

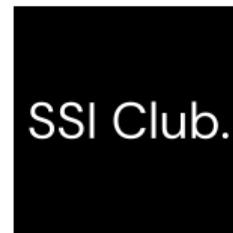
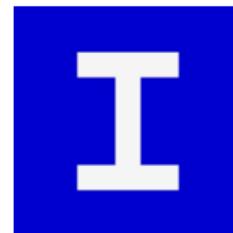
# Large Language models as Markov Chains

Oussama Zekri<sup>1</sup>

AI Paper Fest 2024, SSI Club & Arya.ai

<sup>1</sup> ENS Paris-Saclay and Imperial College London, [oussama.zekri@ens-paris-saclay.fr](mailto:oussama.zekri@ens-paris-saclay.fr)

November 14, 2024



# Co-authors



Ambroise Odonnat



Abdelhakim Benechehab



Linus Bleistein



Nicolas Boullé



Ievgen Redko

from *Huawei Noah's Ark Lab, Inria, Imperial College London.*

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Summary

① Introduction

Large  
Language  
models as  
Markov  
Chains

O. Zekri

② Large Language Models as Markov Chains

Introduction

③ Generalization bounds on pre-training

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

④ In-Context Learning and experiments

In-Context  
Learning and  
experiments

⑤ Take home message

Take home  
message

# Introduction

école  
normale  
supérieure  
paris-saclay

# Introduction

Large  
Language  
models as  
Markov  
Chains

O. Zekri

## Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Background on Autoregressive Language Modeling

**Goal:** Predict the next word based on previous ones.

Large  
Language  
models as  
Markov  
Chains

O. Zekri

## Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Background on Autoregressive Language Modeling

**Goal:** Predict the next word based on previous ones.

**Autoregressive Property:** Each word depends only on past words.



Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

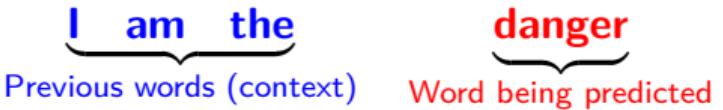
In-Context  
Learning and  
experiments

Take home  
message

# Background on Autoregressive Language Modeling

**Goal:** Predict the next word based on previous ones.

**Autoregressive Property:** Each word depends only on past words.

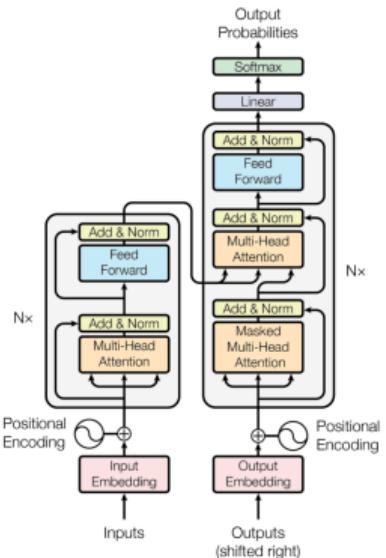


**Modelization:** Probability of a sequence  $(x_1, x_2, \dots, x_N)$ :

$$\begin{aligned}\mathbb{P}(x_1, x_2, \dots, x_N) &= \mathbb{P}(x_1)\mathbb{P}(x_2 \mid x_1) \cdots \mathbb{P}(x_N \mid x_1, x_2, \dots, x_{N-1}) \\ &= \prod_{n=1}^N \mathbb{P}(x_n \mid x_1, x_2, \dots, x_{n-1})\end{aligned}$$

# Background on Autoregressive Models

**Best models so far : Generative Transformers for Autoregressive Modeling  $f_{\Theta}^{T,K}$**



Large  
Language  
models as  
Markov  
Chains

O. Zekri

## Introduction

Large  
Language  
Models as  
Markov Chains

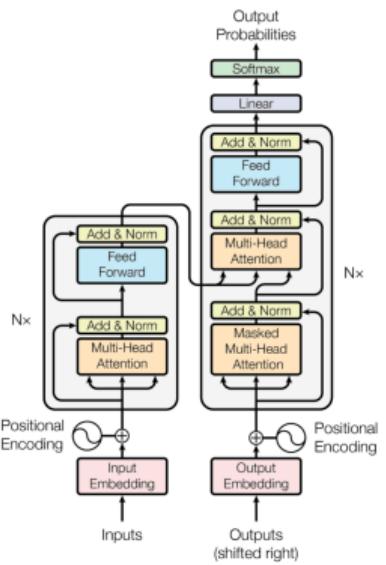
Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Background on Autoregressive Models

**Best models so far :** Generative Transformers for Autoregressive Modeling  $f_{\Theta}^{T,K}$

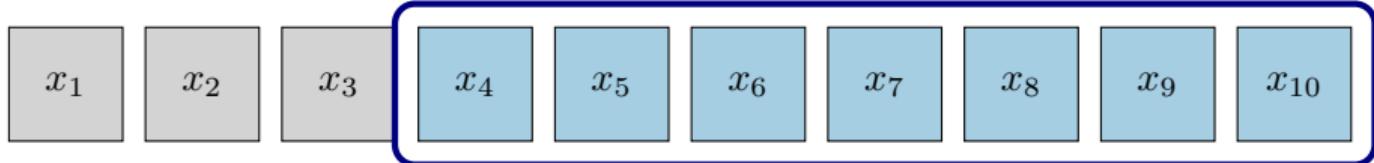
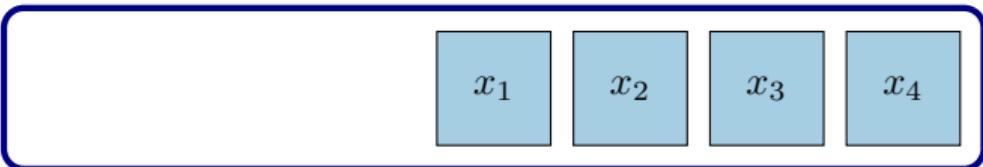


Vocabulary size  $T$ .  
Context window  $K$ .  
Parameter set  $\Theta$ .

**GPT-3** :  $T = 50257$ ,  $K = 2048$  and  
 $|\Theta| \sim 175B$

# Context Window $K$

Context Window  $K = 7$  in navy blue.



**Top.** A sequence of length  $N = 4$ .

**Bottom.** A sequence of length  $N = 10$ .

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Background on Markov Chains

$\Omega$  discrete finite state-space.  $\Omega = \{z_1, \dots, z_{|\Omega|}\}$ .

**Markov Chain:** A sequence of states where each state depends only on the previous one.

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Background on Markov Chains

$\Omega$  discrete finite state-space.  $\Omega = \{z_1, \dots, z_{|\Omega|}\}$ .

**Markov Chain:** A sequence of states where each state depends only on the previous one.

**Mathematical Formulation:** A process  $(Z_n)_{n \geq 0}$  supported on  $\Omega$  is a Markov chain if

$$\mathbb{P}(Z_{n+1} | Z_n, Z_{n-1}, \dots, Z_0) = \mathbb{P}(Z_{n+1} | Z_n)$$

This is called the *Markov property*.

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Are LLMs really Markov chains?

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

**State space:**  $\Omega = \{z_1, \dots, z_{|\Omega|}\}$

**Markov Chain:**  $\mathbb{P}(Z_{n+1} | Z_n, Z_{n-1}, \dots, Z_0) = \mathbb{P}(Z_{n+1} | Z_n)$

  
Previous words (context)      Word being predicted

# Are LLMs really Markov chains?

**State space:**  $\Omega = \{z_1, \dots, z_{|\Omega|}\}$

**Markov Chain:**  $\mathbb{P}(Z_{n+1} | Z_n, Z_{n-1}, \dots, Z_0) = \mathbb{P}(Z_{n+1} | Z_n)$

The diagram shows a sequence of words: "I am the danger". The first three words, "I am the", are underlined with a blue bracket and labeled "Previous words (context)". The last word, "danger", is underlined with a red bracket and labeled "Word being predicted".

- ✗ LLMs are clearly not Markov chains at the token level ( $|\Omega| = T$ ).

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Are LLMs really Markov chains?

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

**State space:**  $\Omega = \{z_1, \dots, z_{|\Omega|}\}$

**Markov Chain:**  $\mathbb{P}(Z_{n+1} | Z_n, Z_{n-1}, \dots, Z_0) = \mathbb{P}(Z_{n+1} | Z_n)$

  
Previous words (context)      Word being predicted

- ✗ LLMs are clearly not Markov chains at the token level ( $|\Omega| = T$ ).
- ✓ Considering the whole sequence, there is something to be done... ( $|\Omega| = ?$ ).

# Large Language Models as Markov Chains

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Large Language Models as Markov Chains

# Correct State Space

Vocabulary space  $\mathcal{V}$  of size  $T$ .

- ▶  $\Omega = \mathcal{V}_K^*$  is the set of all sequences consisting of elements from  $\mathcal{V}$  with up to  $K$  elements.

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

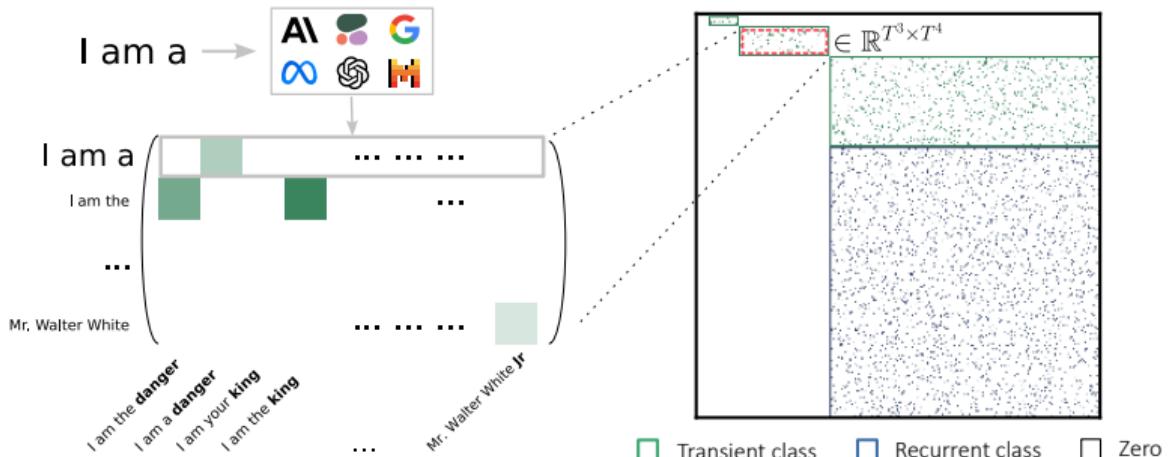
In-Context  
Learning and  
experiments

Take home  
message

# Correct State Space

Vocabulary space  $\mathcal{V}$  of size  $T$ .

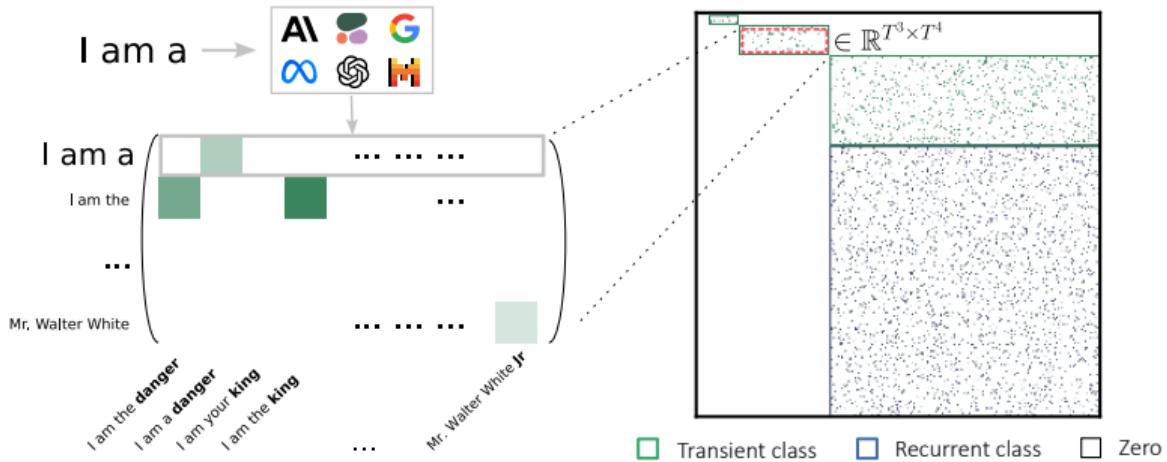
- $\Omega = \mathcal{V}_K^*$  is the set of all sequences consisting of elements from  $\mathcal{V}$  with up to  $K$  elements.



# Correct State Space

Vocabulary space  $\mathcal{V}$  of size  $T$ .

- $\Omega = \mathcal{V}_K^*$  is the set of all sequences consisting of elements from  $\mathcal{V}$  with up to  $K$  elements.



- $|\Omega| = T \left( \frac{T^K - 1}{T - 1} \right)$ . This is  $\sim 50257^{2048} = 1.123456 \times 10^{9628}$  for GPT-3.

# Formalization

Any autoregressive model  $f_{\Theta}^{T,K}$  can be equivalently represented by a **finite** Markov chain, with a sparse transition matrix  $\mathbf{Q}_f \in \mathbb{R}^{|\mathcal{V}_K^*| \times |\mathcal{V}_K^*|}$  defined as:

$$\forall v_i, v_j \in \mathcal{V}_K^*,$$

$$\mathbf{Q}_f(v_i, v_j) = \begin{cases} 0, & \text{if } \exists l \in \{1, \dots, |v_i| - 1\}, \text{s.t. } (v_i)_{l+1} \neq (v_j)_l, \\ \{f_{\Theta}^{T,K}(v_i)\}_j, & \text{otherwise.} \end{cases}$$

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Formalization

Any autoregressive model  $f_{\Theta}^{T,K}$  can be equivalently represented by a **finite** Markov chain, with a sparse transition matrix  $\mathbf{Q}_f \in \mathbb{R}^{|\mathcal{V}_K^*| \times |\mathcal{V}_K^*|}$  defined as:

$$\forall v_i, v_j \in \mathcal{V}_K^*,$$

$$\mathbf{Q}_f(v_i, v_j) = \begin{cases} 0, & \text{if } \exists l \in \{1, \dots, |v_i| - 1\}, \text{s.t. } (v_i)_{l+1} \neq (v_j)_l, \\ \{f_{\Theta}^{T,K}(v_i)\}_j, & \text{otherwise.} \end{cases}$$

The proportion of non-zero elements in  $\mathbf{Q}_f$  is  $\frac{T-1}{T^K - 1}$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Is this Markov Chain point of view is useful?

- ✗  $|\Omega| = T \left( \frac{T^K - 1}{T - 1} \right)$ , grows exponentially with  $K$ .  $|\Omega| \sim 10^{9628}$  for GPT-3.  $\mathbf{Q}_f$  cannot be stored.
- ~, Model weights, a few GPUs and a single forward pass are all you need to access the row you want in the matrix!
- ✓ Connection to the rich theory of **finite** Markov Chain.
- ✓ Gives an insight into the dynamics of LLMs.



Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Toy example

**Toy example "Baby"**  
LLM.

- ▶  $\mathcal{V} = \{0, 1\}$  (i.e.  
 $T = 2$ ),
- ▶  $K = 3$ ,
- ▶  $|\Theta| = 12688$  and
- ▶  $|\Omega| = 14$ .

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Stationary distribution

- ▶ A stationary distribution  $\pi$  represents the long-term behavior of a Markov chain.

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Stationary distribution

- ▶ A stationary distribution  $\pi$  represents the long-term behavior of a Markov chain.
- ▶ It verifies  $\pi \mathbf{Q}_f = \pi$  and each row of  $\mathbf{Q}_f^n$  tends to  $\pi$  when  $n \rightarrow \infty$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Stationary distribution

- ▶ A stationary distribution  $\pi$  represents the long-term behavior of a Markov chain.
- ▶ It verifies  $\pi \mathbf{Q}_f = \pi$  and each row of  $\mathbf{Q}_f^n$  tends to  $\pi$  when  $n \rightarrow \infty$ .
- ▶ A finite state unichain has a unique stationary distribution.

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

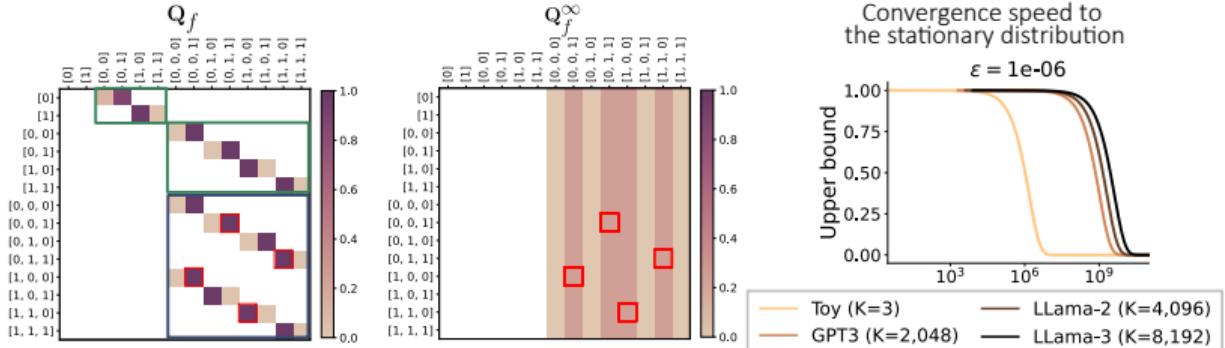
Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Stationary distribution

- ▶ A stationary distribution  $\pi$  represents the long-term behavior of a Markov chain.
- ▶ It verifies  $\pi \mathbf{Q}_f = \pi$  and each row of  $\mathbf{Q}_f^n$  tends to  $\pi$  when  $n \rightarrow \infty$ .
- ▶ A finite state unichain has a unique stationary distribution.



# Speed of convergence

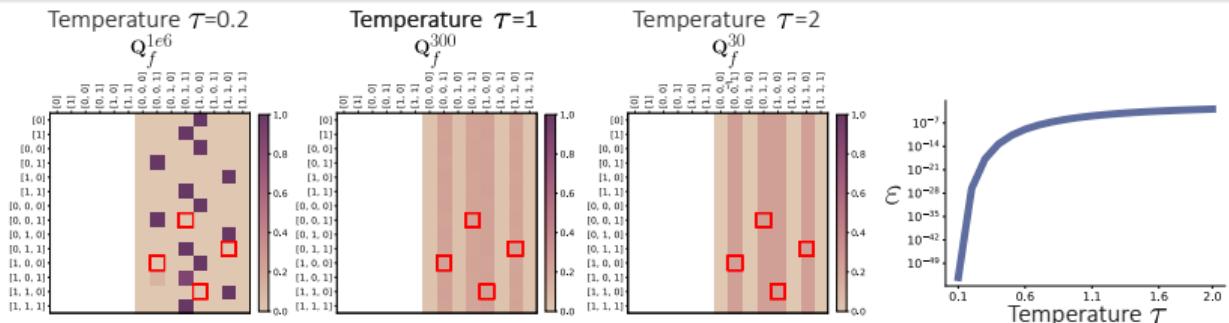
- ▶ Convergence speed to the stationary distribution.
- ▶ Impact of the temperature.

## Proposition

For all  $n \geq K$ ,

$$|(\mathbf{Q}_f^n)_{i,j} - (e\pi)_{i,j}| \leq (1 - 2\varepsilon)^{\lfloor \frac{n}{K} \rfloor - 1},$$

where  $\varepsilon = \min_{i,j \in \mathcal{R}^2} \{(\mathbf{Q}_f^K)_{i,j}\} > 0$ .



# Generalization bounds on pre-training

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Generalization bounds on pre-training

# Setup

**Question :** How far our matrix  $\mathbf{Q}_f$  is from the reference matrix  $\mathbf{Q}^*$ ?

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Setup

**Question :** How far our matrix  $\mathbf{Q}_f$  is from the reference matrix  $\mathbf{Q}^*$ ?

- ▶ For GPT-3, only  $5 \times 10^{11}$  training tokens, but  $T^{K+1} \approx 10^{9632}$  non zero elements in  $\mathbf{Q}_f$ .

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Setup

**Question :** How far our matrix  $\mathbf{Q}_f$  is from the reference matrix  $\mathbf{Q}^*$ ?

- ▶ For GPT-3, only  $5 \times 10^{11}$  training tokens, but  $T^{K+1} \approx 10^{9632}$  non zero elements in  $\mathbf{Q}_f$ .
- ▶ Generalization capacity of the model?

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Setup

**Question :** How far our matrix  $\mathbf{Q}_f$  is from the reference matrix  $\mathbf{Q}^*$ ?

- ▶ For GPT-3, only  $5 \times 10^{11}$  training tokens, but  $T^{K+1} \approx 10^{9632}$  non zero elements in  $\mathbf{Q}_f$ .
- ▶ Generalization capacity of the model?

## Generalization problem

$$\mathcal{R}(\Theta) := \mathbb{E}[\widehat{\mathcal{R}}(\Theta)], \quad \widehat{\mathcal{R}}(\Theta) := \frac{1}{N} \sum_{n=1}^N d_{\text{TV}}(\mathbf{Q}^*(\mathbf{S}_n, \cdot), \mathbf{Q}_f(\mathbf{S}_n, \cdot)), \quad (1)$$

The generalization problem consists of bounding the difference  $\mathcal{R}(\Theta) - \widehat{\mathcal{R}}(\Theta)$ .

# Pre-training generalization bound

## Assumptions

- ✓ Pre-training data  $S = (\mathbf{S}_1, \dots, \mathbf{S}_{N_{\text{train}}})$  a sequence of **dependent** random variables with a Marton coupling matrix  $\Gamma$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Pre-training generalization bound

## Assumptions

- ✓ Pre-training data  $S = (\mathbf{S}_1, \dots, \mathbf{S}_{N_{\text{train}}})$  a sequence of **dependent** random variables with a Marton coupling matrix  $\Gamma$ .
- ✓ Assumption only on the last transformer layer: Bounded unembedding matrix, i.e.  $\|\mathbf{W}_U^\top\|_{2,1} \leq B_U$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Pre-training generalization bound

## Assumptions

- ✓ Pre-training data  $S = (\mathbf{S}_1, \dots, \mathbf{S}_{N_{\text{train}}})$  a sequence of **dependent** random variables with a Marton coupling matrix  $\Gamma$ .
- ✓ Assumption only on the last transformer layer: Bounded unembedding matrix, i.e.  $\|\mathbf{W}_U^\top\|_{2,1} \leq B_U$ .

## Theorem

Let  $0 < \delta < 1$ , then with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\text{pre}}(\boldsymbol{\Theta}) \leq \widehat{\mathcal{R}}_{\text{pre}}(\boldsymbol{\Theta}) + \frac{\bar{B}}{\sqrt{N_{\text{train}}}} \sqrt{\log\left(\frac{2}{\delta}\right)},$$

where  $\bar{B} = 2\|\boldsymbol{\Gamma}\| \max\{\log(T) + 2B_U/\tau, \log(1/c_0)\}^{1/2}$ .

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Fine-grained bound

More fine-grained bound with additional assumption.

$$\widetilde{\mathcal{W}} = \left\{ \Theta \mid \forall \ell \in [L], \begin{array}{l} \|\mathbf{W}_V^{(\ell)}\|_\infty \leq B_V, \quad \|\mathbf{W}_O^{(\ell)}\|_\infty \leq B_O, \\ \|\mathbf{W}_1^{(\ell)}\|_\infty \leq B_1, \quad \|\mathbf{W}_2^{(\ell)}\|_\infty \leq B_2, \quad \|\mathbf{W}_U^\top\|_{2,1} \leq B_U \end{array} \right\}.$$

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Fine-grained bound

More fine-grained bound with additional assumption.

$$\widetilde{\mathcal{W}} = \left\{ \Theta \mid \forall \ell \in [L], \begin{array}{l} \|\mathbf{W}_V^{(\ell)}\|_\infty \leq B_V, \quad \|\mathbf{W}_O^{(\ell)}\|_\infty \leq B_O, \\ \|\mathbf{W}_1^{(\ell)}\|_\infty \leq B_1, \quad \|\mathbf{W}_2^{(\ell)}\|_\infty \leq B_2, \quad \|\mathbf{W}_U^\top\|_{2,1} \leq B_U \end{array} \right\}.$$

## Corollary

Let  $0 < \delta < 1$ , then with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\text{pre}}(\Theta) \leq \widehat{\mathcal{R}}_{\text{pre}}(\Theta) + \frac{\bar{B}}{\sqrt{N_{\text{train}}}} \sqrt{\log\left(\frac{2}{\delta}\right)},$$

where  $\bar{B} = 2\|\boldsymbol{\Gamma}\| \max\{\log(T) + 2(B_\Theta)^L/\tau, \log(1/c_0)\}^{1/2}$ , and  $B_\Theta = [(1 + rmB_1B_2)(1 + \frac{r^3}{H}B_O B_V)](B_{\text{tok}} B_U)^{1/L}$ .

Large Language models as Markov Chains  
O. Zekri

Introduction

Large Language Models as Markov Chains

Generalization bounds on pre-training

In-Context Learning and experiments

Take home message

# Sample complexity

**Question :** How much training data do I need for  $\mathbf{Q}_f$  to be close to  $\mathbf{Q}^*$ ?

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Sample complexity

**Question :** How much training data do I need for  $\mathbf{Q}_f$  to be close to  $\mathbf{Q}^*$ ?

- ✓ Number of sequences that an LLM requires such that  $\mathbf{Q}_f$  is  $\varepsilon$ -close to  $\mathbf{Q}^*$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Sample complexity

**Question :** How much training data do I need for  $\mathbf{Q}_f$  to be close to  $\mathbf{Q}^*$ ?

- ✓ Number of sequences that an LLM requires such that  $\mathbf{Q}_f$  is  $\varepsilon$ -close to  $\mathbf{Q}^*$ .
- ✓ Dependency on model parameters.

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Sample complexity

**Question :** How much training data do I need for  $\mathbf{Q}_f$  to be close to  $\mathbf{Q}^*$ ?

- ✓ Number of sequences that an LLM requires such that  $\mathbf{Q}_f$  is  $\varepsilon$ -close to  $\mathbf{Q}^*$ .
- ✓ Dependency on model parameters.

## Corollary

Let  $\delta \in [0, 1]$  and let  $\epsilon > 0$ . If  $N_{\text{train}} \geq N^* := \lceil \frac{4\bar{B}^2}{\epsilon^2} \log\left(\frac{2}{\delta}\right) \rceil$  and if we assume a perfect pre-training error for  $f_\Theta$ , then we have with probability at least  $1 - \delta$ ,

$$\mathbb{E}_{\mathbf{S} \sim \mathbb{P}_{\mathcal{L}}} \|\mathbf{Q}^*(\mathbf{S}, \cdot) - \mathbf{Q}_f(\mathbf{S}, \cdot)\|_1 \leq \epsilon.$$

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# In-Context Learning and experiments

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# In Context Learning

**In-Context Learning** : Model's ability to learn and adapt to patterns, without updating its internal parameters.

## Demonstrations

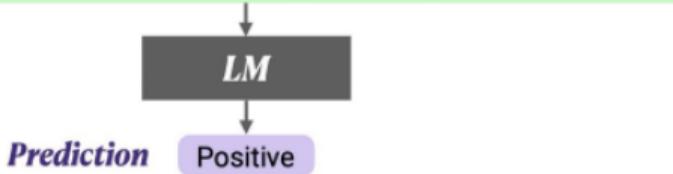
Circulation revenue has increased by 5% in Finland. \n Positive

Panostaja did not disclose the purchase price. \n Neutral

Paying off the national debt will be extremely painful. \n Negative

The acquisition will have an immediate positive impact. \n \_\_\_\_\_

*Test input*



# In Context Learning

**In-Context Learning :** Model's ability to learn and adapt to patterns, without updating its internal parameters.

## Demonstrations

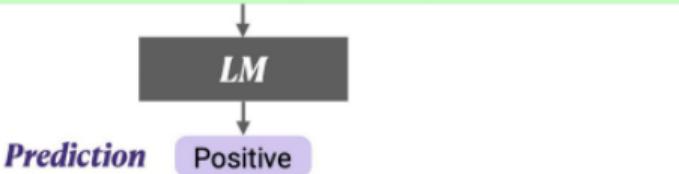
Circulation revenue has increased by 5% in Finland. \n Positive

Panostaja did not disclose the purchase price. \n Neutral

Paying off the national debt will be extremely painful. \n Negative

The acquisition will have an immediate positive impact. \n \_\_\_\_\_

## Test input



- ▶ Not as costly as model pre-training.

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# In Context Learning

**In-Context Learning :** Model's ability to learn and adapt to patterns, without updating its internal parameters.

## Demonstrations

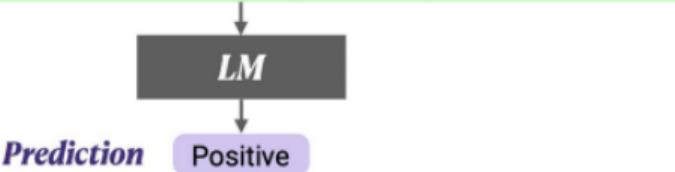
Circulation revenue has increased by 5% in Finland. \n Positive

Panostaja did not disclose the purchase price. \n Neutral

Paying off the national debt will be extremely painful. \n Negative

The acquisition will have an immediate positive impact. \n \_\_\_\_\_

## Test input



- ▶ Not as costly as model pre-training.
- ▶ Access to a ground truth matrix  $\mathbf{Q}^*$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# In Context Learning of Markov chains

## Setup

- ▶ A  $d$ -state Markov chain  $X = (\mathbf{X}_1, \dots, \mathbf{X}_{N_{\text{icl}}})$ , and the sequence of its first  $n$  terms is denoted by  $S = (\mathbf{S}_1, \dots, \mathbf{S}_n)$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# In Context Learning of Markov chains

## Setup

- ▶ A  $d$ -state Markov chain  $X = (\mathbf{X}_1, \dots, \mathbf{X}_{N_{\text{icl}}})$ , and the sequence of its first  $n$  terms is denoted by  $S = (\mathbf{S}_1, \dots, \mathbf{S}_n)$ .
- ▶ Mixing time of  $S$ , denoted as  $t_{\text{mix}}(\varepsilon)$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# In Context Learning of Markov chains

## Setup

- ▶ A  $d$ -state Markov chain  $X = (\mathbf{X}_1, \dots, \mathbf{X}_{N_{\text{icl}}})$ , and the sequence of its first  $n$  terms is denoted by  $S = (\mathbf{S}_1, \dots, \mathbf{S}_n)$ .
- ▶ Mixing time of  $S$ , denoted as  $t_{\text{mix}}(\varepsilon)$ .
- ▶ Almost distance  $K(\Theta_1, \Theta_2) := \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\mathbf{S}_n} [d_{\text{TV}}(\mathbb{P}_{\Theta_1}(\cdot | \mathbf{S}_n), \mathbb{P}_{\Theta_2}(\cdot | \mathbf{S}_n))]$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# In Context Learning of Markov chains

## Setup

- ▶ A  $d$ -state Markov chain  $X = (\mathbf{X}_1, \dots, \mathbf{X}_{N_{\text{icl}}})$ , and the sequence of its first  $n$  terms is denoted by  $S = (\mathbf{S}_1, \dots, \mathbf{S}_n)$ .
- ▶ Mixing time of  $S$ , denoted as  $t_{\text{mix}}(\varepsilon)$ .
- ▶ Almost distance  $K(\Theta_1, \Theta_2) := \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\mathbf{S}_n} [d_{\text{TV}}(\mathbb{P}_{\Theta_1}(\cdot | \mathbf{S}_n), \mathbb{P}_{\Theta_2}(\cdot | \mathbf{S}_n))]$ .

## Theorem

Let  $\delta > 0$ . Then, with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\text{icl}}(\Theta) \leq \inf_{\vartheta \in \mathcal{W}_{\text{mc}}} \{\widehat{\mathcal{R}}_{\text{icl}}(\vartheta) + \mathcal{K}(\vartheta, \Theta)\} + \bar{B} \sqrt{\frac{t_{\min}}{N_{\text{icl}}}} \sqrt{\log\left(\frac{2}{\delta}\right)}, \quad (2)$$

where  $\bar{B} = 2 \max\{\log(d) + 2B_U/\tau, \log(1/p_{\min})\}^{1/2}$ .

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

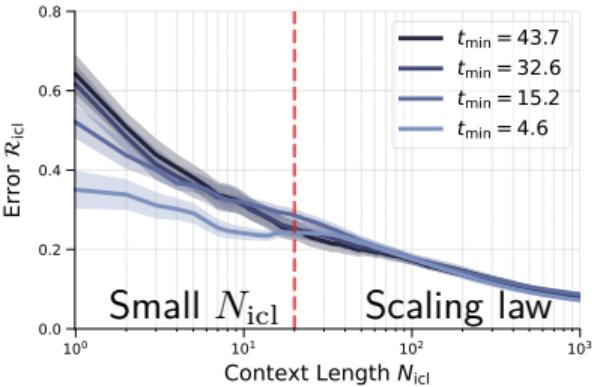
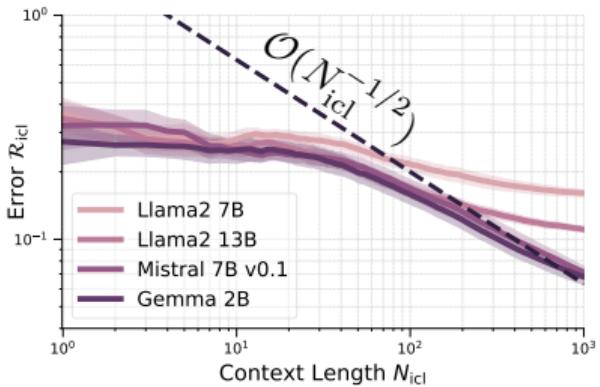
Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

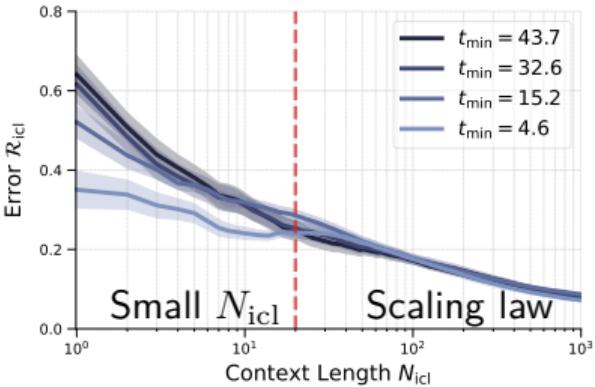
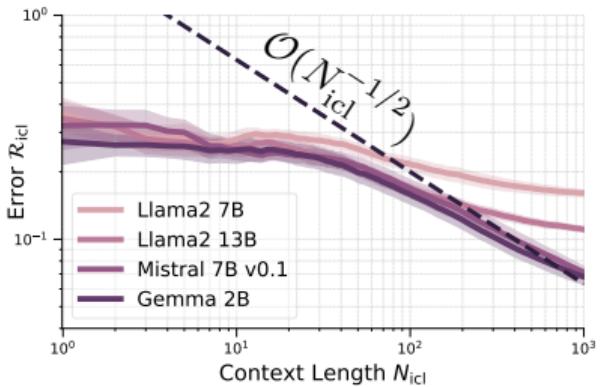
Take home  
message

# Experiments : In-Context Scaling Laws



**Figure: In-context scaling laws.** We plot the risk  $\mathcal{R}_{\text{icl}}$  as functions of  $N_{\text{icl}}$ , with 95% confidence intervals.

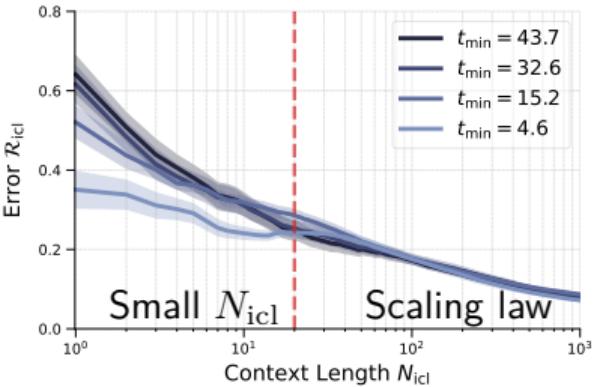
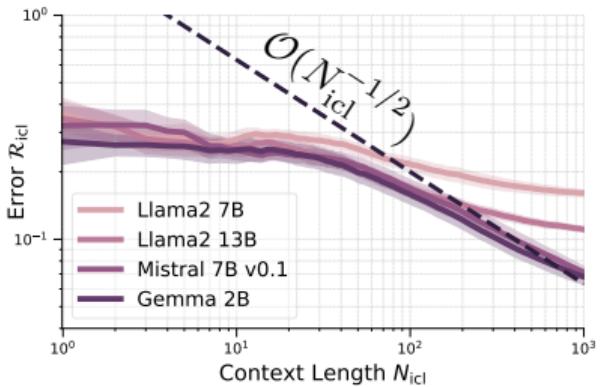
# Experiments : In-Context Scaling Laws



**Figure: In-context scaling laws.** We plot the risk  $\mathcal{R}_{\text{icl}}$  as functions of  $N_{\text{icl}}$ , with 95% confidence intervals.

- ✓ Randomly generated datas : not seen during training.

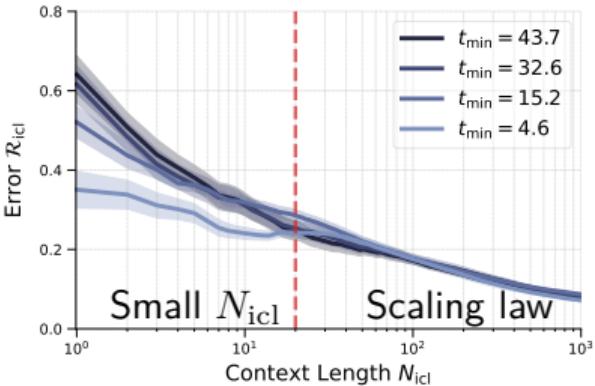
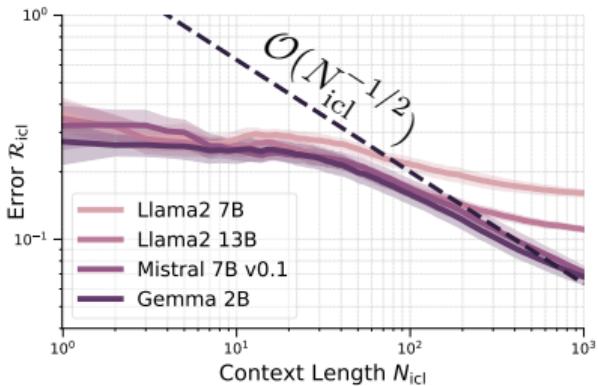
# Experiments : In-Context Scaling Laws



**Figure: In-context scaling laws.** We plot the risk  $\mathcal{R}_{\text{icl}}$  as functions of  $N_{\text{icl}}$ , with 95% confidence intervals.

- ✓ Randomly generated datas : not seen during training.
- ✓  $N_{\text{icl}}$  and  $t_{\min}$  dependence in line with theoretical result.

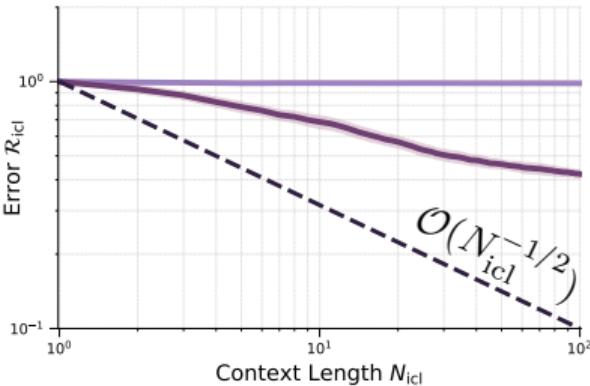
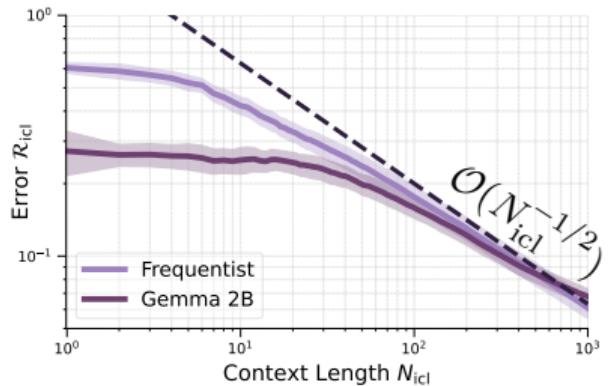
# Experiments : In-Context Scaling Laws



**Figure: In-context scaling laws.** We plot the risk  $\mathcal{R}_{\text{icl}}$  as functions of  $N_{\text{icl}}$ , with 95% confidence intervals.

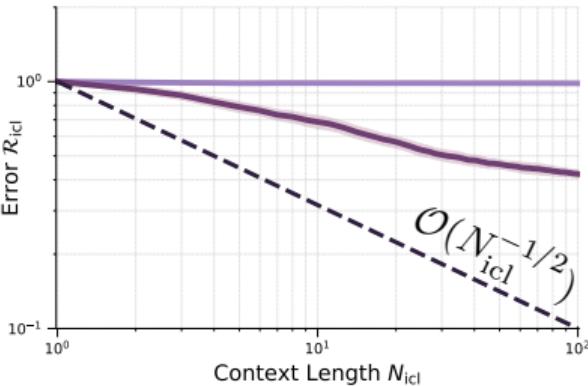
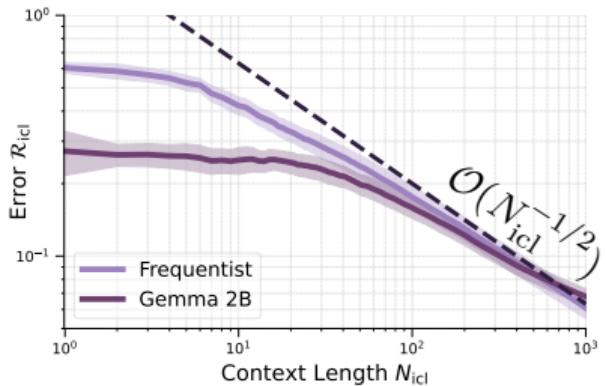
- ✓ Randomly generated datas : not seen during training.
- ✓  $N_{\text{icl}}$  and  $t_{\min}$  dependence in line with theoretical result.
- ✓ Most recent models stay much closer to the theoretical result.

# Experiments : Number of states



**Figure: Impact of the number of states.** **Left.** Random 3-state Markov chain. **Right.** Brownian motion discretized as a 700-state Markov chain.

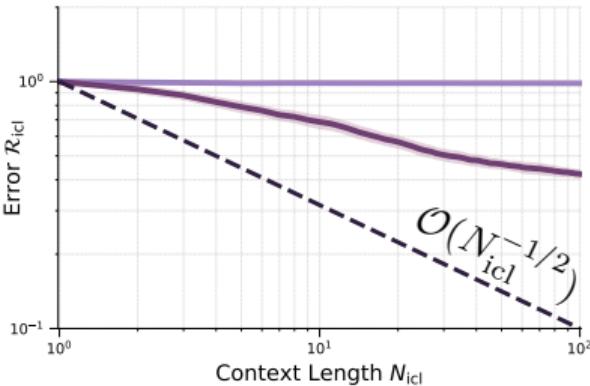
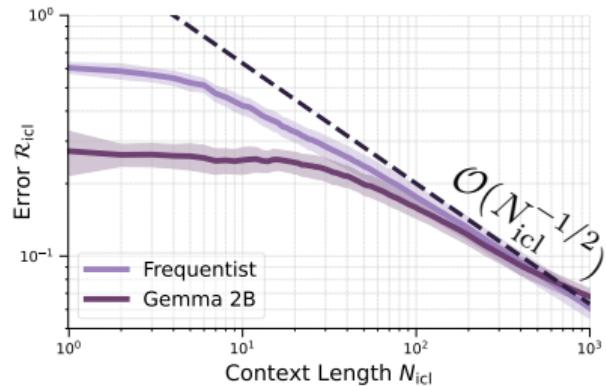
# Experiments : Number of states



**Figure: Impact of the number of states.** **Left.** Random 3-state Markov chain. **Right.** Brownian motion discretized as a 700-state Markov chain.

- ✓ Frequentist (baseline method)  $\mathcal{O}(\sqrt{d/N_{\text{icl}}})$  vs. LLM's bound  $\mathcal{O}(\sqrt{\log(d)/N_{\text{icl}}})$ .

# Experiments : Number of states



**Figure: Impact of the number of states.** **Left.** Random 3-state Markov chain. **Right.** Brownian motion discretized as a 700-state Markov chain.

- ✓ Frequentist (baseline method)  $\mathcal{O}(\sqrt{d/N_{\text{icl}}})$  vs. LLM's bound  $\mathcal{O}(\sqrt{\log(d)/N_{\text{icl}}})$ .
- ✓ As  $d$  grows, frequentist method struggles.

# Take home message

Large  
Language  
models as  
Markov  
Chains

O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Take home message

- ✓ Explicit characterization of the inference mechanism in LLMs through an equivalent finite-state Markov chain.

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Take home message

- ✓ Explicit characterization of the inference mechanism in LLMs through an equivalent finite-state Markov chain.
- ✓ Existence and uniqueness of a stationary distribution. Generalization bounds on pre-training and in-context learning (ICL) phases.

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Take home message

- ✓ Explicit characterization of the inference mechanism in LLMs through an equivalent finite-state Markov chain.
- ✓ Existence and uniqueness of a stationary distribution. Generalization bounds on pre-training and in-context learning (ICL) phases.
- ✓ Experiments validate our theory with Llama2 7B & 13B, Gemma 2B, Mistral 7B, and even Llama 3.2.

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message

# Thank you for your attention !

You can follow me on social networks!

Large  
Language  
models as  
Markov  
Chains  
O. Zekri

Introduction

Large  
Language  
Models as  
Markov Chains

Generalization  
bounds on  
pre-training

In-Context  
Learning and  
experiments

Take home  
message



@oussamazekri\_



My website : [www.oussamazekri.fr](http://www.oussamazekri.fr)