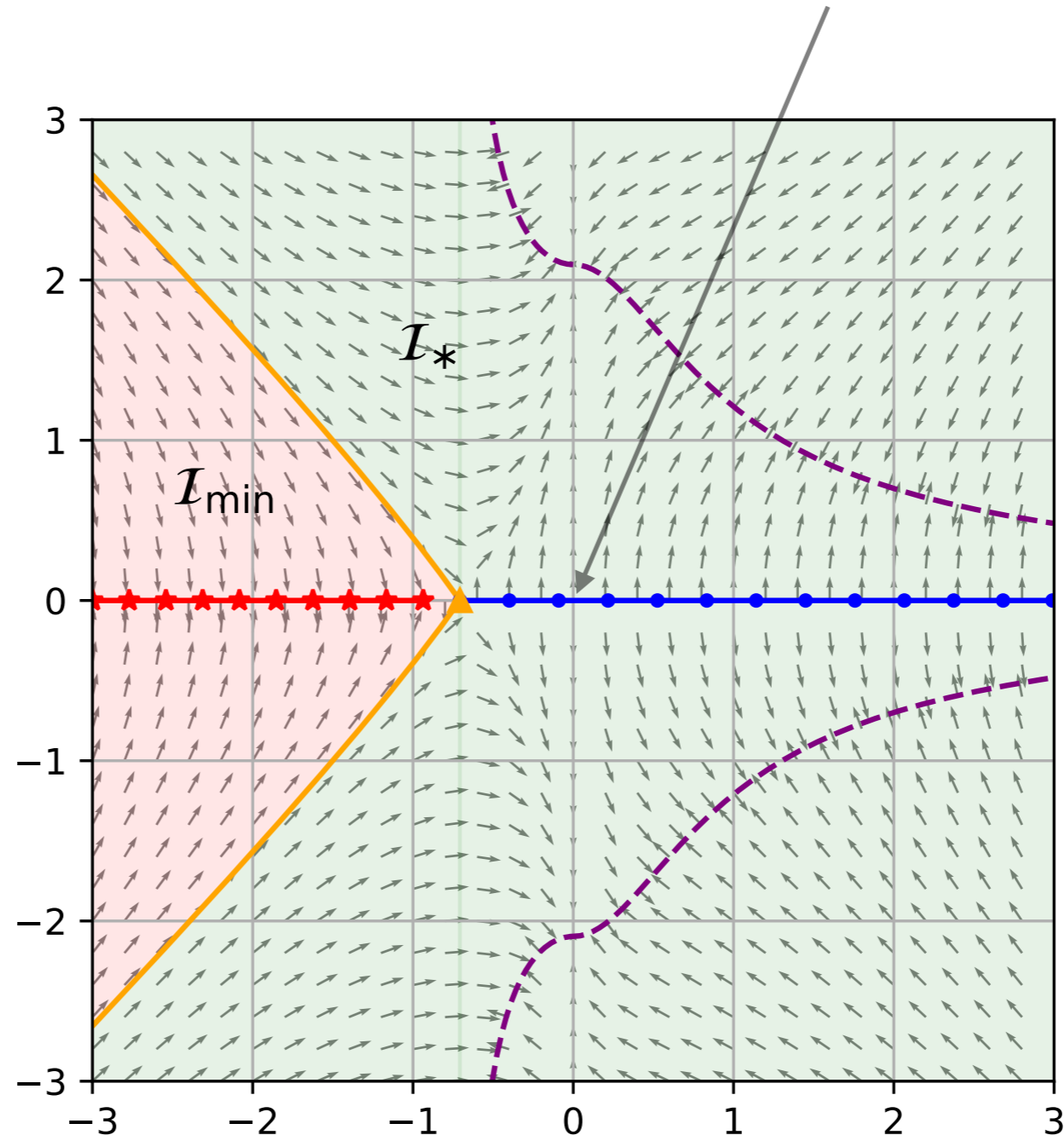
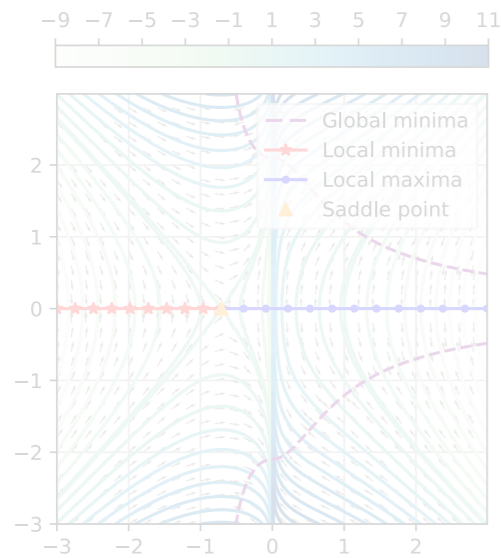


$$p+q < 1$$

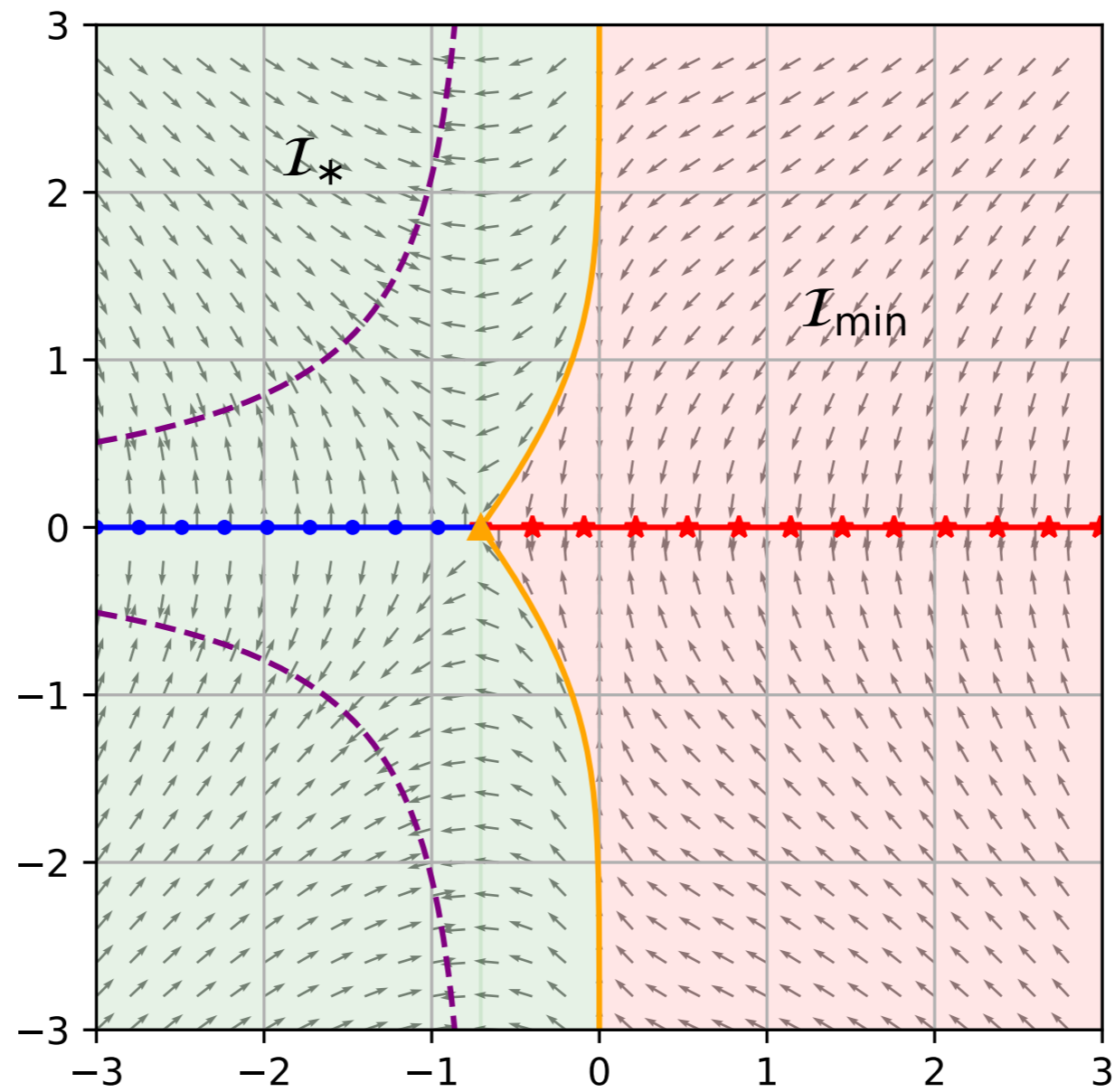


$$p+q < 1$$

Gaussian init. converges to global minima

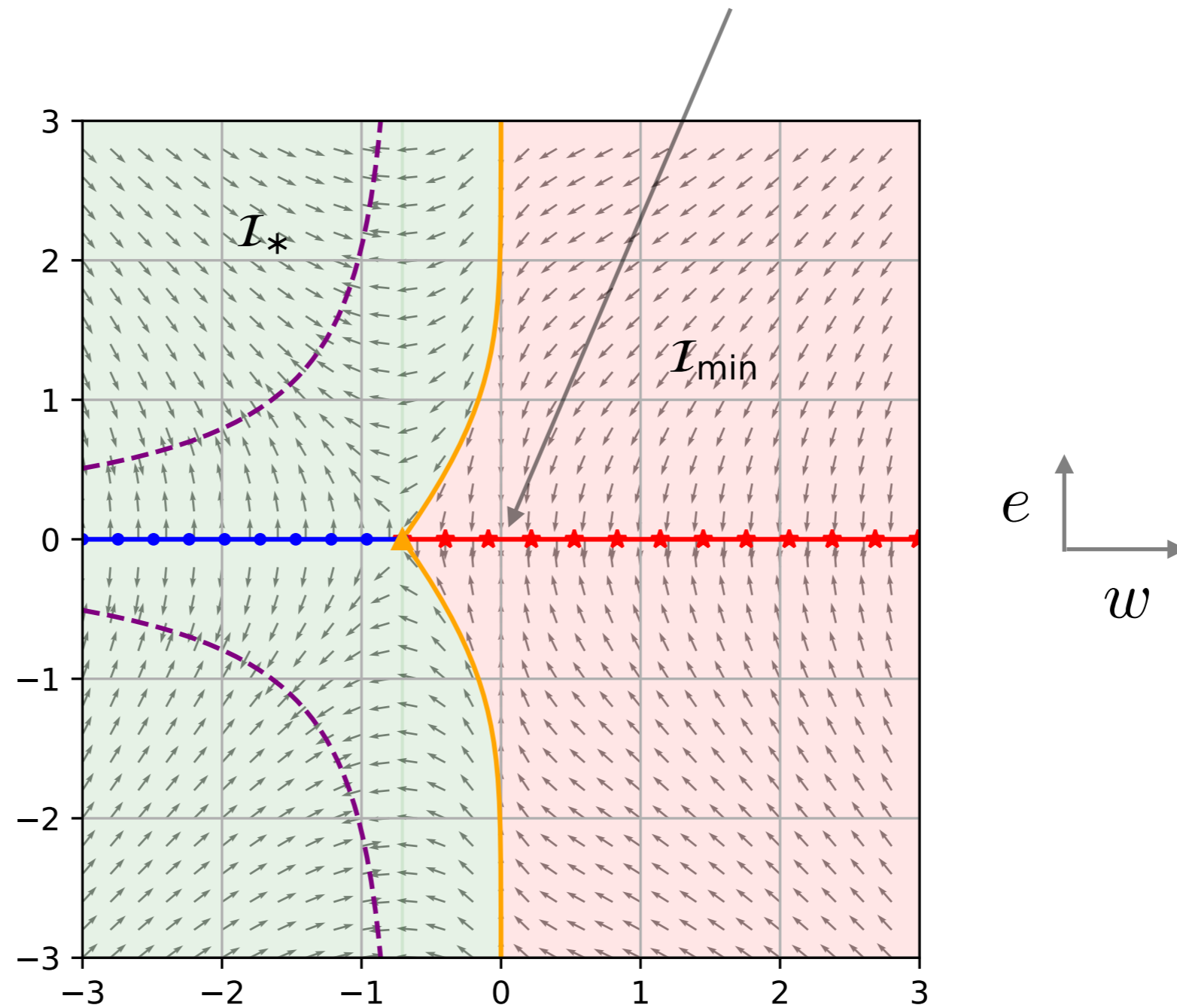


$$p+q > 1$$

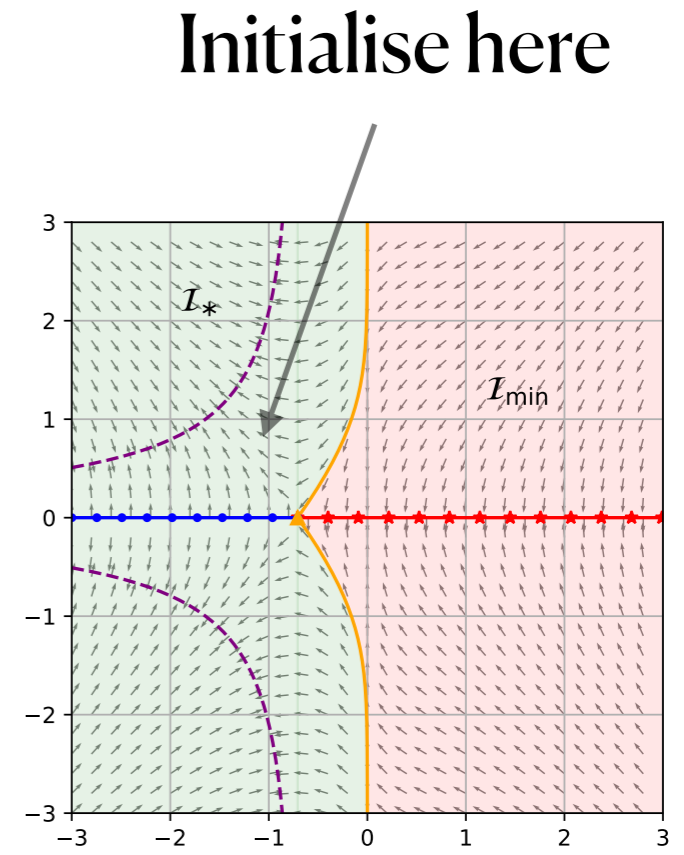
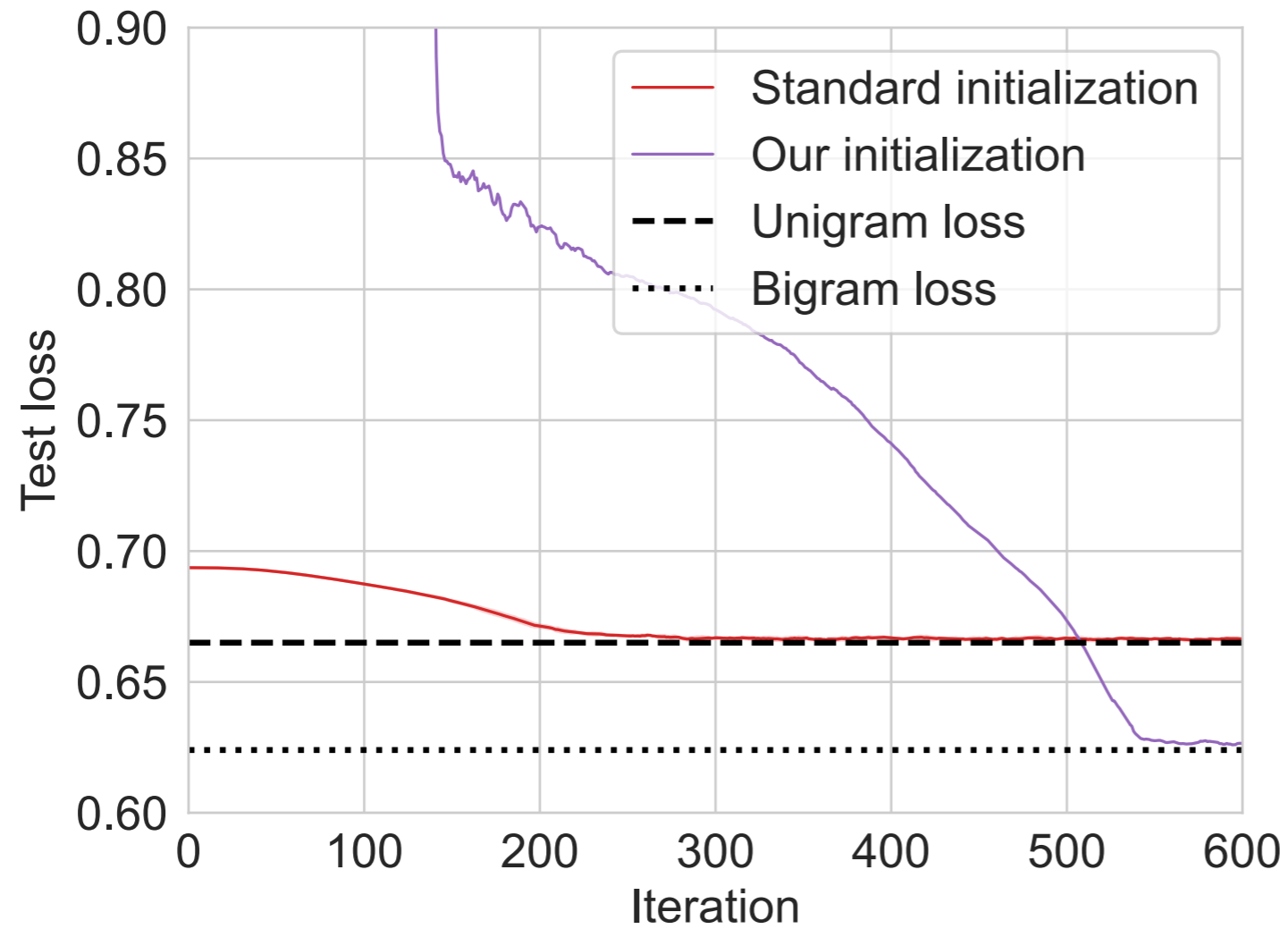


$$p+q > 1$$

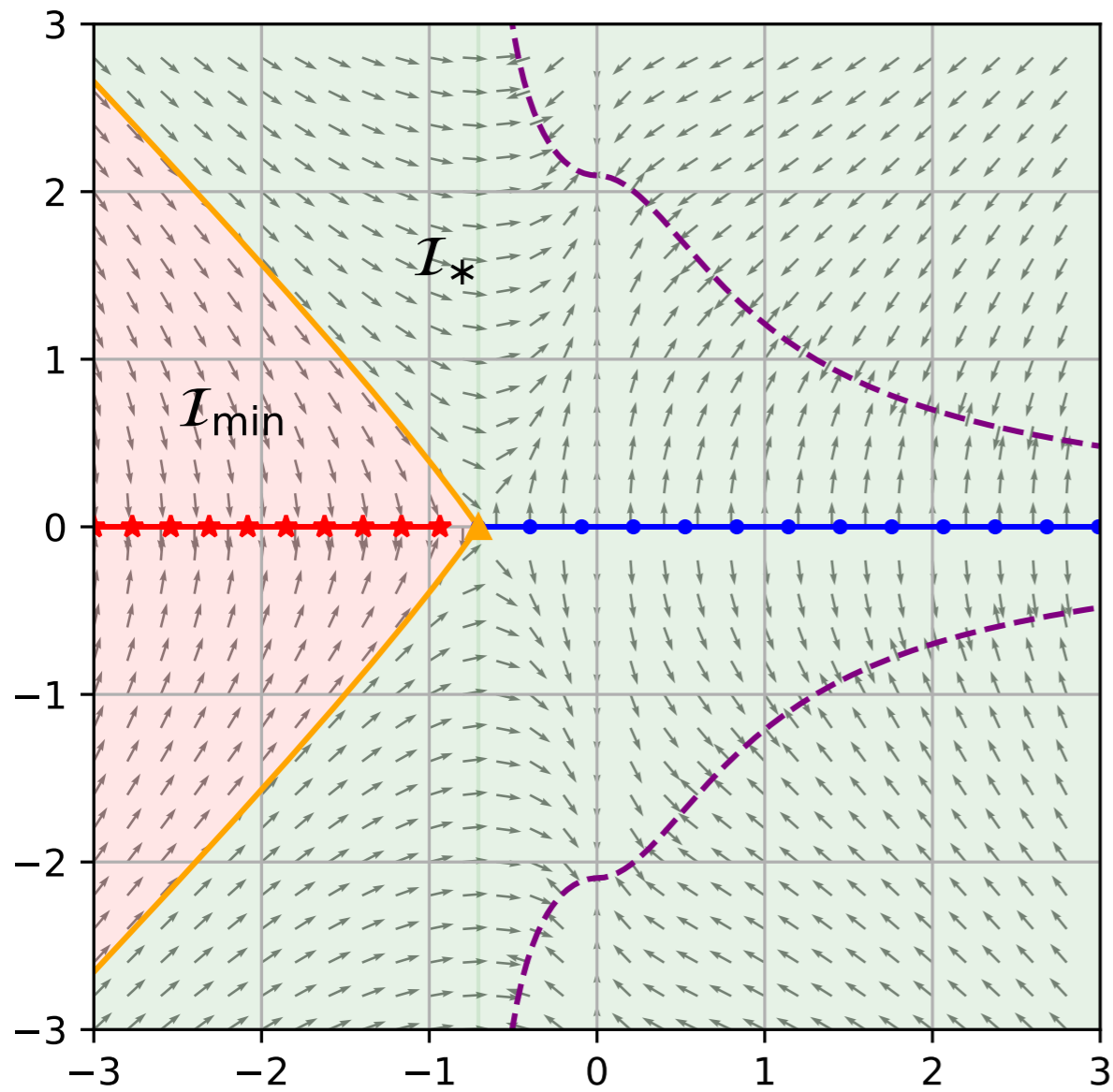
Gets stuck at local minima!



Can we escape it?

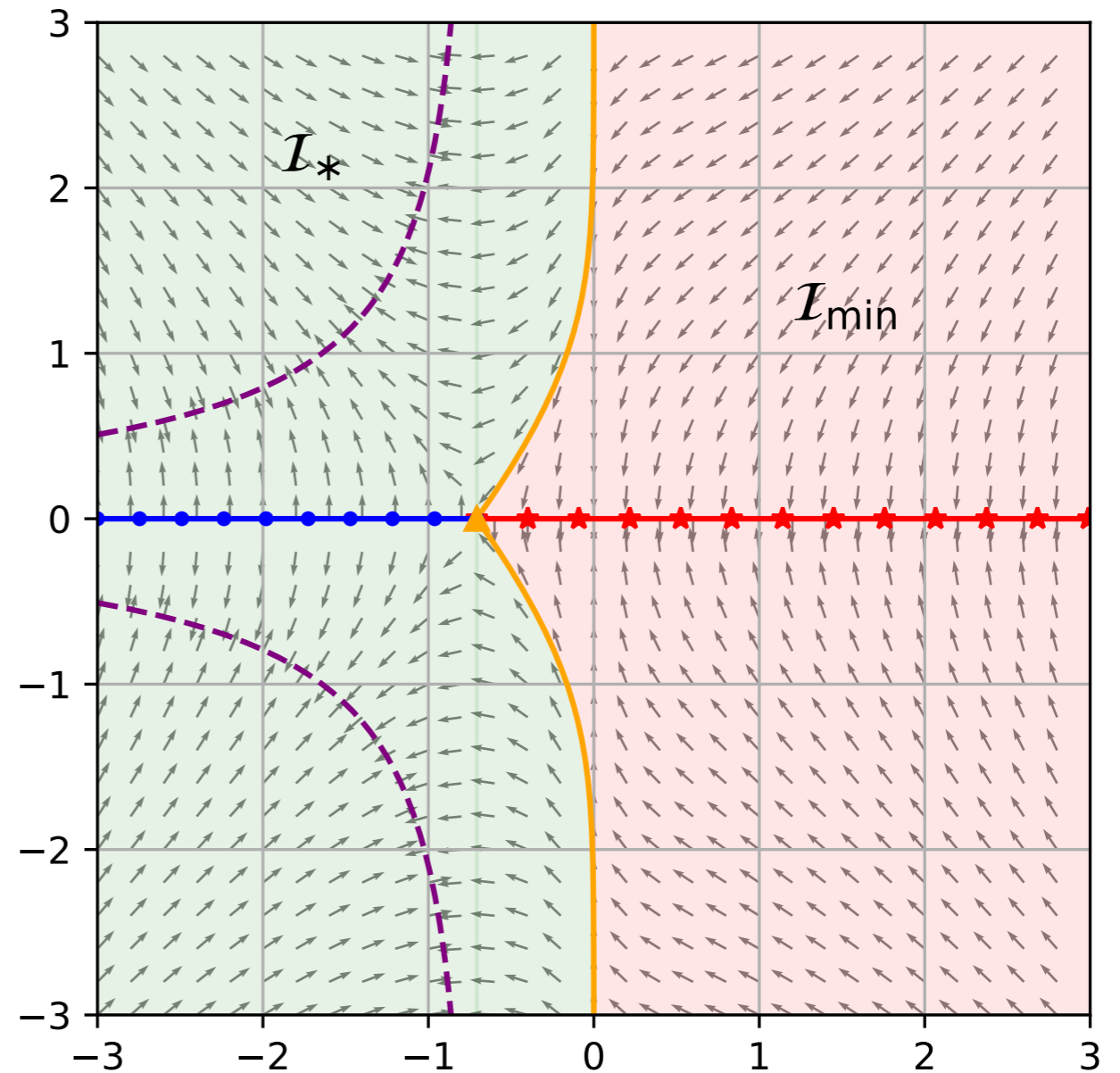


Converges to global minima

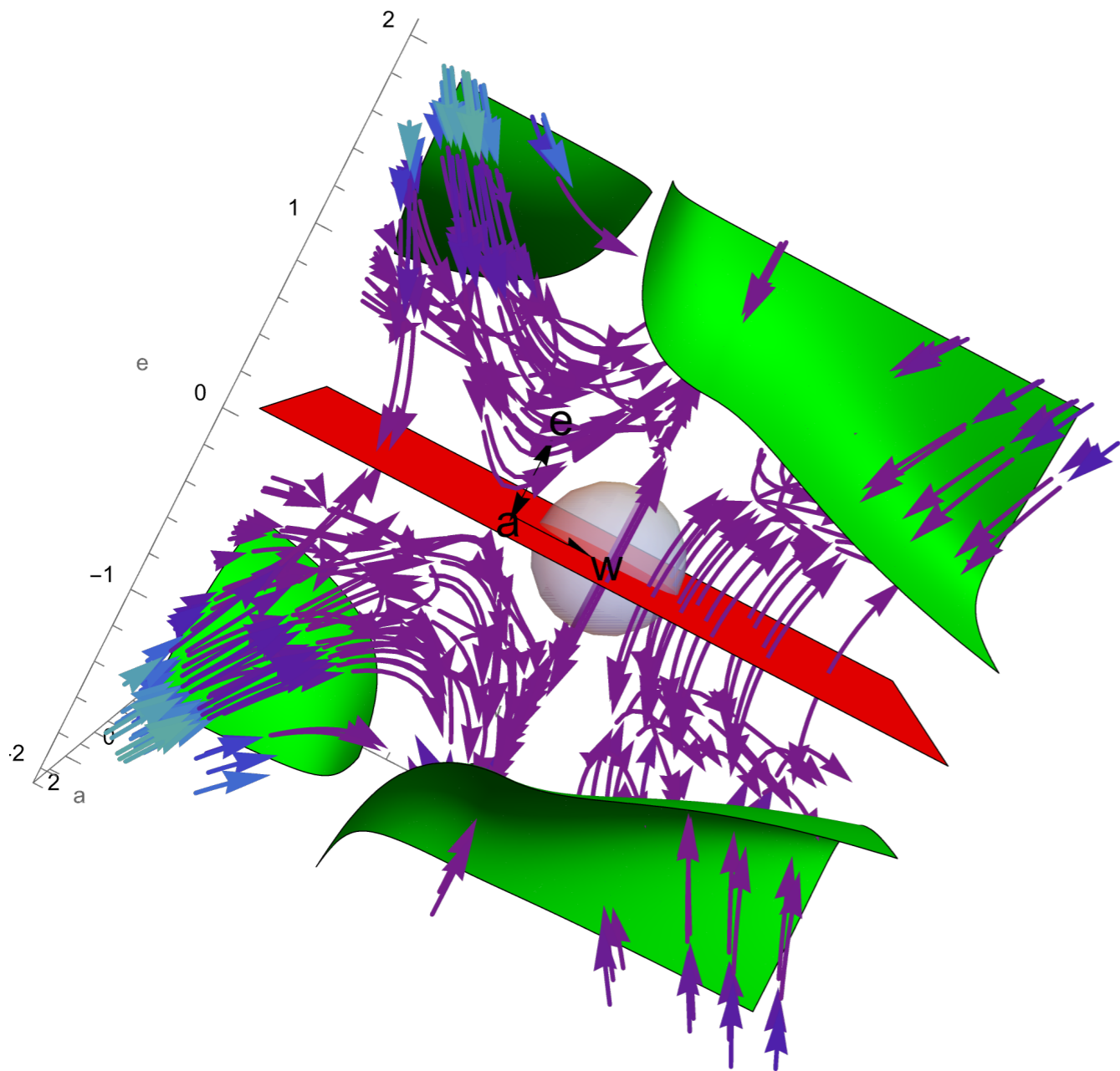


$$p + q < 1$$

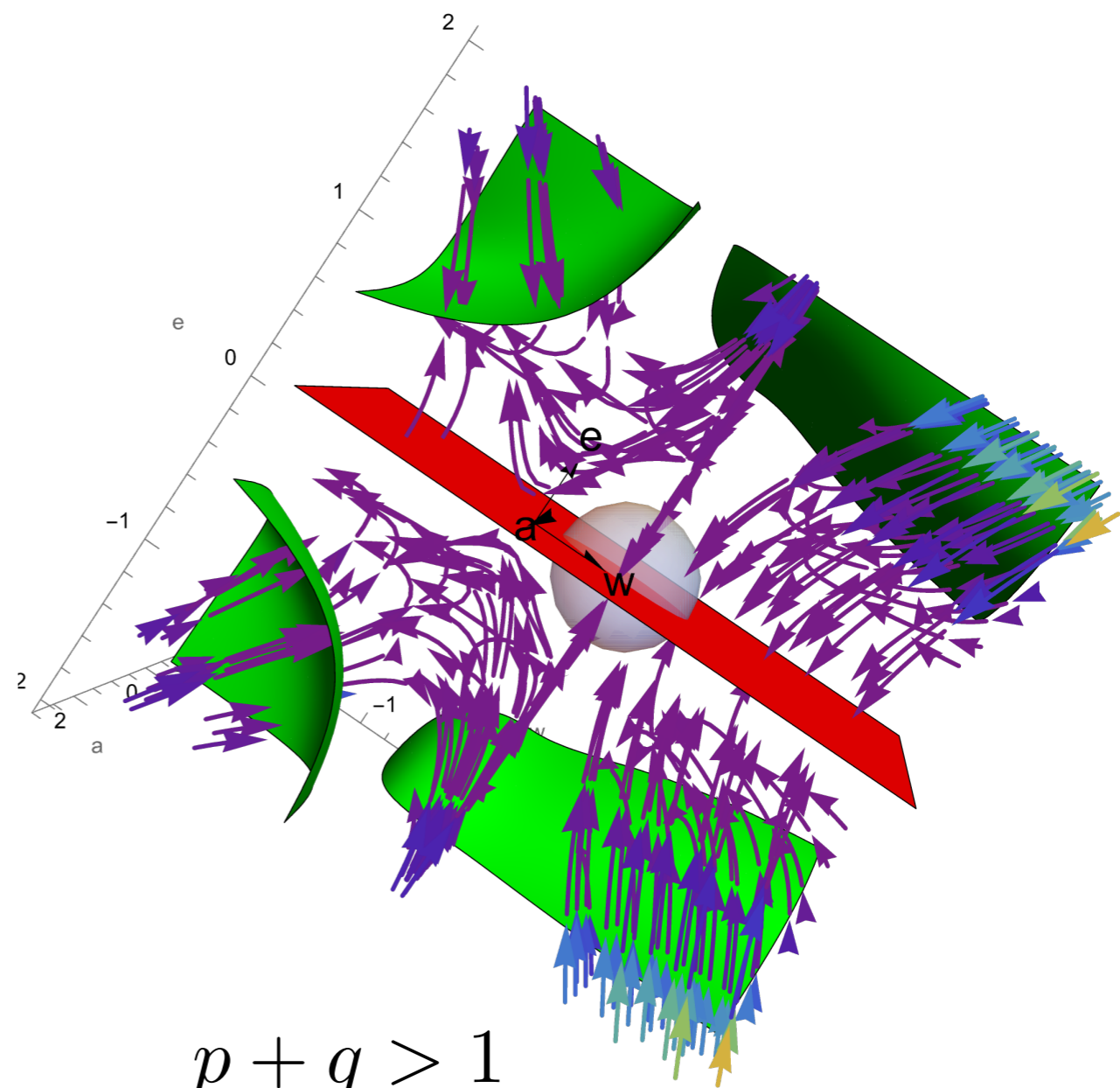
Gets stuck at local minima!



$$p + q > 1$$



$$p + q < 1$$



$$p + q > 1$$

■ Global minima
 ■ Local minima
 ■ Ball around origin

Part I.a



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 I.b



Markovian inputs

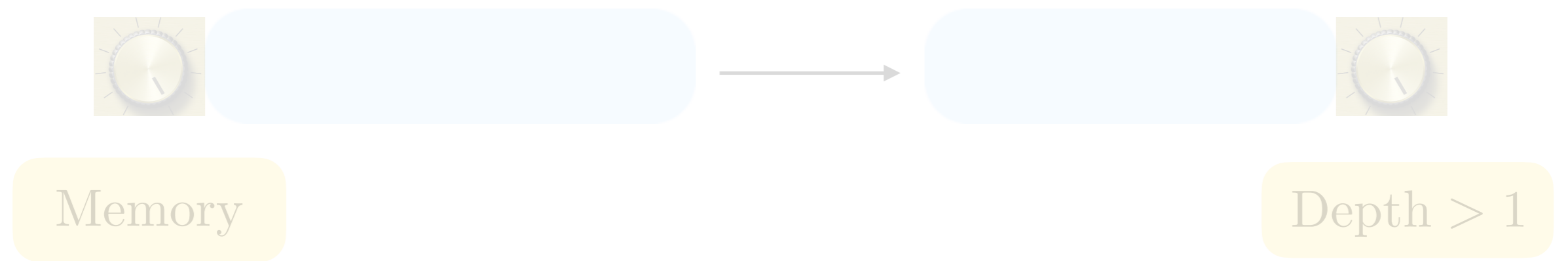


Transformers



Memory

Depth > 1

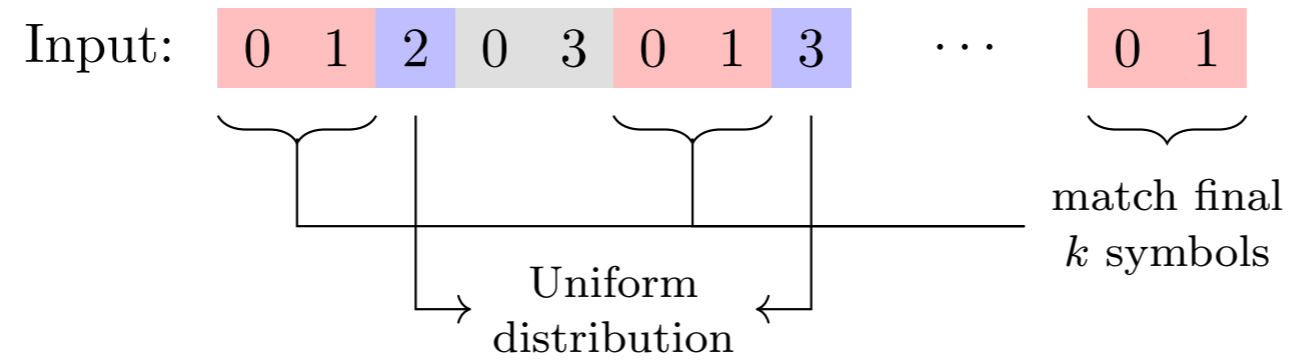


In-context learning (ICL) emerges!

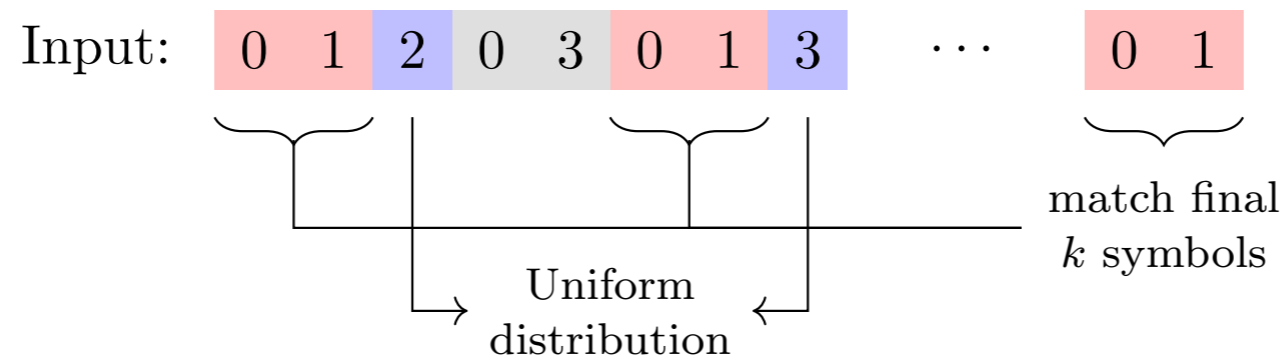
ICL

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

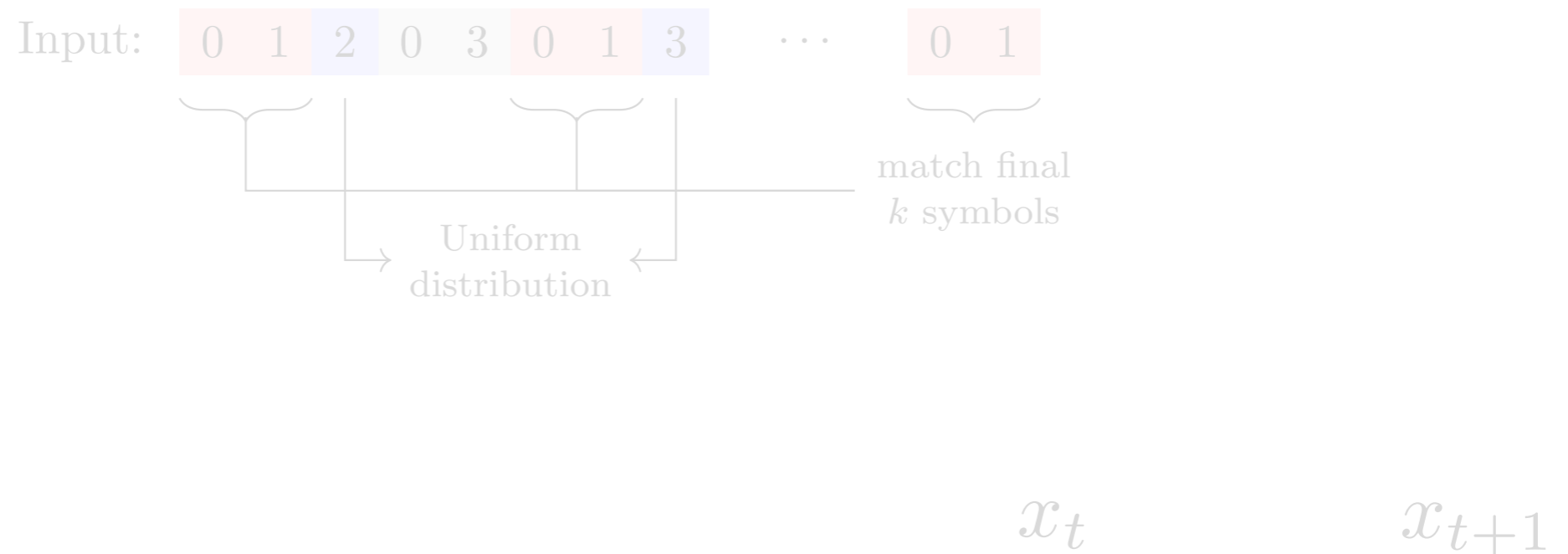


ICL for k th-order Markov



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

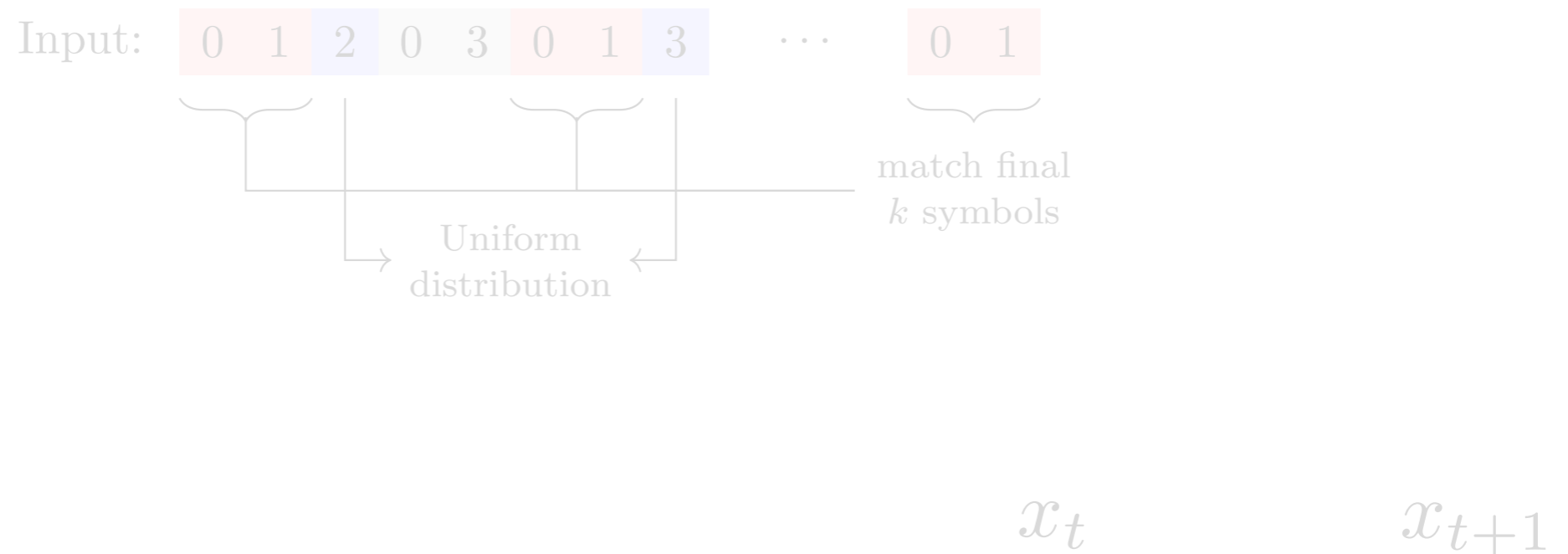
ICL for k th-order Markov



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 k th-order Markov



In-context estimator:

Context-matching

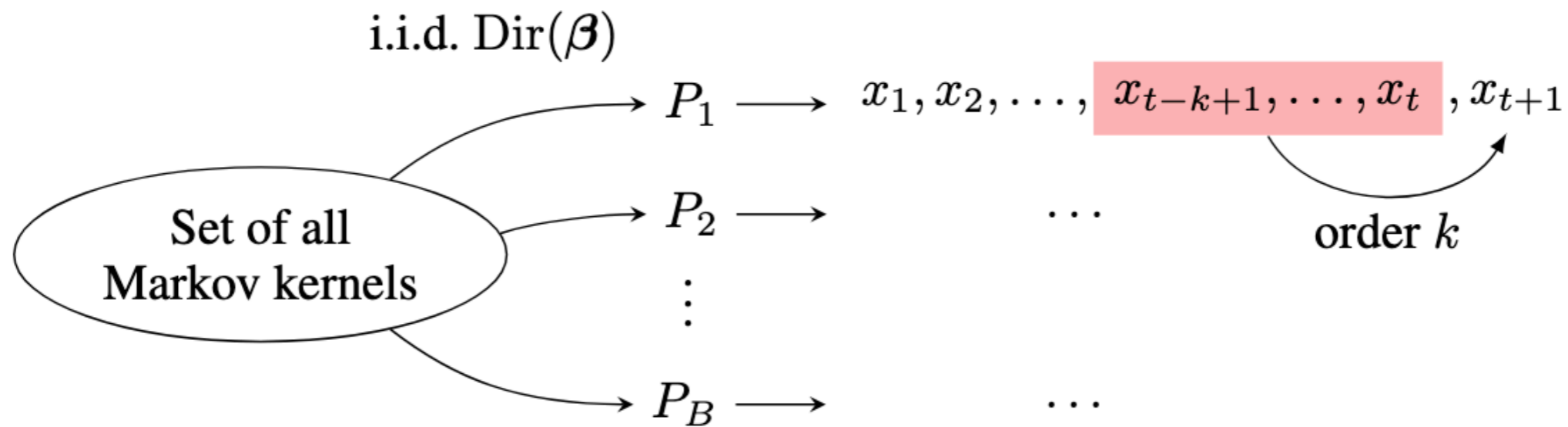
$$\widehat{\text{Pr}}_k(x \mid x_1, \dots, x_t) \triangleq \frac{\sum_{i=k+1}^t \mathbb{I}(x_i = x, x_{i-1} = x_t, \dots, x_{i-k} = x_{t-k+1})}{\sum_{i=k+1}^t \mathbb{I}(x_{i-1} = x_n, \dots, x_{i-k} = x_{t-k+1})}$$

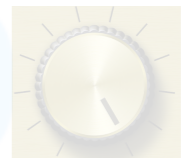
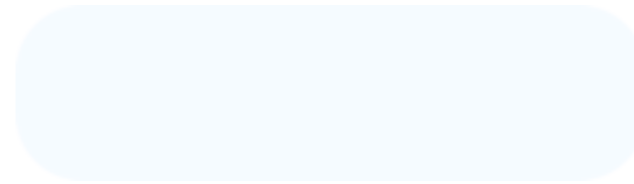
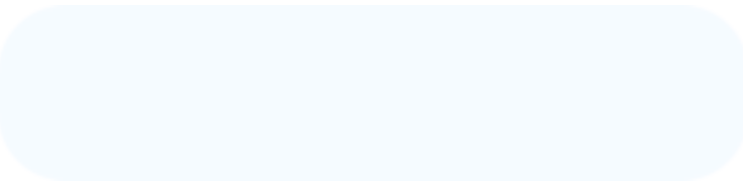
In-context Markov chains



Memory = k

Depth

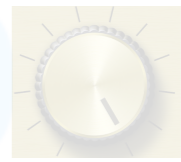
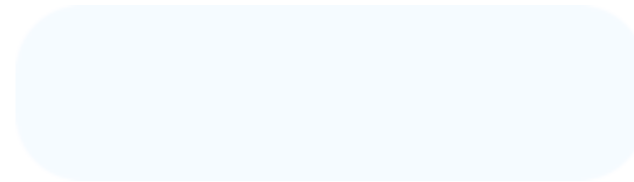
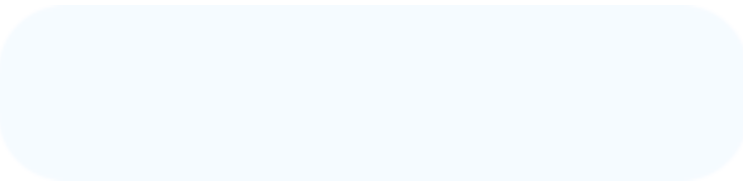




Memory = k

Depth?

What we know so far

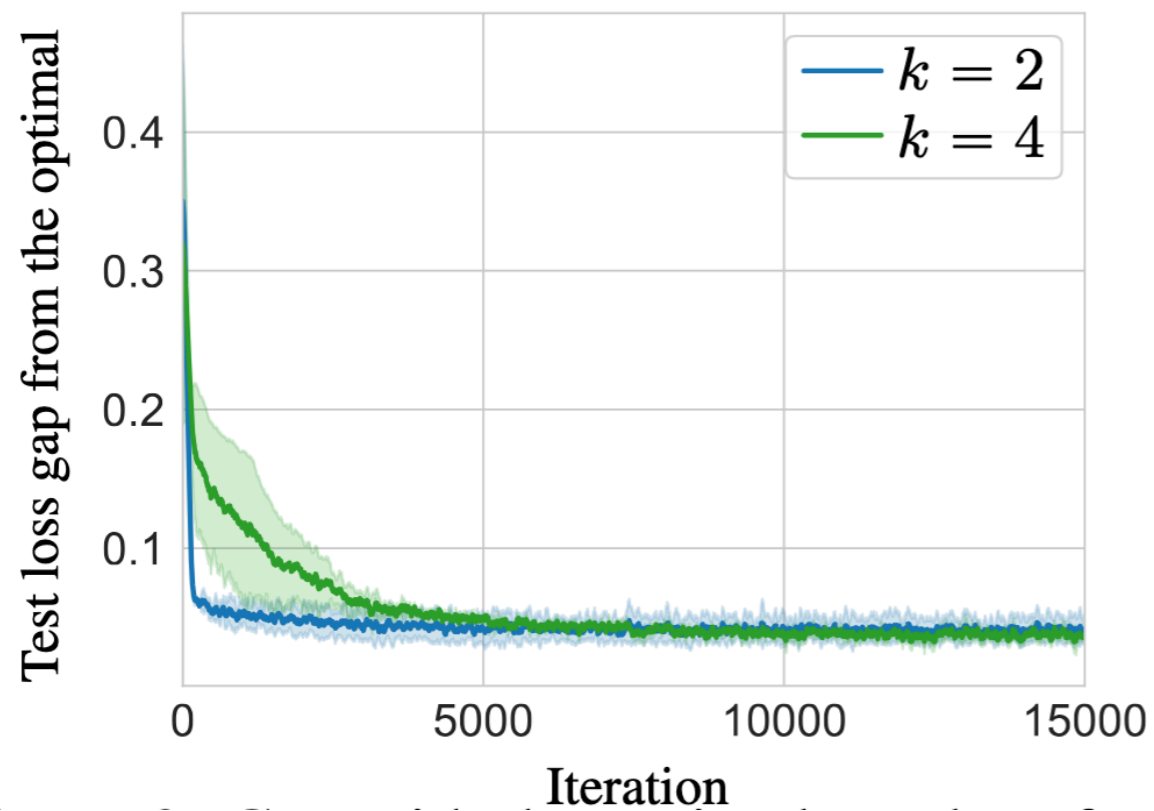


Memory = k

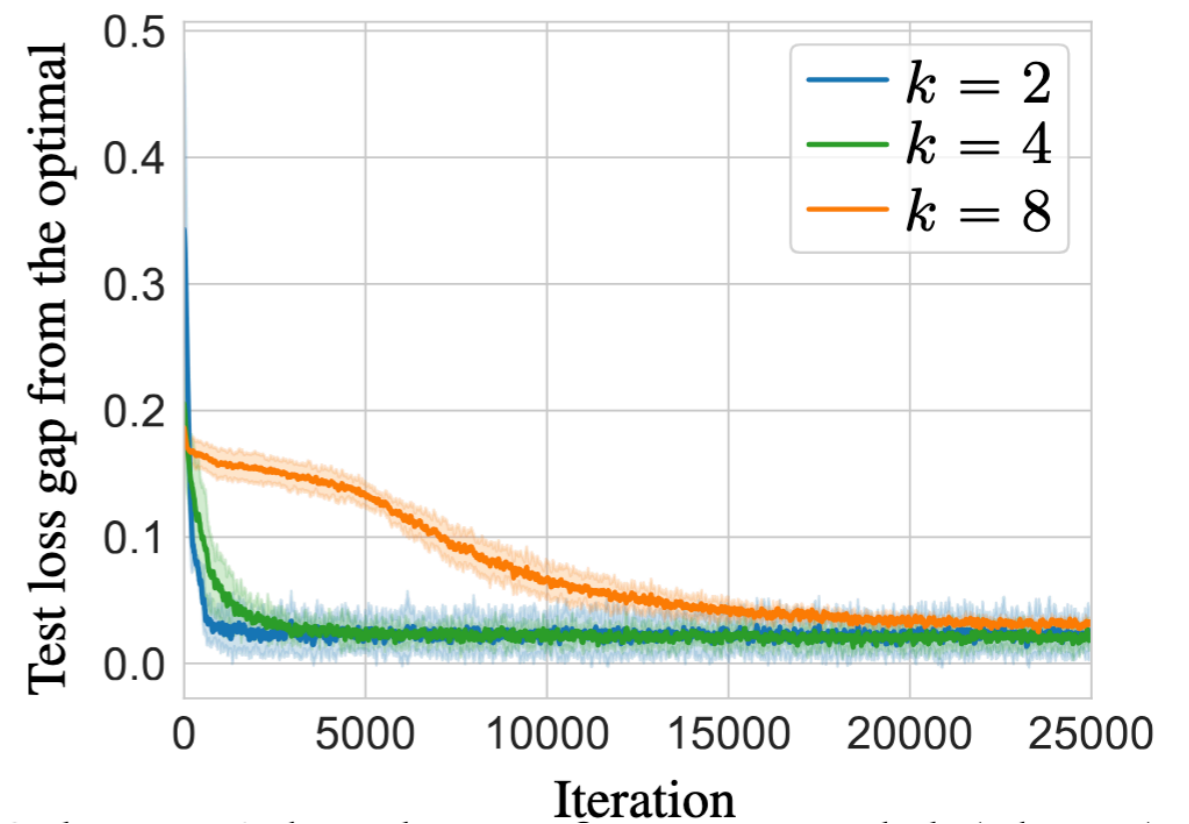
Depth?

- ▶ (Edelman et al., 2024 & Nichani 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 layer norm, you need logarithmic depth

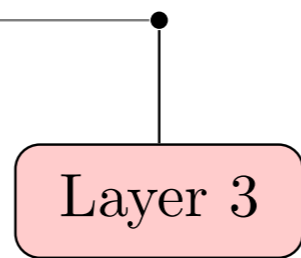
Intuition

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

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

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

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



Match contexts & count

Capture context

x_1, \dots, x_T

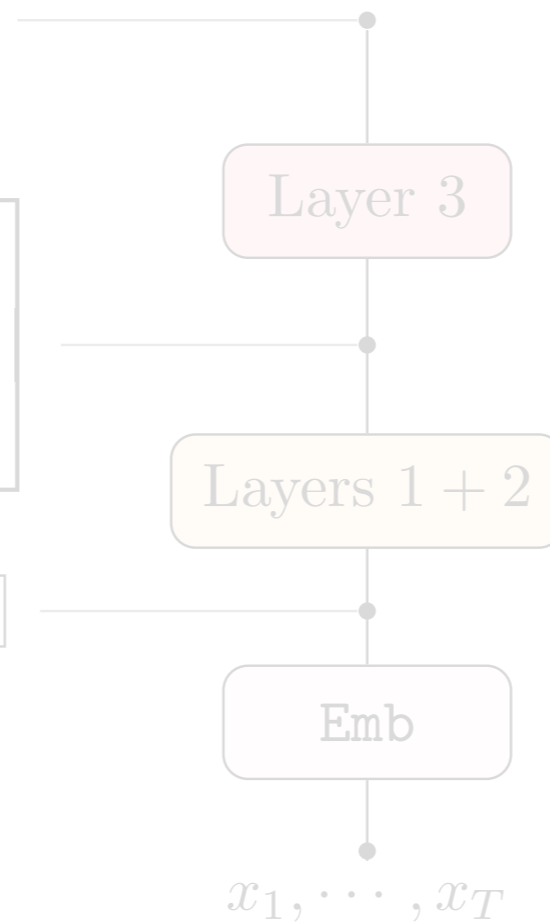
Intuition

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

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

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

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



2 layers suffice!

[Ekbote et al. 2025]

Part I.b



Markovian inputs



Transformers



Memory

Depth > 1

What do they learn?



Part I.b



Markovian inputs

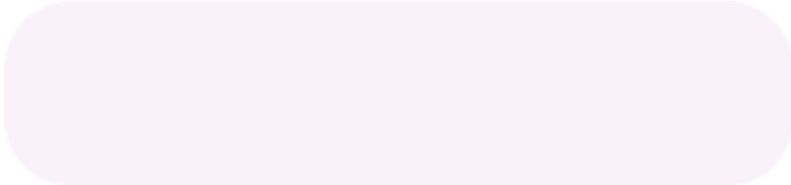


Transformers



Memory

Depth > 1

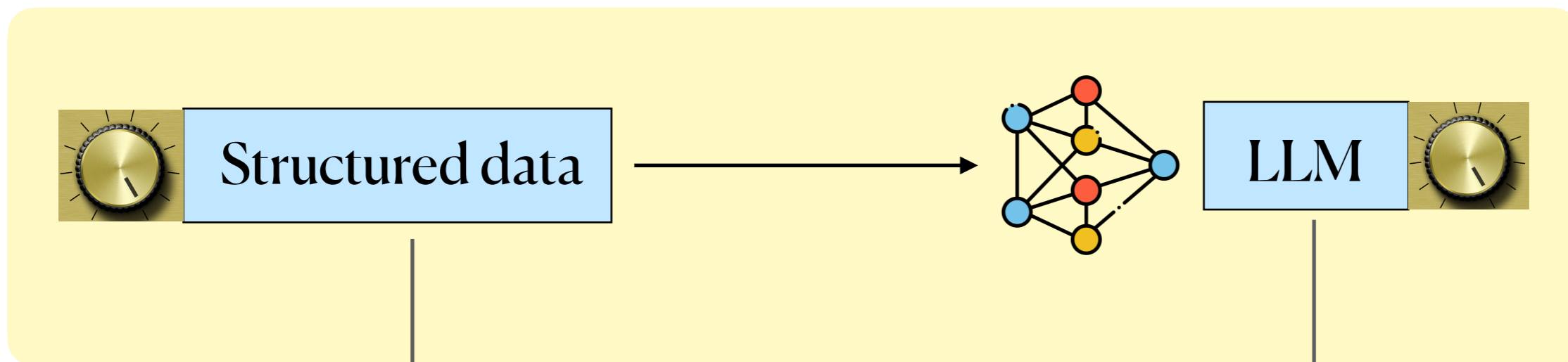


How do they learn?

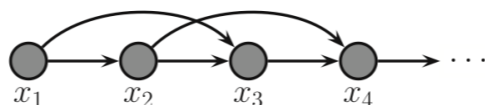


Ongoing..

Part I



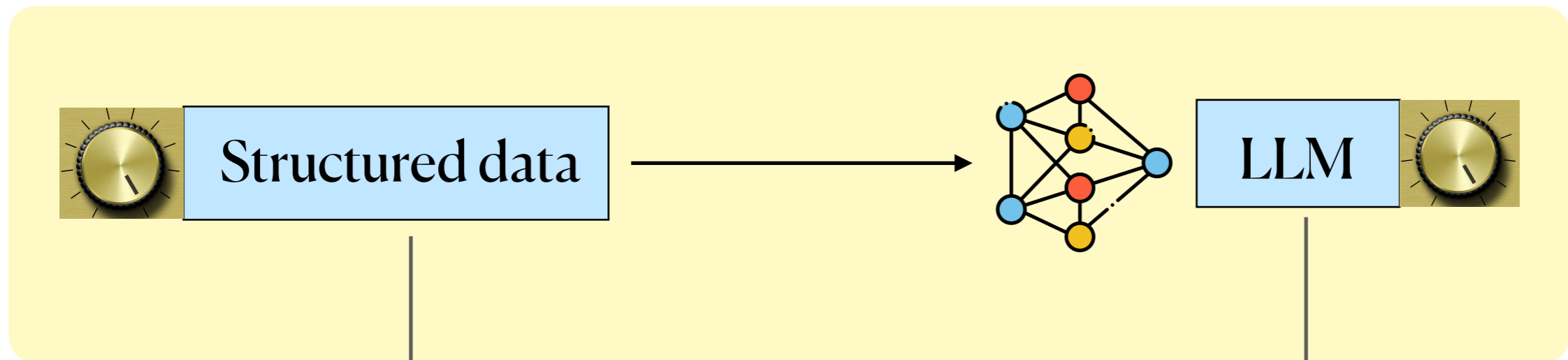
Markov processes



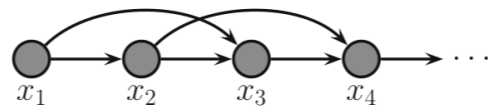
Transformers



Part II



Markov processes



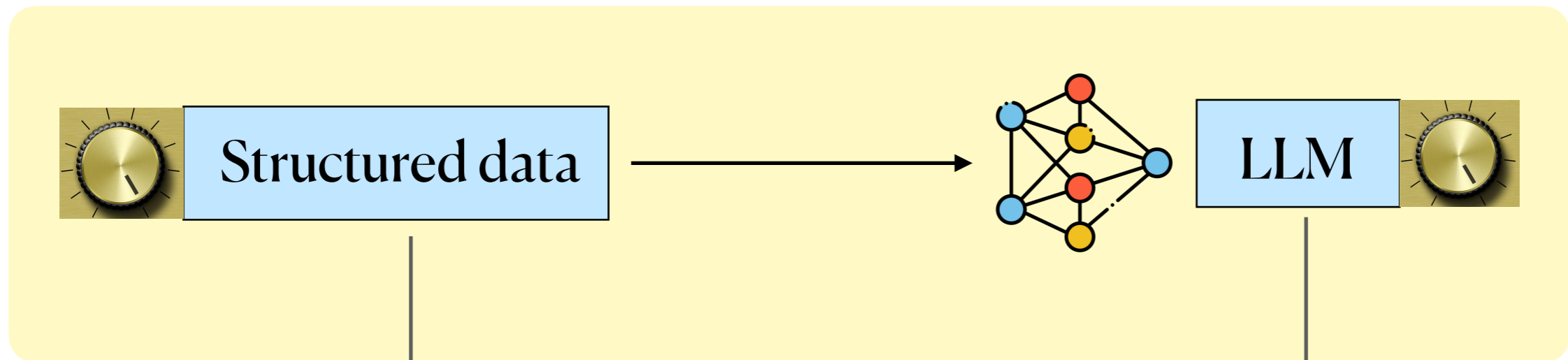
State-Space Models
(Mamba)



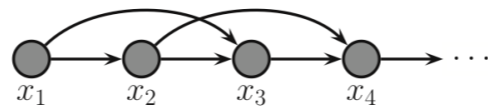
**Even a single-layer Mamba learns the
optimal Laplacian smoothing estimator!**

**Convolution plays a critical role, more so
than gating and non-linearities**

Part II



Markov processes



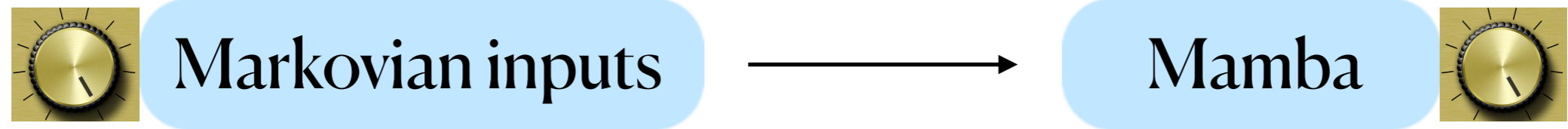
State-Space Models
(Mamba)



**From Markov to Laplace: How Mamba In-Context
Learns Markov chains**

ICLR 2026

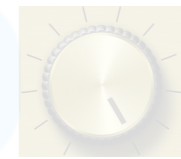
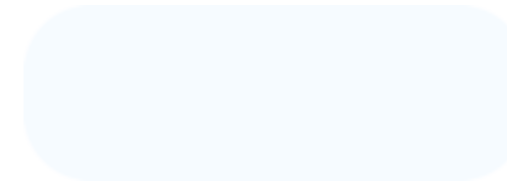
Mamba with Markov



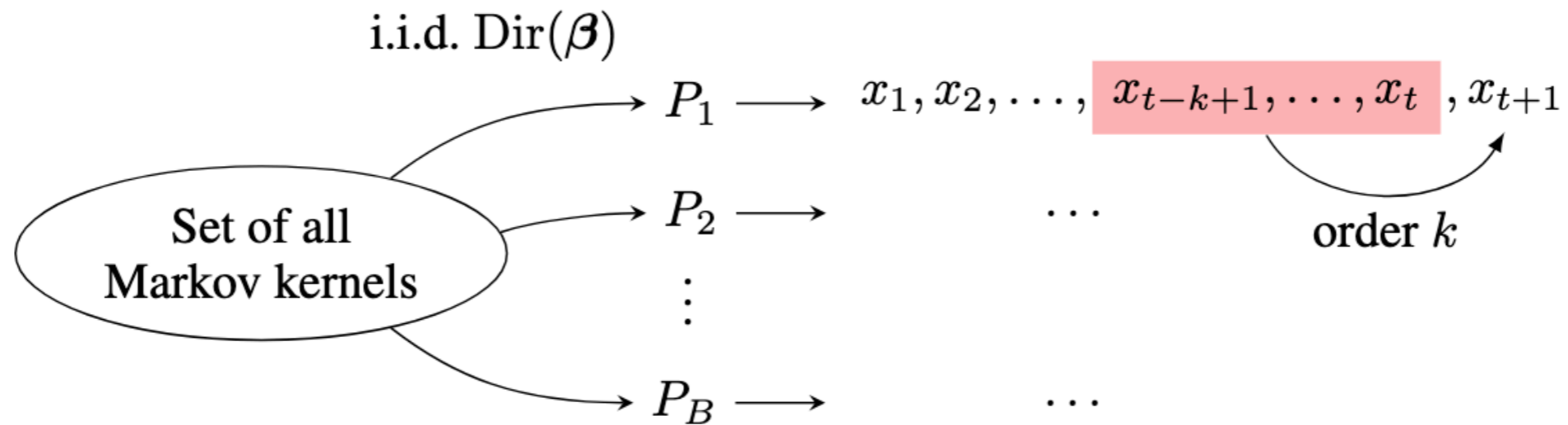
In-context Markov chains

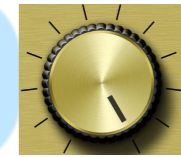
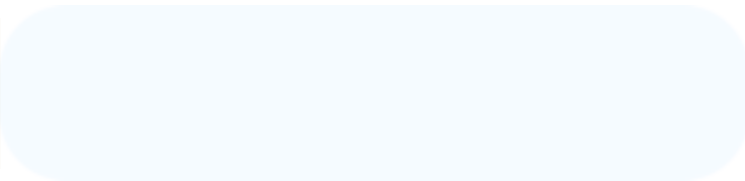
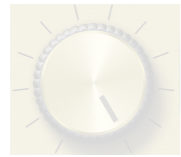


Markovian inputs



Memory = k

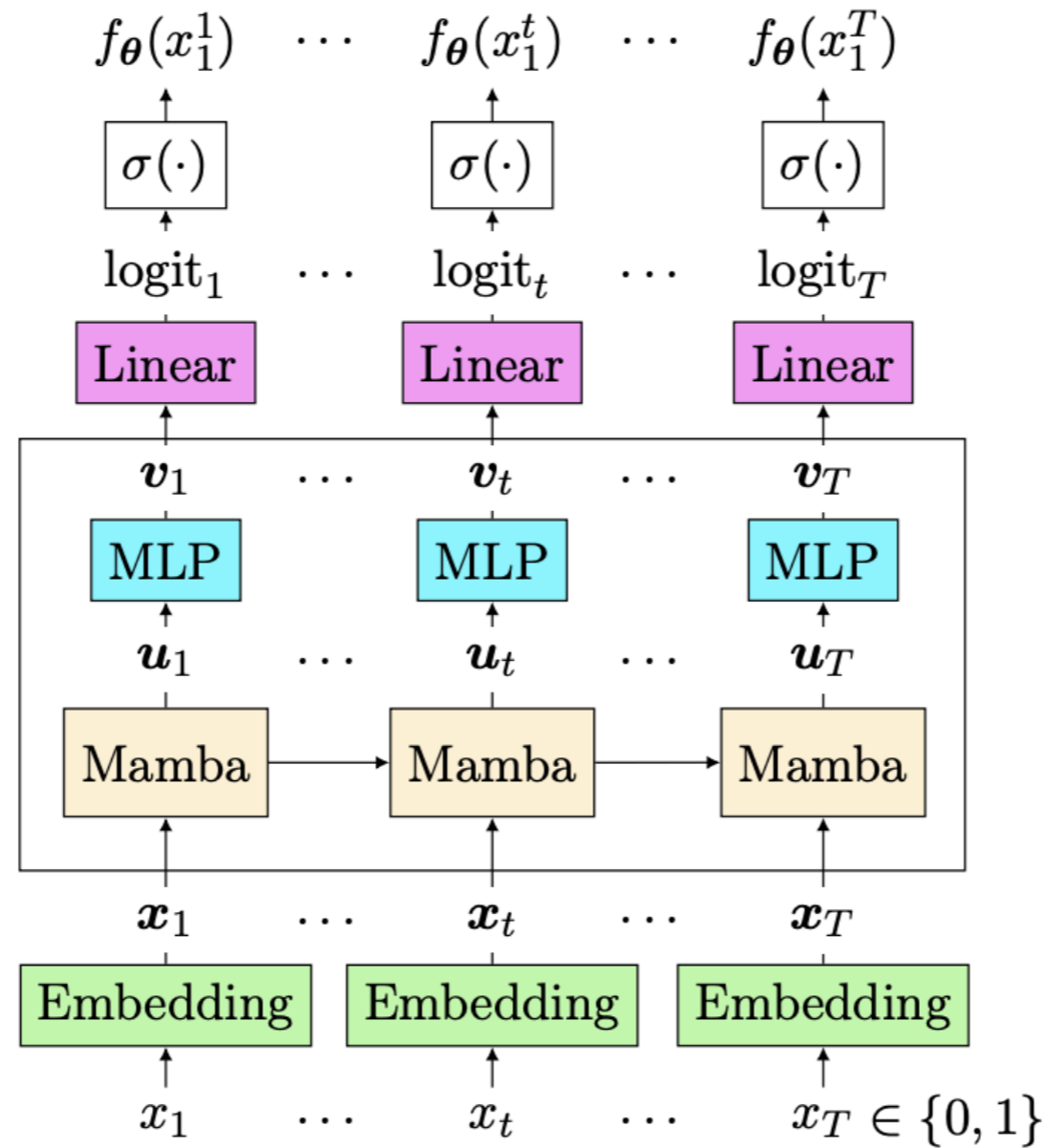




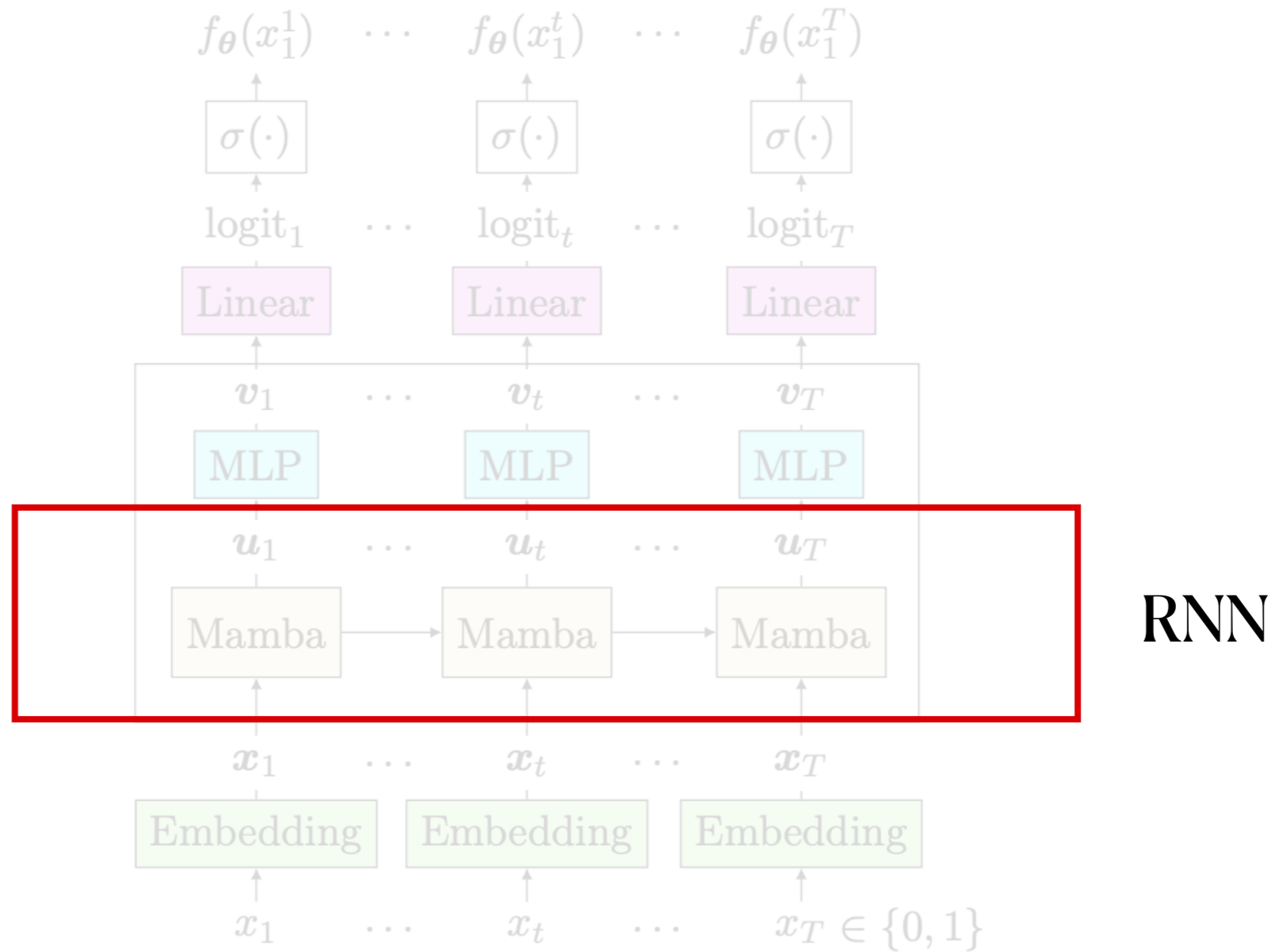
Memory = k



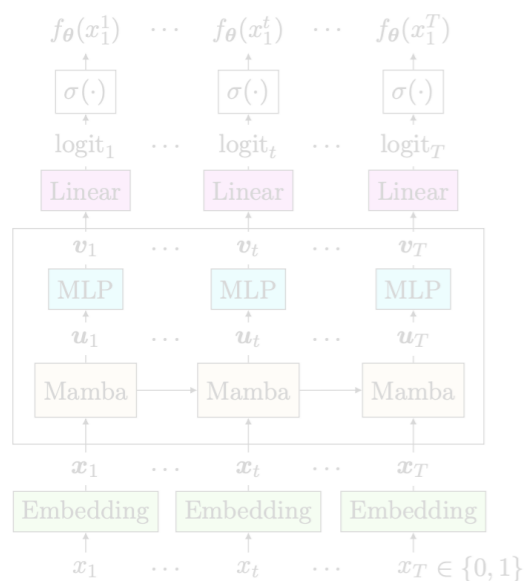
Mamba-based language model



Mamba-based language model



Mamba



$$H_t = a_t H_{t-1} + \tilde{\mathbf{x}}_t \mathbf{b}_t^\top \in \mathbb{R}^{ed \times N},$$

$$\mathbf{y}_t = H_t \mathbf{c}_t \in \mathbb{R}^{ed},$$

$$\mathbf{z}_t = \mathbf{y}_t \odot \text{ReLU}(W_z \mathbf{x}_t) \in \mathbb{R}^{ed},$$

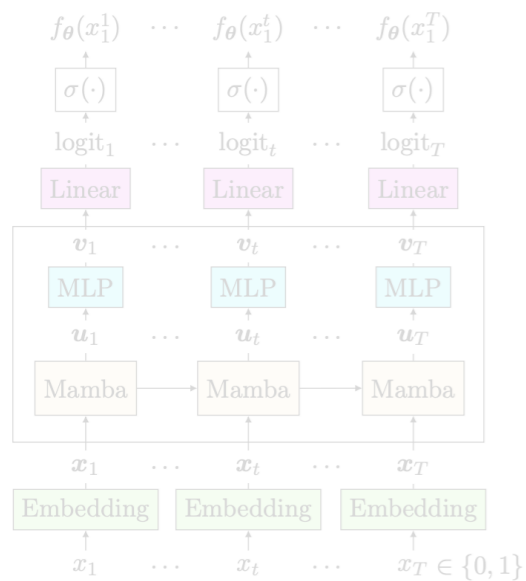
$$\mathbf{o}_t = W_o \mathbf{z}_t \in \mathbb{R}^d,$$

**Input-selective
recurrence**

$$\text{Mamba} : \mathbf{x} = (\mathbf{x}_t)_{t=1}^T \in \mathbb{R}^{d \times T} \mapsto \mathbf{o} = (\mathbf{o}_t)_{t=1}^T \in \mathbb{R}^{d \times T}$$

Mamba

Temporal convolution

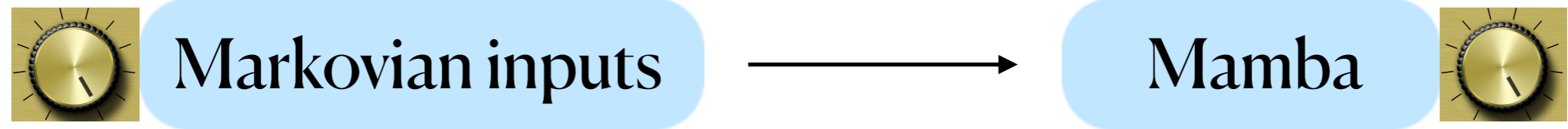


$$\begin{aligned}
 H_t &= a_t H_{t-1} + \tilde{\mathbf{x}}_t \mathbf{b}_t^\top \in \mathbb{R}^{ed \times N}, \\
 \mathbf{y}_t &= H_t \mathbf{c}_t \in \mathbb{R}^{ed}, \\
 \mathbf{z}_t &= \mathbf{y}_t \odot \text{ReLU}(W_z \mathbf{x}_t) \in \mathbb{R}^{ed}, \\
 \mathbf{o}_t &= W_o \mathbf{z}_t \in \mathbb{R}^d,
 \end{aligned}$$

$$\begin{aligned}
 a_t &\triangleq \exp(-a \cdot \Delta_t) \in (0, 1), \\
 \Delta_t &\triangleq \text{softplus}(\langle \mathbf{w}_\Delta, \mathbf{x}_t \rangle + \delta) \in \mathbb{R}, \\
 \tilde{\mathbf{x}}_t &\triangleq \text{ReLU}(\text{conv}_X(W_X \mathbf{x}_{t-w+1}^t)) \cdot \Delta_t, \\
 \mathbf{b}_t &\triangleq \text{ReLU}(\text{conv}_B(W_B \mathbf{x}_{t-w+1}^t)), \\
 \mathbf{c}_t &\triangleq \text{ReLU}(\text{conv}_C(W_C \mathbf{x}_{t-w+1}^t)),
 \end{aligned}$$

$$: \mathbf{x} = (\mathbf{x}_t)_{t=1}^T \in \mathbb{R}^{d \times T} \mapsto \mathbf{o} = (\mathbf{o}_t)_{t=1}^T \in \mathbb{R}^{d \times T}$$

Mamba with Markov



Next-token prediction loss

$$L(\boldsymbol{\theta}) = -\frac{1}{T} \sum_{t \in [T]} \mathbb{E}_P \mathbb{E}_{x_1^{t+1} \sim P} [\log f_{\boldsymbol{\theta}}^{(x_{t+1})}(x_1^t)]$$

Cross-entropy loss between x_{t+1} and prediction probability $f_{\boldsymbol{\theta}}(x_1^t) = \mathbb{P}_{\boldsymbol{\theta}}(\cdot | x_1^t)$

Optimal estimator: Laplacian smoothing

$$\mathbb{P}_\beta(x_{t+1} = j \mid x_1^t) \triangleq \mathbb{E}_P[\mathbb{P}_\beta(x_{t+1} = j \mid x_1^t)] = \frac{n_j + \beta}{n + 2\beta}$$

Prediction probability



Add- β estimator

Optimal estimator: Laplacian smoothing

$$\mathbb{P}_\beta(x_{t+1} = j \mid x_1^t) \triangleq \mathbb{E}_P[\mathbb{P}_\beta(x_{t+1} = j \mid x_1^t)] = \frac{n_j + \beta}{n + 2\beta}$$

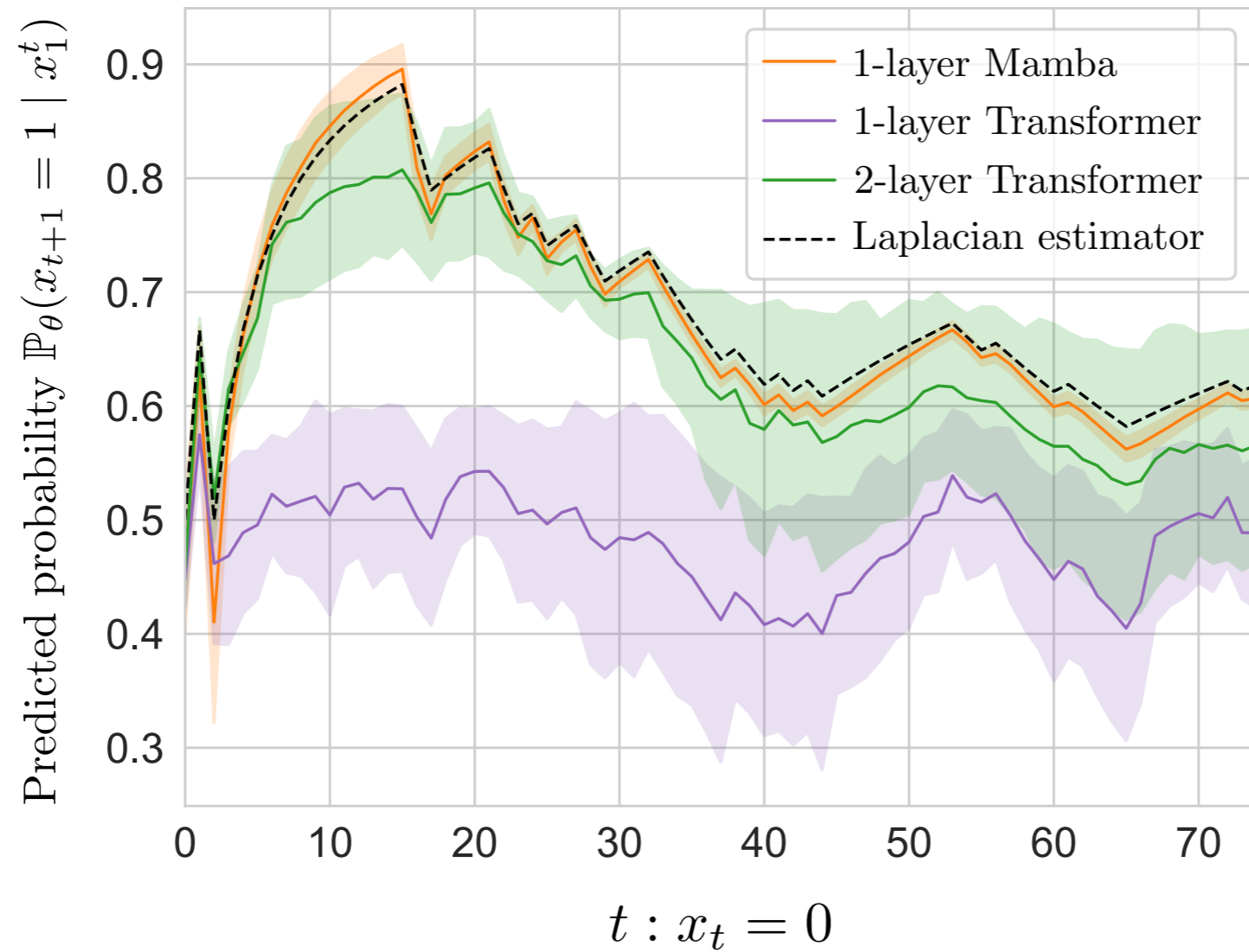
Prediction probability

Add- β estimator

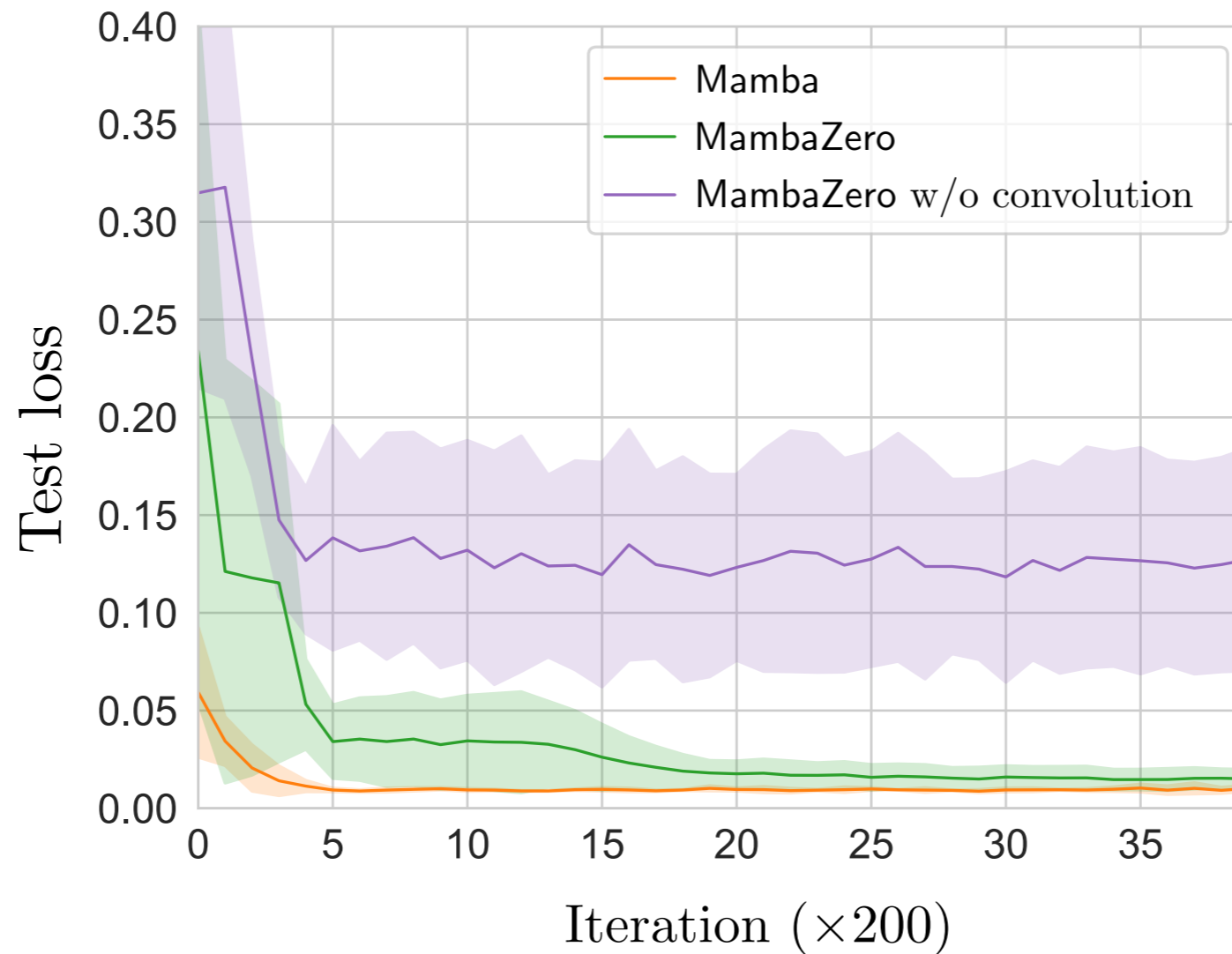
n_j : no. of times token j follows the context x_{t-k+1}^t
 n : no. of times context x_{t-k+1}^t occurs

Can Mamba learn Laplacian smoothing?

Can Mamba learn Laplacian smoothing?



Convolution is the key



**MambaZero =
Mamba with just
the convolution**

Main result

MambaZero realizes Laplacian smoothing for first-order

For first-order Markov processes on any finite state-space, there exists a realization of MambaZero model such that for all sequences $(x_t)_{t \geq 1}$ and all $t \geq 1$:

$$D_{\text{KL}} \left(\mathbb{P}_{\beta}(\cdot \mid x_1^t) \parallel \mathbb{P}_{\theta}(\cdot \mid x_1^t) \right) = 0.$$

Main result

MambaZero realizes Laplacian smoothing for first-order

For first-order Markov processes on any finite state-space, there exists a realization of MambaZero model such that for all sequences $(x_t)_{t \geq 1}$ and all $t \geq 1$:

$$D_{\text{KL}} \left(\mathbb{P}_{\beta}(\cdot \mid x_1^t) \parallel \mathbb{P}_{\theta}(\cdot \mid x_1^t) \right) = 0.$$

Empirically true for higher-orders too

Intuition

$$\begin{aligned} H_t &= a_t H_{t-1} + \tilde{\mathbf{x}}_t \mathbf{b}_t^\top \in \mathbb{R}^{ed \times N}, \\ \mathbf{y}_t &= H_t \mathbf{c}_t \in \mathbb{R}^{ed}, \\ \mathbf{z}_t &= \mathbf{y}_t \odot \text{ReLU}(W_z \mathbf{x}_t) \in \mathbb{R}^{ed}, \\ \mathbf{o}_t &= W_o \mathbf{z}_t \in \mathbb{R}^d, \end{aligned}$$



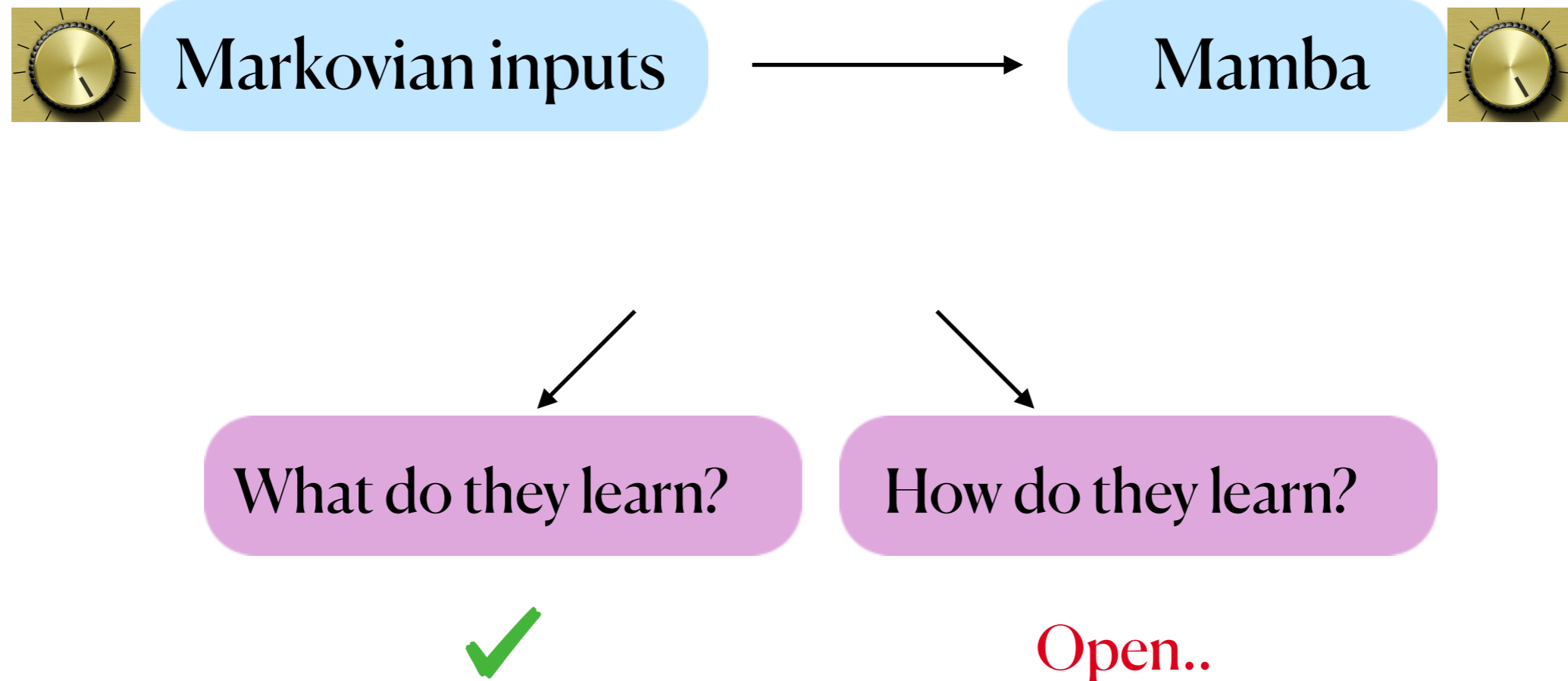
**Capture only context-relevant
counting**

$$\begin{aligned} a_t &\triangleq \exp(-a \cdot \Delta_t) \in (0, 1), \\ \Delta_t &\triangleq \text{softplus}(\langle \mathbf{w}_\Delta, \mathbf{x}_t \rangle + \delta) \in \mathbb{R}, \\ \tilde{\mathbf{x}}_t &\triangleq \text{ReLU}(\text{conv}_X(W_X \mathbf{x}_{t-w+1}^t)) \cdot \Delta_t, \\ \mathbf{b}_t &\triangleq \text{ReLU}(\text{conv}_B(W_B \mathbf{x}_{t-w+1}^t)), \\ \mathbf{c}_t &\triangleq \text{ReLU}(\text{conv}_C(W_C \mathbf{x}_{t-w+1}^t)), \end{aligned}$$

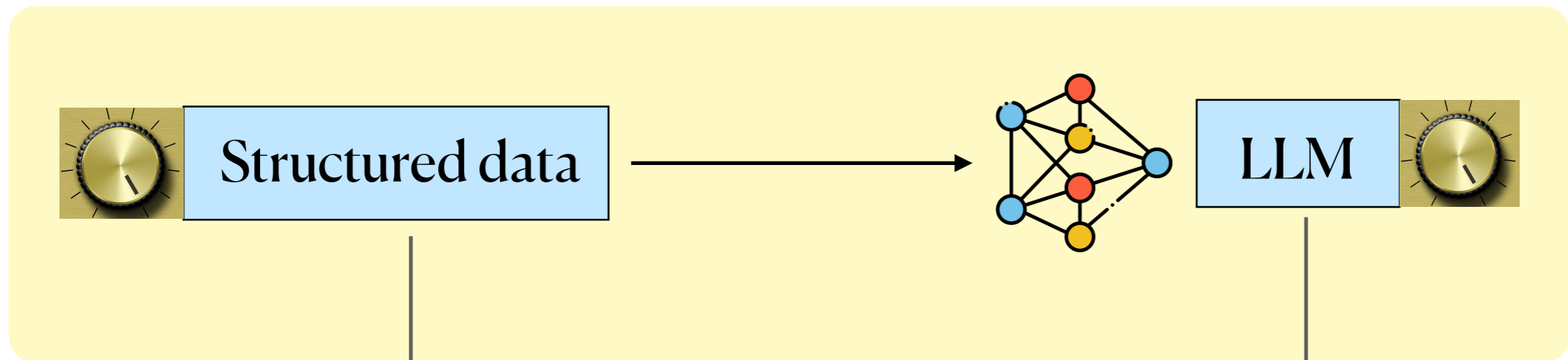


**Appropriate convolution kernel
and matrices**

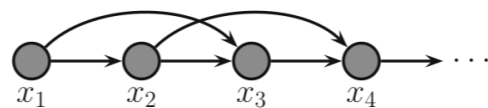
Mamba with Markov



Summary



Markov processes

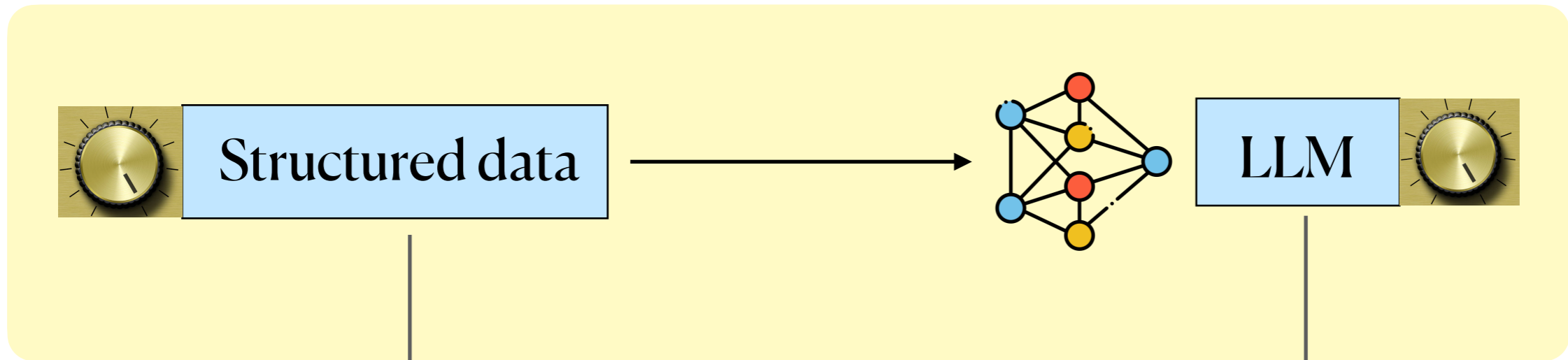


Transformers

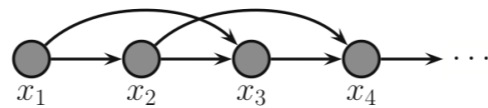
State-Space Models



Summary



Markov processes



Transformers



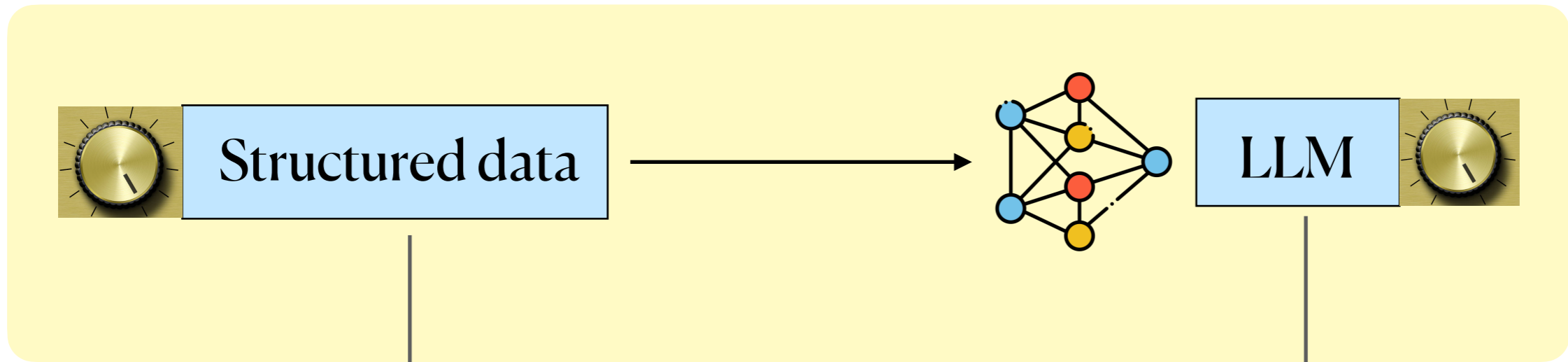
State-Space Models



Non-linearity

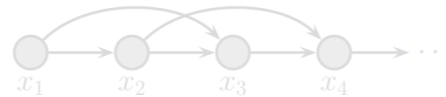
Convolution

Future directions

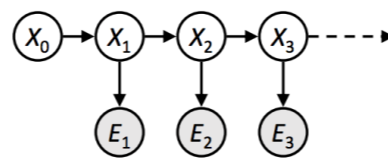


Multi-hop reasoning

baebc**a**bebdea.



Hidden-Markov-Models

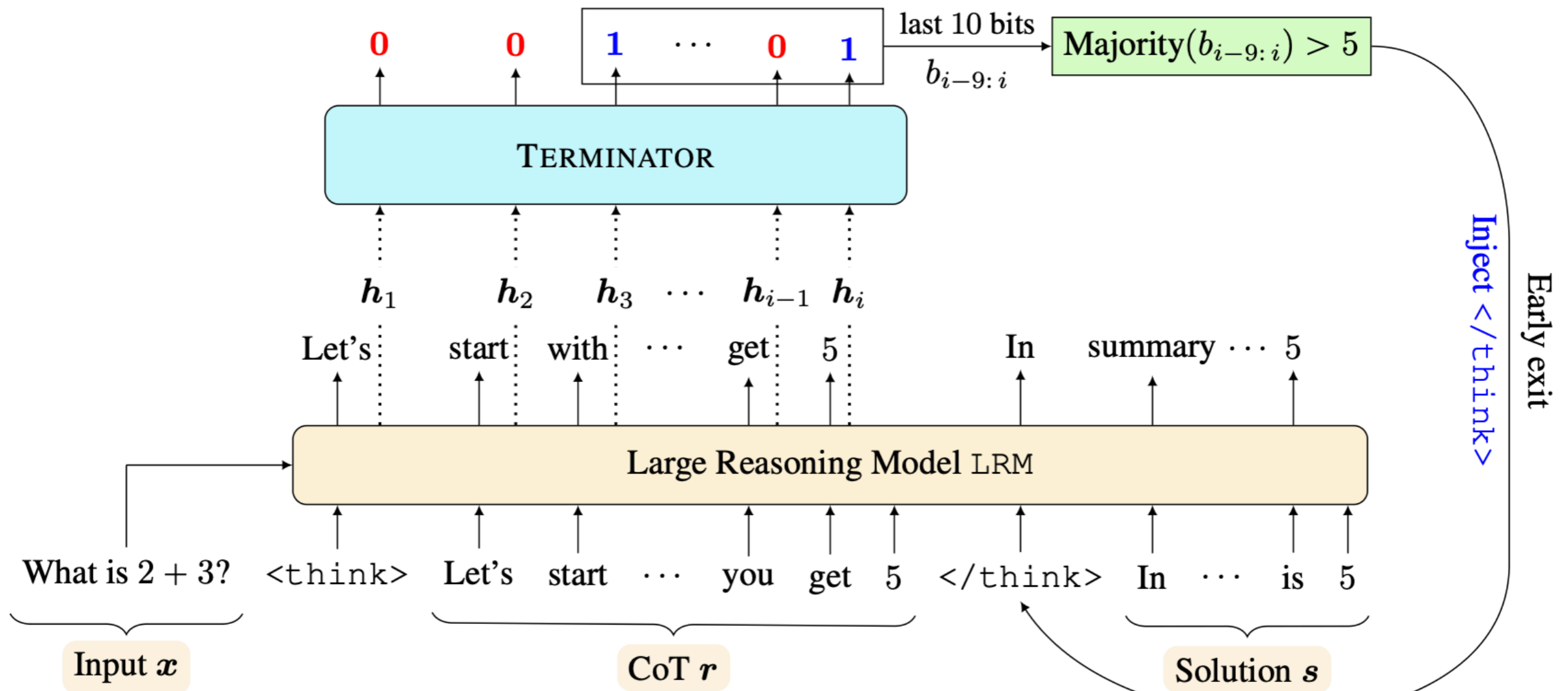


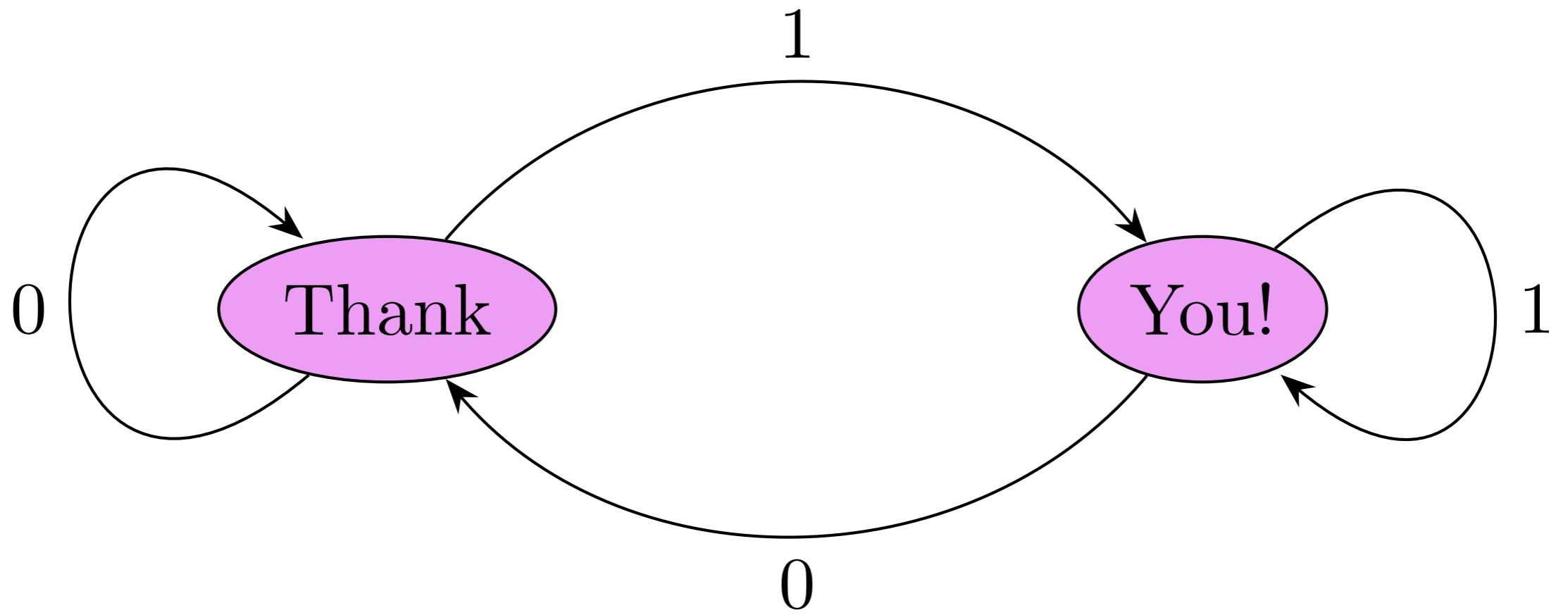
Diffusion Models



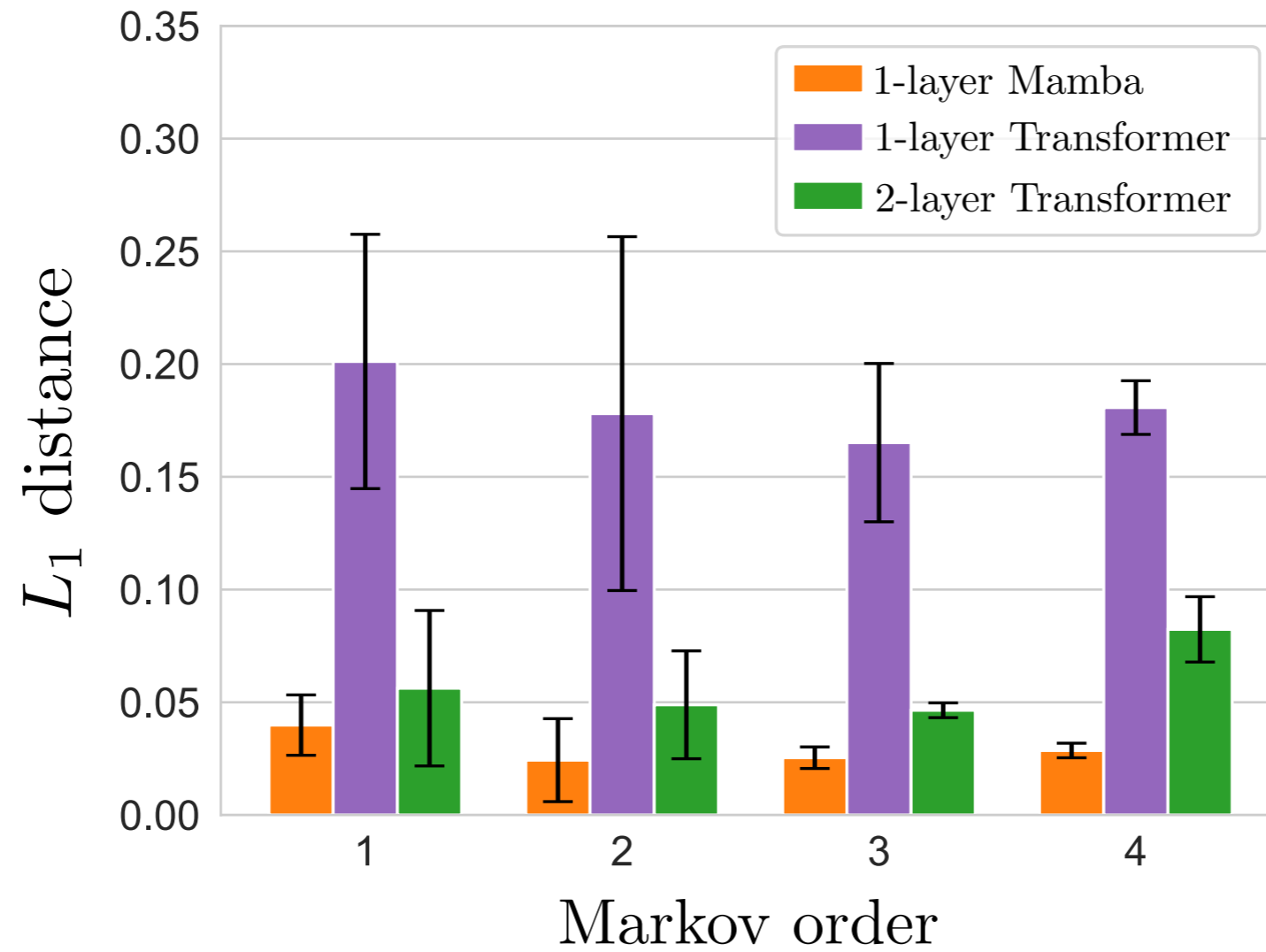
Efficient and hybrid architectures

Terminator: Efficient Reasoning





All Markov orders



Convolution window vs. Markov order

