# LARGE LANGUAGE MODELS AS MARKOV CHAINS

### Anonymous authors

000

001 002 003

004

005 006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

023

Paper under double-blind review

## ABSTRACT

Large language models (LLMs) have proven to be remarkably efficient, both across a wide range of natural language processing tasks and well beyond them. However, a comprehensive theoretical analysis of the origins of their impressive performance remains elusive. In this paper, we approach this challenging task by drawing an equivalence between generic autoregressive language models with vocabulary of size T and context window of size K and Markov chains defined on a finite state space of size  $\mathcal{O}(T^K)$ . We derive several surprising findings related to the existence of a stationary distribution of Markov chains that capture the inference power of LLMs, their speed of convergence to it, and the influence of the temperature on the latter. We then prove pre-training and in-context generalization bounds and show how the drawn equivalence allows us to enrich their interpretation. Finally, we illustrate our theoretical guarantees with experiments on several recent LLMs to highlight how they capture the behavior observed in practice.

## 1 INTRODUCTION

The fields of machine learning and artificial intelligence have recently seen significant progress with the introduction of large language models (LLMs) (Brown et al., 2020; Touvron et al., 2023a), built on the transformer architecture (Vaswani et al., 2017). These models, trained on vast amounts of data, have been applied in many natural language processing tasks, including machine translation (Brown et al., 2020), text generation, question answering (Roberts et al., 2020), and sentiment analysis (Zhang et al., 2023a). Although successful in practice, the origins of the impressive performance of LLMs remain elusive, as there is no widely accepted agreement in the scientific community on how they achieve remarkable reasoning capabilities that go far beyond their training data (Brown et al., 2020).

This work takes a step towards bridging the knowledge gap mentioned above by providing an *explicit* characterization of the LLM's inference capabilities. For this, we adopt an intuitive, yet overlooked, approach that interprets LLMs as Markov chains operating on a finite state space of sequences and tokens (see Fig. 1). A key insight is that despite the seeming infinity of LLMs generating capacity, they have a limited vocabulary and context window making all their possible input and output sequences enumerable. We show that despite the prohibitively large size of the latter set, it exhibits a structure that makes it amenable to theoretical analysis. We further generalize recent theoretical advances on the generalization of LLMs and leverage our proposed point of view to provide a more insightful interpretation of them.

042

Markov chains and Large Language Models. While none of the prior works considered the 043 equivalence between LLMs and Markov chains presented in this work<sup>1</sup>, some used the Markovian gen-044 erative process to better understand the intrinsic capabilities of the transformer architecture. Makkuva et al. (2024) assume that the input data is generated by an unknown high-order Markov chain to 046 analyze the learning dynamics of the self-attention mechanism in a single-layer single-head trans-047 former. Similarly, Edelman et al. (2024) study in-context learning of a transformer model trained on 048 samples drawn from a bi-gram Markov chain. Ildiz et al. (2024) establish an equivalence between context-conditioned Markov chains and the self-attention mechanism in transformers to show that self-attention weights can be learned, under certain conditions, by prompting the model. In contrast 051 to this prior work, we seek to model any transformer-based LLM as a Markov chain. Hence, our

<sup>&</sup>lt;sup>1</sup>We note, however, a recent blog post (Nardo, 2023) discussing a similar idea on a high level without analyzing it at the level of details of this work.



Figure 1: LLM as a Markov chain. A large language model with vocabulary size T and context window K is equivalent to a Markov chain with a sparse and block-structured transition matrix of size  $\sum_{i \le K} T^i \sim \mathcal{O}(T^K)$ . The latter captures all possible outputs of a given LLM for all possible input sequences allowed by its vocabulary and context window.

analysis provides insights into LLMs beyond understanding the self-attention mechanism in simplified transformers.

071 **In-Context Learning (ICL).** ICL is the ability of LLMs to adapt their predictions during inference 072 by leveraging examples or prompts directly without updating their parameters. Xie et al. (2022) 073 provide a theoretical guarantee for ICL by showing its equivalence to implicit Bayesian inference, 074 while Jeon et al. (2024) study ICL in the Bayesian setup by adopting a point of view that relates 075 it to meta-learning from non-independent and identically distributed (i.i.d.) data. Their guarantees introduce a desirable dependence on the length of the input sequence and the number of sequences 076 seen, yet their data generative process assumes that each sequence is outputted by a transformer model 077 drawn from a prior distribution. A follow-up work by Wies et al. (2024) generalizes this analysis within a broader Probably Approximately Correct (PAC) framework and sheds light on the few-shot 079 nature of ICL. Li et al. (2023) studied the generalization error of trained ICL transformers from a stability viewpoint and provided generalization guarantees for temporally dependent prompts that can 081 be seen as Markov chains of different orders. The closest work on ICL to our analysis is (Zhang et al., 082 2023b) upon which we improve by relaxing many assumptions and, most importantly, by providing 083 stronger interpretations and experimental observations. 084

085 Generalization bounds for LLMs. Deriving generalization bounds for neural networks is inherently difficult due to the complexity of the operations performed by the model. This can require 087 expressing it as a continuous process (Marion, 2023). For LLMs, an avenue to obtain such bounds is 880 to rely on the PAC-Bayes framework. Lotfi et al. (2023; 2024) leverage compression techniques in combination with PAC-Bayes type bounds to obtain tight generalization bounds both at the document 089 and token level. These works are connected to the existing literature on compressibility and the 090 intrinsic dimension of neural networks (Aghajanyan et al., 2020; Yaras et al., 2024) and typically 091 focus on fine-tuning LLMs with LoRA (Hu et al., 2022) inspired adapters. The work of Zhang et al. 092 (2023b) considers the Bayesian framework to derive generalization bounds for pre-training where the 093 data is assumed to be a Markov chain. Our theoretical analysis of pre-training is done without relying 094 on Bayesian modeling and under realistic data generating assumption covering many common types 095 of data used to train LLMs.

096 098

101

054

055

056

059 060

061

062

063

064

065

066 067 068

069

**Summary of our contributions.** Our main contributions are summarized as follows.

- 099 1) We provide an explicit characterization of LLM's inference mechanism by showing its equivalence 100 to a finite-state Markov chain. We analyze the transition matrix of the latter and prove the existence and uniqueness of its stationary distribution. We give a rate of convergence to this distribution that depends on the vocabulary and context window sizes, and the model's temperature. 102
- 103 2) By leveraging concentration inequalities for *dependent* random variables, we obtain generalization 104 bounds for LLMs in both pretraining and in-context inference. Our bounds are proved under minimal assumptions on the model and data and depend on the model's depth, dictionary, and 105 dataset sizes, as well as the intrinsic properties of the temporally-dependent sequences it was trained on. We highlight the insights that stem from these bounds by relating them to the minimax 107 bounds of Markov chain learning.

3) We experimentally show that the most recent LLMs dating from 2023-2024 obey the in-context scaling laws predicted by our theoretical results. One highlight is that LLMs are better Markov chains learners than the minimax optimal frequentist approach (Wolfer & Kontorovich, 2019).

We underline that, in this work, the term LLM refers to a deep transformer-based model trained on non-iid data whose inference is based on the next-token prediction principle in an autoregressive fashion. The latter implies that such a model transitions between a sequence of tokens to a sequence of tokens. Hence, in Section 3, the **Markov chain** formalization transitions between states that are **sequences of tokens** (instead of single tokens). The vast majority of existing LLMs fall into our definition suggesting that our results apply to them.

Organization of the paper. Section 2 provides background material on autoregressive models and Markov chains. We formalize an equivalence between these two models in Section 3 and illustrate it on a toy example. In Section 4, we derive generalization bounds for LLMs trained on *non-iid* data and prompted on Markov chains. Our results are empirically verified in Section 5.

122 123 124

125

126

127

133

138 139

145

156 157

158

108

110

111

## 2 BACKGROUND KNOWLEDGE

We recall some elementary facts about Markov chains (Paulin, 2015; Roberts & Rosenthal, 2004) and LLMs. More notations and background materials are available in Appendices A to C.

128 129 130 131 132 Markov chains. Let  $\Omega$  be a discrete finite set of size  $|\Omega|$ . A discrete-time, time-homogeneous Markov chain MC( $\Omega$ ,  $\mathbf{Q}$ ) defined on a state space  $\Omega = \{x_i\}_{i=1}^{|\Omega|}$  with transition matrix  $\mathbf{Q} \in \mathbb{R}^{|\Omega| \times |\Omega|}$ with entries  $\mathbf{Q}_{ij} = \mathbf{Q}(x_i, x_j) \in [0, 1]$  is a sequence of random variables  $(\mathbf{X}_1, \mathbf{X}_2, \ldots)$  taking values in  $\Omega$  such that for any  $n \in \mathbb{N}$  and  $(x_1, \ldots, x_{n+1}) \in \Omega^{n+1}$ , we have

$$\mathbb{P}(\mathbf{X}_{n+1} = x_{n+1} \mid \mathbf{X}_n = x_n, \dots, \mathbf{X}_1 = x_1) = \mathbb{P}(\mathbf{X}_{n+1} = x_{n+1} \mid \mathbf{X}_n = x_n) =: \mathbf{Q}(x_n, x_{n+1}).$$

134 A distribution  $\pi$  on  $\Omega$  is said to be a stationary distribution if  $\mathbf{Q}\pi = \pi$ . Under mild conditions 135 on  $\mathbf{Q}$ ,  $\mathbf{MC}(\Omega, \mathbf{Q})$  has a unique stationary distribution to which it converges, i.e., for any  $x \in \Omega$ , 136  $\lim_{n\to\infty} d_{\mathrm{TV}}(\mathbf{Q}^n(x,\cdot),\pi) = 0$ , where  $\mathbf{Q}^n(x,\cdot)$  denotes the probability of  $\mathbf{X}_n$  conditioned on 137  $\mathbf{X}_1 = x$  and the total variation between two distributions  $\mathbb{P}$  and  $\mathbb{Q}$ , defined on  $(\Omega, \mathcal{F})$ , is

$$d_{\mathrm{TV}}(\mathbb{P}, \mathbb{Q}) \coloneqq \sup_{A \in \mathcal{F}} |\mathbb{P}(A) - \mathbb{Q}(A)|$$

140 We recall that the mixing time  $t_{mix}(\varepsilon)$  of a Markov chain is the minimal time needed to be  $\varepsilon$ -close 141 to its stationary distribution (see Definition C.8). Intuitively, a Markov chain mixes slowly when 142 it remains close to the initial state after a given number of steps and doesn't explore its state space. 143 A Markov chain that exhibits a fast mixing time on the contrary quickly forgets its initial state and 144 transitions more easily to a wider set of states.

**Large language models.** Let  $\mathcal{V}$  denote a dictionary of size T used to encode an arbitrary sequence 146 into a sequence of predefined tokens belonging to  $\mathcal{V}$ . We assume that our model admits a maximum of 147 K tokens as input, referred to as the *context window* of the model. The domain of the autoregressive 148 LLM is the set of all sequences consisting of elements from  $\mathcal{V}$  with up to K elements. We denote this 149 by  $\mathcal{V}_{K}^{*}$ , which represents a restriction of Kleene closure of  $\mathcal{V}$ , i.e.,  $\mathcal{V}_{K}^{*} := \{v \in \mathcal{V}^{*}, |v| \leq K\}$  with |v| the length of v. We define an LLM with trainable parameters  $\Theta$  as a function  $f_{\Theta}^{T,K} : \mathcal{V}_{K}^{*} \to \Delta(\mathcal{V})$ , 150 151 where  $\Delta(\mathcal{V})$  is the probability simplex over  $\mathcal{V}$ , that given a sequence of tokens v outputs a probability 152 distribution over the whole state space indicating the likelihood for each of its elements to appear 153 after v (see Appendix B for more details). We consider a setting where the learner's objective is 154 to approximate the probabilities of sequences over an input vocabulary given by some reference distribution  $\mathbb{P}_{\mathcal{L}}: \mathcal{P}(\mathcal{V}_K^*) \to [0,1]^2$ . 155

3 LARGE LANGUAGE MODELS AS MARKOV CHAINS

We formally define a Markov chain that explicitly captures the full inference capacity of a given LLM  $f_{\Theta}$ . We build upon a high-level idea that associates a tokenized input prompt with a state  $v_i$ , from

 ${}^{2}\mathcal{P}(\mathcal{V}_{K}^{*})$  denotes the powerset of  $\mathcal{V}_{K}^{*}$ .

162 which we transition to a new state  $v_i = [v_i, v]$  by concatenating the token v predicted by an LLM 163 to it. We then provide a theoretical characterization of this Markov chain highlighting its intriguing properties and asymptotic behaviour.

#### 3.1 MARKOV CHAIN FORMALIZATION

We begin by defining the transition matrix associated with a large language model  $f_{\Theta}^{T,K}$ .

**Proposition 3.1.** Any large language model  $f_{\Theta}^{T,K}$  can be equivalently represented by a Markov chain  $MC(\mathcal{V}_{K}^{*}, \mathbf{Q}_{f})$ , with a sparse transition matrix  $\mathbf{Q}_{f} \in \mathbb{R}^{|\mathcal{V}_{K}^{*}| \times |\mathcal{V}_{K}^{*}|}$  defined as:

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

where  $|\mathcal{V}_K^*| = T(T^K - 1)/(T - 1)$ . The proportion of non-zero elements in  $\mathbf{Q}_f$  is  $(T - 1)/(T^K - 1)$ .

179 We discuss the intuition behind the definition of  $\mathbf{Q}_{f}$  pro-181 vided above and illustrate it in Figure 2 for a case of T = 2and K = 3. For this, we first note that given an input se-182 quence  $v_i$  of size  $|v_i| < K$ , a transition to any state  $v_i$ 183 has a probability of 0 if  $v_j \neq [v_i, v]$  for some  $v \in \mathcal{V}$ , 184 i.e., if the state we transition to is not a concatenation 185 of the input sequence with an additional token from the vocabulary (for instance, a state  $\{0\}$  cannot transition to 187  $\{1, 0\}$  in one step). Applying this reasoning for different 188 values of k < K defines green rectangular blocks of size  $T^k \times T^{k+1}$  in the transition matrix portrayed in Figure 2. 189

When one reaches the blue *square* block in the transition

matrix, the input sequence reaches the maximum context



Figure 2: Illustration of Proposition 3.1 with T = 2 and K = 3.

window length  $v_i$ : the model can no longer append tokens to the input sequence and has to delete 192 the first token from it to proceed. This blue block is of size  $T^K \times T^K$ : it captures transitions 193 between all possible sequences of the maximum admissible length. We define similarly the reference transition matrix  $\mathbf{Q}^*$  of the language where the probability of transitions  $\{f_{\Theta}^{T,K}(v_i)\}_j$  are replaced 194 195 by ground-truth probabilities  $\mathbb{P}_{\mathcal{L}}(v_i \mid v_i)$ . In order to use  $\mathbf{Q}_f$  as  $f_{\Theta}$ , it is now sufficient to define an 196 input distribution  $\delta_0$  of the Markov chain based on input prompt v as a one-hot encoding vector of 197 size  $|\mathcal{V}_{K}^{*}|$  with 1 at the position of the state corresponding to v. Then, the transition to the next state simply writes as  $\delta_1 = \mathbf{Q}_f \delta_0$ . The output of  $f_{\Theta}(v)$  for individual tokens in  $\mathcal{V}$  would then correspond 199 exactly to the probabilities in  $\delta_1$  for states that are concatenations of v with T tokens from V. This 200 process is illustrated in Figure 2. 201

We now characterize this Markov chain and note that, since  $\mathcal{V}_{K}^{*}$  is finite, MC( $\mathcal{V}_{K}^{*}, \mathbf{Q}_{f}$ ) admits a 202 stationary distribution. This stationary distribution is unique given the structure of the transition 203 matrix  $\mathbf{Q}_{f}$ , as established in the following result. 204

**Proposition 3.2.** Let  $MC(\mathcal{V}_{K}^{*}, \mathbf{Q}_{f})$  be a Markov chain defined in Proposition 3.1. Then  $MC(\mathcal{V}_{K}^{*}, \mathbf{Q}_{f})$  is an ergodic unichain and has a unique stationary distribution.

209 A unichain is a chain that has at most one recurrent class plus some additional transient states. From 210 Proposition 3.1, we note immediately that green blocks in Fig. 1 represent transient classes, meaning that applying  $\mathbf{Q}_f$  to the input prompt, represented by a one-hot encoding of size  $|\mathcal{V}_K^*|$ , will transition 211 to a state that corresponds to a sequence of length increased by one with an additional, most likely, 212 token appended to it. This process is repeated if the model is called further on: we append tokens 213 until we reach the context window limit K. At this point, we reach the recurrent class, represented in 214 blue, in which the chain stays until it reaches its unique stationary distribution. We now characterize 215 how many times one should apply  $\mathbf{Q}_{f}$  to the input to reach the stationary distribution.

164 165

166

167 168

169 170

171

172 173 174

175 176

177 178

190

191

205

206

**Proposition 3.3.** Given an ergodic finite-state unichain  $MC(\mathcal{V}_K^*, \mathbf{Q}_f)$  and  $e = (1, 1, ..., 1)^\top$ , then  $\lim_{n\to\infty} \mathbf{Q}_f^n = e\pi$  where  $\pi$  is the stationary distribution of the recurrent class  $\mathscr{R}$  of states, expanded by 0's for each transient state of the unichain. Moreover, for all  $n \ge K$ ,

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

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

226

227

228

229

230

231

232

233

264

216

217

218

219 220 221

> The stationary distribution is the long-term equilibrium of the Markov chain defined by the LLM and can be interpreted as a proxy of its understanding of natural language in its token space. It is independent of the initial state (i.e., input prompt) but rather captures the absolute frequencies of occurrences of certain tokens seen during pre-training. For a well-performing model, it is hence likely to be heavy-tailed, meaning that rare states have a non-zero probability of occurring due to language's ambiguity and complexity. Proposition 3.3 shows that reaching the stationary distribution requires more generation steps for models with larger context window K. Additionally, convergence depends on  $\varepsilon$  (that is, the smallest element of the K<sup>th</sup> power of the transition matrix), which is related to the ability of the chain to explore the state space after having forgotten the input prompt.

# 3.2 ILLUSTRATION ON A TOY MODEL

We illustrate the results of Section 3 on a toy model trained on a sequence of 0s and 1s. Here, each 236 subsequent token is 0 if the sum of three previous tokens is even and 0 otherwise. Therefore, T = 2237 and K = 3. We generate a sequence of 40 digits, resulting in 37 distinct supervised examples, 238 and train a small "GPT-like" model (Karpathy, 2023) on it. We extract the logits from the model 239 by prompting it with all possible combinations of 0s and 1s of length less than three to obtain the 240 transition matrix  $\mathbf{Q}_f \in \mathbb{R}^{14 \times 14}$  depicted in Fig. 3(a). The transition matrix's structure (e.g., presence 241 of transient and recurrent classes) matches the one presented in Fig. 1. Fig. 3(b) displays the stationary 242 distribution of the trained model obtained by raising  $\mathbf{Q}_f$  to power 10<sup>5</sup>. We note that it has a strong 243 bias toward seen training samples in accordance with our intuition behind the stationary distribution 244 presented earlier. Finally, Fig. 3(c) illustrates the convergence rate of the toy model, predicted by 245 Proposition 3.3, and compares it to models with larger dictionary size T and context window K. In Fig. 3(c), we set  $\varepsilon = 10^{-6}$  and note that this parameter reflects the ability of the LLM to explore the 246 247 state space.



Figure 3: Markov chain with a small GPT-like model. (a) Transition matrix  $\mathbf{Q}_f$  of the model where denotes the examples from the training set. (b) Stationary distribution of the trained model assigning almost uniform probabilities to the states seen during training. (c) Convergence rate to the stationary distribution for the considered toy model along with three LLMs, highlighting the dependence on K. The y-axis is the upper bound in Proposition 3.3.

**Role of the temperature.** To better illustrate the role of  $\varepsilon$ , we now plot the transition matrix of the studied Markov chain obtained when applying different temperature scaling to the logits returned by the trained model. As the temperature is commonly linked to the ability of LLMs to transition more freely to a large set of states (Chen & Ding, 2023), we expect that lower temperatures should impact negatively the speed of the convergence to the stationary distribution. In Fig. 4(a), we show that for a low temperature (0.2), the Markov chain mixes slowly and is unable to reach its stationary distribution (same line in the transition matrix as in Fig. 3(c)) even after  $10^6$  steps. In the case of a more commonly used temperature equal to 1 (Fig. 4(b)), the model requires only 300 steps to converge. Finally, setting the model's temperature to 2 (Fig. 4(c)) makes the convergence extremely fast, reaching the stationary distribution after only 30 steps. The interplay between  $\varepsilon$  and the model's temperature is displayed in Fig. 4(d), increasing the temperature leads to a drastic improvement in the convergence speed.

276 277 278

279

281

282

283

284

285

287

292

300



Figure 4: Dependence of  $\varepsilon$  on the temperature of the model. (a) For low temperatures,  $\varepsilon$  becomes too small to achieve convergence to the stationary distribution. (b)-(c) Increasing the temperature from 1 to 2 leads to a  $\times 10$  faster convergence. (d)  $\varepsilon$  (log-scale) increase for temperature values in [0.1, 2].

### 4 GENERALIZATION BOUNDS FOR LARGE LANGUAGE MODELS

The inference of any large language model  $f_{\Theta}$  can be fully captured by a Markov chain with a finite transition kernel  $\mathbf{Q}_f$  defined as above. The formalization of Section 3.1 allows us to see and study the generalization of  $f_{\Theta}$  as its capacity to infer correctly all the elements of  $\mathbf{Q}_f$  that approximate the true reference matrix of transition probabilities  $\mathbf{Q}^*$ . The hardness of this task lies in achieving precise inference having observed a negligible amount of  $\mathbf{Q}^*$ 's elements during its pre-training. For GPT-3 (Brown et al., 2020), this represents  $5 \times 10^{11}$  training tokens, which pales in comparison with the number of non-zero elements in  $\mathbf{Q}_f$ , given by  $T^{K+1} \approx 10^{9632}$ .

**Risk definition.** We denote by  $X = (\mathbf{X}_1, \dots, \mathbf{X}_N)$  the tokens in  $\mathcal{V}$  that  $f_{\Theta}$  observes (e.g., during pre-training or at inference time). The training sequences of tokens can be written as  $\mathbf{S}_n = (\mathbf{X}_1, \dots, \mathbf{X}_n)$  if  $n \leq K$  and  $\mathbf{S}_n = (\mathbf{X}_{n-K+1}, \dots, \mathbf{X}_n)$  otherwise due to the deletion process (see Definition B.2). In particular, the  $\mathbf{S}_n$  are elements of  $\mathcal{V}_K^*$ . For any  $n \in [N]$ , the true probability of next token  $\mathbf{X}_{n+1}$  given a past sequence  $\mathbf{S}_n$  is defined as  $\mathbb{P}_{\mathcal{L}}(\cdot | \mathbf{S}_n) \in \Delta_T$  and the probability estimated by the model writes  $\mathbb{P}_{\Theta}(\cdot | \mathbf{S}_n)$ . We assume the existence of a constant  $c_0 > 0$ such that for any  $n \in [N]$  and  $(x_1, \dots, x_{n+1}) \in \Omega^{n+1}$ ,

$$\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} = x_{n+1} \mid \mathbf{X}_n = x_n, \dots, \mathbf{X}_1 = x_1) \ge c_0 > 0.$$
(1)

This is a common assumption used previously in (Hu et al., 2024; Wies et al., 2024; Xie et al., 2022; Zhang et al., 2023b). Following the Markov chain formalization introduced in Section 3.1, we define the theoretical and empirical risks for any  $\Theta \in W$  as<sup>3</sup>

313 314

315

323

308

$$\mathcal{R}(\boldsymbol{\Theta}) \coloneqq \mathbb{E}_{\mathbf{S} \sim \mathbb{P}_{\mathcal{L}}}[d_{\mathrm{TV}}(\mathbf{Q}^{*}(\mathbf{S}, \cdot), \mathbf{Q}_{f}(\mathbf{S}, \cdot))], \, \widehat{\mathcal{R}}(\boldsymbol{\Theta}) \coloneqq \frac{1}{N} \sum_{n=1}^{N} d_{\mathrm{TV}}(\mathbb{P}_{\mathcal{L}}(\cdot \mid \mathbf{S}_{n}), \mathbb{P}_{\boldsymbol{\Theta}}(\cdot \mid \mathbf{S}_{n})).$$
(2)

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

Remark 4.1 (Choice of risk). Our risk definition departs from usual generalization bounds in statistical learning where risks are mostly derived from empirical risk minimization (Bach, 2024; Marion, 2023; Redko et al., 2019; Vapnik, 1999). As we want to assess how well the model estimates the probability distribution of the next token, we rather follow (Hu et al., 2024; Zhang et al., 2023b) and the learning and identity testing of Markov chains literature (Wolfer & Kontorovich, 2019; 2023) and use the total variation distance.

$${}^{3}\mathcal{R}(\boldsymbol{\Theta}) = \mathbb{E}[\widehat{\mathcal{R}}(\boldsymbol{\Theta})] = \mathbb{E}_{\mathbf{S} \sim \mathbb{P}_{\mathcal{L}}}[d_{\mathrm{TV}}(\mathbb{P}_{\mathcal{L}}(\cdot \mid \mathbf{S}), \mathbb{P}_{\boldsymbol{\Theta}}(\cdot \mid \mathbf{S}))] = \mathbb{E}_{\mathbf{S} \sim \mathbb{P}_{\mathcal{L}}}[d_{\mathrm{TV}}(\mathbf{Q}^{*}(\mathbf{S}, \cdot), \mathbf{Q}_{f}(\mathbf{S}, \cdot))].$$

324 **Transformer model.** Without loss of generality,  $f_{\Theta}$  is assumed to be a transformer model with L 325 layers and H heads, consisting of alternating multi-head attention (MHA) and feed-forward blocks 326 (more details in Appendix B). The first layer receives an input  $S^{(0)} = S$  embedded in *r*-dimensional 327 space. To obtain a probability distribution on the vocabulary  $\mathcal{V}$ , the output  $\mathbf{S}^{(L)} \in \mathbb{R}^{r \times T}$  of the final layer is projected back to the vocabulary size by an "unembedding layer"  $\mathbf{W}_U \in \mathbb{R}^{T \times r}$  and averaged 328 over the columns to obtain a vector in  $\mathbb{R}^T$ . A softmax layer is finally applied to obtain the probability distribution of the next token  $\mathbb{P}_{\Theta}(\cdot | \mathbf{S}) \coloneqq \operatorname{softmax}\left(\frac{1}{n\tau}\mathbf{W}_U\mathbf{S}^{(L)}\mathbb{1}_n\right) \in \Delta_T$ , where  $\Theta$  denotes the 330 parameters of the entire network and  $\tau$  is the softmax temperature (Hinton, 2015). Unless otherwise 331 specified, we assume that the unembedding layer is bounded. The classes of parameters and neural 332 networks it generates respectively write  $\mathcal{W} = \{ \Theta \text{ s.t. } \| \mathbf{W}_U^{\top} \|_{2,1} \leq B_U \}$  and  $\mathcal{F} = \{ f_{\Theta} \text{ s.t. } \Theta \in \mathcal{W} \}$ . 333

334 335

336

337

338

339

340

341

342 343

344

345

346 347

348

349

355

356

361

367

377

#### 4.1 PRE-TRAINING THEORETICAL ANALYSIS

We now significantly extend the scope of our theoretical contributions by assuming that the pretraining data  $S = (\mathbf{S}_1, \dots, \mathbf{S}_{N_{\text{train}}})$  is a sequence of *dependent* random variables with a mild coupling structure, namely that a Marton coupling with mixing matrix  $\Gamma$  exists for  $S = (\mathbf{S}_1, \dots, \mathbf{S}_{N_{\text{train}}})$ .<sup>4</sup> This ensures our setting remains very broad as it subsumes the case of independent variables, mdependent variables, language bigrams (Bietti et al., 2023), and the Markov chain setting considered in state-of-the-art ICL analysis of LLMs (Hu et al., 2024; Zhang et al., 2023b).

**Generalization bound.** We denote the risks by  $\mathcal{R}_{pre}(\Theta)$  and  $\widehat{\mathcal{R}}_{pre}(\Theta)$  to indicate that we take  $N = N_{\text{train}}$  in Eq. (2). Below, we state our main result, whose proof is deferred to Appendix D.4, which provides a generalization bound on the estimation risk of pre-training.

**Theorem 4.1** (Pre-training generalization bound). *Consider an LLM*  $f_{\Theta} \in \mathcal{F}$ . We denote by  $\Gamma$  the mixing matrix of the pre-training sequences of tokens  $(\mathbf{S}_1, \ldots, \mathbf{S}_{N_{\text{train}}})$ . Let  $0 < \delta < 1$ , then with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\mathrm{pre}}(\mathbf{\Theta}) \leq \widehat{\mathcal{R}}_{\mathrm{pre}}(\mathbf{\Theta}) + rac{ar{B}}{\sqrt{N_{\mathrm{train}}}} \sqrt{\log\left(rac{2}{\delta}
ight)},$$

where  $\overline{B} = 2 \| \Gamma \| \max\{ \log (T) + 2B_U / \tau, \log (1/c_0) \}^{1/2}$  is a constant depending on the parameters of the problem.

357 The bound in Theorem 4.1 depends on the intrinsic structure of the pre-training data through the norm 358 of the mixing matrix  $\|\mathbf{\Gamma}\|$ . If the pre-training data S is a Markov chain with state space  $\Omega$ , this norm 359 captures exactly the mixing time of the latter, making sequences that mix at a slower pace harder 360 to learn. Secondly, and perhaps most surprisingly, this bound becomes model-independent when  $\max\{\log(T) + 2B_U/\tau, \log(1/c_0)\}$  is dominated by  $\log(1/c_0)$  term. Hence, if  $B_U \approx \mathcal{O}(T\sqrt{r})$ , which happens in practice due to the common normalization of the unembedding layer, then the 362 model's hidden dimension r and vocabulary size T should be large enough to ensure  $\log(T)$  + 363  $2B_U/\tau \ge \log(1/c_0)$  for some unknown reference constant  $c_0$ . Below this threshold, the architecture 364 of  $f_{\Theta}$  is not expressive enough to have any tangible impact on its generalization, although it may 365 affect the training error  $\mathcal{R}_{\text{pre}}(\Theta)$ . 366

**Depth-dependent variation.** We extend Theorem 4.1 to make its dependency on  $f_{\Theta}$  more fine-368 grained. Rather than assuming that only the norm of the embedding layer's matrix is bounded, we 369 follow the setting of prior work (Edelman et al., 2022; Furuya et al., 2024; Marion, 2023; Zhang 370 et al., 2023b) and consider the parameter space defined as follows: 371

$$\widetilde{\mathcal{W}} = \{ \boldsymbol{\Theta} \in \mathcal{W} \mid \forall \ell \in [L], \| \mathbf{W}_{V}^{(\ell)} \|_{\infty} \le B_{V}, \| \mathbf{W}_{O}^{(\ell)} \|_{\infty} \le B_{O}, \| \mathbf{W}_{1}^{(\ell)} \|_{\infty} \le B_{1}, \| \mathbf{W}_{2}^{(\ell)} \|_{\infty} \le B_{2} \}$$

374 The definition of W concerns the query, key, and value matrices of all layers and heads. Similarly 375 to Zhang et al. (2023b, Assumption 5.1), we assume that each token has an  $\ell_1$ -norm bounded by  $B_{\rm tok}$ . We have the following generalization bound, whose proof is deferred to Appendix D.5. 376

 $<sup>{}^{4} \| \</sup>mathbf{\Gamma} \| = 1$  for independent variables and more details on Marton coupling can be found in Appendix C.3.

**Corollary 4.2** (Depth-dependent bound). Consider an LLM  $f_{\Theta} \in \tilde{\mathcal{F}} := \{f_{\Theta} \mid \Theta \in \tilde{\mathcal{W}}\}$ . With the same assumptions as in Theorem 4.1, we have

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

where  $\overline{B} = 2 \|\mathbf{\Gamma}\| \max\{\log{(T)} + 2(B_{\Theta})^L/\tau, \log{(1/c_0)}\}^{1/2}$  is a constant depending on the parameters of the problem, and  $B_{\Theta} = [(1 + rmB_1B_2)(1 + \frac{r^3}{H}B_OB_V)](B_{\text{tok}}B_U)^{1/L}$ .

We note that B exhibits an exponential dependence on the depth of the transformer, which also amplifies the hidden dimensionality (width) of the embedding layer r. This contrasts with the dependency in m, the hidden dimensionality of the MLP block, which is linear. All these factors are commonly associated with higher expressive power of transformers suggesting that they should contribute to a better minimization of  $\hat{\mathcal{R}}_{pre}(\Theta)$  at the expense of requiring more training data. The number of heads H can be used as a counterbalance to increasing the width in the cubic term  $r^3$ , suggesting that a good balance between these parameters may lead to more data-efficient models.

**Sample complexity of LLMs.** Our goal is to show the asymptotic dependence on the number of sequences that an LLM requires such that  $Q_f$  is  $\varepsilon$ -close to the reference transition matrix  $Q^*$ . We then derive a sample complexity bound. The proof is deferred to Appendix D.6.

**Corollary 4.3** (Sample complexity). Let  $\overline{B}$  be the parameter-dependent constant of Theorem 4.1 or Corollary 4.2. Let  $\delta \in [0,1]$  and let  $\epsilon > 0$ . If  $N_{\text{train}} \ge N^* \coloneqq \lfloor \frac{4\overline{B}^2}{\epsilon^2} \log \left(\frac{2}{\delta}\right) \rfloor$  and if we assume a perfect pre-training error for  $f_{\Theta}$ , then we have with probability at least  $1 - \delta$ ,

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

This result allows us to contextualize LLMs' ability to learn Markov chains with respect to the existing 406 literature. To the best of our knowledge, the only existing approach with theoretical guarantees for 407 learning Markov chains is the frequentist method: counting the number of occurrences of different 408 states to fill in the matrix  $\mathbf{Q}_f$ . Wolfer & Kontorovich (2019) show that the sample complexity of ap-409 proximating  $\mathbf{Q}^*$  up to  $\epsilon$  with such approach requires at most  $\mathcal{O}(\max\{|\mathcal{V}_k^*|/\epsilon^2\gamma_s, 1/\gamma_s\pi^*\})$  samples, 410 where  $\gamma_s$  is a (pseudo) spectral gap of the Markov chain and  $\pi^*$  is the smallest element of its stationary 411 distribution. The authors state that the frequentist approach is minimax optimal (up to logarithmic 412 factors). Our bound has a dependence that behaves as  $\bar{B}^2 = \mathcal{O}(\max\{\log T + \frac{2T\sqrt{\tau}}{\tau}, \log(1/c_0)\})$ . 413 Given that in practice T > r, it then simplifies to  $\mathcal{O}(\max\{T/\epsilon^2 \tau, 1/\epsilon^2\})$ . Note that the LLMs' 414 sample complexity is linear in the vocabulary size T, which is remarkable compared to the sample complexity of the frequentist approach, which scales as  $\mathcal{O}(T^K)$ . We show in Section 5 that this is 415 confirmed experimentally: LLM's ability to learn Markov chains exceeds the frequentist approach 416 for Markov chains with a large state space. 417

### 419 4.2 IN-CONTEXT LEARNING OF MARKOV CHAINS

Although insightful, the analysis presented above is related to the pre-training of LLMs – a process 421 that is hard and extremely costly to reproduce in practice. Similarly, we do not have access to the 422 ground-truth matrix  $\mathbf{Q}^*$  to reason about LLM's ability to infer it in practice. To provide theoretical 423 results that can be confirmed experimentally, we now turn our attention to in-context learning of 424 Markov chains: a setup where one provides an LLM with an input sequence formed by a Markov 425 chain of size  $N_{\rm icl}$  defined over a state space  $\Omega$  of size  $d^5$ . Different from the setting of Section 4.1, 426 we now can explicitly use a transition kernel  $\mathbb{P}$  of this Markov chain for the theoretical analysis 427 by replacing  $\mathbb{P}_{\mathcal{L}}$  with it in the definition of  $\mathcal{R}_{icl}(\Theta)$  and  $\mathcal{R}_{icl}(\Theta)$  in Eq. (2) (see Appendix D.7 for 428 details on the problem setup). To relate the generalization error to the pre-training error, we quantify

418

420

378

379

380 381 382

384

385

386 387 388

389

390

391

392

393

394 395

396

397

398 399

400

401 402 403

<sup>429</sup> 430

<sup>&</sup>lt;sup>5</sup>This is different from another variation of ICL where supervised (x,y) pairs are provided in-context. Rather, the supervision is provided from observing transitions between states  $(x_i, x_{i+1} = f(x_i))$  as discussed in (Li et al., 2023, Fig.1).

432 the discrepancy between an LLM pre-trained mostly on textual data, and a hypothetical LLM with 433 parameters in  $\mathcal{W}_{mc}$  that is pre-trained on a dataset of Markov chains with the same data distribution 434 as the Markov chain used as an input during in-context inference. We define the divergence between 435 two estimated transition matrices  $\mathbb{P}_{\Theta_1}, \mathbb{P}_{\Theta_2}$  as

$$\mathcal{K}(\mathbf{\Theta}_1, \mathbf{\Theta}_2) \coloneqq \frac{1}{N} \sum_{n=1}^{N} \mathbb{E}_{\mathbf{S}_n} [d_{\mathrm{TV}}(\mathbb{P}_{\mathbf{\Theta}_1}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\mathbf{\Theta}_2}(\cdot \mid \mathbf{S}_n))].$$
(3)

The operator  $\mathcal{K}$  is akin to a distance (the separation property is only verified almost surely, see Appendix C.4 for more details). The next result, whose proof is deferred to Appendix D.7, provides a generalization bound on the in-context learning phase.

**Theorem 4.4** (In-Context Learning generalization bound). Consider an LLM  $f_{\Theta} \in \mathcal{F}$ . We provide as input of  $f_{\Theta}$  a d-state Markov chain  $X = (\mathbf{X}_1, \ldots, \mathbf{X}_{N_{icl}})$ . The sequence of subsequences of the first n terms is denoted by  $S = (\mathbf{S}_1, \dots, \mathbf{S}_n)$ . S is also a Markov chain, and we denote by  $t_{\min}(\varepsilon)$  its mixing time. Let  $t_{\min} \coloneqq \inf_{0 < \varepsilon < 1} t_{\min}(\frac{\varepsilon}{2})(\frac{2-\varepsilon}{1-\varepsilon})^2$ . Let  $\delta > 0$ . Then, with probability at least  $1 - \delta$ ,

$$\mathcal{R}_{\rm icl}(\boldsymbol{\Theta}) \leq \inf_{\boldsymbol{\vartheta} \in \mathcal{W}_{\rm mc}} \{ \widehat{\mathcal{R}}_{\rm icl}(\boldsymbol{\vartheta}) + \mathcal{K}(\boldsymbol{\vartheta}, \boldsymbol{\Theta}) \} + \bar{B} \sqrt{\frac{t_{\rm min}}{N_{\rm icl}}} \sqrt{\log\left(\frac{2}{\delta}\right)},\tag{4}$$

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

We first note that instead of the norm of the mixing matrix  $\Gamma$  seen before, we now have an explicit 455 dependency on  $t_{\min}$ , which is related to the mixing time of the input Markov chain. This, together 456 with the availability of the ground-truth transition matrix, allows us to use Theorem 4.4 to derive and verify experimentally the scaling laws of ICL for popular LLMs. Theorem 4.4 also suggests that an 458 LLM pre-trained on diverse data sequences different from Markov chains should exhibit a certain degree of invariance to correctly infer the transition probabilities of the latter. This is reminiscent of 460 the domain adaptation bounds (Redko et al., 2019) that also commonly involve a distribution shift (i.e., a distance or a divergence) term that vanishes if the model is invariant to classes of transformations linking the distribution of the input data with that on which it is applied during inference. A recent success of applying LLMs to time series data (Gruver et al., 2023), for instance, suggests that this term is indeed small for certain types of data not used during pre-training.

#### 5 NUMERICAL EXPERIMENTS

468 Theorem 4.4 provides a practically verifiable result which naturally stems from our analysis in Section 4. We then evaluate the ability of recent LLMs, namely Mistral 7Bv0.1 (Jiang et al., 469 2023), Llama2 7B & 13B (Touvron et al., 2023b), and Gemma 2B (Team et al., 2024) to infer 470 transition probabilities of Markov chains in-context. We associate each state in the d-state Markov 471 chain with a token from the set  $\{0, \ldots, d-1\}$ , concatenated to obtain a prompt of length  $N_{\rm icl}$ . 472 Bearing in mind the differences in the tokenization mechanisms of the different models, we add comas 473 whenever necessary to ensure that each state is tokenized separately. More details on the experimental 474 setup and additional experiments with more Markov chains and with Llama3.2 (Dubey et al., 2024) 475 are available in Appendix E.1.

476

483

440

441

442 443

444

445

446

447

453 454

457

459

461

462

463

464 465

466 467

477 **Dependence on**  $N_{icl}$ . We first analyze the effect of  $N_{icl}$  on the risk calculated for a randomly gener-478 ated 3-state Markov transition matrix. From the results presented in Fig. 5(left), we note that Llama2 models deviate from our  $\mathcal{O}(N_{\rm icl}^{-1/2})$  theoretical scaling law, while most recent models (Mistral 479 480 and Gemma) stay much closer to Theorem 4.4, similarly to what was observed by Cabannes et al. 481 (2024). Being randomly generated, the Markov chains provided to the models have not been seen during training, and older (weaker) models naturally struggle to generalize. 482

**Dependence on**  $t_{\min}$ . Theorem 4.4 states that Markov chains with slow mixing (higher  $t_{\min}$ ) are 484 slower to learn. We now plot the true risk for a single model with different values of  $t_{\min}$  highlighting 485 in Fig. 5(right) a two-stage regime of ICL. In a first stage, the bound in Eq. (4) is dominated by 486

487

488 489

490

491

492

493

494

495

500

501

502

503

504

505

506

507 508

509

510

515

516

517

518

519

520 521

522

523

526

527

529

531

Gemma 2B

Context Length Nic

as  $N_{\rm icl}$  increases beyond  $N_{\rm icl} \approx 20$ . = 43.7 (a) (b) (c) in = 32.6 in = 15.2 32 F = 4.6  $\mathcal{R}_{\rm lcl}$ Rici R<sub>icl</sub> Error Error Error Llama2 13B 0.2 Mistral 7B v0.1

Small Nici

0.0

 $\sqrt{t_{\min}/N_{icl}}$  for small  $N_{icl}$ , and depends strongly on  $t_{\min}$ , while the scaling law  $\mathcal{O}(N_{icl}^{-1/2})$  dominates



Context Length N<sub>icl</sub>

Scaling law

Ratio N<sub>icl</sub>/t<sub>min</sub>

**Dependence on** *d*. We now verify Theorem 4.4 for Markov chains with a different state space size (previously d = 3). We also consider a baseline given by the frequentist method mentioned before. We recall that, for the latter, its dependence on d behaves like  $\mathcal{O}(\sqrt{d/N_{icl}})$ , while Theorem 4.4 gives  $\mathcal{O}(\sqrt{\log(d)}/N_{\rm icl})$ . For Markov chains with a small number of states d, there is no clear difference between the frequentist estimator and a LLM. However, as d grows the frequentist estimator struggles to estimate the transition matrix due to the  $\mathcal{O}(\sqrt{d})$  scaling factor. This is verified experimentally in Fig. 6, where we vary the parameter d from 3 (left) to 700 (right). We observe that the LLM follows the theoretical neural scaling law  $\mathcal{O}(N_{\rm icl}^{-1/2})$  and outperforms the frequentist method for d = 700, while being close to it for d = 3. We conclude that our analysis gives theoretical insights on the ICL neural scaling law observed empirically in (Liu et al., 2024). The additional experiments conducted in Appendix E.5 show that our bounds remain valid for large values of d.



Figure 6: Impact of the number of states. We plot the risks  $\mathcal{R}_{icl}$  as functions of  $N_{icl}$  for Gemma 2B and the frequentist approach (Wolfer & Kontorovich, 2019) with 95% confidence intervals. Left. The input sequence is a random 3-state Markov chain. **Right.** The input sequence is a Brownian motion discretized as a 700-state Markov chain, similarly to Liu et al. (2024). 524 525

CONCLUSION 6

528 This paper proposed an explicit characterization of the inference mechanism in large language models through an equivalent finite-state Markov chain. We provided an insightful theoretical analysis 530 based on the established characterization and the ability of the LLM to infer the transition kernel approximating the true transition probabilities of language. We adapted our results to in-context 532 learning where experiments confirm our theoretical insights. In the future, we hope that the proposed 533 equivalence will have far-reaching implications on our understanding of LLMs and allow for a more 534 fine-grained understanding of their expressiveness.

536 REFERENCES 537

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, 538 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.

540 Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity 541 and representation learning of low rank mdps. In Advances in Neural Information Processing 542 Systems, volume 33, pp. 20095–20107, 2020. 543 Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. Intrinsic dimensionality explains the 544 effectiveness of language model fine-tuning. arXiv preprint arXiv:2012.13255, 2020. 546 Mehdi Ali, Michael Fromm, Klaudia Thellmann, et al. Tokenizer Choice For LLM Training: 547 Negligible or Crucial? In Findings of the Association for Computational Linguistics: NAACL 2024, pp. 3907–3924. Association for Computational Linguistics, 2024. 548 549 Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan 550 Schalkwyk, Andrew M Dai, and Anja et al. Hauth. Gemini: a family of highly capable multimodal 551 models. arXiv preprint arXiv:2312.11805, 2023. 552 Francis Bach. Learning Theory from First Principles. MIT Press, 2024. 553 554 Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a 555 Transformer: A Memory Viewpoint. In Advances on Neural Information Processing Systems, 556 2023. Mathieu Blondel, Andre Martins, and Vlad Niculae. Learning classifiers with fenchel-young losses: 558 Generalized entropies, margins, and algorithms. In Kamalika Chaudhuri and Masashi Sugiyama 559 (eds.), Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics, volume 89 of Proceedings of Machine Learning Research, pp. 606–615. PMLR, 16–18 561 Apr 2019. URL https://proceedings.mlr.press/v89/blondel19a.html. 562 563 Tom Brown, Benjamin Mann, Nick Ryder, et al. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pp. 1877–1901, 2020. 564 565 Vivien Cabannes, Elvis Dohmatob, and Alberto Bietti. Scaling Laws for Associative Memories. In 566 International Conference on Learning Representations, 2024. 567 H. Chen and N. Ding. Probing the "creativity" of large language models: Can models produce 568 divergent semantic association? In EMNLP, pp. 12881-12888. ACL, 2023. 569 570 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep 571 bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar 572 Solorio (eds.), Conference of the North American Chapter of the Association for Computational 573 Linguistics: Human Language Technologies, pp. 4171–4186, 2019. 574 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 575 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. 576 arXiv preprint arXiv:2407.21783, 2024. 577 Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable 578 creation in self-attention mechanisms. In International Conference on Machine Learning, pp. 579 5793-5831. PMLR, 2022. 580 581 Benjamin L Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The evolution 582 of statistical induction heads: In-context learning markov chains. arXiv preprint arXiv:2402.11004, 583 2024. 584 Gerald B Folland. Real analysis: modern techniques and their applications, volume 40. John Wiley 585 & Sons, 1999. 586 Takashi Furuya, Maarten V de Hoop, and Gabriel Peyré. Transformers are universal in-context learners. arXiv preprint arXiv:2408.01367, 2024. 588 589 Robert G Gallager. Finite state Markov chains. In Discrete Stochastic Processes, pp. 103-147. 590 Springer, 1996. Siavash Golkar, Mariel Pettee, Michael Eickenberg, Alberto Bietti, Miles Cranmer, Geraud Krawezik, 592 Francois Lanusse, Michael McCabe, Ruben Ohana, Liam Parker, et al. xval: A continuous number

encoding for large language models. arXiv preprint arXiv:2310.02989, 2023.

604

610

624

625

626

Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew G Wilson. Large language models are zero-shot time series forecasters. *Advances in Neural Information Processing Systems*, 36, 2023.

- Yi Hao, Alon Orlitsky, and Venkatadheeraj Pichapati. On learning markov chains. In
  S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.,
  2018. URL https://proceedings.neurips.cc/paper\_files/paper/2018/
  file/d34ab169b70c9dcd35e62896010cd9ff-Paper.pdf.
- Geoffrey Hinton. Distilling the Knowledge in a Neural Network. *arXiv preprint arXiv:1503.02531*, 2015.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
   Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference* on Learning Representations, 2022.
- Xinyang Hu, Fengzhuo Zhang, Siyu Chen, and Zhuoran Yang. Unveiling the Statistical Foundations
   of Chain-of-Thought Prompting Methods. *arXiv preprint arXiv:2408.14511*, 2024.
- M Emrullah Ildiz, Yixiao Huang, Yingcong Li, Ankit Singh Rawat, and Samet Oymak. From Self-Attention to Markov Models: Unveiling the Dynamics of Generative Transformers. *arXiv preprint arXiv:2402.13512*, 2024.
- Hong Jun Jeon, Jason D. Lee, Qi Lei, and Benjamin Van Roy. An Information-Theoretic Analysis
  of In-Context Learning. In *International Conference on Machine Learning*, volume 235, pp. 21522–21554, 2024.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7B. arXiv preprint arXiv:2310.06825, 2023.
- Daniel Jurafsky and James H. Martin. Speech and Language Processing: An Introduction to Natural
   Language Processing, Computational Linguistics, and Speech Recognition with Language Models.
   Pearson, 3rd edition, 2024.
  - Andrej Karpathy. minGPT: A minimal PyTorch re-implementation of the GPT (Generative Pretrained Transformer). https://github.com/karpathy/minGPT, 2023. GitHub repository.
- Juno Kim, Tai Nakamaki, and Taiji Suzuki. Transformers are Minimax Optimal Nonparametric
   In-Context Learners. *arXiv preprint arXiv:2408.12186*, 2024.
- Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*, pp. 19565–19594. PMLR, 2023.
- Toni JB Liu, Nicolas Boullé, Raphaël Sarfati, and Christopher J Earls. LLMs learn governing
   principles of dynamical systems, revealing an in-context neural scaling law. *arXiv preprint arXiv:2402.00795*, 2024.
- Sanae Lotfi, Marc Finzi, Yilun Kuang, Tim GJ Rudner, Micah Goldblum, and Andrew Gordon Wilson. Non-vacuous generalization bounds for large language models. *arXiv preprint arXiv:2312.17173*, 2023.
- Sanae Lotfi, Yilun Kuang, Brandon Amos, Micah Goldblum, Marc Finzi, and Andrew Gordon
   Wilson. Unlocking tokens as data points for generalization bounds on larger language models.
   *arXiv preprint arXiv:2407.18158*, 2024.
- Ashok Vardhan Makkuva, Marco Bondaschi, Adway Girish, Alliot Nagle, Martin Jaggi, Hyeji Kim, and Michael Gastpar. Attention with Markov: A framework for principled analysis of transformers via markov chains. *arXiv preprint arXiv:2402.04161*, 2024.
- 647 Pierre Marion. Generalization bounds for neural ordinary differential equations and deep residual networks. In *Advances on Neural Information Processing Systems*, volume 36, 2023.

648 649 650	Katalin Marton. Measure concentration for Euclidean distance in the case of dependent random variables. <i>Ann. Probab.</i> , 32(3):2526–2544, 2004.
651 652	Cleo Nardo. Remarks 1-18 on GPT (compressed). https://www.lesswrong.com/posts/ 7qSHKYRnqyrumEfbt/, 2023. Accessed: 2024-09-27.
653 654 655	Daniel Paulin. Concentration inequalities for Markov chains by Marton couplings and spectral methods. <i>Electron. J. Probab.</i> , 20:79, 2015.
656 657 658	Ievgen Redko, Amaury Habrard, Emilie Morvant, Marc Sebban, and Younès Bennani. State of the Art of Statistical Learning Theory. In <i>Advances in Domain Adaptation Theory</i> , pp. 1–19. Elsevier, 2019.
659 660 661 662	Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the parameters of a language model? In <i>Empirical Methods in Natural Language Processing</i> , pp. 5418–5426. Association for Computational Linguistics, 2020.
663 664	Gareth O. Roberts and Jeffrey S. Rosenthal. General state space Markov chains and MCMC algorithms. <i>Probab. Surveys</i> , 1:20 – 71, 2004.
665 666 667 668	Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural Machine Translation of Rare Words with Subword Units. In Katrin Erk and Noah A. Smith (eds.), <i>Association for Computational Linguistics</i> , pp. 1715–1725, 2016.
669 670	Aaditya K Singh and DJ Strouse. Tokenization counts: the impact of tokenization on arithmetic in frontier llms. <i>arXiv preprint arXiv:2402.14903</i> , 2024.
671 672 673 674	Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. <i>arXiv preprint arXiv:2403.08295</i> , 2024.
675 676 677	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023a.
678 679 680 681	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023b.
682	Alexandre B. Tsybakov. Introduction to Nonparametric Estimation. Springer, 1st edition, 2008.
683 684 685	V. N. Vapnik. An overview of statistical learning theory. <i>IEEE Trans. Neural Netw.</i> , 10(5):988–999, 1999.
686 687 688	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>Advances in Neural Information Processing Systems</i> , volume 30, 2017.
689 690 691	Noam Wies, Yoav Levine, and Amnon Shashua. The learnability of in-context learning. In Advances in Neural Information Processing Systems, volume 36, 2024.
692 693	Geoffrey Wolfer and Aryeh Kontorovich. Minimax learning of ergodic Markov chains. In <i>Algorithmic Learning Theory</i> , pp. 904–930. PMLR, 2019.
695 696	Geoffrey Wolfer and Aryeh Kontorovich. Learning and identity testing of Markov chains. In <i>Handbook of Statistics</i> , volume 49, pp. 85–102. Elsevier, 2023.
697 698 699	Zhenglong Wu, Qi Qi, Zirui Zhuang, Haifeng Sun, and Jingyu Wang. Pre-tokenization of numbers for large language models. In <i>The Second Tiny Papers Track at ICLR 2024</i> , 2024.
700 701	Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An Explanation of In-context Learning as Implicit Bayesian Inference. In <i>International Conference on Learning Representations</i> , 2022.

Can Yaras, Peng Wang, Laura Balzano, and Qing Qu. Compressible dynamics in deep overparameterized low-rank learning & adaptation. *arXiv preprint arXiv:2406.04112*, 2024.

- Biao Zhang and Rico Sennrich. Root Mean Square Layer Normalization. In Advances in Neural Information Processing Systems, volume 32, 2019.
- Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. Sentiment analysis in the era of large language models: A reality check. *arXiv preprint arXiv:2305.15005*, 2023a.

Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. What and how does in-context learning learn? Bayesian model averaging, parameterization, and generalization. *arXiv preprint arXiv:2305.19420*, 2023b.

	U
709	
710	Yuf
711	10
712	a
713	
714	
715	
716	
717	
718	
719	
720	
721	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	

705

706 707

708

# Appendix

 **Roadmap.** In Appendix A, we first recall our notations. We provide additional details on large language models and transformers in Appendix B. Important notions and definitions related to Markov chains and Marton couplings are given in Appendix C. The detailed proofs of our theoretical results are given in Appendix D. Finally, we provide additional experiments in Appendix E.

# TABLE OF CONTENTS

A	Nota	ations	16
В	Bac	kground on Large Language Models	16
	<b>B</b> .1	Large Language Models	16
	B.2	Transformer Architecture	16
	B.3	Autoregressive Transformer-based LLM	18
С	Bac	kground on Markov Chains	18
	C.1	Basic Notions	18
	C.2	Ergodic Unichains	20
	C.3	Marton Couplings	20
	C.4	An (Almost) Distance between Markov Chains	21
D	Proc	ofs	22
	D.1	Proof of Proposition 3.1	22
	D.2	Proof of Proposition 3.2	23
	D.3	Proof of Proposition 3.3	24
	D.4	Proof of Theorem 4.1	25
	D.5	Proof of Corollary 4.2	33
	D.6	Proof of Corollary 4.3	37
	D.7	Proof of Theorem 4.4	37
E	Add	itional Experiments	44
	E.1	Experimental Setup and Tokenization	44
	E.2	Impact of the Number of States d	44
	E.3	More Structured Markov Chains	45
	E.4	Recent Models: Impact of the Tokenization	49
	E.5	Dynamical Systems	49
F	Exte	ended Results with the KL Divergence	50
-	E.1	Pre-training Generalization Bounds	50
	F.2	Limitations	51

# <sup>810</sup> A NOTATIONS

812 We denote  $\{1, \dots, N\}$  as [N]. We represent scalar values with regular letters (e.g., parameter  $\lambda$ ), 813 vectors with bold lowercase letters (e.g., vector x), and matrices with bold capital letters (e.g., matrix 814 A). The *i*-th row of the matrix A is denoted by  $A_i$ , its *j*-th column is denoted by  $A_{j}$  and its 815 transpose is denoted by by  $\mathbf{A}^{\top}$ . The identity matrix of size *n* is denoted by  $\mathbf{I}_n \in \mathbb{R}^{n \times n}$ . The vector of size n with each entry equal to 1 is denoted by  $\mathbb{1}_n$ . We denote by  $\|\mathbf{A}\|_{p,q}$  the  $L_{p,q}$  matrix norm 816 where the p-norm is over columns and the q-norm is over rows. We denote by  $\|\mathbf{A}\|$  the operator 817 818 norm of A induced by the  $\ell_2$  norm and by  $\|\mathbf{A}\|_{\infty} = \max_{ij} |\mathbf{A}_{ij}|$  the operator norm induced by the  $\ell_{\infty}$ -norm. Similarly,  $\mathbf{x}^{\top}$  is the transpose of the vector  $\mathbf{x}$  and  $\|\mathbf{x}\|_{p}$  is its  $\ell_{p}$ -norm. The total variation 819 between two probability distributions  $\mathbb{P}, \mathbb{Q}$  is denoted by  $d_{TV}(\mathbb{P}, \mathbb{Q})$ . The term "almost surely" is 820 denoted by the notation "a.s." while the term random variable is denoted by the notation "r.v.". 821  $\Delta_n \coloneqq \{\mathbf{p} \in [0,1]^n | \sum_{i=1}^n \mathbf{p}_i = 1\}$  is the probability simplex of  $\mathbb{R}^n$ . 822

823 824

825 826

827

828

829 830

831 832

833

834 835

841

842 843 844

845 846 847

848 849

850

851

852

## **B** BACKGROUND ON LARGE LANGUAGE MODELS

We first recall important notions regarding large language models before focusing on the most widely used ones, namely the transformer-based LLMs. We describe the components of the vanilla transformer architecture before describing the whole network at the heart of such a model and formally defining the class of parameters and neural networks considered in our work.

## B.1 LARGE LANGUAGE MODELS

In this section, we recall how the sequences of tokens are processed by the large language model notably regarding the next token generation and the deletion process.

**Definition B.1** (Generation process). Given an input  $s \in \mathcal{V}_K^*$  of size p, an large language model outputs a probability mass function  $f_{\Theta}^{T,K}(s)$  over the discrete vocabulary space. A next token x is then sampled from  $f_{\Theta}^{T,K}(s)$ , to construct a new sequence (s, x) of size p + 1.

Generation can be repeated by considering (s, x) as new input sequence and iterating this process. Since these models are designed to handle only sequences of size at most K, a deletion process is required.

**Definition B.2** (Deletion process). Given an input s of size p > K, an large language model outputs a probability mass function  $f_{\Theta}^{T,K}(s_K)$  where  $s_K$  is a truncation of K tokens of the sequence s. large language models implement front truncation, which is done by setting  $s_K$  as the last K tokens of s.

As shown in Fig. 7, only the last K tokens of a long input sequence are used. This is why we speak of *deletion*, since we ignore the first tokens.

Note that it is possible to implement other kinds of truncation, but large language models usually do not (Brown et al., 2020; Touvron et al., 2023a), however, in models like BERT (Devlin et al., 2019), which are not autoregressive, back truncation as described in Fig. 8 is also an option.

**B.2** TRANSFORMER ARCHITECTURE

The most popular autoregressive LLMs rely on the transformer architecture (Vaswani et al., 2017) which we describe below following (Brown et al., 2020; Edelman et al., 2022; Zhang et al., 2023b). An autoregressive transformer-based LLM takes as input a sequence of length n, with  $n \le K$  and Kis the context window, tokens with values in a vocabulary  $\mathcal{V}$  of size T. The tokens are embedded into a r-dimensional space and the input can be written as  $\mathbf{S} \in \mathbb{R}^{r \times n}$ . We consider a transformer model with L layers and h heads. The output of the  $\ell$ -th layer writes  $\mathbf{S}^{(\ell)}$  and is fed as input of the  $(\ell + 1)$ -th layer. The input of the whole model is  $\mathbf{S}^{(0)} = \mathbf{S}$ . Below, we describe the operations performed by the model, including the embeddings of the tokens.



Figure 7: Deletion process, front truncation. A large language model with context window K = 7 in navy blue, processing sequences of different lengths. Top. A sequence of length 4. Bottom. Front truncation of a sequence of length 10.



Figure 8: **Back truncation**. A large language model with context window K = 7 in navy blue, processing back truncation of a sequence of a sequence of length 10.

- Token embeddings. The tokens are embedded in a r-dimensional space via an embedding layer W<sub>E</sub> which results in an input of the form S<sup>r×n</sup>;
- **Positional embeddings.** (Learnable) positional embeddings are added to each token depending on its position in the input sequence. This breaks the permutation-invariance of the transformer architecture and leads, by abuse of notation, to an output  $\mathbf{S} \in \mathbb{R}^{r \times n}$ ;
- Multi-head attention (MHA). Given an input sequence  $\mathbf{S} \in \mathbb{R}^{r \times n}$ , query, key, and value matrices  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{r \times r}$  (here the value and output matrices are merged for ease of notations), the self-attention module computes

$$\mathcal{A}(\mathbf{S}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V) \coloneqq \operatorname{softmax} \left( \mathbf{W}_Q \mathbf{S} (\mathbf{W}_K \mathbf{S})^\top / \sqrt{r} \right) (\mathbf{W}_V \mathbf{S}) \in \mathbb{R}^{r \times n},$$

with softmax:  $\mathbf{x} \in \mathbb{R}^n \to \exp(\mathbf{x}) / \sum_i \exp(\mathbf{x})_i \in \Delta_n$ . The operation described below corresponds to single-head self-attention. In practice, multi-head attention (MHA) is used with H heads and the query and key matrices are in  $\mathbb{R}^{\frac{r}{H} \times \frac{r}{H}}$  and the value matrix is in  $\mathbb{R}^{\frac{r}{H} \times r}$  (r, H are taken such that  $\frac{r}{H}$  is an integer). The MHA module concatenates on the row dimension the outputs of  $\mathcal{A}$  for each head and then projects it back to the embedding dimension r with an output matrix  $\mathbf{W}_O \in \mathbb{R}^{r \times r}$ . By abuse of notation, we also denote by  $\mathcal{A}$  this operation which results in an output of dimension  $r \times n$ , and we include the output matrix in the argument of the operator. The  $\ell$ -th layer of the transformer applies attention with layer-specific weight matrices and a residual connection that leads to an output

$$\mathbf{Z}^{(\ell)} = \mathbf{S}^{(\ell-1)} + \mathcal{A}\left(\mathbf{S}^{(\ell-1)}; \mathbf{W}_Q^{(\ell)}, \mathbf{W}_K^{(\ell)}, \mathbf{W}_V^{(\ell)}, \mathbf{W}_Q^{(\ell)}\right).$$

This is followed by a layer normalization (Zhang & Sennrich, 2019) that projects each token into the  $\ell_2$ -unit ball, i.e., each column  $\mathbf{S}_{:,n}^{(\ell)}$  has an  $\ell_2$ -norm lower than 1;

910 • Feed-forward block (FF). Finally, a feed-forward block is applied, consisting of a two-layer
911 MLP with hidden dimension m, layer-specific weight matrices W<sub>1</sub><sup>(ℓ)</sup> ∈ ℝ<sup>m×r</sup>, W<sub>2</sub><sup>(ℓ)</sup> ∈ ℝ<sup>r×m</sup> and
912 ReLU activation denoted by ReLU(x) = max{0, x} and applied entry-wise. The output of this layer reads

$$\mathbf{Y}^{(\ell)} = \mathbf{W}_2^{(\ell)} \operatorname{ReLU}\left(\mathbf{W}_1^{(\ell)} \mathbf{Z}^{(\ell)}\right).$$

916 It is followed by a residual connection to produce the output

 $\mathbf{S}^{(\ell)} =$ 

it is followed by a residual connection to produce the

$$\mathbf{Z}^{(\ell)} + \mathbf{W}_2^{(\ell)} \operatorname{ReLU}\left(\mathbf{W}_1^{(\ell)} \mathbf{Z}^{(\ell)}\right) \in \mathbb{R}^{r \times n}$$

on which layer normalization (Zhang & Sennrich, 2019) is applied ensuring that each column  $\mathbf{S}_{:,n}^{(\ell)}$ has an  $\ell_2$ -norm lower than 1.

• softmax output layer. In the autoregressive setting, the model outputs a probability distribution on the vocabulary  $\mathcal{V}$ . To that end, the output  $\mathbf{S}^{(L)} \in \mathbb{R}^{r \times n}$  of the final layer is projected back to the vocabulary size by an "unembedding layer"  $\mathbf{W}_U \in \mathbb{R}^{T \times r}$  and averaged over the columns to obtain a vector in  $\mathbb{R}^T$ . A softmax layer is finally applied on top of it to obtain the probability distribution of the next token  $\mathbb{P}_{\Theta}(\cdot | \mathbf{S})$ . Formally, we have

$$\mathbb{P}_{\Theta}(\cdot \mid \mathbf{S}) = \operatorname{softmax}\left(\frac{1}{n\tau}\mathbf{W}_{U}\mathbf{S}^{(L)}\mathbb{1}_{n}\right) \in \Delta_{T}$$

*n* is the length (i.e., number of columns) of the input sequence **S** (and thus of the last layer output  $\mathbf{S}^{(L)}$ ),  $\boldsymbol{\Theta}$  denotes the parameters of the whole network that subsume the parameters of each layer and each block and  $\tau$  is the softmax temperature (Hinton, 2015).

### B.3 AUTOREGRESSIVE TRANSFORMER-BASED LLM

The architecture described above is used in most of the transformer-based autoregressive LLM (Anil et al., 2023; Brown et al., 2020; Dubey et al., 2024; Jiang et al., 2023). In the theoretical analysis of Section 4, and unless specified otherwise, we remain faithful to their practical implementation and only make the following mild assumption: we assume that the unembedding layer is bounded. The class of parameters and the class of neural networks it generates respectively writes

$$\mathcal{W} = \{ \mathbf{\Theta} \mid \|\mathbf{W}_U^{\top}\|_{2,1} \leq B_U \} \text{ and } \mathcal{F} = \{ f_{\mathbf{\Theta}} \mid \mathbf{\Theta} \in \mathcal{W} \}.$$

It should be noted that this assumption is significantly weaker than what is usually done in the literature (Edelman et al., 2022; Zhang et al., 2023b).

### C BACKGROUND ON MARKOV CHAINS

We recall below some important notions related to Markov chains based on (Paulin, 2015; Roberts & Rosenthal, 2004) and that will be used in our proofs.

### C.1 BASIC NOTIONS

Consider two distribution probabilities  $\mathbb{P}$  and  $\mathbb{Q}$  defined on a measurable space  $(\Omega, \mathcal{F})$ .

**Definition C.1.** *The total variation between*  $\mathbb{P}$  *and*  $\mathbb{Q}$  *is defined as* 

$$d_{\mathrm{TV}}(\mathbb{P},\mathbb{Q}) \coloneqq \sup_{A \in \mathcal{F}} |\mathbb{P}(A) - \mathbb{Q}(A)|$$

In the setting considered in the main paper, we consider Markov chains with finite discrete state space  $\Omega$ . In this section, we refer to  $\Omega$  as a general Polish space, whose elements are referred to as *states*.

Informally, a discrete-time, time-homogeneous Markov chain with state space  $\Omega$  is a sequence of random variables  $(\mathbf{X}_1, \mathbf{X}_2, \ldots)$  taking values in  $\Omega$  such that the next observation is independent of the past given the present. This property is referred to as the Markov property and is defined below.

**Definition C.2.** A sequence of random variables  $(\mathbf{X}_1, \mathbf{X}_2, ...)$  is said to satisfy the Markov property if for all  $n \ge 1$  and any  $(x_1, ..., x_{n+1}) \in \Omega^{n+1}$ 

$$\mathbb{P}(\mathbf{X}_{n+1} = x_{n+1} \mid \mathbf{X}_n = x_n, \cdots, \mathbf{X}_1 = x_1) = \mathbb{P}(\mathbf{X}_{n+1} = x_{n+1} \mid \mathbf{X}_n = x_n).$$

To a given Markov chain, we associate its *transition kernel*  $\mathbf{Q} : \Omega^2 \to [0, 1]$  which collects the transition probabilities from one state to another

 $\forall n \in \mathbb{N}, (x, y) \in \Omega^2, \quad \mathbf{Q}(x, y) = \mathbb{P}(\mathbf{S}_{n+1} = y \mid \mathbf{S}_n = x).$ 

In the main text, we refer to  $\mathbf{Q}$  as a transition *matrix* as the Markov chains we consider are of finite state space.

**Definition C.3.** A distribution  $\pi$  on  $\Omega$  is said to be a stationary distribution if the action of the transition kernel leaves  $\pi$  unchanged, that is

$$(\mathbf{Q}\pi)(A) := \int_{y \in A} \mathbf{Q}(x, y) d\pi(x) = \pi(A)$$

for all  $A \in \mathcal{F}$ .

 A natural question is whether such a distribution exists for a generic Markov chain. Before stating an existence theorem, we introduce a classification of states below.

Class of states. All definitions bellow are borrowed from (Gallager, 1996)

**Definition C.4** (Accessibility and communication). A state x is accessible from y (abbreviated as  $x \to y$ ) if there exists n > 0 such that  $\mathbf{Q}^n(x, y) > 0$ . Two distinct states x and y communicate (abbreviated  $x \leftrightarrow y$ ) if x is accessible from y and y is accessible from x.

Accessibility and communication concepts define how states can reach each other within a Markov chain. This leads to an important classification of states into transient and recurrent categories.

**Definition C.5** (Recurrent and transient states). For finite-state Markov chains, a recurrent state is a state *i* that is accessible from all states that are accessible from *i* (*i* is recurrent if  $i \rightarrow j$  implies that  $j \rightarrow i$ ). A transient state is a state that is not recurrent.

With the distinction between recurrent and transient states established, we can now group states into classes based on their communication properties.

**Definition C.6** (Class of states). A class C of states is a non-empty set of states such that each  $i \in C$  communicates with every other state  $j \in C$  and communicates with no  $j \notin C$ 

**Aperiodicity and Ergodicity.** Aperiodicity ensures that the system does not exhibit cyclic behavior, which is a key condition for understanding the asymptotic behavior of states.

**Definition C.7** (Aperiodicity). The period of a state *i*, denoted d(i), is the greatest common divisor (gcd) of those values of *n* for which  $\mathbf{Q}^n(i, i) > 0$ . If the period is 1, the state is aperiodic.

Under some conditions on the Markov chain (aperiodicity and irreducibility (Roberts & Rosenthal, 2004)), it can be proven that the chain converges to its stationary distribution i.e. for any  $x \in \Omega$ ,  $\lim_{n\to\infty} d_{\rm TV}(Q^n(x,\cdot),\pi) = 0$ , where  $Q^n(x,\cdot)$  denotes the probability of  $S_n$  conditioned on  $S_1 = x$ .

1013 We recall below the notion of mixing time that assesses the time taken by the Markov chain to be  $\varepsilon$ -close to its stationary distribution (see Definition C.8).

**Definition C.8** (Mixing time for time-homogeneous Markov chains (Paulin, 2015)). Let  $X := (\mathbf{X}_1, \mathbf{X}_2, ...)$  be a time-homogeneous Markov chain with a state space  $\Omega$ , a transition kernel Q, and a stationary distribution  $\pi$ . Its mixing time is defined for any  $\varepsilon \in [0, 1]$  as

$$t_{\min}(\varepsilon) \coloneqq \min \left\{ t \mid d(t) \le \varepsilon \right\} \text{ where } d(t) \coloneqq \sup_{x \in \Omega} d_{\mathrm{TV}} \left( Q^t(x, \cdot), \pi \right).$$

We also introduce the quantity

$$t_{\min} \coloneqq \inf_{0 \le \varepsilon < 1} t_{\min} \left(\frac{\varepsilon}{2}\right) \cdot \left(\frac{2 - \varepsilon}{1 - \varepsilon}\right)^2$$

1026	which will be useful later on.
1027	<b>Remark C.1</b> (Well-posedness of $t_{min}$ ). As we only consider finite state-space Markov chains in
1028	our work, we know that a stationary distribution always exists. However, its uniqueness and the
1029	convergence to it require additional assumptions (see Appendix C.2). In particular, not all Markov
1030	chains admit a finite $t_{\min}(\varepsilon)$ , $t_{\min}$ for some $\varepsilon < \frac{1}{2}$ . In such case, $t_{\min}$ can be infinite. In our practical
1031	experimentation, this is never the case despite considering various Markov chains.
1032	
1033	C.2 ERGODIC UNICHAINS
1034	
1035	We are now ready to state the following theorem, which formalizes the classification of states into
1036	recurrent, transient, and aperiodic classes.
1037	
1038	Theorem C.9 (Recurrent and transient classes). For finite state Markov chains, either all states
1039	in a class are transient or all are recurrent. We refer to these classes as transient and recurrent,
1040	respectively.
1041	For any Markov chain all states in a class have the same period. If the period is 1, the class is
1042	said to be aperiodic
1043	
1044	Having categorized states into recurrent, transient, and aperiodic classes, we can now define ergodicity.
1045	
1046	Definition C 10 (Encodicity) Encodering of the state Manhamatania and an elicated a factoria
1047	<b>Definition C.10</b> (Ergodicity). For a finite-state Markov chain, an ergodic class of states is a class that is both recurrent and apariodic. A Markov chain consisting antirally of one areadic
1048	class that is both recurrent and aperiodic. A Markov chain consisting entirety of one ergodic
1049	ciuss is cuitea an ergoaic chain.
1050	
1051	Unichains. We now introduce the concept of unichains.
1052	
1053	Definition C.11 (Unichains and ergodic unichains). A unichain is a finite-state Markov chain
1054	containing a single recurrent class and transient states. An ergodic unichain is a unichain for
1055	which the recurrent class is ergodic.
1056	
1057	C.3 MARTON COUPLINGS
1058	
1059	While we consider Markov chain inputs in Section 4.2, we consider less structured inputs during the
1060	pre-training phase Section 4.1.
1061	More specifically, we model the sequences of tokens used during pre-training as generic dependent
1062	random variables. To derive meaningful results, we rely on the notion of Marton couplings introduced
1063	by Marton (2004). A Marton coupling can be seen as a weak dependency structure between random
1064	variables. The associated notion of the mixing matrix, analogous to the mixing time of a Markov
1065	chain, is used to assess the strength of the dependence between those variables.
1066	This minimal modeling choice is made to remain as faithful as possible to the pre-training considered
1067	in practical applications of LLMs, for which the pre-training data is not public and may contain
1068	arbitrary data points such as video, code snippets, text and images (Achiam et al., 2023; Anil et al.,
1069	2023; Brown et al., 2020; Dubey et al., 2024; Jiang et al., 2023; Touvron et al., 2023a).
1070	As shown in Paulin (2015, Pamark 2.2.), considering sequences of random variables linked through a
1071	As shown in raumi (2013, Kemaik 2.2.), considering sequences of random variables inked infolginal Marton coupling is a weaker assumption than what is usually done in the literature on generalization
1072	bounds which typically relies on independent random variables and Markov chains (Hu et al. 2024).
1073	Marion, 2023: Wolfer & Kontorovich. 2019: Zhang et al., 2023b).
1074	
1075	In particular, the results stated in Section 4.1 encompass the case where the pre-training input
1076	sequences of lokens are independent random variables (Kim et al., 2024) or Markov chains (Zhang et al. 2022b) We also note that Markov chains can model bigging and the restriction of the second second second
1077	et al., 20230), we also note that ivial to chains can model bigrams used in natural language (Bletti et al. 2023; Jurafsky & Martin 2024)
1078	$(1 a_1, 2023, Juraisky & Waltin, 2024).$
1070	We do not provide an exhaustive review of Marton couplings. We will simply recall its definition and

(2015). Consider a sequence of dependent random variables  $S = (\mathbf{S}_1, \dots, \mathbf{S}_N)$  taking values in a polish space  $\Omega = \Omega_1 \times \dots \times \Omega_N$ . We will denote by  $\mathbb{P}(\mathbf{S}_1, \dots, \mathbf{S}_N)$  the distribution of S.

**Definition C.12** (Marton coupling). We define a Marton coupling for S as a set of couplings

$$\left(S^{(s_1,\ldots,s_i,s'_i)}, S^{\prime(s_1,\ldots,s_i,s'_i)}\right) \in \Omega \times \Omega,$$

for every  $i \in [N]$ , every  $s_1 \in \Omega_1, \ldots, s_i \in \Omega_i, s'_i \in \Omega_i$ , satisfying the following conditions.

(i) 
$$\mathbf{S}_{1}^{(s_{1},\ldots,s_{i},s_{i}')} = s_{1}, \qquad \dots, \qquad \mathbf{S}_{i}^{(s_{1},\ldots,s_{i},s_{i}')} = s_{i}, \\ \mathbf{S}_{1}^{\prime(s_{1},\ldots,s_{i},s_{i}')} = s_{1}, \qquad \dots, \qquad \mathbf{S}_{i-1}^{\prime(s_{1},\ldots,s_{i},s_{i}')} = s_{i-1}, \qquad \mathbf{S}_{i}^{\prime(s_{1},\ldots,s_{i},s_{i}')} = s_{i}'.$$

(ii) 
$$\left(\mathbf{S}_{i+1}^{(s_1,\ldots,s_i,s'_i)},\ldots,\mathbf{S}_N^{(s_1,\ldots,s_i,s'_i)}\right)$$
  
 $\sim \mathbb{P}(\mathbf{S}_{i+1},\ldots,\mathbf{S}_N \mid \mathbf{S}_1 = s_1,\ldots,\mathbf{S}_i = s_i),$   
 $\left(\mathbf{S}_{i+1}^{\prime(x_1,\ldots,x_i,x'_i)},\ldots,\mathbf{S}_N^{\prime(x_1,\ldots,x_i,x'_i)}\right)$   
 $\sim \mathbb{P}(\mathbf{S}_{i+1},\ldots,\mathbf{S}_N \mid \mathbf{S}_1 = x_1,\ldots,\mathbf{S}_{i-1} = x_{i-1},\mathbf{S}_i = x'_i).$ 

(iii) If 
$$x_i = x'_i$$
, then  $S^{(x_1,...,x_i,x'_i)} = S'^{(x_1,...,x_i,x'_i)}$ 

**Definition C.13** (Mixing matrix (Paulin, 2015)). For a Marton coupling, we define the mixing matrix  $\Gamma \in \mathbb{R}^{N \times N}$  as an upper diagonal matrix with

$$\forall 1 \leq i < j \leq N, \quad \begin{cases} & \Gamma_{i,i} := 1, \\ & \Gamma_{j,i} := 0 \\ & & \Gamma_{i,j} := \operatorname{sup}_{s_1,\dots,s_i,s'_i} \mathbb{P}\left[\mathbf{S}_j^{(s_1,\dots,s_i,s'_i)} \neq \mathbf{S'}_j^{(s_1,\dots,s_i,s'_i)}\right]. \end{cases}$$

For independent random variables, one can define a Marton coupling with a mixing matrix equal to the identity (see Paulin, 2015, Remark 2.2). In particular, it means that for independent variables, we have the operator norm of the mixing matrix equal to 1, i.e.,  $\|\Gamma\| = 1$ .

1115 C.4 AN (ALMOST) DISTANCE BETWEEN MARKOV CHAINS

In Theorem D.23, We state elementary properties of  $\mathcal{K}$  in the proposition below.

**Proposition C.14** (Properties of  $\mathcal{K}$ ).  $\mathcal{K}$  is an almost-distance between transition matrices in the sense that it satisfies the properties below:

1. Non-negativity. For any  $\Theta_1, \Theta_2, \mathcal{K}(\Theta_1, \Theta_2) \geq 0$ .

2. Almost sure positivity.  $\mathcal{K}(\Theta_1, \Theta_2) = 0 \iff \forall n \in [N], \mathbb{P}_{\Theta_1}(\cdot \mid \mathbf{S}_n) = \mathbb{P}_{\Theta_2}(\cdot \mid \mathbf{S}_n) a.s.$ 

3. Symmetry. For any  $\Theta_1, \Theta_2, \mathcal{K}(\Theta_1, \Theta_2) = \mathcal{K}(\Theta_1, \Theta_2).$ 

4. Triangular inequality. For any  $\Theta_1, \Theta_2, \Theta_3, \mathcal{K}(\Theta_1, \Theta_3) \leq \mathcal{K}(\Theta_1, \Theta_2) + \mathcal{K}(\Theta_2, \Theta_3).$ 

<sup>1128</sup> *Proof of Proposition C.14.* We first recall the following technical lemma.

**Lemma C.15** (Proposition 2.16 in Folland (1999)). Let Y be a non-negative random variable defined on a probability space  $\Omega$  with probability function  $\mathbb{P}$ . If  $\mathbb{E}[Y] = 0$ , then Y = 0 almost

1134 surely, i.e., 1135  $\mathbb{P}(\{\omega \in \Omega \mid Y(\omega) = 0\}) = 1$ 1136 1137 The non-negativity and symmetry of  $\mathcal{K}$  directly come from the symmetry and non-negativity of the 1138 total variation distance. The triangular inequality follows from the fact that the total variation is a 1139 distance and that the expectation respects inequalities. For the almost positivity, consider  $\Theta_1, \Theta_2$ such that  $\mathcal{K}(\Theta_1, \Theta_2) = 0$ . By non-negativity of all the terms in the sum, it means that for all  $n \in [N]$ , 1140 we have 1141  $\mathbb{E}_{\mathbf{S}_n}[d_{\mathrm{TV}}(\mathbb{P}_{\Theta_1}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\Theta_2}(\cdot \mid \mathbf{S}_n))] = 0.$ 1142 As the total variation is a distance, we know that the random variable under the expectation is 1143 non-negative. Applying Lemma C.15 leads to 1144 1145  $d_{\mathrm{TV}}(\mathbb{P}_{\Theta_1}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\Theta_2}(\cdot \mid \mathbf{S}_n)) = 0$  almost surely. 1146 On the probability space, deprived of the set where the distance is non-zero (which is of null measure), 1147 the total variation is equal to zero and as a distance between probability distributions, it means that 1148 on this subset of the probability space, the probabilities are equal. Again, as the set on which they are 1149 not equal is of null measure, we have 1150  $\mathbb{P}_{\Theta_1}(\cdot \mid \mathbf{S}_n) = \mathbb{P}_{\Theta_2}(\cdot \mid \mathbf{S}_n)$  almost surely. 1151 1152 Putting everything together, we have 1153  $\forall n \in [N], \mathbb{P}_{\Theta_1}(\cdot \mid \mathbf{S}_n) = \mathbb{P}_{\Theta_2}(\cdot \mid \mathbf{S}_n)$  a.s., (5)1154 which concludes the direct sense. The converse sense is proved by assuming that Eq. (5) holds and 1155 using the distance properties of the total variation. This concludes the proof. 1156 1157 1158 D PROOFS 1159 1160 D.1 PROOF OF PROPOSITION 3.1 1161 We detail below the proof of Proposition 3.1. 1162 1163 Proof of Proposition 3.1. Step 1: large language models as Markov chains. Given an input 1164  $v_i \in \mathcal{V}_K^*$  of p tokens, an large language model outputs a probability mass function  $f_{\Theta}^{T,K}(v_i)$  over 1165 the discrete vocabulary space. As the temperature is positive, i.e.,  $\tau > 0$ , and as the exponential is 1166 positive, we know that all the tokens in the vocabulary will be given a positive mass. 1167 1168 A next sequence  $v_j \in \mathcal{V}_K^*$  is then sampled according to  $f_{\Theta}^{T,K}(v_i)$ . But the  $v_j$  sequences that fit 1169 necessarily contain the  $v_i$  sequence (except possibly the first element of  $v_i$ , thanks to Definition B.2), 1170 i.e.  $\forall l, (v_i)_l = (v_i)_{l+1}$ . Note also the size of  $v_i$  is p+1 when p < k and k when p = k. All other 1171 sequences  $v_i$  that do not satisfy this condition are not suitable. 1172 In that sense,  $f_{\Theta}^{T,K}$  can be represented by a Markov chain  $MC(\mathcal{V}_K^*, \mathbf{Q}_f)$  with transition kernel  $\mathbf{Q}_f \in \mathbb{R}^{|\mathcal{V}_K^*| \times |\mathcal{V}_K^*|}$ , as defined in Proposition 3.1. 1173 1174 1175 Step 2: Proportion of non-zero elements. We denote by  $\mathcal{R}$  the set of states of length K. The set 1176 of states of length strictly less than equal K is denoted by  $\mathscr{T}$ . We can construct a transition matrix 1177  $P_{\mathscr{R}} \in \mathbb{R}^{T^{K} \times T^{K}}$  with the states of this class, containing the probabilities of moving from one state of

1179 1180 1181

1178

$$\mathbf{Q}_f = \left( \begin{array}{c|c} P_{\mathscr{T}} & P_{\mathscr{T}\mathscr{R}} \\ \hline 0 & P_{\mathscr{R}} \end{array} \right). \tag{6}$$

1182 1183

Now, let us count the number of non-zero elements in each of these 4 large blocks.

 $\frac{1186}{1187} \qquad \frac{P_{\mathscr{T}} \text{ block : }}{\text{green blocks contained in } P_{\mathscr{T}}. \text{ The block number } i \in [K-2] \text{ is of size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ is of size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ is of size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ is of size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ is of size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ is of size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ is of size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ is of size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ so for a size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ so for a size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ so for a size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ so for a size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ so for a size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ so for a size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ so for a size } T^i \times T^{i+1}. \text{ Since each } i \in [K-2] \text{ so for a size } T^i \times T^i \text{ so for a$ 

 $\mathscr{R}$  to another.  $P_{\mathscr{R}}$  corresponds to the blue block in Fig. 1 while green rectangle blocks correspond to

part of  $P_{\mathcal{T}}$  and  $P_{\mathcal{TR}}$  in the following description of large language models as Markov chains,

sentence of size *i* can be completed with non-zero probability, by any other token, there are a total of  $\sum_{p=1}^{T^i} T = T^{i+1}$  non-zero elements. There are therefore  $\sum_{i=1}^{K-2} T^{i+1}$  non-zero elements in the entire  $P_{\mathcal{T}}$  block.

1192  $\frac{P_{\mathcal{TR}} \text{ block}:}{P_{\mathcal{TR}} \text{ block}:}$  The size of this block is  $\left[\frac{T}{T-1}(T^{K-1}-1)\right] \times T^{K}$ . The green block contained in  $P_{\mathcal{TR}}$  that contains non-zero elements is of size  $T^{K-1} \times T^{K}$ . Since each sentence of size K-1 can be completed with non-zero probability, by any other token, there are a total of  $\sum_{p=1}^{T^{K-1}} T = T^{K}$  non-zero elements.

1197 1198 1199 1199 1200 1201  $\frac{P_{\mathscr{R}} \text{ block}:}{P_{\mathscr{R}} \text{ block}:} \text{ The size of this block is } T^K \times T^K. \text{ Each sentence } v = (v_1, \ldots, v_K) \text{ of size } K \text{ is mapped to another sentence } v' = (v'_1, \ldots, v'_K) \text{ of size } K \text{ with non-zero probability, if and only if } v'_1 = v_2, v'_2 = v_3, \ldots, v'_{k-1} = v_K. \text{ The final token } v'_K \text{ can by any other token in the vocabulary. It means that there are a total of } \sum_{p=1}^{T^K} T = T^{K+1} \text{ non-zero elements.}$ 

 $\frac{1202}{1203} \qquad \underline{0's \text{ block : }} \text{ Obviously, there are no non-zero elements in this block.}$ 

1204 Finally, there are

1205

1207

1210 1211

1212 1213

1215 1216 1217

1218

$$\sum_{i=1}^{K-2} T^{i+1} + T^K + T^{K+1} = \sum_{i=1}^{K} T^{i+1} = T^2 \left( \frac{T^K - 1}{T - 1} \right)$$

non-zero elements. This means that the proportion of non-zero elements in the matrix is exactly

$$\frac{T^2\left(\frac{T^K-1}{T-1}\right)}{\left(T\left(\frac{T^K-1}{T-1}\right)\right)^2} = \frac{T-1}{T^K-1}.$$

1214 Note that for large T and K we have that

$$\frac{T-1}{T^K-1} \sim \frac{1}{T^{K-1}}.$$

#### 1219 1220 D.2 PROOF OF PROPOSITION 3.2

We begin with a preliminary lemma.

**Lemma D.1** (Powers of  $\mathbf{Q}_f$  greater than K). For any initial state *i*, the following hold:

1225 1226 1227

1223

1224

1228 1229

•  $\forall k \ge K, \forall j \in \mathscr{T}, (\mathbf{Q}_f^k)_{i,j} = 0,$ 

• 
$$\forall k \geq K, \forall j \in \mathscr{R}, (\mathbf{Q}_f^k)_{i,j} > 0.$$

*Proof.* By considering  $\mathbf{Q}_f$  as defined in (6), we can compute its powers. For any  $k \geq 1$ ,

$$\mathbf{Q}_{f}^{k} = \frac{\begin{pmatrix} P_{\mathscr{T}}^{k} & B_{k} \\ \\ 0 & P_{\mathscr{R}}^{k} \end{pmatrix},$$

where  $B_k = \sum_{m=0}^{k-1} P_{\mathscr{T}}^m P_{\mathscr{T}} P_{\mathscr{T}} P_{\mathscr{T}}^{k-1-m}$ .

To prove the first item, we focus on the blocks on the left of  $\mathbf{Q}_f$ . Since the lower left block is zero, we have that  $\forall k \geq 1, \forall i \in \mathscr{R}, \forall j \in \mathscr{T}, (\mathbf{Q}_f^k)_{i,j} = 0$ . In the upper left block, the element  $(P_{\mathscr{T}}^k)_{i,j}$ designates the probability of moving from one transient state  $i \in \mathscr{T}$  to another transient state  $j \in \mathscr{T}$ after k iterations. According to Definition B.1, if state i is a sequence of  $p \geq 1$  tokens, state j is necessarily a sequence of min $\{K, p + K\} = K$  elements. Thus,  $P_{\mathscr{T}}$  is a nilpotent matrix and  $\forall k \geq K, \forall i, j, (P_{\mathscr{T}}^k)_{i,j} = 0$ . This proves that  $\forall k \geq K, \forall j \in \mathscr{T}, (\mathbf{Q}_f^k)_{i,j} = 0$ . 1242 We now move on to the second item. From the above,  $\forall k \geq K, B_k = \sum_{m=0}^{K-1} P_{\mathscr{T}}^m P_{\mathscr{T}\mathscr{R}} P_{\mathscr{R}}^{k-1-m}$ . 1243 Note that this sum is finite, but there is still a dependence on k, in the powers of the matrix  $P_{\mathscr{R}}$ . In 1244 the lower right block, the element  $(P_{\mathscr{R}}^k)_{i,j}$  designates the probability of moving from one recurrent 1245 state  $i \in \mathscr{R}$  to another recurrent state  $j \in \mathscr{R}$  after k iterations. According to the definition of  $\mathbf{Q}_f$  in 1246 Proposition 3.1 and Definition B.2,  $\forall k \geq K, \forall i, j \in \mathscr{R}^2, (P_{\mathscr{R}}^k)_{i,j}$  is nonzero. Exploiting this also in 1247  $B_k$ , we obtain the result, i.e.  $\forall k \geq K, \forall j \in \mathscr{R}, (\mathbf{Q}_f^k)_{i,j} > 0$ .

We are now ready to prove Proposition 3.2, which is inspired by Gallager (1996).

1251 *Proof of Proposition 3.2.* The states of length strictly less than equal K (elements of  $\mathscr{T}$ ) are transient, 1252 because of Definition B.1. To discuss the nature of states of length K (elements of  $\mathcal{R}$ ), let us introduce 1253 a result regarding the powers of the  $Q_f$  matrix as defined in (6). Thanks to Lemma D.1, the set of states  $\mathscr{R}$  (i.e. the states of length K) form a class. Lemma D.1 gives us also that  $\mathscr{R}$  is ergodic. In 1254 fact, every state in  $\mathscr{R}$  only communicates with all the other states in  $\mathscr{R}$ , which proves the recurrence. 1255 Since  $\forall i, j \in \mathscr{R}^2$ ,  $(\mathbf{Q}_F^K)_{i,j} > 0$ , we can move between any two states in exactly K steps, regardless 1256 of the initial position. This ensures that  $\mathscr{R}$  is aperiodic because the transition probabilities do not 1257 depend on a specific cycle, and states can be revisited at various time steps, not just multiples of 1258 a particular number. More simply, by considering a token x, the state defined as  $i = xx \dots x$  has 1259 K times 1260

period 1, i.e.  $(\mathbf{Q}_f)_{i,i} > 0$ . This is a consequence of Definition B.2 and Proposition 3.1. Thanks to Theorem C.9, it means that the whole class  $\mathscr{R}$  is aperiodic. Finally, this means that  $MC(\mathcal{V}_K^*, \mathbf{Q}_f)$  are ergodic unichains, in the sense of Definition C.11.

# 1264 D.3 PROOF OF PROPOSITION 3.3

We start by introducing three technical lemmas that will be useful in the proof of Proposition 3.3. We start with the Chapman–Kolmogorov equation.

Lemma D.2 (Chapman-Kolmogorov equation). Let P be a matrix of size d. Then, P satisfies

$$\forall i, j \in [d]^2, \forall n, n' \in \mathbb{N}^2, (P^{n+n'})_{i,j} = \sum_{k=1}^d (P^n)_{i,k} (P^{n'})_{k,j}$$

1275 *Proof.* The result follows from the fact that  $\forall n, n' \in \mathbb{N}^2$ ,  $P^{n+n'} = P^n P^{n'}$ .

Then, we introduce a simple but useful result of monotonicity.

**Lemma D.3** (Lemma 3.3.1. in Gallager (1996)). Let the transition matrix P of a finite state Markov chain. Then, for all states j and  $n \ge 1$ , we have

$$\max_{i} (P^{n+1})_{i,j} \le \max_{i} (P^n)_{i,j}, \quad and \quad \min_{i} (P^{n+1})_{i,j} \ge \min_{i} (P^n)_{i,j}$$

We now refer to a result on Markov chains with positive transition matrices.

**Lemma D.4** (Lemma 3.3.2. in Gallager (1996)). Let the transition matrix P of a finite state Markov chain satisfy  $\forall i, j, P_{i,j} > 0$ , and let  $\alpha = \min_{i,j} P_{i,j} > 0$ . Then, for all states j and  $n \ge 1$ ,

$$\max_{i} (P^{n})_{i,j} - \min_{i} (P^{n})_{i,j} \le (1 - 2\alpha) \left( \max_{i} (P^{n})_{i,j} - \min_{i} (P^{n})_{i,j} \right), \max_{i} (P^{n})_{i,j} - \min_{i} (P^{n})_{i,j} \le (1 - 2\alpha)^{n}, \lim_{n \to \infty} \max_{i} (P^{n})_{i,j} = \lim_{n \to \infty} \min_{i} (P^{n})_{i,j} > 0.$$

1295

1268 1269

1270

1272 1273 1274

1277

1278 1279

1280

1285 1286

1287

We are now ready to prove Proposition 3.3 using a similar argument as in Gallager (1996).

*Proof of Proposition 3.3.* Let  $\mathcal{T}$  and  $\mathcal{R}$  denote respectively the sets of transient and recurrent states. For any state j, we define  $\pi_j := \lim_{n \to \infty} \max_i (\mathbf{Q}_f^n)_{i,j} = \lim_{n \to \infty} \min_i (\mathbf{Q}_f^n)_{i,j}$ . Then  $\pi =$  $(\pi_i)_{i \in \Omega}$  is the stationary distribution for  $\mathbf{Q}_f$ . 

**Step 1: Stationary distribution for transient states.** Lemma D.1 gives us that  $\forall i, \forall k \ge K, \forall j \in K, j \in K, \forall j \in K, \forall j \in K, \forall j \in K, j \in K, j \in K, j \in$  $\mathscr{T}, (\mathbf{Q}_{f}^{k})_{i,j} = 0$ . This means that  $\forall j \in \mathscr{T}, \pi_{j} = 0$  and hence the limit is reached at most after K iteration. 

Step 2: Stationary distribution for recurrent states. Lemma D.1 gives us  $\forall i, j \in \mathscr{R}^2, (\mathbf{Q}_f^K)_{i,j} >$ 0. By defining  $\varepsilon := \min_{i,j \in \mathscr{R}^2} (\mathbf{Q}_f^K)_{i,j}$ , Lemma D.4 shows that for any integer  $\ell \ge 1$ , 

 $\max_{i \in \mathscr{R}} \left( \mathbf{Q}_f^{\ell K} \right)_{i,j} - \min_{i \in \mathscr{R}} \left( \mathbf{Q}_f^{\ell K} \right)_{i,j} \le (1 - 2\varepsilon) \left( \max_{i \in \mathscr{R}} \left( \mathbf{Q}_f^{\ell K} \right)_{i,j} - \min_{i \in \mathscr{R}} \left( \mathbf{Q}_f^{\ell K} \right)_{i,j} \right),$ (7)

$$\max_{i \in \mathscr{R}} \left( \mathbf{Q}_f^{\ell K} \right)_{i,j} - \min_{i \in \mathscr{R}} \left( \mathbf{Q}_f^{\ell K} \right)_{i,j} \le (1 - 2\varepsilon)^\ell, \tag{8}$$

$$\lim_{\ell \to \infty} \max_{i \in \mathscr{R}} \left( \mathbf{Q}_f^{\ell K} \right)_{i,j} = \lim_{\ell \to \infty} \min_{i \in \mathscr{R}} \left( \mathbf{Q}_f^{\ell K} \right)_{i,j} > 0.$$
(9)

Thanks to Lemma D.3,  $\max_{i} (\mathbf{Q}_{f}^{n+1})_{i,j}$  is non-decreasing in n, so the limit on the left in Eq. (9) can be replaced with a limit in n. The same argument for the limit on the right gives that,  $\forall j \in \mathscr{R}$ , 

$$\max_{i \in \mathscr{R}} (\mathbf{Q}_f^n)_{i,j} - \min_{i \in \mathscr{R}} (\mathbf{Q}_f^n)_{i,j} \le (1 - 2\varepsilon)^{\lfloor n/K \rfloor},$$

$$\lim_{n \to \infty} \max_{i \in \mathscr{R}} (\mathbf{Q}_f^n)_{i,j} = \lim_{n \to \infty} \min_{i \in \mathscr{R}} (\mathbf{Q}_f^n)_{i,j} > 0,$$

where we have taken the floor function to also convert the result of (8). Since  $\pi_i$  lies between the minimum and maximum  $(\mathbf{Q}_{f}^{n})_{i,j}$  for each n, we have that  $\forall i, j \in \mathscr{R}^{2}$ , 

$$\left| (\mathbf{Q}_{f}^{n})_{i,j} - \pi_{j} \right| \leq (1 - 2\varepsilon)^{\lfloor \frac{n}{K} \rfloor}$$

It means that  $\forall i, j \in \mathscr{R}^2, \pi_j = \lim_{n \to \infty} (\mathbf{Q}_f^n)_{i,j}$ . This also gives us the convergence rate when the initial state i is recurrent. In the next step, we consider the general convergence rate, regardless of the nature of the initial state *i*. 

Step 3: Convergence bound. We proceed to the remaining case, i.e. the case where the initial state  $i \in \mathscr{T}$  and the final state  $j \in \mathscr{R}$ . Lemma D.2 says that  $\forall n \geq K$ , 

$$(\mathbf{Q}_f^n)_{i,j} = \sum_{k \in \mathscr{T}} (\mathbf{Q}_f^K)_{i,k} (\mathbf{Q}_f^{n-K})_{k,j} + \sum_{k \in \mathscr{R}} (\mathbf{Q}_f^K)_{i,k} (\mathbf{Q}_f^{n-K})_{k,j}.$$

We then have that  $\forall i \in \mathscr{T}, \forall n \in \mathbb{N}$ , 

$$\begin{aligned} |(\mathbf{Q}_{f}^{n})_{i,j} - \pi_{j}| &\leq \Big| \sum_{k \in \mathscr{T}} (\mathbf{Q}_{f}^{K})_{i,k} \big[ (\mathbf{Q}_{f}^{n-K})_{k,j} - \pi_{j} \big] + \sum_{k \in \mathscr{R}} (\mathbf{Q}_{f}^{K})_{i,k} \big[ (\mathbf{Q}_{f}^{n-K})_{k,j} - \pi_{j} \big] \Big| \\ &\leq \sum_{k \in \mathscr{T}} (\mathbf{Q}_{f}^{K})_{i,k} \big| (\mathbf{Q}_{f}^{n-K})_{k,j} - \pi_{j} \big| + \sum_{k \in \mathscr{R}} (\mathbf{Q}_{f}^{K})_{i,k} \big| (\mathbf{Q}_{f}^{n-K})_{k,j} - \pi_{j} \big| \end{aligned}$$

$$\leq \sum_{k \in \mathscr{T}} (\mathbf{Q}_f^K)_{i,k} + \sum_{k \in \mathscr{R}} (\mathbf{Q}_f^K)_{i,k} \big| (\mathbf{Q}_f^{n-K})_{k,j} - \pi_j \big|$$

$$\leq (1-2\varepsilon)^{\lfloor \frac{n-\kappa}{K} \rfloor},$$

where the first sum vanishes and  $\sum_{k \in \mathscr{R}} (\mathbf{Q}_f^K)_{i,k} \leq 1$ . Finally,  $\forall i \in \mathscr{T}, \forall n \geq K$ , 

 $\left| (\mathbf{Q}_{f}^{n})_{i,j} - \pi_{j} \right| \leq (1 - 2\varepsilon)^{\left\lfloor \frac{n}{K} \right\rfloor - 1}.$ 

Combining this with the result of Step 2 concludes the proof. 

#### D.4 PROOF OF THEOREM 4.1

In this section, we detail the proof of Theorem 4.1. We provide below an overview of the proof before detailing it.

**Overview of the proof.** We are going to use McDiarmid's inequality for dependent random variables of Paulin (2015, Theorem 2.9). To adapt the arguments of Paulin (2015, Theorem 2.9) to our setting, we bound the total variation between the true probability of the next token and the one estimated by the LLM. The rest of this section is organized as follows. First in Appendix D.4.1, we adapt the concentration inequality of Paulin (2015, Theorem 2.9). Then in Appendix D.4.2, we show how to bound the total variation between the true and the estimated probability of the next token. , in Appendix D.4.3, we restate Theorem 4.1 and conclude the proof.

1358 D.4.1 CONCENTRATION INEQUALITIES FOR DEPENDENT RANDOM VARIABLES

We first state a concentration inequality for dependent random variables that will be used to obtain our final bound.

**Proposition D.5** (McDiarmid's inequality for dependent random variables). Let  $S := (\mathbf{S}_1, \ldots, \mathbf{S}_N)$  be a sequence of random variables that take values in  $\Omega = \Omega_1 \times \ldots \times \Omega_N$ . Assume there exists a Marton coupling for S with mixing matrix  $\mathbf{\Gamma}$ . Let  $\|\mathbf{\Gamma}\|$  be the operator norm of  $\mathbf{\Gamma}$ . If  $f: \Omega \to \mathbb{R}$  is such that there exists  $\mathbf{c} \in \mathbb{R}^N$  satisfying

$$\forall \mathbf{x}, \mathbf{y} \in \Omega, \quad f(\mathbf{x}) - f(\mathbf{y}) \le \sum_{i=1}^{N} \mathbf{c}_{i} \mathbb{1}_{\{\mathbf{x}_{i} \neq \mathbf{y}_{i}\}},$$

then we have for any  $u \ge 0$ ,

1357

1363

1364

1365

1367

1369 1370

1371 1372

1373 1374 1375

1376

1393 1394

1400 1401 1402

$$\mathbb{P}(|f(S) - \mathbb{E}_S[f(S)]| \ge u) \le 2 \exp\left(\frac{-2u^2}{\|\mathbf{\Gamma}\|^2 \|\mathbf{c}\|_2^2}\right)$$

*Proof.* Consider a function f verifying the properties of Proposition D.5. Paulin (2015, Theorem 2.9) ensures that for a partition  $\hat{S}$  of S (see Paulin, 2015, Definition 2.3) the following inequality holds

$$\forall u \ge 0, \quad \mathbb{P}\Big(|f(\hat{S}) - \mathbb{E}\Big[f(\hat{S})\Big]| \ge u\Big) \le 2\exp\left(\frac{-2u^2}{\|\mathbf{\Gamma} \cdot C(\mathbf{c})\|_2^2}\right),\tag{10}$$

1381 where  $C(\mathbf{c})$  is a vector of  $\mathbb{R}^N$  whose *i*-th element is the sum of the  $\mathbf{c}_j$  such that *j* is an index of 1382 the elements of  $\hat{\mathbf{S}}_i$ . Taking the trivial partition  $\hat{S} = S$  implies that the index of the elements in  $\hat{\mathbf{S}}_i$ 1383 are reduced to  $\{i\}$ . Hence the *i*-th entry of  $C(\mathbf{c})$  is equal to  $\mathbf{c}_i$  and  $C(\mathbf{c}) = \mathbf{c}$ . By definition of the 1384 operator norm (naturally induced by the  $\ell_2$ -norm), we have

$$\|\mathbf{\Gamma} \cdot \mathbf{c}\|_2 = \frac{\|\mathbf{\Gamma}\mathbf{c}\|_2}{\|\mathbf{c}\|_2} \cdot \|\mathbf{c}\|_2 \leq \underbrace{\sup_{\mathbf{x}\neq 0} \frac{\|\mathbf{\Gamma}\mathbf{x}\|_2}{\|\mathbf{x}\|_2}}_{=\|\mathbf{\Gamma}\|} \cdot \|\mathbf{c}\|_2 \leq \|\mathbf{\Gamma}\| \cdot \|\mathbf{c}\|_2,$$

where the first inequality comes from the fact that c is non-zero (otherwise the only possible f is the zero function which is not of great interest). Using the fact that the function  $x \to \exp\left(-\frac{2u^2}{x}\right)$  is increasing, we obtain

 $\exp\left(\frac{-2u^2}{\|\mathbf{\Gamma}\cdot\mathbf{c}\|_2^2}\right) \le \exp\left(\frac{-2u^2}{\|\mathbf{\Gamma}\|^2 \cdot \|\mathbf{c}\|_2^2}\right),$ 

which concludes the proof.

By looking at the definition of the risk  $\widehat{\mathcal{R}}_{pre}(\Theta)$ , we can see that applying Proposition D.5 to the function

$$f \colon (\mathbf{S}_1, \dots, \mathbf{S}_{N_{\mathrm{train}}}) = \frac{1}{N_{\mathrm{train}}} \sum_{n=1}^{N_{\mathrm{train}}} d_{\mathrm{TV}}(\mathbb{P}_{\mathcal{L}}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\mathbf{\Theta}}(\cdot \mid \mathbf{S}_n)),$$

would lead to the desired bound as we already know S admits a Marton coupling with mixing matrix  $\Gamma$ . We investigate in the next section how to find the bounding vector c to apply Proposition D.5.

#### D.4.2 FINDING THE BOUNDING VECTOR

**Technical lemmas.** We first recall the following important notions from (Tsybakov, 2008). Let  $(\Omega, \mathcal{F})$  be a measure space and consider two probability distributions  $\mathbb{P}, \mathbb{Q}$  defined on  $(\Omega, \mathcal{F})$ . For any  $\sigma$ -finite measure  $\nu$  on  $(\Omega, \mathcal{F})$  such that  $\mathbb{P}, \mathbb{Q}$  are absolutely continuous with respect to  $\nu$ , we can define  $p = \frac{d\mathbb{P}}{d\nu}$ ,  $q = \frac{d\mathbb{Q}}{d\nu}$  which can also be written as  $\mathbb{P}(d\omega) = q(\omega)\nu(d\omega)$  and  $\mathbb{Q}(d\omega) = p(\omega)\nu(d\omega)$ . We will adopt both notations interchangeably. It should be noted that there always exists at least one such measure  $\nu$  as one can take  $\nu = \mathbb{P} + \mathbb{Q}$ . With these notations, the squared Hellinger distance between  $\mathbb{P}$  and  $\mathbb{Q}$  is defined as 

$$H(\mathbb{P},\mathbb{Q})^2 \coloneqq \int_{\omega \in \Omega} \left(\sqrt{p(\omega)} - \sqrt{q(\omega)}\right)^2 \nu(d\omega) = \int_{\omega \in \Omega} \left(\sqrt{\mathbb{P}(d\omega)} - \sqrt{\mathbb{Q}(d\omega)}\right)^2.$$

The lemma below shows that the total variation between two probability distributions is controlled from above by the absolute value of the logarithm of their ratio. 

**Lemma D.6.** Consider two probability distributions  $\mathbb{P}, \mathbb{Q}$  defined on a measure space  $(\Omega, \mathcal{F})$ and a  $\sigma$ -finite measure  $\nu$  on  $(\Omega, \mathcal{F})$ . Let p, q be the corresponding probabilities densities, i.e., we have  $\mathbb{P}(d\omega) = q(\omega)\nu(d\omega)$  and  $\mathbb{Q}(d\omega) = p(\omega)\nu(d\omega)$ , the total variation between  $\mathbb{P}$  and  $\mathbb{Q}$ satisfies

$$d_{\mathrm{TV}}(\mathbb{P}, \mathbb{Q}) \le \left(2\int_{\omega\in\Omega} \left|\log\sqrt{\frac{\mathbb{P}(d\omega)}{\mathbb{Q}(d\omega)}}\right| q(\omega)d\nu(d\omega)\right)^{1/2}$$

If there exists a non-negative constant B such that for any  $z \in \Omega$ ,  $\left| \log \sqrt{\frac{\mathbb{P}(z)}{\mathbb{O}(z)}} \right| \leq B$ , then we have  $d_{\mathrm{TV}}(\mathbb{P} \ \mathbb{O}) < \sqrt{2B}$ 

$$a_{1V}(\mathbb{I}, \mathbb{Q}) \leq \sqrt{2D}.$$

*Proof.* We have the following relation between the total variation and the Hellinger distance (cf. Tsybakov, 2008, Lem. 2.3, Chapt. 2, p. 86):

$$d_{\mathrm{TV}}(\mathbb{P},\mathbb{Q})^2 \le H(\mathbb{P},\mathbb{Q})^2 \cdot \left(1 - \underbrace{H(\mathbb{P},\mathbb{Q})^2/4}_{\ge 0}\right) \le H(\mathbb{P},\mathbb{Q})^2,\tag{11}$$

where the last inequality uses the positivity of the Hellinger distance. Inspired by the decomposition of the Hellinger distance in (Agarwal et al., 2020, Lem. 25), we have 

$$H(\mathbb{P}, \mathbb{Q})^{2} = \int_{\omega \in \Omega} \left( \sqrt{\mathbb{P}(d\omega)} - \sqrt{\mathbb{Q}(d\omega)} \right)^{2} = \int_{\omega \in \Omega} \left( \mathbb{P}(d\omega) + \mathbb{Q}(d\omega) - 2\sqrt{\mathbb{P}(d\omega)}\sqrt{\mathbb{Q}(d\omega)} \right)$$
$$= 2 \cdot \left( 1 - \int_{\omega \in \Omega} \sqrt{\mathbb{P}(d\omega)}\sqrt{\mathbb{Q}(d\omega)} \right) = 2 \cdot \left( 1 - \int_{\omega \in \Omega} \sqrt{\frac{\mathbb{P}(d\omega)}{\mathbb{Q}(d\omega)}} \mathbb{Q}(d\omega) \right)$$
$$= 2 \cdot \left( 1 - \int_{\omega \in \Omega} \sqrt{\frac{\mathbb{P}(d\omega)}{\mathbb{Q}(d\omega)}} q(\omega) d\nu(d\omega) \right) \qquad \text{(by definition of } \mathbb{Q}(d\omega)$$

$$\leq -2\log\left(\int_{\omega\in\Omega}\sqrt{\frac{\mathbb{P}(d\omega)}{\mathbb{Q}(d\omega)}}q(\omega)d\nu(d\omega)\right)$$

It follows using Eq. (11)

$$d_{\mathrm{TV}}(\mathbb{P},\mathbb{Q})^2 \le H(\mathbb{P},\mathbb{Q})^2$$

$$\leq 2 \int_{\mathbb{T} \subseteq \Omega} -\log\left(\sqrt{\frac{\mathbb{P}(d\omega)}{\mathbb{Q}(d\omega)}}\right) q(\omega) d\nu(d\omega)$$

$$\leq 2 \int_{\omega \in \Omega} -$$

(by Jensen as 
$$-\log$$
 is convex)

(by definition of  $\mathbb{Q}(d\omega)$ )

 $(\text{using } 1 - x < -\log(x))$ 

1455  
1456  
1457
$$\leq 2 \left| \int_{\omega \in \Omega} -\log\left(\sqrt{\frac{\mathbb{P}(d\omega)}{\mathbb{Q}(d\omega)}}\right) q(\omega) d\nu(d\omega) \right|$$

$$\begin{aligned} & 1458 \\ & 1459 \\ & 1460 \\ & 1460 \\ & 1461 \\ & 1462 \\ & 1462 \\ & 1463 \end{aligned} \qquad & \leq 2 \int_{\omega \in \Omega} \left| \log \left( \sqrt{\frac{\mathbb{P}(d\omega)}{\mathbb{Q}(d\omega)}} \right) \right| q(\omega) d\nu(d\omega) \qquad \text{(by Jensen as } |\cdot| \text{ is convex}) \\ & \text{(by Jensen as } |\cdot| \text{ is convex}) \\ & \text{(first part of Lemma D.6)} \end{aligned}$$

$$\leq 2 \int_{\omega \in \Omega} \underbrace{\left| \log \left( \sqrt{\frac{\mathbb{P}(d\omega)}{\mathbb{Q}(d\omega)}} \right) \right|}_{\leq B} q(\omega) d\nu(d\omega) \qquad \text{(first part of Lemma D.6)}$$
$$\leq 2B \underbrace{\int_{\omega \in \Omega} q(\omega) d\nu(d\omega)}_{\leq 2B} \leq 2B. \qquad \text{(second part of Lemma D.6)}$$

This concludes both parts of the proof. 

> The next lemma provides a lower bound on the softmax output if its input is upper-bounded (in  $\ell_1$ -norm).

**Lemma D.7.** Let 
$$\mathbf{x} \in \mathbb{R}^m$$
 be such that  $\|\mathbf{x}\|_1 \leq c_1$  for some  $c_1 > 0$ . Then, we have

$$\operatorname{softmax}(\mathbf{u}) \ge \frac{1}{m \exp\left(2c_1\right)}$$

where the inequality holds for each component of  $softmax(\mathbf{u})$ .

*Proof.* Using the fact that

$$\|\mathbf{x}\|_1 = \sum_{i=1}^m |\mathbf{x}_i| \le c_1$$

we know that for any  $i \in [m]$ , we have 

$$-c_1 \leq \mathbf{x}_i \leq c_1.$$

Hence, using the fact that the exponential is increasing, we have for any  $i \in [m]$ 

e

$$\exp\left(-c_{1}\right) \leq \exp\left(\mathbf{x}_{i}\right) \leq \exp\left(c_{1}\right).$$
(12)

Summing and taking the inverse leads to

$$\sum_{i=1}^{m} \exp\left(-c_{1}\right) \leq \sum_{i=1}^{m} \exp\left(\mathbf{x}_{j}\right) \leq \sum_{i=1}^{m} \exp\left(c_{1}\right)$$

$$\iff \frac{1}{\sum_{j=1}^{m} \exp\left(c_{1}\right)} \leq \frac{1}{\sum_{j=1}^{m} \exp\left(\mathbf{x}_{j}\right)} \leq \frac{1}{\sum_{j=1}^{m} \exp\left(-c_{1}\right)}.$$
(13)

Combining Eq. (12) and Eq. (13) yields 

1498  
1499  
1500
$$\frac{\exp(-c_1)}{\sum_{j=1}^{m} \exp(c_1)} \le \frac{\exp(\mathbf{x}_i)}{\sum_{j=1}^{m} \exp(\mathbf{x}_j)} \le \frac{\exp(c_1)}{\sum_{j=1}^{m} \exp(-c_1)}$$

As we desire a lower bound, we only focus on the left-hand side of the previous inequality. Multiplying the numerator and denominator by  $\exp(c_1)$  leads to 

$$\forall i \in [m], \quad \text{softmax}(\mathbf{x})_i = \frac{\exp(\mathbf{x}_i)}{\sum_{j=1}^m \exp(\mathbf{x}_j)} \ge \frac{1}{m \exp(2c_1)},$$

which concludes the proof.

**Upper-bounding the total variation.** We now proceed with finding an upper bound on the total variation between the true probability of the next token and the one estimated by the LLM  $f_{\Theta}$ . It will enable us to find the bounding vector c. The next lemma shows that the input of the softmax layer of the model is bounded. 

**Lemma D.8.** Consider an LLM  $f_{\Theta} \in \mathcal{F}$ . For any input sequence  $\mathbf{S} \in \mathbb{R}^{r \times n}$ , the following inequality holds

$$\|\frac{1}{n\tau}\mathbf{W}_U\mathbf{S}^{(L)}\mathbb{1}_n\|_1 \le \frac{1}{\tau}\|\mathbf{W}_U^{\top}\|_{2,1}$$

where  $\tau$  is the temperature,  $\mathbf{W}_U$  is the unembedding matrix (which is bounded as stated in the definition of the parameters space W), and  $\mathbf{S}^{(L)}$  is the output of the last transformer layer.

*Proof.* We recall that the layer normalization ensures that at each layer, the tokens are in the unit  $\ell_2$ -ball. This is, in particular, the case for the output of the last layer  $\mathbf{S}^{(L)}$ . It means that the columns of  $\mathbf{S}^{(L)}$  verifies ( )

$$\forall k \in [n], \quad \|\mathbf{S}_{\cdot,k}^{(L)}\|_2 \le 1,$$
(14)

which implies

$$\max_{1 \le k \le n} \|\mathbf{S}_{\cdot,i}^{(L)}\|_2 \le 1.$$
(15)

Recalling that the  $L_{p,q}$ -norm of a matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$  can be rewritten as

$$\|\mathbf{A}\|_{p,q} \coloneqq \left(\sum_{j=1}^{m} \left(\sum_{i=1}^{n} |\mathbf{A}_{ij}|^{p}\right)^{\frac{q}{p}}\right)^{\frac{1}{q}} = \|(\|\mathbf{A}_{\cdot,j}\|_{p})_{j=1}^{m}\|_{q},$$
(16)

the  $\ell_1$ -norm of the last layer before the softmax layer satisfies 

$$\begin{aligned} \|\frac{1}{n\tau} \mathbf{W}_{U} \mathbf{S}^{(L)} \mathbf{1}_{n} \|_{1} &= \frac{1}{n\tau} \sum_{i=1}^{T} \left| \sum_{j=1}^{r} \mathbf{W}_{ij} \sum_{k=1}^{n} \mathbf{S}_{jk} \right| = \frac{1}{n\tau} \sum_{i=1}^{T} \left| \sum_{j=1}^{r} \sum_{k=1}^{n} \mathbf{W}_{ij} \mathbf{S}_{jk} \right| \\ &\leq \frac{1}{n\tau} \sum_{i=1}^{T} \sum_{j=1}^{r} \sum_{k=1}^{n} |\mathbf{W}_{ij} \mathbf{S}_{jk}| \\ &\leq \frac{1}{n\tau} \sum_{i=1}^{T} \sum_{k=1}^{T} |\mathbf{W}_{i}^{\mathsf{T}} \mathbf{S}_{\cdot,k}| \\ &\leq \frac{1}{n\tau} \sum_{i=1}^{T} \sum_{k=1}^{n} |\mathbf{W}_{i} \|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{n\tau} \sum_{i=1}^{T} \sum_{k=1}^{n} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{n\tau} \sum_{i=1}^{T} \sum_{k=1}^{n} \|\mathbf{W}_{i}\|_{2} \max_{1 \leq k \leq n} \|\mathbf{S}_{\cdot,k}\|_{2} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \|\mathbf{S}_{\cdot,k}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \leq \frac{1}{\tau} \|\mathbf{W}_{i}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T} \|\mathbf{W}_{i}\|_{2} \\ &\leq \frac{1}{\tau} \sum_{i=1}^{T$$

where we dropped the subscript and superscript on  $\mathbf{W}_U$  and  $\mathbf{S}^{(L)}$  to ease the notations. This concludes the proof. 

The previous lemma can be used to show that the logarithm of the ratio between the true probability of the next token and the one estimated by the LLM  $f_{\Theta}$  is upper bounded as a function of the vocabulary size T, the temperature, the upper-bound on  $\mathbf{W}_U$  and some constant related to the ambiguity of language (see Eq. (1)). 

**Proposition D.9** (Upper-bound on the logarithm). Consider an LLM  $f_{\Theta} \in \mathcal{F}$  with vocabulary size T. We recall that  $B_U$  is the upper bound on the norm of  $\mathbf{W}_U$  in the definition of parameter space W,  $\tau$  is the softmax temperature and  $c_0$  is the constant related to the ambiguity of language

(see Eq. (1)). We have  

$$\forall n \in [N], \quad \left| \log \left( \frac{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \right) \right| \leq \bar{B} = \max\{ \log (T) + \frac{2B_U}{\tau}, \log \left( \frac{1}{c_0} \right) \}.$$

*Proof.* The main idea of the proof is to bound the probability ratio and use the fact that log is 1572 non-decreasing. Let  $n \in [N]$ . The model  $f_{\Theta}$  receives as input sequences of tokens  $\mathbf{S}_n$  of size  $n \leq K$ . 1573 We first lower-bound each term of the probability ratio. From Eq. (1), we have

$$\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n) \ge c_0. \tag{17}$$

We want to obtain a similar inequality for  $\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} | \mathbf{S}_n)$ . As the parameters  $\Theta$  of the LLM are in  $\mathcal{W}$ , we know that  $\|\mathbf{W}_U^{\top}\|_{2,1} \leq B_U$ . Lemma D.8 ensures that

$$\|\frac{1}{n\tau}\mathbf{W}_U\mathbf{S}^{(L)}\mathbb{1}_n\|_1 \le \frac{1}{\tau}\|\mathbf{W}_U^{\top}\|_{2,1} \le \frac{B_U}{\tau}$$

We can then apply Lemma D.7 with  $c_1 = \frac{B_U}{\tau}$  and given that  $\frac{1}{n\tau} \mathbf{W}_U \mathbf{S}^{(L)} \mathbb{1}_n \in \mathbb{R}^T$ , it leads to

$$\mathbb{P}_{\Theta}(\cdot \mid \mathbf{S}_n) = \operatorname{softmax}\left(\frac{1}{n\tau}\mathbf{W}_U\mathbf{S}^{(L)}\mathbb{1}_n\right) \ge \frac{1}{T\exp\left(2B_U/\tau\right)}$$

where the inequality holds for each component of  $\mathbb{P}_{\Theta}(\cdot | \mathbf{S}_n)$ . This is in particular the case for  $\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} | \mathbf{S}_n)$  which is the entry we are interested in, i.e., we have

$$\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n) \ge \frac{1}{T \exp\left(2B_U/\tau\right)}.$$
(18)

1588 Going back to the ratio of probability, consider the situation where we have

$$\frac{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\boldsymbol{\Theta}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \ge 1.$$

1591 Then, using Eq. (18), we have 1592  $\mathbb{D}_{(18)}$ 

$$1 \leq \frac{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \leq \frac{1}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \leq T \exp\left(2B_U/\tau\right),$$

1595 which implies, as the log is non-decreasing monotonically,

$$0 \le \log\left(\frac{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\mathbf{\Theta}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}\right) \le \log\left(T \exp\left(2B_U/\tau\right)\right) = \log\left(T\right) + \frac{2B_U}{\tau}.$$
(19)

<sup>1598</sup> Similarly, consider the case where we have

$$\frac{\mathbb{P}_{\boldsymbol{\Theta}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \geq 1$$

and similarly to above, we can use Eq. (17) to obtain

$$1 \leq \frac{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \leq \frac{1}{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \leq \frac{1}{c_0}$$

 $\frac{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \le 1.$ 

1608 This implies

$$0 \le \log\left(\frac{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}\right) \le \log\left(\frac{1}{c_0}\right)$$

<sup>1611</sup> which also rewrites

$$0 \leq -\log\left(\frac{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}\right) \leq \log\left(\frac{1}{c_0}\right).$$
(20)

By definition of the absolute value, combining Eq. (19) and Eq. (20) leads to  $|P_{2}(\mathbf{X} \to |\mathbf{S}|)| = 2B_{2} = (11)$ 

$$\left| \log \left( \frac{\mathbb{P}_{\mathcal{L}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\mathbf{\Theta}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \right) \right| \le \max\{ \log \left(T\right) + \frac{2B_U}{\tau}, \log \left(\frac{1}{c_0}\right) \}.$$

This concludes the proof.

We are now ready to upper-bound the total variation.

**Corollary D.10** (Upper-bound on the total variation). Consider an LLM  $f_{\Theta} \in \mathcal{F}$  with vocabulary size T. We recall that  $B_U$  is the upper bound on the norm of  $\mathbf{W}_U$  in the definition of parameter space  $W, \tau$  is the softmax temperature and  $c_0$  is the constant related to the ambiguity of language (see Eq. (1)). For  $n \in [N]$ , we have

$$d_{\mathrm{TV}}(\mathbb{P}_{\mathcal{L}}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\Theta}(\cdot \mid \mathbf{S}_n)) \le \sqrt{2\max\{\log\left(T\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{c_0}\right)\}} \coloneqq c_2.$$
(21)

*Proof.* Using Proposition D.9, we can directly apply Lemma D.6 with  $B = \max\{\log(T) +$  $\frac{2B_U}{\tau}, \log\left(\frac{1}{c_0}\right)$  for any  $n \in [N]$ . This leads to

$$\forall n \in [N], \quad d_{\mathrm{TV}}(\mathbb{P}_{\mathcal{L}}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\mathbf{\Theta}}(\cdot \mid \mathbf{S}_n)) \le \sqrt{2 \max\{\log\left(T\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{c_0}\right)\}}.$$

This concludes the proof.

D.4.3 CONCLUDING THE PROOF 

We are now ready to state our main result.

**Theorem D.11** (Restatement of Theorem 4.1). Consider an LLM  $f_{\Theta} \in \mathcal{F}$  with vocabulary size T. We denote by  $\Gamma$  the mixing matrix of the pretraining sequences of tokens  $(\mathbf{S}_1, \ldots, \mathbf{S}_{N_{\text{train}}})$ . *Let*  $\delta > 0$ . *Then, with probability at least*  $1 - \delta$ *,* 

$$\mathcal{R}_{ ext{pre}}(\mathbf{\Theta}) \leq \widehat{\mathcal{R}}_{ ext{pre}}(\mathbf{\Theta}) + rac{ar{B}}{\sqrt{N_{ ext{train}}}} \sqrt{\log\left(rac{2}{\delta}
ight)},$$

where B is a constant depending on the parameters of the problem. More precisely,

$$\bar{B} = 2 \| \mathbf{\Gamma} \| \sqrt{\max\{\log\left(T\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{c_0}\right)\}}$$

Proof of Theorem 4.1. By definition of the risk, we have 

$$\widehat{\mathcal{R}}_{\text{pre}}(\Theta) = \frac{1}{N_{\text{train}}} \sum_{n=1}^{N_{\text{train}}} \underbrace{d_{\text{TV}}(\mathbb{P}_{\mathcal{L}}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\Theta}(\cdot \mid \mathbf{S}_n))}_{=g_n(\mathbf{S}_n)} = \frac{1}{N_{\text{train}}} \sum_{n=1}^{N_{\text{train}}} g_n(\mathbf{S}_n)$$
$$= f(\mathbf{S}_1, \dots, \mathbf{S}_{N_{\text{train}}}) = f(S).$$

Using Corollary D.10, we know that 

$$|g_n(\mathbf{S}_n)| \le \sqrt{2\max\{\log{(T)} + \frac{2B_U}{\tau}, \log{\left(\frac{1}{c_0}\right)}\}} \coloneqq c_2.$$

By definition, each sequence of tokens  $S_n$  takes its values in  $\mathcal{V}^n$  (again by abuse of notation,  $n = \min\{n, K\}$ ) and  $\tilde{S}$  takes its values in  $\mathcal{V}^1 \times \ldots \times \mathcal{V}^{N_{\text{train}}}$ . For any two sequences  $\zeta, \Sigma$  with values in  $\mathcal{V}^1 \times \ldots \times \mathcal{V}^{N_{\text{train}}}$ , we have 

$$f(\zeta) - f(\Sigma) = \frac{1}{N_{\text{train}}} \sum_{n=1}^{N_{\text{train}}} \left( \underbrace{d_{\text{TV}}(\mathbb{P}_{\mathcal{L}}(\cdot \mid \boldsymbol{\zeta}_n), \mathbb{P}_{\Theta}(\cdot \mid \boldsymbol{\zeta}_n))}_{=g_n(\boldsymbol{\zeta}_n)} - \underbrace{d_{\text{TV}}(\mathbb{P}_{\mathcal{L}}(\cdot \mid \boldsymbol{\Sigma}_n), \mathbb{P}_{\Theta}(\cdot \mid \boldsymbol{\Sigma}_n))}_{=g_n(\boldsymbol{\Sigma}_n)} \right)$$

1672  
1673 
$$= \frac{1}{N_{\text{train}}} \sum_{n=1}^{N_{\text{train}}} (g_n(\boldsymbol{\zeta}_n) - g_n(\boldsymbol{\Sigma}_n))$$

\

$$\begin{array}{ll} & \begin{array}{l} 1674\\ 1675\\ 1676\\ 1676\\ 1676\\ 1677\\ 1678\\ 1678\\ 1679\\ 16$$

Let  $u \ge 0$ . We have the following events ordering

$$(\mathbb{E}_S[f(S)] - f(S) \ge u) \subseteq (\mathbb{E}_S[f(S)] - f(S) \ge u) \cup (f(S) - \mathbb{E}_S[f(S)] \ge u)$$
$$= (|f(S) - \mathbb{E}_S[f(S)]| \ge u).$$

Hence, as u was taken arbitrary and using Eq. (22), we have 

$$\forall u \ge 0, \quad \mathbb{P}(\mathbb{E}_S[f(S)] - f(S) \ge u) \le 2 \exp\left(\frac{-2u^2}{\|\mathbf{\Gamma}\|^2 \|\mathbf{c}\|_2^2}\right)$$

We recall that by definition

$$f(S) = \widehat{\mathcal{R}}_{\text{pre}}(\Theta) \text{ and } \mathcal{R}_{\text{pre}}(\Theta) = \mathbb{E}_{S} \Big[ \widehat{\mathcal{R}}_{\text{pre}}(\Theta) \Big].$$

Since the previous inequality holds for any  $u \ge 0$ , we can hence choose u such that 

$$\delta = 2 \exp\left(\frac{-2u^2}{\|\mathbf{\Gamma}\|^2 \|\mathbf{c}\|_2^2}\right) \iff \frac{-2u^2}{\|\mathbf{\Gamma}\|^2 \|\mathbf{c}\|_2^2} = \log\left(\frac{\delta}{2}\right) \iff u^2 = \frac{1}{2} \|\mathbf{\Gamma}\|^2 \|\mathbf{c}\|_2^2 \log\left(\frac{2}{\delta}\right)$$
$$\iff u = \frac{1}{\sqrt{2}} \|\mathbf{\Gamma}\| \|\mathbf{c}\|_2 \sqrt{\log\left(\frac{2}{\delta}\right)}.$$

Using the fact that

$$\|\mathbf{c}\|_{2} = \sqrt{\sum_{n=1}^{N_{\text{train}}} \mathbf{c}_{n}^{2}} = \sqrt{\sum_{n=1}^{N_{\text{train}}} \left(\frac{2c_{2}}{N_{\text{train}}}\right)^{2}} = \sqrt{\sum_{n=1}^{N_{\text{train}}} \frac{4c_{2}^{2}}{N_{\text{train}}^{2}}} = \sqrt{\frac{4c_{2}^{2}}{N_{\text{train}}}} = \frac{2c_{2}}{\sqrt{N_{\text{train}}}}$$

and using the fact that 
$$c_2 = \sqrt{2 \max\{\log(T) + \frac{2B_U}{\tau}, \log\left(\frac{1}{c_0}\right)\}}$$
 from Corollary D.10, we obtain  
1730

1731  
1732  
1733  
1734  

$$u = \frac{1}{\sqrt{2}} \frac{2c_2}{\sqrt{N_{\text{train}}}} \|\mathbf{\Gamma}\| \sqrt{\log\left(\frac{2}{\delta}\right)} = \frac{\sqrt{2}c_2}{\sqrt{N_{\text{train}}}} \|\mathbf{\Gamma}\| \sqrt{\log\left(\frac{2}{\delta}\right)}$$
1734  

$$u = \frac{1}{\sqrt{2}} \frac{2c_2}{\sqrt{N_{\text{train}}}} \|\mathbf{\Gamma}\| \sqrt{\log\left(\frac{2}{\delta}\right)} = \frac{\sqrt{2}c_2}{\sqrt{N_{\text{train}}}} \|\mathbf{\Gamma}\| \sqrt{\log\left(\frac{2}{\delta}\right)}$$

$$=\frac{2\|\mathbf{\Gamma}\|\sqrt{\max\{\log\left(T\right)+\frac{2B_U}{\tau},\log\left(\frac{1}{c_0}\right)\}}}{\sqrt{N_{\text{train}}}}\sqrt{\log\left(\frac{2}{\delta}\right)}=\frac{\bar{B}}{\sqrt{N_{\text{train}}}}\sqrt{\log\left(\frac{2}{\delta}\right)},$$

where we define

$$\bar{B} = 2 \|\mathbf{\Gamma}\| \sqrt{\max\{\log\left(T\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{c_0}\right)\}}.$$

1741 Putting everything together, we have 1742

$$\mathbb{P}\left(\mathcal{R}_{\rm pre}(\boldsymbol{\Theta}) - \widehat{\mathcal{R}}_{\rm pre}(\boldsymbol{\Theta}) \geq \frac{\bar{B}}{\sqrt{N_{\rm train}}} \sqrt{\log\left(\frac{2}{\delta}\right)}\right) \leq \delta.$$

Taking the opposite event leads to the following inequality with probability at least  $1 - \delta$ 

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

which concludes the proof.

1752 D.5 PROOF OF COROLLARY 4.2

1754 As the layer norm is not applied anymore, each token is no longer in the  $\ell_2$ -unit ball, and Lemma D.8 1755 does not hold anymore. We want to provide an analogous lemma for our setting. We first prove the 1756 following technical lemmas. 

**Lemma D.12.** The ReLU is a norm-decreasing operator, i.e., we have

 $\forall \mathbf{A} \in \mathbb{R}^{n \times m}, \quad \|\text{ReLU}(\mathbf{A})\|_{1,1} \le \|\mathbf{A}\|_{1,1},$ 

where the ReLU is applied entry-wise.

1764 Proof. Recalling that  $\operatorname{ReLU}(x) = \max\{0, x\}$  is applied entry-wise, using the fact that 1765  $|\max\{0, x\}| \le |x|$  and considering A and  $\tilde{A} = \operatorname{ReLU}(A)$ , we have

$$\|\tilde{\mathbf{A}}\|_{1,1} = \sum_{i,j} |\tilde{\mathbf{A}}_{i,j}| = \sum_{i,j} |\max\{0, \tilde{\mathbf{A}}_{i,j}\}| \le \sum_{i,j} |\mathbf{A}_{i,j}| \le \|\mathbf{A}\|_{1,1}$$

which concludes the proof.

**Lemma D.13.** The  $L_{1,1}$ -norm verifies the following property:

$$\forall \mathbf{A} \in \mathbb{R}^{n \times m}, \mathbf{B} \in \mathbb{R}^{m \times p}, \quad \|\mathbf{A}\mathbf{B}\|_{1,1} \le n \|\mathbf{A}\|_{\infty} \|\mathbf{B}\|_{1,1}.$$

Proof. We have

$$\begin{aligned} \|\mathbf{AB}\|_{1,1} &= \sum_{j=1}^{p} \sum_{i=1}^{n} |(\mathbf{AB})_{ij}| = \sum_{j=1}^{p} \sum_{i=1}^{n} |\sum_{k=1}^{m} \mathbf{A}_{ik} \mathbf{B}_{kj}| \le \sum_{j=1}^{p} \sum_{i=1}^{n} \sum_{k=1}^{m} |\mathbf{A}_{ik} \mathbf{B}_{kj}| \\ &\le \sum_{j=1}^{p} \sum_{i=1}^{n} \sum_{k=1}^{m} |\mathbf{A}_{ik}| |\mathbf{B}_{kj}| \le \max_{ik} |\mathbf{A}_{ik}| \sum_{j=1}^{p} \sum_{i=1}^{n} \sum_{k=1}^{m} |\mathbf{B}_{kj}| \end{aligned}$$

$$\leq n \|\mathbf{A}\|_{\infty} \sum_{j=1}^{p} \sum_{k=1}^{m} |\mathbf{B}_{kj}| \leq n \|\mathbf{A}\|_{\infty} \|\mathbf{B}\|_{1,1},$$

1785 which concludes the proof.

**Lemma D.14.** The 
$$L_{2,1}$$
 and  $L_{\infty,1}$ -norms verify the following relation

 $\forall \mathbf{A} \in \mathbb{R}^{n \times m}, \quad \|\mathbf{A}\|_{\infty,1} \le \|\mathbf{A}\|_{2,1}.$ 

*Proof.* By definition of the  $L_{p,q}$ -norm, we have

$$\begin{split} |\mathbf{A}\|_{\infty,1} &= \sum_{j=1}^{M} \max_{1 \le i \le n} |\mathbf{A}_{ij}| = \sum_{j=1}^{M} \sqrt{\max_{1 \le i \le n} |\mathbf{A}_{ij}^2|} \quad (\text{as } x \to x^2 \text{ is increasing}) \\ &\leq \sum_{j=1}^{M} \sqrt{\sum_{i=1}^{n} |\mathbf{A}_{ij}^2|} \le \sum_{j=1}^{M} \|\mathbf{A}_{\cdot,j}\|_2 \le \|\mathbf{A}\|_{2,1}, \end{split}$$

where the first inequality comes from adding non-negative terms.

We are now ready to state the lemma analogous to Lemma D.8.

**Lemma D.15.** Consider an LLM  $f_{\Theta} \in \tilde{F}$  with L layers. For any input sequence  $\mathbf{S} \in \mathbb{R}^{r \times n}$ , the following inequality holds

$$\|\frac{1}{n\tau}\mathbf{W}_U\mathbf{S}^{(L)}\mathbb{1}_n\|_1 \le \frac{c_3}{\tau}\|\mathbf{W}_U^\top\|_{2,1}$$

where  $\tau$  is the temperature and  $c_3$  is a constant depending on the parameters upper-bound. More precisely,

$$c_3 = \left[ \left(1 + rmB_1B_2\right) \cdot \left(1 + \frac{r^3}{H}B_OB_V\right) \right]^L \cdot B_{\text{tok}}.$$

 $\mathbf{W}_U$  is the unembedding matrix (which is bounded as stated in the definition of the parameters space  $\mathcal{W}$ ), and  $\mathbf{S}^{(L)}$  is the output of the last transformer layer.

*Proof of Lemma D.15.* Our model  $f_{\Theta} \in \tilde{\mathcal{F}}$  is given as input a sequence  $\mathbf{S} \in \mathbb{R}^{r \times n}$ . With similar computations than in Lemma D.8, we have

$$\frac{1}{n\tau} \|\mathbf{W}_U \mathbf{S}^{(L)} \mathbb{1}_n\|_1 = \frac{1}{n\tau} \sum_{i=1}^T \left| \sum_{j=1}^r \mathbf{W}_{ij} \sum_{k=1}^n \mathbf{S}_{jk} \right| = \frac{1}{n\tau} \sum_{i=1}^T \left| \sum_{j=1}^r \sum_{k=1}^n \mathbf{W}_{ij} \mathbf{S}_{jk} \right|$$
$$\leq \frac{1}{n\tau} \sum_{i=1}^T \sum_{j=1}^r \sum_{k=1}^n |\mathbf{W}_{ij} \mathbf{S}_{jk}| \qquad (\text{triangle})$$

$$\leq \frac{1}{n\tau} \sum_{i=1}^{T} \sum_{k=1}^{n} \left| \mathbf{W}_{i}^{\top} \mathbf{S}_{\cdot,k} \right| \leq \frac{1}{n\tau} \sum_{i=1}^{T} \sum_{k=1}^{n} \|\mathbf{W}_{i}\|_{\infty} \|\mathbf{S}_{\cdot,k}\|_{1} \quad \text{(Hölder inequality)}$$

$$\sum_{\substack{i=1\\ 1829\\ 1830}} \leq \frac{1}{n\tau} \left( \sum_{i=1}^{T} \|\mathbf{W}_i\|_{\infty} \right) \cdot \left( \sum_{k=1}^{n} \|\mathbf{S}_{\cdot,k}\|_{1} \right) \leq \frac{1}{n\tau} \|\mathbf{W}_U^{\top}\|_{\infty,1} \|\mathbf{S}^{(L)}\|_{1,1}$$

$$\sum_{\substack{i=1\\ 1831}} \sum_{\substack{i=1\\ 1831}} \sum$$

1831  
1832 
$$\leq \frac{1}{n\tau} \| \mathbf{W}_U^{\top} \|_{2,1} \| \mathbf{S}^{(L)} \|_{1,1},$$
 (Lemma D.14)

where, again, we dropped the subscript and superscript on  $\mathbf{W}_U$  and  $\mathbf{S}^{(L)}$  to ease the notations. We obtain

$$\|\frac{1}{n\tau}\mathbf{W}_{U}\mathbf{S}^{(L)}\mathbb{1}_{n}\|_{1} \leq \frac{1}{n\tau}\|\mathbf{W}_{U}^{\top}\|_{2,1}\|\mathbf{S}^{(L)}\|_{1,1}.$$
(23)

(triangular inequality)

As we do not use layer normalization, we want to find another way to bound  $S^{(L)}$ . To that end, we will first express  $\mathbf{S}^{(\ell)}$ , the output of the  $(\ell)$ -th layer of the transformer, as a function of  $\mathbf{S}^{(\ell-1)}$ , the output of the  $(\ell - 1)$ -th layer. Using the definition of the transformer model (see Appendix B), we have

$$\begin{cases} \mathbf{Z}^{(\ell)} = \mathbf{S}^{(\ell-1)} + \mathcal{A}\left(\mathbf{S}^{(\ell-1)}; \mathbf{W}_Q^{(\ell)}, \mathbf{W}_K^{(\ell)}, \mathbf{W}_V^{(\ell)}, \mathbf{W}_Q^{(\ell)}\right), \\ \mathbf{Y}^{(\ell)} = \mathbf{W}_2^{(\ell)} \operatorname{ReLU}\left(\mathbf{W}_1^{(\ell)} \mathbf{Z}^{(\ell)}\right), \end{cases}$$

$$\mathbf{S}^{(\ell)} = \mathbf{Z}^{(\ell)} + \mathbf{Y}^{(\ell)}.$$

We will compute each layer's  $L_{1,1}$ -norm. 

Step 1: MHA. By definition, denoting the number of heads by H, we know that  $\mathcal{A}\left(\mathbf{S}^{(\ell-1)}; \mathbf{W}_Q^{(\ell)}, \mathbf{W}_K^{(\ell)}, \mathbf{W}_V^{(\ell)}, \mathbf{W}_O^{(\ell)}\right) \in \mathbb{R}^{r \times n} \text{ multiplies } \mathbf{W}^{(\ell)} \in \mathbb{R}^{r \times r} \text{ with the concatenation}$ on the rows of the H softmax layers that each writes 

softmax 
$$\left( \mathbf{W}_{Q}^{(\ell)} \mathbf{S}^{(\ell)} \left( \mathbf{W}_{K}^{(\ell)} \mathbf{S}^{(\ell-1)} \right)^{\top} / \sqrt{r} \right) \left( \mathbf{W}_{V}^{(\ell)} \mathbf{S}^{(\ell-1)} \right) \in \mathbb{R}^{\frac{r}{H} \times n}$$

We keep the notations  $\ell$  without explicating the index of the head to ease notations. Denoting the concatenation on the rows by  $\mathbf{C}^{(\ell)} \in \mathbb{R}^{r \times n}$ , we have 

$$\begin{aligned} \|\mathcal{A}\Big(\mathbf{S}^{(\ell-1)}; \mathbf{W}_{Q}^{(\ell)}, \mathbf{W}_{K}^{(\ell)}, \mathbf{W}_{V}^{(\ell)}, \mathbf{W}_{O}^{(\ell)}\Big)\|_{1,1} &= \|\mathbf{W}_{O}^{(\ell)} \mathbf{C}^{(\ell)}\|_{1,1} \leq r \cdot \|\mathbf{W}_{O}^{(\ell)}\|_{\infty} \|\mathbf{C}^{(\ell)}\|_{1,1} \\ &\leq rB_{O} \|\mathbf{C}^{(\ell)}\|_{1,1}. \end{aligned}$$
 (definition of  $\tilde{\mathcal{W}}$ )

Moreover, by definition of  $\mathbf{C}^{(\ell)}$ , we have 

$$\|\mathbf{C}^{(\ell)}\|_{1,1} = \sum_{j=1}^{r} \sum_{i=1}^{r} |\mathbf{C}_{ij}^{(\ell)}| = \sum_{j=1}^{r} \sum_{i=1}^{r/H} \sum_{h=1}^{H} |\mathbf{C}_{ij}^{(\ell,h)}| = \sum_{h=1}^{H} \|\mathbf{C}^{(\ell,h)}\|_{1,1},$$
(24)

where  $\mathbf{C}^{(\ell,h)} \in \mathbb{R}^{\frac{T}{H} \times n}$  is the softmax matrix of the *h*-th layer. We recall that the softmax matrix is a row-stochastic matrix of  $\mathbb{R}^{\frac{r}{H} \times r}$  so it has all values lower than 1. In the next computations, we drop the h index on the query, key, and value matrices to ease the notations. Using Lemma D.13 on the softmax matrix and on the value matrix  $\mathbf{W}_{V}^{(\ell)} \in \mathbb{R}^{\frac{r}{H} \times r}$ , we have 

$$\begin{aligned} \|\mathbf{C}^{(\ell,h)}\|_{1,1} &= \|\operatorname{softmax}\left(\mathbf{W}_{Q}^{(\ell)}\mathbf{S}^{(\ell)}\left(\mathbf{W}_{K}^{(\ell)}\mathbf{S}^{(\ell-1)}\right)^{\top}/\sqrt{r}\right)\left(\mathbf{W}_{V}^{(\ell)}\mathbf{S}^{(\ell-1)}\right)\|_{1,1} \\ &\leq \frac{r}{H} \cdot \|\operatorname{softmax}\left(\mathbf{W}_{Q}^{(\ell)}\mathbf{S}^{(\ell)}\left(\mathbf{W}_{K}^{(\ell)}\mathbf{S}^{(\ell-1)}\right)^{\top}/\sqrt{r}\right)\|_{\infty} \cdot \|\left(\mathbf{W}_{V}^{(\ell)}\mathbf{S}^{(\ell-1)}\right)\|_{1,1} \end{aligned}$$

1871  
1872 
$$\leq \frac{r}{H} \cdot \| \left( \mathbf{W}_V^{(\ell)} \mathbf{S}^{(\ell-1)} \right) \|_{1,1}$$
 (the softmax matrix is row-stochastic)

Combining the previous inequality with Eq. (24) leads to 

1877  
1878 
$$\|\mathbf{C}^{(\ell)}\|_{1,1} \le \frac{r^2}{H} B_V \|\mathbf{S}^{(\ell-1)}\|_{1,1}$$

In the end, the multi-head attention norm verifies 

$$\|\mathcal{A}\left(\mathbf{S}^{(\ell-1)}; \mathbf{W}_{Q}^{(\ell)}, \mathbf{W}_{K}^{(\ell)}, \mathbf{W}_{V}^{(\ell)}, \mathbf{W}_{O}^{(\ell)}\right)\|_{1,1} \le \frac{r^{3}}{H} B_{O} B_{V} \|\mathbf{S}^{(\ell-1)}\|_{1,1}.$$

Using the triangular inequality, we obtain 

$$\|\mathbf{Z}^{(\ell)}\|_{1,1} \le \left(1 + \frac{r^3}{H} B_O B_V\right) \cdot \|\mathbf{S}^{(\ell-1)}\|_{1,1}.$$
(25)

**Step 2: FF.** We recall that  $\mathbf{W}_1 \in \mathbb{R}^{m \times r}$  and  $\mathbf{W}_2 \in \mathbb{R}^{r \times m}$ . Using similar arguments to the above, we have 

$$\|\mathbf{Y}^{(\ell)}\|_{1,1} = \|\mathbf{W}_2^{(\ell)} \operatorname{ReLU}\left(\mathbf{W}_1^{(\ell)} \mathbf{Z}^{(\ell)}\right)\|_{1,1}$$

$$\leq r \cdot \|\mathbf{W}_{2}^{(\ell)}\|_{\infty} \|\operatorname{ReLU}\left(\mathbf{W}_{1}^{(\ell)}\mathbf{Z}^{(\ell)}\right)\|_{1,1} \qquad (\operatorname{Lemma D.13})$$

1892 
$$\leq r \cdot \|\mathbf{W}_{2}^{(\ell)}\|_{\infty} \|\mathbf{W}_{1}^{(\ell)} \mathbf{Z}^{(\ell)}\|_{1,1}$$
 (Lemma D.12)

$$\leq r \cdot m \cdot \|\mathbf{W}_{2}^{(\ell)}\|_{\infty} \|\mathbf{W}_{1}^{(\ell)}\|_{\infty} \|\mathbf{Z}^{(\ell)}\|_{1,1}$$
 (Lemma D.13)

$$\leq rmB_1B_2 \|\mathbf{Z}^{(\ell)}\|_{1,1}.$$
 (definition of  $\tilde{\mathcal{W}}$ )

**Step 3: output layer.** Again, applying the triangular inequality and using the previous inequality and Eq. (25), we have

$$\begin{aligned} \|\mathbf{S}^{(\ell)}\|_{1,1} &\leq \|\mathbf{Z}^{(\ell)}\|_{1,1} + \|\mathbf{Y}^{(\ell)}\|_{1,1} \leq (1 + rmB_1B_2)\|\mathbf{Z}^{(\ell)}\|_{1,1} \\ &\leq (1 + rmB_1B_2) \left(1 + \frac{r^3}{H}B_OB_V\right) \|\mathbf{S}^{(\ell-1)}\|_{1,1}. \end{aligned}$$

1904 Iterating through the layers, recalling that  $S^{(0)} = S$ , we finally obtain 

$$\|\mathbf{S}^{(L)}\|_{1,1} \le \left[ (1 + rmB_1B_2) \left( 1 + \frac{r^3}{H}B_OB_V \right) \right]^L \|\mathbf{S}\|_{1,1},$$

where S is the input sequence. Combining this inequality with Eq. (23) leads to

$$\|\frac{1}{n\tau}\mathbf{W}_{U}\mathbf{S}^{(L)}\mathbb{1}_{n}\|_{1} \leq \left[(1+rmB_{1}B_{2})\left(1+\frac{r^{3}}{H}B_{O}B_{V}\right)\right]^{L}\frac{\|\mathbf{S}\|_{1,1}}{n}\left(\frac{1}{\tau}\|\mathbf{W}_{U}^{\top}\|_{2,1}\right).$$

Using the fact that each token has a  $\ell_1$ -norm bounded by  $B_{tok}$ . Hence, each column of S is too and we have

$$\frac{1}{n} \|\mathbf{S}\|_{1,1} = \frac{1}{n} \sum_{j=1}^{n} \sum_{i=1}^{r} |\mathbf{S}_{ij}| = \frac{1}{n} \sum_{j=1}^{n} \underbrace{\|\mathbf{S}_{\cdot,j}\|_{1}}_{\leq B_{\text{tok}}} \leq B_{\text{tok}}.$$

1919 Combining the last two inequalities concludes the proof.

We can now restate Corollary 4.2.

**Corollary D.16** (Restatement of Corollary 4.2). Consider an LLM  $f_{\Theta} \in \mathcal{F}$  with vocabulary size T composed of L transformer blocks and H attention heads. We denote by  $\Gamma$  the mixing matrix of the pretraining sequences of tokens  $(\mathbf{S}_1, \ldots, \mathbf{S}_{N_{\text{train}}})$ . Let  $\delta > 0$ . Then, with probability at least  $1 - \delta$ ,

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

where  $\overline{B}$  is a constant depending on the parameters of the problem. More precisely,

$$\bar{B} = 2 \|\mathbf{\Gamma}\| \sqrt{\max\{\log\left(T\right) + \frac{2(B_{\Theta})^L}{\tau}, \log\left(\frac{1}{c_0}\right)\}},$$

with 
$$B_{\Theta} = \left[ (1 + rmB_1B_2) \left( 1 + \frac{r^3}{H} B_O B_V \right) \right] (B_{\text{tok}} B_U)^{1/L}.$$

Proof of Corollary 4.2. We first note that the only change from Lemma D.8 to Lemma D.15 is the multiplicative constant  $c_3 = \left[(1 + rmB_1B_2)\left(1 + \frac{r^3}{H}B_OB_V\right)\right]^L B_{\text{tok}}$  in front of  $\frac{1}{\tau} \|\mathbf{W}_U^{\mathsf{T}}\|_{2,1}$ . In particular, as we know that  $\tilde{\mathcal{W}} \subset \mathcal{W}$ , we also have  $\|\mathbf{W}_U^{\mathsf{T}}\|_{2,1} \leq B_U$ . Hence, we can apply the proof of Theorem 4.1 in a straightforward manner by changing  $\frac{B_U}{\tau}$  by  $c_3 \cdot \frac{B_U}{\tau}$ . This concludes the proof.

<sup>1944</sup> D.6 PROOF OF COROLLARY 4.3

<sup>1946</sup> We detail the proof of Corollary 4.3 below.

*Proof.* We first note that by definition of the total variation distance (Wolfer & Kontorovich, 2019), we have

$$\begin{split} \mathbb{E}_{\mathbf{S} \sim \mathbb{P}_{\mathcal{L}}} \| \mathbf{Q}^{*}(\mathbf{S}, \cdot) - \mathbf{Q}_{f}(\mathbf{S}, \cdot) \|_{1} &= \mathbb{E}_{\mathbf{S} \sim \mathbb{P}_{\mathcal{L}}} [2 \cdot d_{\mathrm{TV}}(\mathbf{Q}^{*}(\mathbf{S}, \cdot), \mathbf{Q}_{f}(\mathbf{S}, \cdot))] \\ &= 2 \cdot \mathbb{E}_{\mathbf{S} \sim \mathbb{P}_{\mathcal{L}}} [d_{\mathrm{TV}}(\mathbf{Q}^{*}(\mathbf{S}, \cdot), \mathbf{Q}_{f}(\mathbf{S}, \cdot))] \\ &= 2 \cdot \mathcal{R}_{\mathrm{pre}}(\mathbf{\Theta}). \qquad \text{(by definition of the risk Eq. (2))} \end{split}$$

Applying Theorem 4.1 (or similarly Corollary 4.2), we know that

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

where  $\overline{B}$  is formally defined in Theorem 4.1 (respectively Corollary 4.2). Assuming a perfect pretraining error amounts to consider  $\widehat{\mathcal{R}}_{\text{pre}}(\Theta) = 0$ . We denote by  $N^*$  the integer such that the error is equal to  $\frac{\epsilon}{2}$ , i.e.,

1962 1963

1964 1965

1966

1957 1958

1947

$$\frac{\bar{B}}{\sqrt{N^*}}\sqrt{\log\left(\frac{2}{\delta}\right)} = \frac{\epsilon}{2} \iff \frac{\bar{B}^2}{N^*}\log\left(\frac{2}{\delta}\right) = \frac{\epsilon^2}{4} \iff N^* = \left(\frac{2\bar{B}}{\epsilon}\right)^2\log\left(\frac{2}{\delta}\right).$$

Taking the ceiling function ensures that  $N^*$  is an integer. Hence, taking  $N_{\text{train}} \geq N^* = \left\lceil \left(\frac{2\bar{B}}{\epsilon}\right)^2 \log\left(\frac{2}{\delta}\right) \right\rceil$  ensures that

1967 1968 1969

1970 1971

1973 1974

1976

$$\frac{\bar{B}}{/N_{\text{train}}}\sqrt{\log\left(\frac{2}{\delta}\right)} \le \frac{\bar{B}}{\sqrt{N^*}}\sqrt{\log\left(\frac{2}{\delta}\right)} = \frac{\epsilon}{2}.$$

1972 Putting everything together, taking  $N_{\text{train}} \ge N^*$  leads to

$$\mathbb{E}_{\mathbf{S}\sim\mathbb{P}_{\mathcal{L}}}\|\mathbf{Q}^{*}(\mathbf{S},\cdot)-\mathbf{Q}_{f}(\mathbf{S},\cdot)\|_{1}\leq 2\cdot\mathcal{R}_{\text{pre}}(\mathbf{\Theta})\leq 2\cdot\frac{\epsilon}{2}=\epsilon$$

1975 which concludes the proof.

1979 In this section, we detail the proof of Theorem 4.4. We first recall the problem setup.

Markov chains inputs. In this section, we give as input of the model a single Markov chain 1981  $X = (\mathbf{X}_1, \dots, \mathbf{X}_{N_{\text{icl}}})$  with finite, discrete state space  $\Omega$  of size d with transition probability  $\mathbb{P}$ . We 1982 assume the  $X_n$  are already tokenized and thus we have  $\Omega \subset \mathcal{V}$ . We denote the sequence of tokens the LLM receives by  $\mathbf{S}_n = (\mathbf{X}_1, \dots, \mathbf{X}_n)$  if  $n \leq K$  and  $\mathbf{S}_n = (\mathbf{X}_{n-K+1}, \dots, \mathbf{X}_n)$  otherwise due 1984 to the deletion process (see Definition B.2). In particular, the  $S_n$  are elements of  $\mathcal{V}_K^*$ . We note that  $S = (\mathbf{S}_1, \dots, \mathbf{S}_{N_{\text{icl}}})$  is also a Markov chain (see Appendix D.7.1). By definition of  $\mathbb{P}$ , we know that 1986 for any  $n \in [N_{icl}]$ , the next token  $\mathbf{X}_{n+1}$  follows the distribution  $\mathbb{P}(\cdot \mid \mathbf{S}_n)$ . We assume that there 1987 exists a positive constant  $p_{\min}$  that lower bounds all the transition probability between states, i.e., 1988  $\forall n \in [N_{icl}], \forall x, y \in \Omega, \quad \mathbb{P}(\mathbf{X}_{n+1} = y \mid \mathbf{X}_n = x) \ge p_{\min} > 0.$  This is akin to the ambiguity of 1989 language constant  $c_0$  considered in the previous section and in Hu et al. (2024); Wies et al. (2024); Xie et al. (2022); Zhang et al. (2023b).

199<sup>-</sup>

**Next token probability distribution.** An important difference with the setting considered in Theorem 4.1 is that here, we predict a probability distribution on the state space  $\Omega$  of the Markov chain and not on the vocabulary of the LLM  $\mathcal{V}$ . To that end, we restrict the predicted probability given the past tokens  $\mathbf{S}_n$  to the state space  $\Omega$ . Formally, denoting the output of the last layer of  $f_{\Theta}$  by  $\mathbf{S}^{(L)}$ , the last layer before the softmax outputs a vector  $\mathbf{u} = \frac{1}{n\tau} \mathbf{W}_U \mathbf{S}^{(L)} \mathbb{1}_n \in \mathbb{R}^T$ . We first extract the entries of  $\mathbf{u}$  whose index i are such that the i-th element of the vocabulary space  $\mathcal{V}$  is in  $\Omega$ . This can be formalized as follows. We denote by  $\mathbb{I}_d = (i_1 \leq i_2 \leq \ldots \leq i_d) \in [T]^d$  the subset of d distinct

1998 elements of [T] and consider the matrix  $\mathbf{M}_j = \mathbf{e}_{i_j}^{\top}$ , where  $\mathbf{e}_{i_j} \in \mathbb{R}^T$  has value 1 at entry  $i_j \in \mathbb{I}$  and 0 elsewhere. Extracting only the d entries of u that corresponds to the state space yields a vector in  $\mathbb{R}^d$ 2000 that writes  $\mathbf{v} = \frac{1}{n\tau} \mathbf{M} \mathbf{W}_U \mathbf{S}^{(L)} \mathbb{1}_n \in \mathbb{R}^T$ . Similarly to Appendix B, the probability distribution of 2001 next token  $\mathbf{X}_{n+1}$  provided by the LLM  $f_{\Theta}$  now writes 2002

$$\mathbb{P}_{\Theta}(\cdot \mid \mathbf{S}_n) = \operatorname{softmax}\left(\frac{1}{n\tau}\mathbf{M}\mathbf{W}_U\mathbf{S}^{(L)}\mathbb{1}_n\right) \in \Delta_d.$$

2005 We aim to obtain a similar generalization bound than in Theorem 4.1 where the reference probability 2006 distribution is the Markov chain transition probability P instead of the probability distribution of 2007 language  $\mathbb{P}_{\mathcal{L}}$ . In particular,  $\mathbb{P}$  will replace  $\mathbb{P}_{\mathcal{L}}$  in the definition of the risks in Eq. (2). We provide 2008 below an overview of the proof before detailing it. 2009

2010 **Overview of the proof.** We are going to use McDiarmid's inequality for Markov chains of Paulin 2011 (2015, Corollary 2.11). To adapt their arguments to our setting, we bound the total variation between 2012 the true probability of the next token and the one estimated by the LLM. The rest of this section is organized as follows. First, in Appendix D.7.1, we show that  $S = (\mathbf{S}_1, \dots, \mathbf{S}_{N_{\text{icl}}})$  is a Markov chain. 2013 Then in Appendix D.7.2, we adapt the concentration inequality of Paulin (2015, Corollary 2.11). 2014 Afterwards in Appendix D.7.3, we show how to bound the total variation between the true and the 2015 estimated probability of the next token. Finally Appendix D.7.4 concludes the proof. 2016

D.7.1 CONNECTION BETWEEN TOKENS AND SEQUENCES OF TOKENS MARKOV CHAINS 2018

We first show that  $S = (\mathbf{S}_1, \dots, \mathbf{S}_{N_{icl}})$  is also a Markov chain. 2020

> **Lemma D.17.** Consider a sequence (not necessarily a Markov chain)  $X = (\mathbf{X}_1, \ldots, \mathbf{X}_N)$  with values in  $\Omega$  and let  $\mathbf{S}_n = (\mathbf{X}_1, \dots, \mathbf{X}_n)$  if n < K and  $\mathbf{S}_n = (\mathbf{X}_{n-K+1}, \dots, \mathbf{X}_n)$  otherwise. Then, the sequence  $S = (\mathbf{S}_1, \dots, \mathbf{S}_N)$  is a Markov chain with state space  $\Omega_K^*$  that contains the sequence of elements in  $\Omega$  of length smaller than K.

*Proof.* By definition of the  $S_n$ , we know that they take values in  $\Omega_K^*$ . Let  $x_1, \ldots, x_{n+1} \in \Omega$ . We 2027 first assume that n > K and denote  $s_i = (x_{n-K+1}, \ldots, x_i)$ . We have 2028

$$\mathbb{P}(\mathbf{S}_{n+1} = s_{n+1} | \mathbf{S}_n = s_n, \dots, \mathbf{S}_{n-K+1} = s_{n-K+1})$$
  
=  $\mathbb{P}(\mathbf{S}_{n+1} = s_{n+1} | \mathbf{X}_n = x_n, \dots, \mathbf{X}_{n-K+1} = x_{n-K+1})$   
=  $\mathbb{P}(\mathbf{S}_{n+1} = s_{n+1} | \mathbf{S}_n = s_n).$  (by definition of  $\mathbf{S}_n$ )

2033 Similarly, we assume n < K and denote  $s_i = (x_1, \ldots, x_i)$ . We have

$$\mathbb{P}(\mathbf{S}_{n+1} = s_{n+1} | \mathbf{S}_n = s_n, \dots, \mathbf{S}_1 = s_1)$$
  
=  $\mathbb{P}(\mathbf{S}_{n+1} = s_{n+1} | \mathbf{X}_n = x_n, \dots, \mathbf{X}_1 = x_1)$   
=  $\mathbb{P}(\mathbf{S}_{n+1} = s_{n+1} | \mathbf{S}_n = s_n).$  (by definition of  $\mathbf{S}_n$ )

Finally, for n = K, we denote  $s_i = (x_1, \ldots, x_i)$  for  $i \leq K$  and  $s_{K+1} = (x_2, \ldots, x_{K+1})$ . We have  $\mathbb{P}(\mathbf{S}_{K+1} = s_{K+1} \mid \mathbf{S}_n = s_{m+1})$  $\mathbf{S}_2 = \mathbf{S}_2$ 

2040 
$$\mathbb{P}(\mathbf{S}_{K+1} = s_{K+1} | \mathbf{S}_n = s_n, \dots, \mathbf{S}_2 = s_2)$$
2041 
$$= \mathbb{P}(\mathbf{S}_{K+1} = s_{K+1} | \mathbf{X}_K = x_K, \dots, \mathbf{X}_1 = x_1)$$
2042 
$$= \mathbb{P}(\mathbf{S}_{K+1} = s_{K+1} | \mathbf{S}_K = s_K).$$
(by definition of  $\mathbf{S}_K$ )

This establishes the Markov property for **S**. 2045

#### D.7.2 CONCENTRATION INEQUALITIES FOR MARKOV CHAINS

2048 We first state a concentration inequality for time-homogeneous Markov chains that will be used to 2049 obtain our final bound. 2050

2046

2047

2003 2004

2017

2019

2021

2022

2023

2024

2025 2026

2034 2035 2036

2037 2038

**Proposition D.18** (McDiarmid's inequality for time-homogeneous Markov chains). Let  $S := (\mathbf{S}_1, \ldots, \mathbf{S}_N)$  be a Markov chain with value in a discrete, finite state space  $\Omega$  and mixing time  $t_{\min}(\varepsilon)$ . Let  $t_{\min} := \inf_{0 \le \varepsilon < 1} t_{\min}(\frac{\varepsilon}{2}) \cdot (\frac{2-\varepsilon}{1-\varepsilon})^2$ . If  $f : \Omega \to \mathbb{R}$  is such that there exists  $\mathbf{c} \in \mathbb{R}^N$  satisfying

$$\forall \mathbf{x}, \mathbf{y} \in \Omega, \quad f(\mathbf{x}) - f(\mathbf{y}) \le \sum_{i=1}^{N} \mathbf{c}_{i} \mathbb{1}_{\{\mathbf{x}_{i} \neq \mathbf{y}_{i}\}},$$

then we have for any  $u \ge 0$ ,

$$\mathbb{P}(|f(S) - \mathbb{E}_S[f(S)]| \ge u) \le 2 \exp\left(\frac{-2u^2}{\|\mathbf{c}\|_2^2 \cdot t_{\min}}\right).$$

*Proof.* We recall that Corollary 2.11 of Paulin (2015) ensures that for such a function f, we have

$$\mathbb{P}(|f(S) - \mathbb{E}[f(S)]| \ge u) \le 2 \exp\left(\frac{-2u^2}{\|\mathbf{c}\|_2^2 \cdot \tau_{\min}}\right),\tag{26}$$

where  $au_{\min}$  is defined as

$$\min \coloneqq \inf_{0 \le \varepsilon < 1} \tau(\varepsilon) \left( \frac{2 - \varepsilon}{1 - \varepsilon} \right)^2,$$

with  $\tau(\varepsilon)$  being the mixing time of a Markov chain *without assuming time homogeneity* (see Paulin (2015, Definition 1.4)). As in our case, we assume the time homogeneity, this inequality in Eq. (26) has to be adapted. Following Remark 1.5 of Paulin (2015), we notice that

$$\forall \varepsilon \in [0,1], \quad \tau(2\varepsilon) \le t_{\min}(\varepsilon) \le \tau(\varepsilon).$$

2078 Let  $0 \le \varepsilon < 1$ . Using the fact that  $\left(\frac{2-\varepsilon}{1-\varepsilon}\right)^2 > 0$ , the previous inequality ensures 

τ

$$\tau(\varepsilon) \le t_{\min}\left(\frac{\varepsilon}{2}\right) \iff \tau(\varepsilon)\left(\frac{2-\varepsilon}{1-\varepsilon}\right)^2 \le t_{\min}\left(\frac{\varepsilon}{2}\right)\left(\frac{2-\varepsilon}{1-\varepsilon}\right)^2.$$

Taking the infimum on the left-hand side leads to

$$\tau_{\min} = \inf_{0 \le \varepsilon < 1} \tau(\varepsilon) \left(\frac{2-\varepsilon}{1-\varepsilon}\right)^2 \le t_{\min}\left(\frac{\varepsilon}{2}\right) \left(\frac{2-\varepsilon}{1-\varepsilon}\right)^2.$$

As we took  $\varepsilon$  arbitrary in [0, 1), we can take the infimum on the right-hand side, which leads to

$$\tau_{\min} \leq t_{\min}.$$

2089 As the function  $x \to \exp\left(\frac{-2u^2}{\|\mathbf{c}\|_2^2 x}\right)$  is decreasing, we finally obtain 2091

$$\exp\left(\frac{-2u^2}{\|\mathbf{c}\|_2^2 \tau_{\min}}\right) \le \exp\left(\frac{-2u^2}{\|\mathbf{c}\|_2^2 t_{\min}}\right).$$
(27)

2094 Combining Eqs. (26) and (27) concludes the proof.

Similarly to Theorem 4.1, we want to apply Proposition D.18 to a function f that consists of sums of total variation. We investigate in the next section how to find the bounding vector **c** to apply Proposition D.18.

### 2100 D.7.3 FINDING THE BOUNDING VECTOR

2102 We want to apply the same arguments as in the proof of Theorem 4.1 to find the bounding vector c. 2103 The only difference in terms of setting is the definition of the probability of the next token. Indeed, in 2104 our case, we apply an extraction matrix  $\mathbf{M} \in \mathbb{R}^{d \times T}$  to recover the *d* states of the input Markov chain. 2105 We first prove the following technical lemma.

**2107 Lemma D.19.** Let  $d \leq T$  and consider a subset of d distinct elements of [T] that writes **2108**  $\mathbb{I}_d = (i_1 \leq i_2 \leq \ldots \leq i_d) \in [T]^d$ . We denote by  $\mathbf{M} \in \mathbb{R}^{d \times T}$  the matrix with rows  $\mathbf{M}_j = \mathbf{e}_{i_j}^{\top}$ , **2109** where  $\mathbf{e}_{i_j} \in \mathbb{R}^T$  has value 1 at entry  $i_j \in \mathbb{I}$  and 0 elsewhere. For any vector  $\mathbf{u} \in \mathbb{R}^T$ , we have

$$\|\mathbf{M}\mathbf{u}\|_{1} \le \|\mathbf{u}\|_{1}$$

2113 *Proof.* By definition of the  $\ell_1$ -norm, we have

2106

2110 2111 2112

2114 2115 2116

2124

2125

2126

2127 2128

2135 2136

2137

2143

2144

2145

2146

2155

2159

$$\|\mathbf{M}\mathbf{u}\|_{1} = \sum_{k=1}^{d} |\sum_{l=1}^{T} \mathbf{M}_{kl} \mathbf{u}_{l}| \le \sum_{k=1}^{d} \sum_{l=1}^{T} |\mathbf{M}_{kl} \mathbf{u}_{l}| \le \sum_{l=1}^{T} |\mathbf{u}_{l}| \sum_{k=1}^{d} |\mathbf{M}_{kl}|.$$

2117 Moreover, each column of M contains at most one non-zero entry (with value 1). Otherwise, it means 2118 that two  $\mathbf{e}_{i_j}$  are identical (as they only have one non-zero entry with value 1, having it at the same 2119 position ensures their equality) which contradicts the fact that the  $i_j$  where taken distinct. Hence, for 2120 all l, we have  $\sum_{k=1}^{d} |\mathbf{M}_{kl}| \leq 1$ , which concludes the proof.

We now prove a lemma analogous to Lemma D.8.

**Lemma D.20.** Let  $\mathbf{S} \in \mathbb{R}^{r \times n}$  denote the entry of the LLM  $f_{\Theta}$  and  $\mathbf{S}^{(L)}$  denote the output of the last layer before the softmax. Let  $d \leq T$  and consider a subset of d distinct elements of [T] that writes  $\mathbb{I}_d = (i_1 \leq i_2 \leq \ldots \leq i_d) \in [T]^d$ . We denote by  $\mathbf{M} \in \mathbb{R}^{d \times T}$  the matrix with rows  $\mathbf{M}_j = \mathbf{e}_{i_j}^{\top}$ , where  $\mathbf{e}_{i_j} \in \mathbb{R}^T$  has value 1 at entry  $i_j \in \mathbb{I}$  and 0 elsewhere. Then, the following inequality holds

$$\frac{1}{n\tau} \|\mathbf{M}\mathbf{W}_U\mathbf{S}^{(L)}\mathbb{1}_n\|_1 \le \frac{1}{\tau} \|\mathbf{W}_U^\top\|_{2,1}$$

2133 *Proof.* Applying Lemma D.19 with the matrix  $\mathbf{M} \in \mathbb{R}^d$  and the vector  $\frac{1}{n\tau} \mathbf{W}_U \mathbf{X}^{(L)} \mathbb{1}_n \in \mathbb{R}^T$  leads to

$$\frac{1}{n\tau} \|\mathbf{M}\mathbf{W}_U \mathbf{S}^{(L)} \mathbb{1}_n\|_1 \le \frac{1}{n\tau} \|\mathbf{W}_U \mathbf{X}^{(L)} \mathbb{1}_n\|_1$$

Applying Lemma D.8 concludes the proof.

The previous lemma can be used to show that the logarithm of the ratio between the true probability of the next token and the one estimated by the LLM  $f_{\Theta}$  is upper bounded as a function of the number of states of the Markov chain *d*, the temperature  $\tau$ , the upper-bound on  $W_U$  and some constant related to the ambiguity of language (see Eq. (1)).

**Proposition D.21** (Upper-bound on the logarithm). Consider an LLM  $f_{\Theta} \in \mathcal{F}$  and an input Markov chain  $X = (\mathbf{X}_1, \ldots, \mathbf{X}_{N_{icl}})$  with d states. We recall that  $B_U$  is the upper bound on the norm of  $\mathbf{W}_U$  in the definition of parameter space  $\mathcal{W}$ ,  $\tau$  is the softmax temperature, and  $p_{\min}$  is the constant related to the minimal transition probability between states. We have

$$\forall n \in [N], \quad \left| \log \left( \frac{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \right) \right| \le \bar{B} = \max\{ \log \left(d\right) + \frac{2B_U}{\tau}, \log \left(\frac{1}{p_{\min}}\right) \}.$$

**Proof.** The main idea of the proof is to bound the probability ratio and use the non-decreasing monotonicity of the log. Let  $n \in [N]$ . The model  $f_{\Theta}$  receives as input sequences of tokens  $\mathbf{S}_n$  of size  $n \leq K$ . We first lower-bound each term of the probability ratio. By definition of  $p_{\min}$ , we have

 $\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n) = \mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{X}_n) \ge p_{\min} > 0,$ (28)

where we used the Markov property for the first equality. We want to obtain a similar inequality for  $\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} | \mathbf{S}_n)$ . As the parameters  $\Theta$  of the LLM are in  $\mathcal{W}$ , we know that  $\|\mathbf{W}_U^{\top}\|_{2,1} \leq B_U$ . Lemma D.20 ensures that

$$\|\frac{1}{n\tau}\mathbf{M}\mathbf{W}_U\mathbf{S}^{(L)}\mathbb{1}_T\|_1 \le \frac{1}{\tau}\|\mathbf{W}_U^\top\|_{2,1} \le \frac{B_U}{\tau}.$$

We can then apply Lemma D.7 with  $c_1 = \frac{B_U}{\tau}$  and given that  $\frac{1}{T\tau} \mathbf{M} \mathbf{W}_U \mathbf{S}^{(L)} \mathbb{1}_T \in \mathbb{R}^d$ , it leads to

$$\mathbb{P}_{\Theta}(\cdot \mid \mathbf{S}_n) = \operatorname{softmax}\left(\frac{1}{n\tau}\mathbf{M}\mathbf{W}_U\mathbf{S}^{(L)}\mathbb{1}_n\right) \geq \frac{1}{d\exp\left(2B_U/\tau\right)},$$

where the inequality holds for each component of  $\mathbb{P}_{\Theta}(\cdot | \mathbf{S}_n)$ . This is in particular the case  $\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} | \mathbf{S}_n)$  which is the entry we are interested in, i.e., we have

$$\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n) \ge \frac{1}{d \exp\left(2B_U/\tau\right)}.$$
(29)

2169 Going back to the ratio of probability, consider the situation where we have 2170

$$\frac{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\mathbf{\Theta}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \ge 1.$$

2173 Then, using Eq. (29), we have

$$1 \leq \frac{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\mathbf{\Theta}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \leq \frac{1}{\mathbb{P}_{\mathbf{\Theta}}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \leq d \exp{(2B_U/\tau)},$$

2177 which implies, as the log is non-decreasing monotonically,

$$0 \le \log\left(\frac{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}\right) \le \log\left(d\exp\left(2B_U/\tau\right)\right) = \log\left(d\right) + \frac{2B_U}{\tau}.$$
(30)

2181 Similarly, consider the case where we have

$$\frac{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \ge 1,$$

and similarly to above, we can use Eq. (28) to obtain

$$1 \le \frac{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \le \frac{1}{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \le \frac{1}{p_{\min}}$$

 $\frac{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \le 1.$ 

2192 This implies

$$0 \le \log \left( \frac{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \right) \le \log \left( \frac{1}{p_{\min}} \right)$$

2195 which also rewrites

Then, we have

$$0 \le -\log\left(\frac{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}\right) \le \log\left(\frac{1}{p_{\min}}\right).$$
(31)

 $(\square \Theta(\Lambda_{n+1} | S_n))$   $(P_{\min})$ 2198 By definition of the absolute value, combining Eq. (30) and Eq. (31) leads to

$$\left| \log \left( \frac{\mathbb{P}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)}{\mathbb{P}_{\Theta}(\mathbf{X}_{n+1} \mid \mathbf{S}_n)} \right) \right| \le \max\{ \log \left(d\right) + \frac{2B_U}{\tau}, \log \left(\frac{1}{p_{\min}}\right) \}.$$

This concludes the proof.

We are now ready to upper-bound the total variation.

**Corollary D.22** (Upper-bound on the total variation). Consider an LLM  $f_{\Theta} \in \mathcal{F}$  and an input Markov chain  $X = (\mathbf{X}_1, \ldots, \mathbf{X}_{N_{icl}})$  with d states. We recall that  $B_U$  is the upper bound on the norm of  $\mathbf{W}_U$  in the definition of parameter space  $\mathcal{W}$ ,  $\tau$  is the softmax temperature, and  $p_{\min}$  is the constant related to the minimal transition probability between states. We have

$$\forall n \in [N], \quad d_{\mathrm{TV}}(\mathbb{P}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\mathbf{\Theta}}(\cdot \mid \mathbf{S}_n)) \le \sqrt{2 \max\{\log\left(d\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{p_{\min}}\right)\}} \coloneqq c_4.$$
(32)

*Proof.* Using Proposition D.21, we can directly apply Lemma D.6 with  $B = \max\{\log(d) + \frac{2B_U}{\tau}, \log\left(\frac{1}{p_{\min}}\right)\}$  for any  $n \in [N]$ . It leads to

$$\forall n \in [N], \quad d_{\mathrm{TV}}(\mathbb{P}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\mathbf{\Theta}}(\cdot \mid \mathbf{S}_n)) \le \sqrt{2 \max\{\log\left(d\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{p_{\min}}\right)\}}.$$

2222 This concludes the proof.

2224 D.7.4 CONCLUDING THE PROOF

We are now ready to state our main result.

**Theorem D.23** (Restatement of Theorem 4.4). Consider an LLM  $f_{\Theta} \in \mathcal{F}$ . We provide as input of  $f_{\Theta}$  a d-state Markov chain  $X = (\mathbf{X}_1, \ldots, \mathbf{X}_{N_{icl}})$ . The sequence of subsequences of the first n terms is denoted by  $S = (\mathbf{S}_1, \ldots, \mathbf{S}_{N_{icl}})$ . S is also a Markov chain, and we denote by  $t_{mix}(\varepsilon)$  its mixing time. Let  $t_{min} \coloneqq \inf_{0 \le \varepsilon < 1} t_{mix}(\frac{\varepsilon}{2}) \cdot (\frac{2-\varepsilon}{1-\varepsilon})^2$ . Let  $\delta > 0$ . Then, with probability at least  $1 - \delta$ ,

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

where  $\overline{B}$  is a constant depending on the parameters of the problem. More precisely,

$$\bar{B} = 2\sqrt{\max\{\log\left(d\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{p_{\min}}\right)\}}.$$

*Proof.* Let  $\vartheta \in W_{mc}$ . We first benefit from the metric properties of the total variation to decompose the risk.

$$\mathcal{R}_{icl}(\boldsymbol{\Theta}) = \frac{1}{N_{icl}} \sum_{n=1}^{N_{icl}} \mathbb{E}_{\mathbf{S}_n} [d_{\mathrm{TV}}(\mathbb{P}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\boldsymbol{\Theta}}(\cdot \mid \mathbf{S}_n))]$$

$$\leq \frac{1}{N_{icl}} \sum_{n=1}^{N_{icl}} \mathbb{E}_{\mathbf{S}_n} [d_{\mathrm{TV}}(\mathbb{P}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\boldsymbol{\vartheta}}(\cdot \mid \mathbf{S}_n)) + d_{\mathrm{TV}}(\mathbb{P}_{\boldsymbol{\vartheta}}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\boldsymbol{\Theta}}(\cdot \mid \mathbf{S}_n))]$$

$$\leq \frac{1}{N_{icl}} \sum_{n=1}^{N_{icl}} \mathbb{E}_{\mathbf{S}_n} [d_{\mathrm{TV}}(\mathbb{P}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\boldsymbol{\vartheta}}(\cdot \mid \mathbf{S}_n))]$$

$$+ \frac{1}{N_{icl}} \sum_{n=1}^{N_{icl}} \mathbb{E}_{\mathbf{S}_n} [d_{\mathrm{TV}}(\mathbb{P}_{\boldsymbol{\vartheta}}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\boldsymbol{\Theta}}(\cdot \mid \mathbf{S}_n))]$$

$$\leq \mathcal{R}_{icl}(\boldsymbol{\vartheta}) + K(\boldsymbol{\vartheta}, \boldsymbol{\Theta}). \tag{33}$$

By definition of the risk, we have

$$\widehat{\mathcal{R}}_{\text{icl}}(\boldsymbol{\vartheta}) = \frac{1}{N_{\text{icl}}} \sum_{n=1}^{N_{\text{icl}}} \underbrace{d_{\text{TV}}(\mathbb{P}(\cdot \mid \mathbf{S}_n), \mathbb{P}_{\boldsymbol{\vartheta}}(\cdot \mid \mathbf{S}_n))}_{=g_n(\mathbf{S}_n)} = \frac{1}{N_{\text{icl}}} \sum_{n=1}^{N_{\text{train}}} g_n(\mathbf{S}_n) = f(\mathbf{S}_1, \dots, \mathbf{S}_{N_{\text{icl}}}) = f(S).$$

Using Corollary D.22, we know that

$$|g_n(\mathbf{S}_n)| \le \sqrt{2\max\{\log\left(d\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{p_{\min}}\right)\}} \coloneqq c_4.$$

2268 Similarly to Theorem 4.1, and using the fact that  $S = (\mathbf{S}_1, \dots, \mathbf{S}_{N_{icl}})$  is a Markov chain, we can show that choosing  $\mathbf{c} \in \mathbb{R}^{N_{icl}}$  with all entries equal to  $\frac{2c_4}{N_{icl}}$  ensures that f verifies the condition in Proposition D.5, i.e.,

$$\forall S, \Sigma, \quad f(S) - f(\Sigma) \le \sum_{n=1}^{N_{\text{icl}}} \mathbf{c}_n \mathbb{1}_{\{\mathbf{S}_n \neq \mathbf{\Sigma}_n\}}.$$

Putting everything together, we can apply Proposition D.18 which leads to

$$\forall u \ge 0, \quad \mathbb{P}(|f(S) - \mathbb{E}_S[f(S)]| \ge u) \le 2 \exp\left(\frac{-2u^2}{t_{\min} \|\mathbf{c}\|_2^2}\right). \tag{34}$$

2278 Let  $u \ge 0$ . We have the following events ordering

$$(\mathbb{E}_S[f(S)] - f(S) \ge u) \subseteq (\mathbb{E}_S[f(S)] - f(S) \ge u) \cup (f(S) - \mathbb{E}_S[f(S)] \ge u)$$
$$= (|f(S) - \mathbb{E}_S[f(S)]| \ge u).$$

Hence, as u was taken arbitrary and using Eq. (34), we have

$$\forall u \ge 0, \quad \mathbb{P}(\mathbb{E}_S[f(S)] - f(S) \ge u) \le 2 \exp\left(\frac{-2u^2}{t_{\min} \|\mathbf{c}\|_2^2}\right).$$

We recall that by definition

$$f(S) = \widehat{\mathcal{R}}_{icl}(\boldsymbol{\vartheta}) \text{ and } \mathcal{R}_{icl}(\boldsymbol{\vartheta}) = \mathbb{E}_{S} \Big[ \widehat{\mathcal{R}}_{icl}(\boldsymbol{\vartheta}) \Big]$$

2290 Moreover, the inequality on the probability holds for any  $u \ge 0$ , we can choose u such that

$$\delta = 2 \exp\left(\frac{-2u^2}{t_{\min}\mathbf{c}\|_2^2}\right) \iff \frac{-2u^2}{t_{\min}\|\mathbf{c}\|_2^2} = \log\left(\frac{\delta}{2}\right) \iff u^2 = \frac{1}{2}t_{\min}\|\mathbf{c}\|_2^2 \log\left(\frac{2}{\delta}\right)$$
$$\iff u = \frac{1}{\sqrt{2}}\sqrt{t_{\min}}\|\mathbf{c}\|_2 \sqrt{\log\left(\frac{2}{\delta}\right)}.$$

Using the fact that

$$\|\mathbf{c}\|_{2} = \sqrt{\sum_{n=1}^{N_{\rm icl}} \mathbf{c}_{n}^{2}} = \sqrt{\sum_{n=1}^{N_{\rm icl}} \left(\frac{2c_{4}}{N_{\rm icl}}\right)^{2}} = \sqrt{\sum_{n=1}^{N_{\rm icl}} \frac{4c_{4}^{2}}{N_{\rm icl}^{2}}} = \sqrt{\frac{4c_{4}^{2}}{N_{\rm icl}}} = \frac{2c_{4}}{\sqrt{N_{\rm icl}}}$$

Using the fact that  $c_4 = \sqrt{2 \max\{\log(d) + \frac{2B_U}{\tau}, \log\left(\frac{1}{p_{\min}}\right)\}}$  (Corollary D.22), we obtain

$$\begin{split} u &= \frac{1}{\sqrt{2}} \frac{2c_4}{\sqrt{N_{\rm icl}}} \sqrt{\log\left(\frac{2}{\delta}\right)} = \frac{\sqrt{2}c_4}{\sqrt{N_{\rm icl}}} \sqrt{\log\left(\frac{2}{\delta}\right)} \\ &= \frac{2\sqrt{t_{\rm min}}}{\sqrt{\max\{\log\left(d\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{p_{\rm min}}\right)\}}}{\sqrt{N_{\rm train}}} \sqrt{\log\left(\frac{2}{\delta}\right)} \\ &= \bar{B} \sqrt{\frac{t_{\rm min}}{N_{\rm icl}}} \sqrt{\log\left(\frac{2}{\delta}\right)}, \end{split}$$

where we define

$$\bar{B} = 2\sqrt{\max\{\log\left(d\right) + \frac{2B_U}{\tau}, \log\left(\frac{1}{p_{\min}}\right)\}}.$$

2319 Putting everything together, we have

2320  
2321 
$$\mathbb{P}\left(\mathcal{R}_{\rm icl}(\vartheta) - \widehat{\mathcal{R}}_{\rm icl}(\vartheta) \ge \bar{B}\sqrt{\frac{t_{\rm min}}{N_{\rm icl}}}\sqrt{\log\left(\frac{2}{\delta}\right)}\right) \le \delta.$$

Taking the opposite event leads to the following inequality with probability at least  $1 - \delta$ 

$$\mathcal{R}_{\rm icl}(\boldsymbol{\vartheta}) \leq \widehat{\mathcal{R}}_{\rm icl}(\boldsymbol{\vartheta}) + \bar{B} \frac{\sqrt{t_{\rm min}}}{\sqrt{N_{\rm icl}}} \sqrt{\log\left(\frac{2}{\delta}\right)}$$

2327 Going back to the decomposition of the risk in Eq. (33) and rearranging the terms, we obtain

$$\mathcal{R}_{\rm icl}(\boldsymbol{\Theta}) \leq \widehat{\mathcal{R}}_{\rm icl}(\boldsymbol{\vartheta}) + K(\boldsymbol{\Theta}, \boldsymbol{\vartheta}) + \bar{B} \frac{\sqrt{t_{\min}}}{\sqrt{N_{\rm icl}}} \sqrt{\log\left(\frac{2}{\delta}\right)}.$$

<sup>2331</sup> As the left-hand side and the bound function of  $\overline{B}$  do not depend on  $\vartheta$ , we can put them both on the left side of the inequality and then take the infimum on  $\vartheta$ . Rearranging the terms to keep only  $\widehat{\mathcal{R}}_{icl}(\Theta)$  on the left side of the inequality leads to

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

which concludes the proof.

# 2340 E ADDITIONAL EXPERIMENTS

# 2342 E.1 EXPERIMENTAL SETUP AND TOKENIZATION

Experimental setup. To ensure a fair validation of our theoretical results, we conduct our experiments on some of the most recent and widely used LLMs: Gemma 2B (Team et al., 2024), Llama2
7B & 13B (Touvron et al., 2023b), Llama3 8B, Llama3.2 1B & 3B (Dubey et al., 2024), and Mistral 7Bv0.1 (Jiang et al., 2023).

**Tokenization.** As the models we consider have different tokenizations, we need to do this step 2349 with extra care as it is a crucial part of the experimental procedure. Indeed, LLMs' ability to handle 2350 numerical values has been proved to be dependent on the tokenization algorithm (Ali et al., 2024; 2351 Gruver et al., 2023; Singh & Strouse, 2024). The most widely used tokenization algorithm to-date, 2352 BPE (Sennrich et al., 2016), tends to assign tokens to arbitrary 3-digits numbers based on their 2353 occurrences in large-scale corpora, and the tokenizer's vocabulary size. As highlighted by (Gruver 2354 et al., 2023), this artifact severely hinders LLMs' ability to predict numerical values in-context. This is 2355 the case for popular LLMs such as GPT-3 (Brown et al., 2020). Newer models (LLama3, GPT-3.5, 2356 GPT-4) however, tend to have hard-coded rules on top of *BPE*, making them able to encode all 2357 3-digits numbers with their own token. Although this feature would accelerate the ICL procedure by eliminating the need for the Hierarchy-PDF algorithm in (Liu et al., 2024), the under-representability 2359 of larger numbers in the training data could be an issue. Other tokenization techniques that are numerical values-focused has been presented in the literature (Golkar et al., 2023; Wu et al., 2024), 2360 paving the way for another research direction that may benefit our method. 2361

2362 2363 2364

2365

2366 2367

2324

2325 2326

2328

2330

2335 2336 2337

2339

2348

Rodmap. In the rest of this section, we extend our experiments to study the following setups:

- In Appendix E.2: impact of the number of states *d*;
- In Appendix E.3: extension to Markov chains with  $p_{\min} = 0$ ;
- In Appendix E.4: impact of the tokenization;
- In Appendix E.5: extension to dynamical systems.
- 2368 2369
- 2370 E.2 IMPACT OF THE NUMBER OF STATES *d*

We further analyze the effect of the number of states *d* on the risk and consider randomly generated *d*-state transition matrices in Fig. 9. After a first stage of stagnation, the risk tends to take the correct scaling law coefficient. As in (Liu et al., 2024), we notice that considering randomly generated transition matrices seems to be difficult for an LLM to learn when there are more than 9 states. We interpret this behavior as the distribution shift term in Theorem 4.4. Indeed, the lack of structure in

these transition matrices can hinder the correct decay of this term. Note also that the increase in dtends to implicitly increase  $t_{\min}$ , which could have an impact on the upper bound on  $\mathcal{R}_{icl}$  (both in the generalization term and in the distribution shift term). We will now consider more structured Markov chains, and look at their impact on decay.



Figure 9: Impact of the number of states *d*. We plot the risk  $\mathcal{R}_{icl}$  as functions of  $N_{icl}$ , with 95% confidence intervals. Upper Left. 2–states Markov transition matrices. Upper Right. 4–states Markov transition matrices. Lower Left. 6–states Markov transition matrices. Lower Right. 8–states Markov transition matrices.

2406

2411

2417 2418

2419 2420

2421 2422

### 2407 E.3 More Structured Markov Chains 2408

In this section, we empirically verify our theoretical results on more general Markov chains that do not verify  $p_{\min} > 0$ .

### 2412 E.3.1 RANDOM WALKS

Random walks are a simple example of more structured Markov chains. Although we still have the possibility of discretizing the kernel of Markov chains with infinite state spaces as it is done in (Liu et al., 2024), we consider two types of random walks on finite state spaces.



Figure 10: Constrained random walk with d = 3.

**Constrained random walk.** We define the transition matrix P of a constrained random walk of d states as in Eq. (35). We draw the probabilistic graph in Fig. 10 for the case d = 3.

2425  
2426  
2427  
2428  
2429  

$$P_{ij} = \begin{cases} 1, & \text{if } i = 0 \text{ and } j = 1, \\ 1, & \text{if } i = d - 1 \text{ and } j = d - 2, \\ 0.5, & \text{if } 1 \le i \le d - 2 \text{ and } j = i - 1, \\ 0.5, & \text{if } 1 \le i \le d - 2 \text{ and } j = i + 1, \\ 0, & \text{otherwise.} \end{cases}$$
(35)



Fig. 11 highlights the scaling laws of Theorem 4.4, as well as the log(d) dependency. As before, the best-performing models generalize almost perfectly.

Figure 11: Constrained random walks. We plot the risk  $\mathcal{R}_{icl}$  as functions of  $N_{icl}$ , with 99% confidence intervals. We consider different size d. Upper Left. Llama2 7B Upper Right. Llama2 13B Lower Left. Mistral 7Bv0.1 Lower Right. Gemma 2B

**Polygonal random walk.** We define the transition matrix  $\mathbb{P}$  of a polygonal random walk of d states as in Eq. (36). We draw the probabilistic graph in Fig. 12 for the case d = 4.

$$P_{ij} = \begin{cases} 0.5, & \text{if } j = (i+1) \mod d \text{ (clockwise transition)}, \\ 0.5, & \text{if } j = (i-1) \mod d \text{ (counterclockwise transition)}, \\ 0, & \text{otherwise}. \end{cases}$$
(36)

We draw the same conclusions as above for this second type of random walk, in Fig. 13.



Figure 12: Polygonal random walk with d = 4.



Figure 13: Polygonal random walks. We plot the risk  $\mathcal{R}_{icl}$  as functions of  $N_{icl}$ , with 99% confidence intervals. We consider different size d. Upper Left. Llama2 7B Upper Right. Llama2 13B Lower Left. Mistral 7Bv0.1 Lower Right. Gemma 2B

### 2511 E.3.2 INNER CLIQUES AND OUTER RIMS

**Inner Cliques and Outer Rims.** We also want to test our method on the class of Markov chain put forward in (Wolfer & Kontorovich, 2019) to derive their lower bound. Let  $\eta > 0$  and d = 3k for some  $k \in \mathbb{N}$ , and define the collection of Markov matrices  $\mathcal{H}_{\eta} = \{M_{\eta,\tau} : \tau \in \{0,1\}^{d/3}\}$ . Every element of this set consists of an *inner clique* and an *outer rim*.  $M_{\eta,\tau}$  is the block matrix defined as follows,

2509 2510

$$\boldsymbol{M}_{\eta,\boldsymbol{\tau}} = \begin{pmatrix} C_{\eta} & R_{\boldsymbol{\tau}} \\ R_{\boldsymbol{\tau}}^{\mathsf{T}} & L_{\boldsymbol{\tau}} \end{pmatrix}$$

where  $C_{\eta} \in \mathbb{R}^{d/3 \times d/3}$ ,  $L_{\tau} \in \mathbb{R}^{2d/3 \times 2d/3}$ , and  $R_{\tau} \in \mathbb{R}^{d/3 \times 2d/3}$  are given by

$$L_{\boldsymbol{\tau}} = \frac{1}{8} \operatorname{diag} \left( 7 - 4\tau_1 \varepsilon, 7 + 4\tau_1 \varepsilon, \dots, 7 - 4\tau_{d/3} \varepsilon, 7 + 4\tau_{d/3} \varepsilon \right),$$

 $C_{\eta} = \begin{pmatrix} \frac{3}{4} - \eta & \frac{\eta}{d/3 - 1} & \cdots & \frac{\eta}{d/3 - 1} \\ \frac{\eta}{d/3 - 1} & \frac{3}{4} - \eta & \ddots & \vdots \\ \vdots & \ddots & \ddots & \frac{\eta}{d/3 - 1} \\ \frac{\eta}{d/3 - 1} & \eta & \eta & \eta \end{pmatrix},$ 

2523 2524 2525

2521 2522

2526

2527

$$R_{\tau} = \frac{1}{8} \begin{pmatrix} 1 + 4\tau_1 \varepsilon & 1 - 4\tau_1 \varepsilon & 0 & \dots & \dots & 0 \\ 0 & 0 & 1 + 4\tau_2 \varepsilon & 1 - 4\tau_2 \varepsilon & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & \dots & 0 & 1 + 4\tau_{d/3} \varepsilon & 1 - 4\tau_{d/3} \varepsilon \end{pmatrix}.$$

2532 2533

2531

2534 We provide in Fig. 14 a probabilistic graph of the case  $M_{\eta,0}$  and d = 9.

Fig. 15 compares different LLMs with the frequentist method, on the case depicted in Fig. 15 with  $\eta = 0.02$ . Although the frequentist method achieves a lower loss, the power laws seem to be the same with LLMs.



Figure 15: We plot the risk  $\mathcal{R}_{icl}$  as functions of  $N_{icl}$ , with 95% confidence intervals. Upper Left. Llama2 7B Upper Right. Llama2 13B Lower Left. Mistral 7Bv0.1 Lower Right. Gemma 2B

2590

# 2592 E.4 RECENT MODELS: IMPACT OF THE TOKENIZATION

As explained in Appendix E.1, models like Llama 3 tokenize 3-digit numbers with a single token. This saves a lot of inference compute time, but not necessarily in terms of performance when considering Markov chains with a few number of states d, since we have to separate the states by a comma to force tokenization into a single digit (e.g. the transitions  $1 \rightarrow 0 \rightarrow 1$  will be prompted as 1, 0, 1 (5 tokens) instead of 101 (1 token). In Fig. 16, we reproduce the same experiment as in Fig. 5(left), but with Llama 3 models. The scaling laws are quite good, but much less so than those obtained with Gemma 2B and Mistral 7Bv0.1 on the same inputs. On the other hand, with these models, it can be extremely interesting to consider Markov chains with many states, as we did in Fig. 6(right). In the next section, we will use LLama3 to learn other dynamic systems presented in Liu et al. (2024). 



Figure 16: In-context scaling laws for LLama 3 herd of models. We plot the risk  $\mathcal{R}_{icl}$  as functions of  $N_{icl}$ , with 95% confidence intervals.

2617 E 4

E.5 DYNAMICAL SYSTEMS

We consider four of the dynamic systems highlighted in (Liu et al., 2024) : a geometric Brownian motion, a correlated Gaussian, an uncorrelated Gaussian and an uncorrelated uniform processes. We display in Fig. 17 the risks of LLama 3 8B and the frequentist method, which once again highlight the emerging capacity of in-context learning.



Figure 17: LLama 3 8B on dynamical systems. We plot the risks  $\mathcal{R}_{icl}$  as functions of  $N_{icl}$  for LLama 3 8B and the frequentist approach (Wolfer & Kontorovich, 2019) with 95% confidence intervals. Upper Left. Geometric Brownian motion. Upper Right. Correlated Gaussian. Lower Left. Uncorrelated Gaussian. Lower Right. Uncorrelated Uniform.

#### F EXTENDED RESULTS WITH THE KL DIVERGENCE

As explained in Remark 4.1, the total variation is the natural choice to define the risks in Eq. (2). Another possibility in the Markov chain literature is to use the KL divergence to compare probability distributions (Hao et al., 2018). This is an interesting candidate as the KL divergence is naturally connected to the cross-entropy loss commonly used to train neural networks (the cross-entropy corresponds to the KL divergence between the true distribution and the predicted softmax distribution (Blondel et al., 2019). In this section, we discuss the extension of the theoretical results of Section 4 by replacing the TV distance with the KL divergence in the risks' definition, i.e.,

$$\mathcal{R}(\mathbf{\Theta}) \coloneqq \mathbb{E}_{\mathbf{S} \sim \mathbb{P}_{\mathcal{L}}}[d_{\mathrm{KL}}(\mathbf{Q}^{*}(\mathbf{S}, \cdot) || \mathbf{Q}_{f}(\mathbf{S}, \cdot))], \, \widehat{\mathcal{R}}(\mathbf{\Theta}) \coloneqq \frac{1}{N} \sum_{n=1}^{N} d_{\mathrm{KL}}(\mathbb{P}_{\mathcal{L}}(\cdot \mid \mathbf{S}_{n}) || \mathbb{P}_{\mathbf{\Theta}}(\cdot \mid \mathbf{S}_{n})). \, (37)$$

#### **PRE-TRAINING GENERALIZATION BOUNDS** F.1

Theorem 4.1, Corollary 4.2 and Corollary 4.3 related to the pre-training phase in Section 4.1 can be obtained similarly if the risks are defined with the KL divergence following Eq. (37). Indeed, the key step to derive the proofs is to obtain a similar result to Lemma D.6 but with the KL divergence. The next lemma provides this result.

**Lemma F.1.** Consider two probability distributions  $\mathbb{P}, \mathbb{Q}$  defined on a measure space  $(\Omega, \mathcal{F})$ and a  $\sigma$ -finite measure  $\nu$  on  $(\Omega, \mathcal{F})$ . Let p, q be the corresponding probabilities densities, i.e., we have  $\mathbb{P}(d\omega) = q(\omega)\nu(d\omega)$  and  $\mathbb{Q}(d\omega) = p(\omega)\nu(d\omega)$ . If there exists a non-negative constant *B* such that for any  $z \in \Omega$ ,  $\left| \log \sqrt{\frac{\mathbb{P}(z)}{\mathbb{Q}(z)}} \right| \le B$ , then we have

$$d_{\mathrm{KL}}(\mathbb{P}||\mathbb{Q}) \leq B$$

Proof. We have

$$0 \leq d_{\mathrm{KL}}(\mathbb{P}||\mathbb{Q}) = |d_{\mathrm{KL}}(\mathbb{P}||\mathbb{Q})|$$
$$= \left| \int \mathbb{P}(z) \log(\frac{\mathbb{P}(z)}{\mathbb{Q}(z)}) dz \right|$$
$$\leq \int |\mathbb{P}(z)| |\log(\frac{\mathbb{P}(z)}{\mathbb{Q}(z)})| dz$$
$$\leq B \int |\mathbb{P}(z)| dz$$
$$= B \int \mathbb{P}(z) dz$$
$$= B,$$

which concludes the proof.

We can now state the results similar to Theorem 4.1, Corollary 4.2 and Corollary 4.3 from the pre-training phase when the risk is defined according to Eq. (37).

**Theorem F.2** (Pre-training generalization bound). *Consider an LLM*  $f_{\Theta} \in \mathcal{F}$ . We denote by  $\Gamma$  the mixing matrix of the pre-training sequences of tokens  $(\mathbf{S}_1, \ldots, \mathbf{S}_{N_{\text{train}}})$ . Let  $0 < \delta < 1$ , then with probability at least  $1 - \delta$ ,

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

where  $\overline{B} = \sqrt{2} \| \mathbf{\Gamma} \| \max\{ \log (T) + 2B_U / \tau, \log (1/c_0) \}$  is a constant depending on the parameters of the problem.

*Proof.* The proof simply follows from the proof of Theorem 4.1 by replacing the upper bound  $\sqrt{2B}$  by *B* (with the appropriate upper-bound *B*) when Lemma D.6 is used in the proof.

**Corollary F.3** (Depth-dependent bound). Consider an LLM  $f_{\Theta} \in \tilde{\mathcal{F}} := \{f_{\Theta} \mid \Theta \in \tilde{\mathcal{W}}\}$ . With the same assumptions as in Theorem 4.1, we have

$$\mathcal{R}_{ ext{pre}}(\mathbf{\Theta}) \leq \widehat{\mathcal{R}}_{ ext{pre}}(\mathbf{\Theta}) + rac{ar{B}}{\sqrt{N_{ ext{train}}}} \sqrt{\log\left(rac{2}{\delta}
ight)},$$

where  $\bar{B} = \sqrt{2} \| \mathbf{\Gamma} \| \max\{ \log{(T)} + 2(B_{\Theta})^L / \tau, \log{(1/c_0)} \}$  is a constant depending on the parameters of the problem, and  $B_{\Theta} = [(1 + rmB_1B_2)(1 + \frac{r^3}{H}B_OB_V)](B_{\text{tok}}B_U)^{1/L}$ .

*Proof.* The proof simply follows from the proof of Theorem 4.1 by replacing the upper bound  $\sqrt{2B}$  by *B* (with the appropriate upper-bound *B*) when Lemma D.6 is used in the proof.

**Corollary F.4** (Sample complexity). Let *B* be the parameter-dependent constant of Theorem F.2 or Corollary F.3. Let  $\delta \in [0,1]$  and let  $\epsilon > 0$ . If  $N_{\text{train}} \ge N^* := \left\lceil \frac{4\overline{B}^2}{\epsilon^2} \log\left(\frac{2}{\delta}\right) \right\rceil$  and if we assume a perfect pre-training error for  $f_{\Theta}$ , then we have with probability at least  $1 - \delta$ ,

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

*Proof.* The proof simply follows from the proof of Theorem 4.1 by replacing the upper bound  $\sqrt{2B}$  by *B* (with the appropriate upper-bound *B*) when Lemma D.6 is used in the proof.

2731 F.2 LIMITATIONS

We recall from Remark 4.1 that the TV distance is a natural choice to compare transition matrices in the Markov chain literature. In addition, while the KL divergence can be used to compare probability distributions, it does not define a metric space. Hence, we cannot straightforwardly extend Theorem 4.4 with the KL divergence because the proof relies on the use of the triangular inequality. As Theorem 4.4 is one of our main results and enables us to show that the theory and the practice align (Section 5), this also contributed to our preference for the TV distance instead of the KL divergence.