

# On the In-context Generation of Language Models

Anonymous ACL submission

## Abstract

Large language models (LLMs) are found to have the ability of in-context generation (ICG): when they are fed with an in-context prompt containing a somehow similar examples, they can implicitly discover the pattern of them and then complete the prompt in the same pattern. ICG is curious, since language models are not completely trained in the way same as the in-context prompt, and the distribution of examples in the prompt differs from that of sequences in the pretrained corpora. This paper provides a systematic study of the ICG ability of language models, covering discussions about its source and influential factors, in the view of both theory and empirical experiments. Concretely, we first propose a plausible latent variable model to describe the distribution of the pretrained corpora, and then formalize ICG as a problem of next topic prediction. With this framework, we can prove that the repetition nature of a few topics ensures the ICG ability on them theoretically. Then, we use this controllable pretrained distribution to generate several medium-scale synthetic datasets (token scale: 2.1B~3.9B) and experiment with different settings of Transformer architectures (parameter scale: 4M~234M). Our experimental results further offer insights into how factors of data and model architectures influence ICG.

## 1 Introduction

As the data and parameter scale continue to increase, large language models (LLMs) have shown strikingly emergent abilities (Wei et al., 2022a), where one of the most exciting ones is in-context learning (ICL) (Brown et al., 2020). Given an *in-context prompt* that concatenates a few *in-context examples* and a query input, LLMs can somehow implicitly guess the "topic" of those examples and complete the query input in the desired way. Furthermore, LLMs can imitate those examples using the topic learned in context (Meyerson et al., 2023).

For instance, Llama2-13B (Touvron et al., 2023) is able to generate plausible sequences of the topic of in-context examples, as shown in Figure 1. This in-context generation (ICG) ability forms the foundation of multiple few-shot prompting methods like ICL and its variants like Chain-of-thoughts (Wei et al., 2022b).

Intuitively, one might comment that LLMs learn the ICG ability from data in the *repetition mode*, which roughly refers to a type of text concatenated with sequences under the same topic. This is true to some extent. As known, typical pretrained corpora contain (e.g. CommonCrawl<sup>1</sup>) internet data which has an unneglectable portion of array-page data such as IMDB review pages<sup>2</sup>. After preprocessing, these pages are converted to repetition mode data, as shown in Figure 1a. However, this isn't enough to explain the ICG ability, since LLMs can also generate sequences of in-context learned topics that don't appear to repeat and even are unseen in the pretrained corpora. For example, Figure 1 shows sampled completions of Llama2-13B given in-context prompts of different types of topics:

1. The first one is a *repeated topic* called "movie review" (Figure 1a), where Llama2-13B naturally has the ICG ability on it since this topic appears to repeat in the pretrained corpora as mentioned.

2. The second type *nonrepeated topic* refers to those that appear in the pretrained corpora but never repeat, e.g., forward method in any class inherited from nn.Module of Pytorch (Paszke et al., 2019) code (Figure 1b). However, Llama2-13B can also generate plausible code of forward method when prompting a few ones.

3. The last type *unseen topic* includes those that never appear in the pretrained corpora. For example, "unnatural addition" generates 2-digit arith-

<sup>1</sup><https://commoncrawl.org>

<sup>2</sup><https://www.imdb.com>

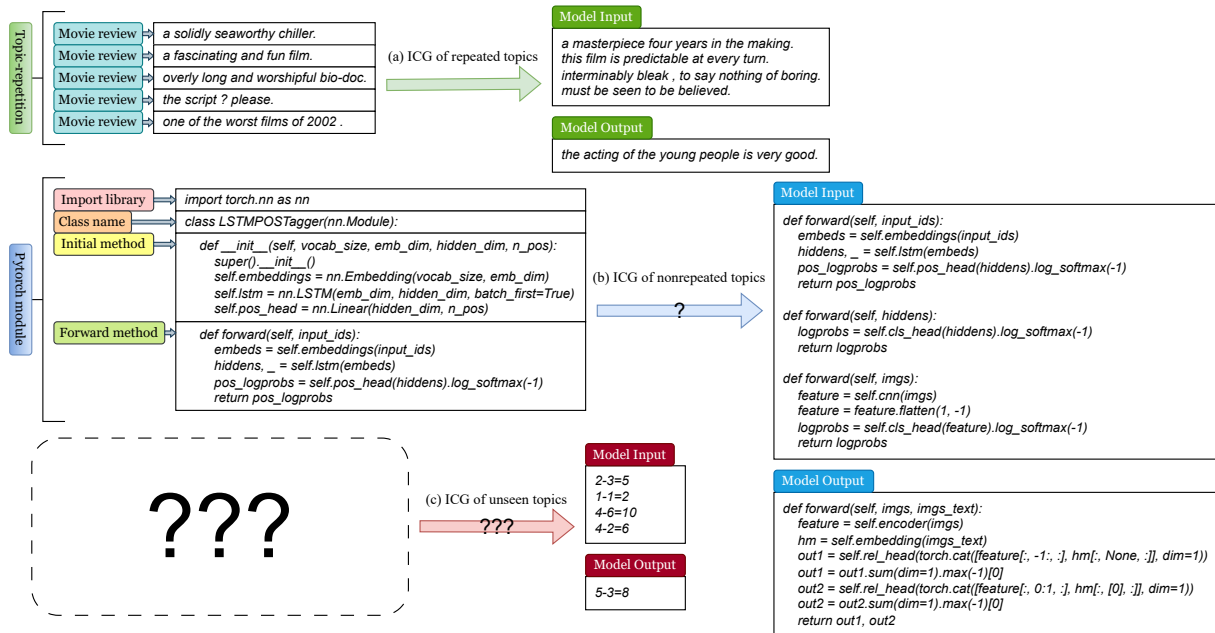


Figure 1: ICG examples (generated from Llama2-13B) of different kinds of topics.

079 metic expressions that input subtraction but expect  
 080 addition (like "1-1=2"), which is intuitively be-  
 081 lieved to never be seen in the pretrained corpora  
 082 (Rong, 2021). However, Llama2-13B can also recog-  
 083 nize this topic and generate plausible sequences  
 084 in context, as shown in Figure 1c.

085 The above results show that LLMs can general-  
 086 ize the repetition mode to nonrepeated and un-  
 087 seen topics. We term this phenomenon as the  
 088 topic generalization of ICG, abbreviated as ICG-  
 089 generalization. ICG-generalization is curious be-  
 090 cause LLMs are not explicitly trained in the way  
 091 they test. The biggest challenge of studying ICG  
 092 and its generalization is that the true pretrained  
 093 distribution is not accessible. Thus, we don't know  
 094 the topic of a span or whether it appears to repeat,  
 095 making it difficult to evaluate the ICG abilities of  
 096 LLMs. To address this, we turn to synthetic data  
 097 generated from a known and controlled pretrained  
 098 distribution (Bowman et al., 2015; McCoy et al.,  
 099 2018; White and Cotterell, 2021; Xie et al., 2021;  
 100 Papadimitriou and Jurafsky, 2023; Jumelet and  
 101 Zuidema, 2023). The distribution is a hierarchical  
 102 latent variable model (LVM) as shown in Figure 2,  
 103 where a document is guided by two kinds of latent  
 104 variables. The distribution is not only plausible to  
 105 explain true pretrained data but also convenient  
 106 for analysis since it decouples different levels of  
 107 uncertainties.

108 Through the proposed pretrained distribution, we  
 109 can naturally formalize ICG as a problem of next  
 topic prediction, and then conduct mathematical

110 analysis. We first theoretically prove that (Theorem  
 111 1), under some mild assumptions, if the language  
 112 model fits the pretrained distribution well, then  
 113 it's guaranteed to have the ICG ability on repeated  
 114 topics in terms of convergence in probability. As  
 115 a result, the ICG distribution (i.e., the generative  
 116 distribution conditioned on the in-context prompt)  
 117 converges to the true topic-paragraph distribution  
 118 in probability. Next, we study ICG-generalization  
 119 via exhaustive experiments, revealing that ICG-  
 120 generalization is caused by both factors of data  
 121 and models. Concretely, we use the controllable  
 122 pretrained distribution to generate several synthetic  
 123 datasets (token scale: 2.1B~3.9B), and train Trans-  
 124 former (Vaswani et al., 2017) language models with  
 125 different settings (parameter scale: 4M~234M). Ex-  
 126 periments show that data compositionality, propor-  
 127 tion of repeated topics, Transformer's parameter  
 128 scale, and window size play crucial roles in en-  
 129 abling ICG-generalization, while the data topic  
 130 uncertainty and Transformer's attention head size  
 131 have few influences<sup>3</sup>. Our study provides insights  
 132 to better understanding the ICG ability and LLMs.

## 2 Settings 133

### 2.1 Pretrained Distribution 134

135 We assume the pretrained distribution is a hierarchi-  
 136 cal LVM as shown in Figure 2, where a document is

<sup>3</sup>These results are consistent with previous works about  
 attention head pruning (Michel et al., 2019; Voita et al., 2019)  
 and the importance of large attention size (Ratner et al., 2023).

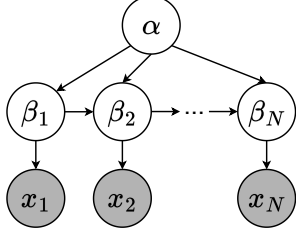


Figure 2: Bayesian network of the pretrained distribution, where the non-shaded nodes are latent variables.

generated via the following steps: 1) Draw a latent mode  $\alpha \in A$  from the mode prior  $p(\alpha)$ . 2) Draw a latent outline  $\beta_{1:N} \in B^N$  containing topics of different paragraphs from the Markov mode-outline distribution  $p(\beta_{1:N}|\alpha)$  parameterized by the mode  $\alpha$ . 3) Sample each paragraph  $x_i \in \Sigma^*$  ( $\Sigma$  is the vocabulary) individually from the topic-paragraph distribution  $p(x|\beta_i)$ , and concatenate them with delimiters. The joint distribution of this LVM is:

$$p(\alpha, \beta_{1:N}, x_{1:N}) = p(\alpha)p(\beta_{1:N}|\alpha) \prod_{i=1}^N p(x_i|\beta_i) \quad (1)$$

This distribution is plausible because: 1) It has a clear realistic interpretation of how humans write documents. Generally, humans would first determine the literature genre (e.g., narrative, letter, and so on), and then plan a specific structure of that genre before writing, as shown in Figure 1. Such a process is modeled via the mode prior  $p(\alpha)$  and the mode-outline distribution  $p(\beta_{1:N}|\alpha)$ . 2) It is capable of describing any language marginal distribution via the marginalization over latent variables. Also, it is convenient to analyze because of disentanglement: two kinds of uncertainties, topic-transition and generation of paragraphs are handled by two separated models  $p(\beta_n|\beta_{1:n-1}, \alpha)$  and  $p(x_n|\beta_n)$ , respectively, but not the entangled marginal  $p(x_{1:N})$ .

### 2.1.1 Assumptions

The pretrained distribution has three additional assumptions. Firstly, as mentioned, typical pretrained distributions for LLMs include the repetition mode  $\hat{\alpha} \in A$  that only generates repeated outlines  $\beta^N$  ( $\beta \in B$ ) ( $\beta^N$  represents a  $N$ -length outline that each topic within is  $\beta$ ). This formally raises the following:

**Assumption 1.** *There exists a mode  $\hat{\alpha} \in A$  called repetition mode such that  $p(\beta_{n+1}|\beta_n, \hat{\alpha}) = 1(\beta_{n+1} = \beta_n)$  for all timesteps  $n$ . Other modes*

$\alpha \in A/\hat{\alpha}$  are called continuous modes, since the outline under them seems to shift gradually and continuously.

Secondly, we have to ensure that different modes and topics are different to get rid of redundancy. That is, they should be distinguished in terms of distance measure of distribution:

**Assumption 2.** *For two different modes  $\alpha, \alpha' \in A$  and an arbitrary context  $x_{1:n}$ , define:*

$$\text{KL}_n(\alpha||\alpha') := \sum_x p(x|x_{1:n}, \alpha) \log \frac{p(x|x_{1:n}, \alpha)}{p(x|x_{1:n}, \alpha')} \quad (2)$$

We assume that  $\text{KL}_n(\alpha||\alpha') \geq \log c_1 > 0$ . Likewise, for two different topics  $\beta, \beta' \in B$ , define:

$$\text{KL}(\beta||\beta') := \sum_x p(x|\beta) \log \frac{p(x|\beta)}{p(x|\beta')} \quad (3)$$

We assume that  $\text{KL}(\beta||\beta') \geq \log c_2 > 0$ .

Thirdly, for convenience and without loss of plausibility, we assume that:

**Assumption 3.** *For each paragraph  $x \in \Sigma^*$ , its support from any topic  $\beta \in B$  is bounded:  $0 < c_3 \leq p(x|\beta) \leq c_4 < 1$ .*

### 2.1.2 Topic Types

With Assumption 1, the likelihood of any repeated outline  $\beta^N$  under the repetition mode  $\hat{\alpha}$  only depends on the topic itself:

$$p(\beta^N|\hat{\alpha}) = p(\beta_1 = \beta|\hat{\alpha}) := p(\beta|\hat{\alpha}) \quad (4)$$

where  $p(\beta|\hat{\alpha})$  is the repetition prior measuring how often the topic  $\beta$  is chosen to repeat under mode  $\hat{\alpha}$ . Analogously, let  $p(\beta)$  be the topic prior assessing the frequency of the topic  $\beta$ :

$$p(\beta) := \sum_{\alpha \in A} p(\beta|\alpha)p(\alpha) \quad (5)$$

According to the appearance, we can formally group topics  $\beta \in B$  into three mutually exclusive sets, as shown in Figure 1:

1. Repeated set  $R$ .  $\forall \beta \in R, p(\beta|\hat{\alpha}) > 0$ . That is, each topic within appears to repeat in the pretrained distribution. By intuition, repeated topics account for a very small proportion of all topics in realistic data, i.e.,  $r_R = |R|/|B|$  is small.

2. Nonrepeated set  $C$ .  $\forall \beta \in C, p(\beta|\hat{\alpha}) = 0, p(\beta) > 0$ . In other words, this set contains topics that don't repeat but appear in the pretrained corpora.

3. Unseen set  $U$ .  $\forall \beta \in U$ ,  $p(\beta) = 0$ . Topics in this set are never seen in the pretrained corpora.

## 2.2 Problem Formalization

Consider a language model  $p_{\text{LM}}$  trained on samples of the above pretrained distribution  $p$ . The ICG ability could be formalized as:

**Hypothesis 1.** *Given a language model  $p_{\text{LM}}$  trained on the pretrained distribution  $p$  and an in-context prompt  $x_{1:N}$ , where each sample  $x_n \sim p(x|\hat{\beta})$ , the in-context topic-repetition rate (ICTR), i.e., the probability that the language model generates a paragraph belong to topic  $\hat{\beta}$  when prompting with  $x_{1:N}$  is somehow close to 1:*

$$p_{\text{LM}}(\hat{\beta}|x_{1:N}) := p_{\text{LM}}(\beta_{N+1} = \hat{\beta}|x_{1:N}) \approx 1 \quad (6)$$

Accordingly, the model ICG distribution  $p_{\text{LM}}(x|x_{1:N})$  is somehow closed to the true topic-paragraph distribution  $p(x|\hat{\beta})$ :

$$p_{\text{LM}}(x|x_{1:N}) \approx p(x|\hat{\beta}) \quad (7)$$

Thus, we formalize ICG as next topic prediction, where language models seem to implicitly choose the topic of in-context examples as the next topic. Our goal is to find support for this hypothesis from the perspective of both theory and empirical experiments.

## 3 Theoretical Support

Intuitively, the pretrained distribution itself ensures the ICG ability for repeated topics  $R$ . This can be explicitly formalized by the following theorem:

**Theorem 1.** *Given an in-context prompt  $x_{1:N}$ , where each sample  $x_n \sim p(x|\hat{\beta})$  and  $\hat{\beta} \in R$ , the pretrained distribution have the following properties:*

1. The data ICTR<sup>4</sup> converges to 1 in probability (corollary 4):

$$\text{plim}_{N \rightarrow \infty} p(\hat{\beta}|x_{1:N}) = 1 \quad (8)$$

where we denote  $p(\beta_{N+1} = \beta|x_{1:N}) := p(\beta|x_{1:N})$ .

2. For any candidate paragraph  $x \in \Sigma^*$ , the data ICG distribution  $p(x|x_{1:N})$  converges to true topic-paragraph  $p(x|\hat{\beta})$  in probability (corollary 5):

$$\text{plim}_{N \rightarrow \infty} p(x|x_{1:N}) = p(x|\hat{\beta}) \quad (9)$$

<sup>4</sup>Note that we use the prefix "data" to distinguish values from pretrained distribution and language model distribution.

If the language model is expressive enough, it would gradually approach the pretrained distribution with the increase of the number of training examples<sup>5</sup>. As a result, it would exhibit the same properties as shown in Theorem 1. Therefore, the ICG ability for repeated topics directly originates from the pretrained corpora.

Detailed theoretical results are provided in Appendix B, and here, we only present a proof sketch. *Proof Sketch.* According to Section 2.1,  $\forall x \in \Sigma^*$ , the data ICG distribution is:

$$p(x|x_{1:N}) = \sum_{\beta \in B} p(\beta|x_{1:N})p(x|\beta) \quad (10)$$

Therefore, the data ICG distribution  $p(x|x_{1:N})$  is dominated by the topic predictive distribution  $p(\beta|x_{1:N})$ , i.e., ICTR.  $p(\beta|x_{1:N})$  can be further decomposed as the mixture of modes:

$$p(\beta|x_{1:N}) = \sum_{\alpha \in A} p(\alpha|x_{1:N})p(\beta|x_{1:N}, \alpha) \quad (11)$$

Firstly, we can prove that if  $\hat{\beta} \in R$ , then  $\text{plim}_{N \rightarrow \infty} p(\hat{\alpha}|x_{1:N}) = 1$  (corollary 1). Therefore, the mixture in formula (11) focuses on the component of repetition mode  $p(\beta|x_{1:N}, \hat{\alpha})$  when  $N$  is large:

$$\begin{aligned} p(\beta|x_{1:N}) &\approx p(\beta|x_{1:N}, \hat{\alpha}) \\ &= \frac{p(\beta|\hat{\alpha}) \prod_{n=1}^N p(x_n|\beta)}{p(x_{1:N}|\hat{\alpha})} \end{aligned} \quad (12)$$

This form is exactly the Bayesian posterior distribution, which is in accord with previous works connecting ICL and Bayesian statistics (Xie et al., 2021; Wang et al., 2023b; Hahn and Goyal, 2023). Likewise, it turns out that the if  $\hat{\beta} \in R$ , then  $\text{plim}_{N \rightarrow \infty} p(\hat{\beta}|x_{1:N}, \hat{\alpha}) = 1$  (corollary 3), thus establishing the first point of theorem 1. Since the data ICG distribution  $p(x|x_{1:N})$  depends on the topic predictive distribution  $p(\beta|x_{1:N})$ , we can prove the second point of theorem 1 analogously<sup>6</sup>. In Appendix B and C, we also present a detailed formula of the convergence, in which the convergence speed depends on the distinguishment of different modes and topics.

<sup>5</sup>Previous works (Xie et al., 2021; Hahn and Goyal, 2023) typically take this as the null hypothesis.

<sup>6</sup>Based on of theorem 1, for regular in-context learning scenario where each example in the prompt is a tuple  $(x_n, y_n)$  consisting with an input  $x_n$  and an output  $y_n$ , we can also obtain similar theoretical conclusions about the ICL ability. Details are shown in proposition 5 and corollary 6.

## 4 Experiments

Theory 1 can't ensure the ICG ability for nonrepeated and unseen topics  $\beta \in C \cup U$  because they have a zero repetition prior  $p(\beta|\hat{\alpha}) = 0$  and so the posterior under repetition mode is also zero:  $p(\beta|x_{1:N}, \hat{\alpha}) = 0$ . Then, the correct component  $p(x|\beta)$  would never be selected under the repetition mode, preventing the ICG/ICL ability as a consequence.

However, this is contrary to the real case, where LLMs have the ICG-generalization ability: they are able to generalize ICG/ICL abilities to nonrepeated and unseen topics  $\beta \in C \cup U$ . We speculate that this might be caused by factors in both data and model side:

- Data side: The compositionality of natural language (Grandy, 1990) and the proportion of repeated topics  $r_R$ . Compositionality considers the meaning of a linguistic unit results from the individual meanings of its sub-parts, and how they are combined (Anderson, 2018). Thus, nonrepeated and unseen topics might share the same "sub-topics" with repeated topics. The bigger the proportion of repeated topics, the more frequently those sub-topics are shared. Therefore, LLMs may be able to recombine those sub-topics to recognize those out-of-distribution topics in the repetition mode and exhibit generalization.

- Model side: The Transformer (Vaswani et al., 2017) structure. As the mainstream architecture of NLP, the success of Transformer is believed to originate from its strong generalization ability (Hupkes et al., 2023).

We conduct rich experiments to verify above arguments.

### 4.1 Synthetic Data

We conduct the experiments on synthetic data generated via the controllable pretrained distribution. As mentioned, the distribution has three components:

1. Mode prior  $p(\alpha)$ . We set the mode prior to be uniform:  $p(\alpha) = 1/|A|$ .

2. Mode-outline distribution  $p(\beta_{1:N}|\alpha)$ . For continuous modes  $\alpha \in A/\hat{\alpha}$ , Since we don't exactly care the outline, we set  $p(\beta_{1:N}|\alpha) = \prod_{n=1}^N p(\beta_n|\alpha)$  for convenience, where  $p(\beta_n|\alpha)$  is

a categorical distribution and its parameter is sampled from a Dirichlet distribution. The Dirichlet parameters are 0 for unseen topics (so that  $p(\beta) = 0$  for  $\beta \in U$ ) and 5 for others. We set the repetition prior to be uniform:  $p(\beta|\hat{\alpha}) = 1/|R| = 1/|B|r_R$  ( $\beta \in R$ ).

3. Topic-paragraph distribution  $p(x|\beta)$ . In order to simulate the compositionality, each topic  $\beta \in B$  is a tuple containing  $M$  subtopics  $\rho^{1:M}$ , where  $\rho^m \in B_*(m \in [M])$  and  $B = B_*^M$ . Accordingly, the paragraph  $x$  also contains  $M$  sub-paragraphs  $s^{1:M}$ , where each sub-paragraph is generated individually:

$$p(x|\beta) = \prod_{m=1}^M p(s^m|\rho^m) \quad (13)$$

The composition arity  $M$  controls the data compositionality. Given a fix number of topics  $|B|$ , the number of subtopics  $|B_*| = \sqrt[M]{|B|}$  decreases when composition arity  $M$  increases, and different topics are more likely to share structures as a result. Here, each sub-paragraph distribution  $p(s^m|\rho^m)$  is a Markov model whose initial probability vector  $\pi_{\rho^m}$  and transition matrix  $\mathbf{A}_{\rho^m}$  are both sampled from  $\text{Dir}(\gamma\mathbf{1})$ , where  $\mathbf{1}$  is an one vector.  $\gamma$  actually controls the uncertainty of different topics, where a lower value is expected to raise the KL divergence between different topic-paragraph models, making them easier to be distinguished, as shown in Appendix D.

#### 4.1.1 Data Parameter Settings

We set the number of modes  $|A| = 32$ , the number of topics  $|B| = 531441^7$ , where 95% of topics are unseen ( $|U| = 504868$ ). We set the vocab size  $|\Sigma| = 324$ , the length of sub-paragraph  $|s^m| = 3$ , and the number of paragraphs in a document  $N = 30$ . Thus, each document contains  $30(3M + 1)$  tokens. For other parameters of pretrained distribution including composition arity  $M$ , the ratio of repeated topics  $r_R$ , and topic uncertainty  $\gamma$ , we adjust their values to study the effects of data properties. In specific, we experiment with  $M \in \{2, 3, 4\}$ ,  $r_R \in \{2^{-d} | d = \{6, 7, \dots, 13\}\}$ , and  $\gamma \in \{0.01, 0.02, \dots, 0.05\}$ .

For each configuration of the pretrained distribution, we generate 10M documents for training. Therefore, the number of tokens in the synthetic dataset ranges from 2.1B to 3.9B. Examples of the synthetic dataset are shown in Figure 6.

<sup>7</sup>Its square, cube and fourth root are all integers.

Models	$L$	$H$	$D$	# params
X <sup>2</sup> S	3	6	384	4M
XS	4	8	448	8M
S	5	8	448	9M
M	6	8	512	15M
L	9	12	768	48M
XL	12	16	1024	114M
X <sup>2</sup> L	16	20	1280	234M

Table 1: Configurations of different models, where  $L$  is the number of layers,  $H$  is the number of attention heads,  $D$  is the hidden dimension. For parameter efficiency, we use grouped query attention (Ainslie et al., 2023) and set the number of key-value heads to be  $H/2$ .

## 4.2 Models

We study the effect of model size, attention window size, and the number of attention heads of Transformer. Table 1 shows configurations of different experimental models, where the parameters scales from 4M to 237M. The models are based on the Transformers (Wolf et al., 2020) implementation of Mistral (Jiang et al., 2023a). We train each model for 1 epoch on one NVIDIA A100 (40GB).

## 4.3 Evaluation Metrics

We aim to evaluate the overall ICG performance and the ICG-generalization ability of models using ICTR. Firstly, we define topic-wise ICTR as the expectation of prompt-wise ICTR:

$$\pi_N^\beta = \mathbb{E}_{p(x_{1:N}|\beta^N)} [p_{\text{LM}}(\beta|x_{1:N})] \quad (14)$$

Then, we can obtain the average ICTR of different kinds of topics:

$$\begin{aligned} \text{ICTR}_N^B &= \frac{1}{|B|} \sum_{\beta \in B} \pi_N^\beta, & \text{ICTR}_N^R &= \frac{1}{|R|} \sum_{\beta \in R} \pi_N^\beta \\ \text{ICTR}_N^C &= \frac{1}{|N|} \sum_{\beta \in C} \pi_N^\beta, & \text{ICTR}_N^U &= \frac{1}{|U|} \sum_{\beta \in U} \pi_N^\beta \end{aligned} \quad (15)$$

Here,  $\text{ICTR}_N^B$  measures the overall ICG ability, while  $\text{ICTR}_N^C$  and  $\text{ICTR}_N^U$  reflect the ICG-generalization ability, where higher values suggest better generalizations. In the experiments, since each pretrained document has 30 paragraphs, the trained model at most supports 29-shot in-context prompts. So by default, we reported  $\text{ICTR}_{29}^{B/R/C/U}$ , which is short of  $\text{ICTR}^{B/R/C/U}$ .

According to the values of the above ICTRs, we further define the following four statuses of a trained model by thresholding:

1. *Underfit*:  $\text{ICTR}^R < 0.65$ .

2. *Overfit*:  $\text{ICTR}^R \geq 0.65$ ,  $\text{ICTR}^C < 0.65$ , and  $\text{ICTR}^U < 0.65$ .

3. *C-Generalization*:  $\text{ICTR}^R \geq 0.65$ ,  $\text{ICTR}^C \geq 0.65$ , and  $\text{ICTR}^U < 0.65$ .

4. *U-Generalization*:  $\text{ICTR}^R \geq 0.65$ ,  $\text{ICTR}^C \geq 0.65$ , and  $\text{ICTR}^U \geq 0.65$ .

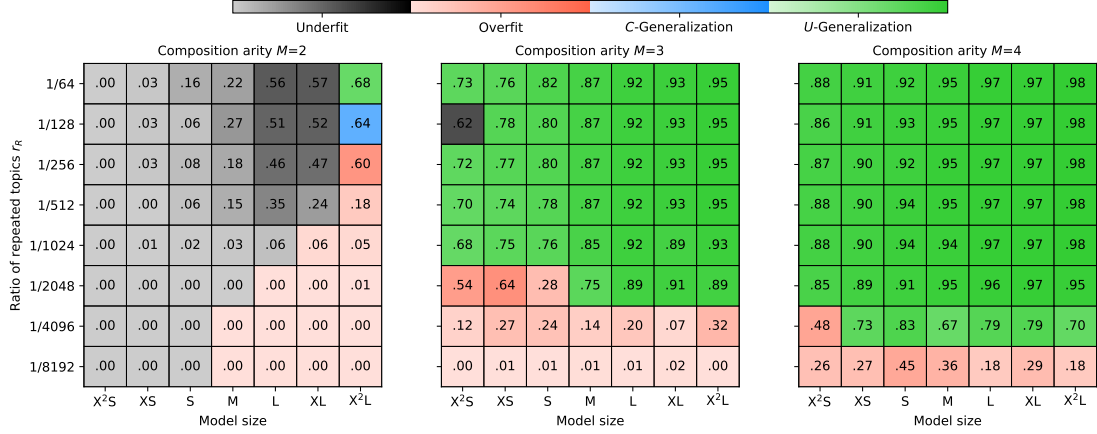
The computation of prompt and topic-wise ICTR is nontrivial, so we present it in Appendix F.

## 4.4 Results & Discussions

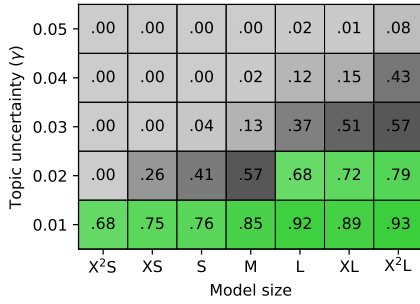
Our experimental results suggest the following arguments.

**Data compositionality enables both ICG and ICG-generalization.** Figure 3a shows the results of different composition arities. Clearly, we can see that data compositionality enables ICG and ICG-generalization, specifically: 1) As the composition arity  $M$  increases, the overall ICG performance consistently improves for models in any sizes trained on the pretrained distribution with different repeated topic proportions  $r_R$ . Notably, the improvement is especially significant when we increase  $M$  from 2 to 3. For example, for all  $r_R$ , the  $\text{ICTR}_{29}^B$  value nears 0 for many small models when  $M = 2$ , but is lifted to a considerable level when  $M = 3$ . 2) The models are easier to generalize on ICG when  $M$  is higher. When  $M = 2$ , most models are even hard to overfit on repeated topics, and only model X<sup>2</sup>L can generalize ICG to both non-repeated and unseen topics only when  $r_R = 1/64$ . On the contrary, when  $M = 3$  or  $M = 4$ , models in all sizes exhibit the ICG-generalization ability with much smaller  $r_R$ .

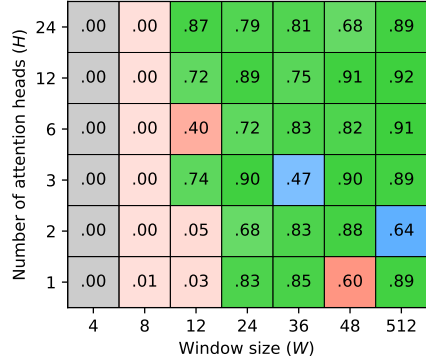
**The model emerges the ICG-generalization as the proportion of repeated topics rises.** As shown in Figure 3a, the model typically tends to overfit only on repeated topics when  $r_R$  is small, and then suddenly emerges the ICG-generalization ability when  $r_R$  hits the threshold. The threshold mainly corresponds to the data compositionality, where a higher composition arity  $M$  leads to a lower threshold and so makes the model easier to generalize. For example, for model X<sup>2</sup>L, the generalization threshold of  $r_R$  is  $1/64$  when  $M = 2$ , and decreases to  $1/2048$  when  $M = 3$ . We speculate this is because the more compositionality of the data, the more likely that nonrepeated and unseen



(a) ICG-generalization results of models in different sizes trained on pretrained distribution with different composition arities  $M$  and proportions of repeated topics  $r_R$ , where the topic uncertainty  $\gamma$  is set to 0.01.



(b) ICG-generalization results of models in different sizes trained on pretrained distribution with different topic uncertainties  $\gamma$ , where we set  $M = 3$  and  $r_R = 1/1024$ .



(c) ICG-generalization results of model L with different window sizes and numbers of attention heads, where we set  $M = 3$ ,  $r_R = 1/1024$ , and  $\gamma = 0.01$ .

Figure 3: ICG-generalization results, where the color suggests the status of the corresponding model, and the number in the cell shows the corresponding  $\text{ICTR}_{29}^B$ .

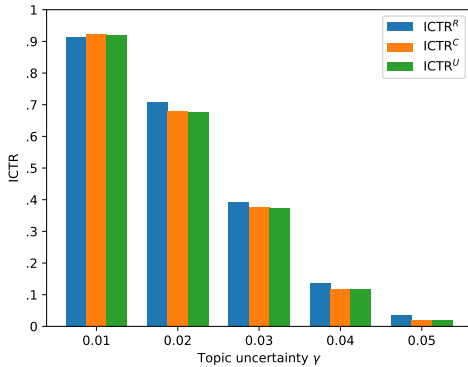


Figure 4:  $\text{ICTR}_{29}^*$  of different topics for model L trained on the pretrained distribution with different topic uncertainty  $\gamma$ , where the other parameters in the pretrained distribution are:  $M = 3$ ,  $r_R = 1/1024$ .

463 topics share sub-topics with repeated ones, there-  
464 fore the less proportion of repeated topics is needed  
465 for generalization.

466 **Topic uncertainty doesn't affect ICG-general-**  
467 **ization.** As shown in Figure 4, Topic uncertainty  
468

468 mainly affects the fitting difficulty of the data rather  
469 than the ICG-generalization ability: As the topic  
470 uncertainty  $\gamma$  increases, the  $\text{ICTR}_{29}$  of model L for  
471 all kinds of topics decreases consistently. However,  
472 we don't observe apparent ICG performance gaps  
473 between those topics.

474 **Larger models do better on ICG and ICG-gen-**  
475 **eralization.** Model size is considered to be a great  
476 factor impacting the ability of language models  
477 (Wei et al., 2022a). This is also verified in our ex-  
478 periments, which we find: 1) As shown in Figure  
479 3a, obviously, larger models not only have better  
480  $\text{ICTR}_{29}^B$ , but also require less repeated topics to  
481 generalize to nonrepeated and unseen topics. 2) As  
482 shown in Figure 3b, larger models are able to deal  
483 with topics with more uncertainties, i.e., bigger  $\gamma$ ,  
484 where models larger than model M are capable of  
485 ICG-generalization when  $\gamma = 0.02$  but smaller  
486 models pose underfit. Especially for model  $X^2S$ ,  
487 whose  $\text{ICTR}_{29}^B$  is 0. 3) As shown in Figure 5a, in  
488 most cases, larger models achieve better  $\text{ICTR}^B$

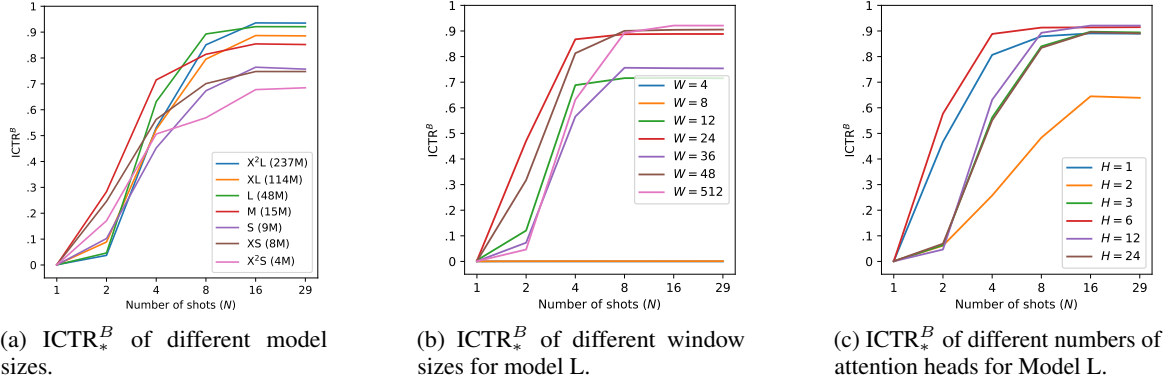


Figure 5:  $\text{ICTR}_*^B$  of different model configurations, where we set  $M = 3$ ,  $\gamma = 0.01$ , and  $r_R = 1/1024$ .

489 given fewer demonstrations. However, curiously, 526  
 490 this does not hold when the number of shots  $N$  is 527  
 491 too small. For example,  $\text{ICTR}_2^B$  of model  $X^2S$ ,  $XS$ , 528  
 492  $S$ , and  $M$  are typically greater than that of model  $L$ , 529  
 493  $XL$ , and  $X^2L$ . We speculate this might be because 530  
 494 when  $N$  is small, larger models are more cautious 531  
 495 in identifying the repetition mode. 532

496 **Big window size is necessary for ICG and ICG-** 533  
 497 **generalization.** Recently, Wang et al. (2023a) 534  
 498 show that LLMs conduct ICL by collecting informa- 535  
 499 tion of demonstrations in the prompt from pre- 536  
 500 vious label words. Specifically, the hidden states 537  
 501 of previous label words are good summarizations 538  
 502 of corresponding demonstrations. Thus, the model 539  
 503 needs to attend to all those previous "anchors" to 540  
 504 conduct ICL, which hints that a small window size 541  
 505 might harm the ICL performance. For example, in 542  
 506 the experimental results of Jiang et al. (2023b), we 543  
 507 can find that the ICL performance of RWKV (Peng 544  
 508 et al., 2023) series is generally inferior to that of 545  
 509 Transformer structures. Our experiments also sup- 546  
 510 port this argument. As shown in Figure 3c and 5b, 547  
 511 when the number of attention heads is fixed, a low 548  
 512 window size would cause underfit. In most cases, 549  
 513 as we increase the window size, the model is shifted 550  
 514 to overfit and finally U-Generalization, the overall 551  
 515  $\text{ICTR}_{29}^B$  also rises at the same time. Note that there 552  
 516 also exists the emergent phenomenon, where the 553  
 517 model suddenly learns ICG and ICG-generalization 554  
 518 when its window size hits a threshold. 555

519 **Big number of heads is not necessary for ICG** 556  
 520 **and ICG-generalization.** Multi-head/group atten- 557  
 521 tion is always believed to be the core driving state- 558  
 522 of-the-art Transformer models. By intuition, dif- 559  
 523 ferent heads can potentially attend onto different 560  
 524 parts of the text, making the model more expressive.  
 525 However, our experiments show this mechanism is

not very important for ICG and ICG-generalization. As shown in Figure 3c, reducing the number of attention heads  $H$  for XL model hardly change the model status. Also, as shown in Figure 3c, at the same size (L), the model with the highest overall ICG performance does not necessarily have the most attention heads. We speculate that this is because the attention pattern for ICG is relatively simple, so different heads are actually functional equivalent. This is consistent with Michel et al. (2019), which finds that the performance of many tasks including machine translation and natural language inference is insensitive to the number of attention heads.

**Generalizations towards nonrepeated and unseen topics are almost the same.** As shown in Figure 3, in most cases, no matter how pretrained distributions and models are configured, the models generally result as either underfit, overfit, or U-Generalization, but hardly in the status of C-Generalization. This suggests that nonrepeated topics, though appear in the pretrained distribution, are not easier for models to generalize.

## 5 Conclusions

This paper provides a systematic study of ICG ability of language models. Firstly, we propose a plausible latent variable pretrained distribution, formalizing ICG as a problem of next topic prediction. Then, we prove that the repetition nature of a few topics ensures the ICG ability on them theoretically. We also conduct rich experiments to study the effects of different factors of data and model architectures on ICG and ICG-generalization. We believe this paper is beneficial to a better understanding of the ICG ability, as well as large language models.



## 561 **Limitations**

562 The major limitation of this work is that we  
563 don't provide a theoretical support for ICG-  
564 generalization, while doing so is non-trivial. Now  
565 we can only speculate the ICG-generalization re-  
566 sults from the smoothing effects of neural probabil-  
567 ity approximator (e.g. Transformer), where unseen  
568 inputs would have non-zero probabilities (Xie et al.,  
569 2017). Therefore, nonrepeated and unseen topics  
570 might have a non-zero repetition prior, thus mak-  
571 ing them possible to be chosen as the topic of the  
572 next paragraph. This phenomenon might be es-  
573 pecially obvious when these topics are similar to  
574 repeated ones according to our experimental results.  
575 Further work on the theoretical understanding of  
576 ICG-generalization might take similarities between  
577 topics into account.

## 578 **References**

579 Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury  
580 Zemlyanskiy, Federico Lebrón, and Sumit Sanghai.  
581 2023. Gqa: Training generalized multi-query trans-  
582 former models from multi-head checkpoints. *arXiv*  
583 *preprint arXiv:2305.13245*.

584 Catherine Anderson. 2018. *Essentials of linguistics*.  
585 McMaster University.

586 Kazuoki Azuma. 1967. Weighted sums of certain de-  
587 pendent random variables. *Tohoku Mathematical*  
588 *Journal, Second Series*, 19(3):357–367.

589 Christopher M Bishop and Nasser M Nasrabadi. 2006.  
590 *Pattern recognition and machine learning*, volume 4.  
591 Springer.

592 Samuel R Bowman, Christopher D Manning, and  
593 Christopher Potts. 2015. Tree-structured composi-  
594 tion in neural networks without tree-structured archi-  
595 tectures. *arXiv preprint arXiv:1506.04834*.

596 Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
597 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
598 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
599 Askell, et al. 2020. Language models are few-shot  
600 learners. *Advances in neural information processing*  
601 *systems*, 33:1877–1901.

602 Richard E Grandy. 1990. Understanding and the princi-  
603 ple of compositionality. *Philosophical Perspectives*,  
604 4:557–572.

605 Michael Hahn and Navin Goyal. 2023. A theory of  
606 emergent in-context learning as implicit structure  
607 induction. *arXiv preprint arXiv:2303.07971*.

608 Wassily Hoeffding. 1994. Probability inequalities for  
609 sums of bounded random variables. *The collected*  
610 *works of Wassily Hoeffding*, pages 409–426.

Dieuwke Hupkes, Mario Giulianelli, Verna Dankers,  
Mikel Artetxe, Yanai Elazar, Tiago Pimentel, Chris-  
tos Christodoulopoulos, Karim Lasri, Naomi Saphra,  
Arabella Sinclair, et al. 2023. A taxonomy and review  
of generalization research in nlp. *Nature Machine*  
*Intelligence*, 5(10):1161–1174. 611 612 613 614 615 616

Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-  
sch, Chris Bamford, Devendra Singh Chaplot, Diego  
de las Casas, Florian Bressand, Gianna Lengyel, Guil-  
laume Lample, Lucile Saulnier, et al. 2023a. Mistral  
7b. *arXiv preprint arXiv:2310.06825*. 617 618 619 620 621

Zhongtao Jiang, Yuanzhe Zhang, Cao Liu, Jun Zhao,  
and Kang Liu. 2023b. Generative calibration for in-  
context learning. *arXiv preprint arXiv:2310.10266*. 622 623 624

Jaap Jumelet and Willem Zuidema. 2023. Transparency  
at the source: Evaluating and interpreting language  
models with access to the true distribution. *arXiv*  
*preprint arXiv:2310.14840*. 625 626 627 628

R Thomas McCoy, Robert Frank, and Tal Linzen. 2018.  
Revisiting the poverty of the stimulus: Hierarchical  
generalization without a hierarchical bias in recurrent  
neural networks. *arXiv preprint arXiv:1802.09091*. 629 630 631 632

Elliot Meyerson, Mark J Nelson, Herbie Bradley, Arash  
Moradi, Amy K Hoover, and Joel Lehman. 2023.  
Language model crossover: Variation through few-  
shot prompting. *arXiv preprint arXiv:2302.12170*. 633 634 635 636

Paul Michel, Omer Levy, and Graham Neubig. 2019.  
Are sixteen heads really better than one? *Advances*  
*in neural information processing systems*, 32. 637 638 639

Isabel Papadimitriou and Dan Jurafsky. 2023. Inject-  
ing structural hints: Using language models to study  
inductive biases in language learning. In *Findings*  
*of the Association for Computational Linguistics:*  
*EMNLP 2023*, pages 8402–8413. 640 641 642 643 644

Adam Paszke, Sam Gross, Francisco Massa, Adam  
Lerer, James Bradbury, Gregory Chanan, Trevor  
Killeen, Zeming Lin, Natalia Gimelshein, Luca  
Antiga, et al. 2019. Pytorch: An imperative style,  
high-performance deep learning library. *Advances in*  
*neural information processing systems*, 32. 645 646 647 648 649 650

Bo Peng, Eric Alcaide, Quentin Anthony, Alon Al-  
balak, Samuel Arcadinho, Huanqi Cao, Xin Cheng,  
Michael Chung, Matteo Grella, Kranthi Kiran GV,  
et al. 2023. Rwkv: Reinventing rnns for the trans-  
former era. *arXiv preprint arXiv:2305.13048*. 651 652 653 654 655

Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram,  
Inbal Magar, Omri Abend, Ehud Karpas, Amnon  
Shashua, Kevin Leyton-Brown, and Yoav Shoham.  
2023. Parallel context windows for large language  
models. In *Proceedings of the 61st Annual Meet-*  
*ing of the Association for Computational Linguistics*  
*(Volume 1: Long Papers)*, pages 6383–6402. 656 657 658 659 660 661 662

Frieda Rong. 2021. [Extrapolating to unnatural lan-  
guage processing with gpt-3's in-context learning:  
The good, the bad, and the mysterious.](#) 663 664 665

666 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-  
667 bert, Amjad Almahairi, Yasmine Babaei, Nikolay  
668 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti  
669 Bhosale, et al. 2023. Llama 2: Open founda-  
670 tion and fine-tuned chat models. *arXiv preprint*  
671 *arXiv:2307.09288*.

672 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob  
673 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
674 Kaiser, and Illia Polosukhin. 2017. Attention is all  
675 you need. *Advances in neural information processing*  
676 *systems*, 30.

677 Elena Voita, David Talbot, Fedor Moiseev, Rico Sen-  
678 nrich, and Ivan Titov. 2019. Analyzing multi-head  
679 self-attention: Specialized heads do the heavy lifting,  
680 the rest can be pruned. In *Proceedings of the 57th*  
681 *Annual Meeting of the Association for Computational*  
682 *Linguistics*, pages 5797–5808.

683 Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou,  
684 Fandong Meng, Jie Zhou, and Xu Sun. 2023a. Label  
685 words are anchors: An information flow perspective  
686 for understanding in-context learning. *arXiv preprint*  
687 *arXiv:2305.14160*.

688 Xinyi Wang, Wanrong Zhu, and William Yang Wang.  
689 2023b. Large language models are implicitly  
690 topic models: Explaining and finding good demon-  
691 strations for in-context learning. *arXiv preprint*  
692 *arXiv:2301.11916*.

693 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,  
694 Barret Zoph, Sebastian Borgeaud, Dani Yogatama,  
695 Maarten Bosma, Denny Zhou, Donald Metzler, et al.  
696 2022a. Emergent abilities of large language models.  
697 *arXiv preprint arXiv:2206.07682*.

698 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten  
699 Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,  
700 et al. 2022b. Chain-of-thought prompting elicits rea-  
701 soning in large language models. *Advances in neural*  
702 *information processing systems*, 35:24824–24837.

703 Jennifer C White and Ryan Cotterell. 2021. Examining  
704 the inductive bias of neural language models with ar-  
705 tificial languages. *arXiv preprint arXiv:2106.01044*.

706 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
707 Chaumond, Clement Delangue, Anthony Moi, Pier-  
708 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,  
709 et al. 2020. Transformers: State-of-the-art natural  
710 language processing. In *Proceedings of the 2020 con-*  
711 *ference on empirical methods in natural language*  
712 *processing: system demonstrations*, pages 38–45.

713 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and  
714 Tengyu Ma. 2021. An explanation of in-context learn-  
715 ing as implicit bayesian inference. *arXiv preprint*  
716 *arXiv:2111.02080*.

717 Ziang Xie, Sida I Wang, Jiwei Li, Daniel Lévy, Aiming  
718 Nie, Dan Jurafsky, and Andrew Y Ng. 2017. Data  
719 noising as smoothing in neural network language  
720 models. *arXiv preprint arXiv:1703.02573*.

## A Lemmas

To access the theoretical results in Appendix B, the following lemmas are useful.

**Lemma 1.** For an arbitrary continuous mode  $\alpha \in A/\hat{\alpha}$ , let

$$s_n = \sum_{i=1}^n \log \frac{p(x_i|x_{1:i-1}, \alpha)}{p(x_i|x_{1:i-1}, \hat{\alpha})} + \text{KL}_{i-1}(\hat{\alpha}||\alpha) \quad (16)$$

where

$$\text{KL}_{i-1}(\hat{\alpha}||\alpha) = \mathbb{E}_{p(x|x_{1:i-1}, \hat{\alpha})} \left[ \log \frac{p(x|x_{1:i-1}, \hat{\alpha})}{p(x|x_{1:i-1}, \alpha)} \right] \quad (17)$$

Then,  $s_n$  is a martingale about  $x_{1:n}$ .

*Proof.* This lemma is easy to prove according to the definition of martingale so we omit it.  $\square$

**Lemma 2.** Let  $z_n$  ( $n \in [N]$ ) be a series of positive random variables,  $\forall t \geq 0$ ,

$$\text{P} \left( \sum_{n=1}^N z_n \geq t \right) \leq \sum_{n=1}^N \text{P} \left( z_n \geq \frac{t}{N} \right) \quad (18)$$

*Proof.* Firstly, we have:

$$\begin{aligned} \text{P} \left( \sum_{n=1}^N z_n \geq t \right) &= \text{P} \left( \sum_{n=1}^N z_n \geq t, z_N \geq \frac{t}{N} \right) \\ &+ \text{P} \left( \sum_{n=1}^{N-1} z_n \geq \frac{N-1}{N}t, z_N \geq \frac{t}{N} \right) \\ &+ \text{P} \left( \sum_{n=1}^N z_n \geq t, \sum_{n=1}^{N-1} z_n \geq \frac{N-1}{N}t \right) \\ &\leq \text{P} \left( \sum_{n=1}^{N-1} z_n \leq \frac{N-1}{N}t, z_N \geq \frac{t}{N} \right) \\ &+ \text{P} \left( \sum_{n=1}^{N-1} z_n \geq \frac{N-1}{N}t, z_N \leq \frac{t}{N} \right) \\ &+ 2\text{P} \left( \sum_{n=1}^{N-1} z_n \geq \frac{N-1}{N}t, z_N \geq \frac{t}{N} \right) \\ &= \text{P} \left( \sum_{n=1}^{N-1} z_n \geq \frac{N-1}{N}t \right) + \text{P} \left( z_N \geq \frac{t}{N} \right) \end{aligned} \quad (19)$$

Then, according to this recursion,

$$\begin{aligned}
& \mathbb{P}\left(\sum_{n=1}^N z_n \geq t\right) \\
& \leq \mathbb{P}\left(\sum_{n=1}^{N-1} z_n \geq \frac{N-1}{N}t\right) + \mathbb{P}\left(z_N \geq \frac{t}{N}\right) \\
& \leq \mathbb{P}\left(\sum_{n=1}^{N-2} z_n \geq \frac{N-2}{N}t\right) + \mathbb{P}\left(z_{N-1} \geq \frac{t}{N}\right) \\
& \quad + \mathbb{P}\left(z_N \geq \frac{t}{N}\right) \\
& \dots \\
& \leq \sum_{n=1}^N \mathbb{P}\left(z_n \geq \frac{t}{N}\right)
\end{aligned} \tag{20}$$

So the result is proved.  $\square$

## B Complete Theoretical Results

We analyze the data ICG distribution  $p(x|x_{1:N})$ , where  $x_{1:N}$  are independent and identical distributed with PDF  $p(x|\hat{\beta})$  and  $x$  is an arbitrary value in the domain of paragraph. As shown in Section 2.1,  $x$  depends on its topic:

$$p(x|x_{1:N}) = \sum_{\beta \in B} p(\beta|x_{1:N})p(x|\beta) \tag{21}$$

where the topic predictive distribution  $p(\beta|x_{1:N}) := p(\beta_{1:N} = \beta|x_{1:N})$  controls the strength of each topic for the  $N + 1$ -th paragraph. We then study the property of this distribution.

Note that the topic predictive distribution can also analogously be factorized as the mixture of modes:

$$p(\beta|x_{1:N}) = \sum_{\alpha \in A} p(\alpha|x_{1:N})p(\beta|x_{1:N}, \alpha) \tag{22}$$

where the mode posterior  $p(\alpha|x_{1:N})$  controls the strength of each mode.

### B.1 Property of mode posterior

Firstly, we study the property of the mode posterior  $p(\alpha|x_{1:N})$ .

**Proposition 1.** *Let:*

$$p_{\max}(\hat{\alpha}) = \max_{\alpha \in A/\hat{\alpha}} p(\alpha) \tag{23}$$

If  $t$  satisfies:

$$\frac{|A|p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}} \leq t < 1 \tag{24}$$

and  $\hat{\beta} \in R$ , for repetition mode  $\hat{\alpha}$ , we have:

$$\begin{aligned}
& \mathbb{P}(1 - p(\hat{\alpha}|x_{1:N}) \geq t) \\
& \leq |A|e^{-\frac{(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(1-t)p_{\max}(\hat{\alpha})})^2}{8N \log^2(c_4/c_3)}}
\end{aligned} \tag{25}$$

For any continuous mode  $\alpha \in A/\hat{\alpha}$ , we also have:

$$\begin{aligned}
& \mathbb{P}(p(\alpha|x_{1:N}) \geq t) \\
& \leq |A|e^{-\frac{(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(1-t)p_{\max}(\hat{\alpha})})^2}{8N \log^2(c_4/c_3)}}
\end{aligned} \tag{26}$$

*Proof.* Firstly, note that the absolute martingale residual difference of  $s_n$  in formula (17) is bounded:

$$\begin{aligned}
& |s_n - s_{n-1}| \\
& = \left| \log \frac{p(x_n|x_{1:n-1}, \alpha)}{p(x_n|x_{1:n-1}, \hat{\alpha})} + \text{KL}_{n-1}(\hat{\alpha}||\alpha) \right| \\
& \leq \left| \log \frac{p(x_n|x_{1:n-1}, \alpha)}{p(x_n|x_{1:n-1}, \hat{\alpha})} \right| + |\text{KL}_{n-1}(\hat{\alpha}||\alpha)| \\
& \leq 2 \log \frac{c_4}{c_3}
\end{aligned} \tag{27}$$

Then, according to Azuma's inequity (Azuma, 1967),  $\forall \epsilon > 0$ , we have:

$$\begin{aligned}
& \mathbb{P}\left(\sum_{n=1}^N \log \frac{p(x_n|x_{1:n-1}, \alpha)}{p(x_n|x_{1:n-1}, \hat{\alpha})} + \text{KL}_{n-1}(\hat{\alpha}||\alpha) \geq \epsilon\right) \\
& \leq e^{-\frac{\epsilon^2}{8N \log^2(c_4/c_3)}}
\end{aligned} \tag{28}$$

Since  $\text{KL}_{i-1}(\hat{\alpha}||\alpha) \geq \log c_1$ , we can rewrite formula (28) as:

$$\begin{aligned}
& \mathbb{P}\left(\sum_{i=1}^N \log \frac{p(x_n|x_{1:n-1}, \alpha)}{p(x_n|x_{1:n-1}, \hat{\alpha})} \geq \epsilon - N \log c_1\right) \\
& \leq e^{-\frac{\epsilon^2}{8N \log^2(c_4/c_3)}}
\end{aligned} \tag{29}$$

Let  $t = e^{\epsilon - N \log c_1} \in [c_1^{-N}, 1)$  and rearrange the formula, we can obtain the following inequality about the ratio of mode likelihoods:

$$\mathbb{P}\left(\frac{p(x_{1:N}|\alpha)}{p(x_{1:N}|\hat{\alpha})} \geq t\right) \leq e^{-\frac{(N \log c_1 + \log t)^2}{8N \log^2(c_4/c_3)}} \tag{30}$$

The ratio of mode likelihoods has a direct impact to the mode posterior. First, for repetition mode  $\hat{\alpha}$ ,

785  $\forall 0 < t < 1$ , we have:

$$\begin{aligned}
\text{P}(1 - p(\hat{\alpha}|x_{1:N}) \geq t) &= \text{P}\left(\frac{1}{p(\hat{\alpha}|x_{1:N})} \geq \frac{1}{1-t}\right) \\
&= \text{P}\left(\sum_{\alpha \in A/\hat{\alpha}} \frac{p(\alpha)p(x_{1:N}|\alpha)}{p(\hat{\alpha})p(x_{1:N}|\hat{\alpha})} \geq \frac{t}{1-t}\right) \\
&\leq \sum_{\alpha \in A/\hat{\alpha}} \text{P}\left(\frac{p(x_{1:N}|\alpha)}{p(x_{1:N}|\hat{\alpha})} \geq \frac{tp(\hat{\alpha})}{(|A|-1)(1-t)p(\alpha)}\right) \\
&\leq \sum_{\alpha \in A/\hat{\alpha}} \text{P}\left(\frac{p(x_{1:N}|\alpha)}{p(x_{1:N}|\hat{\alpha})} \geq \frac{tp(\hat{\alpha})}{|A|(1-t)p_{\max}(\hat{\alpha})}\right)
\end{aligned} \tag{31}$$

786 where we unpack the probability in the third line  
787 using lemma 2. Now, if  
788

$$\begin{aligned}
\frac{tp(\hat{\alpha})}{|A|(1-t)p_{\max}(\hat{\alpha})} &\geq c_1^{-N} \\
\Rightarrow t &\geq \frac{|A|p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}}
\end{aligned} \tag{32}$$

789 then we can apply formula (30):

$$\begin{aligned}
\text{P}(1 - p(\hat{\alpha}|x_{1:N}) \geq t) \\
\leq |A|e^{-\frac{(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(1-t)p_{\max}(\hat{\alpha})})^2}{8N \log^2(c_4/c_3)}}
\end{aligned} \tag{33}$$

792 As for continuous modes  $\alpha \in A/\hat{\alpha}$ , note that:

$$\begin{aligned}
\text{P}(p(\alpha|x_{1:N}) \geq t) &\leq \text{P}\left(\sum_{\alpha \in A/\hat{\alpha}} p(\alpha|x_{1:N}) \geq t\right) \\
&= \text{P}(1 - p(\hat{\alpha}|x_{1:N}) \geq t) \\
&\leq |A|e^{-\frac{(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(1-t)p_{\max}(\hat{\alpha})})^2}{8N \log^2(c_4/c_3)}}
\end{aligned} \tag{34}$$

793  
794  $\square$   
795 Based on proposition 1, we can immediately  
796 obtain the following two corollaries:

797 **Corollary 1.** If  $\hat{\beta} \in R$ ,  $\text{plim}_{N \rightarrow \infty} p(\hat{\alpha}|x_{1:N}) = 1$

798 *Proof.* To prove the results, we need to prove that,  
799  $\forall \epsilon > 0, \delta > 0$ , there exists  $N_0$  such that when  
800  $N \geq N_0$ ,

$$\text{P}(1 - p(\hat{\alpha}|x_{1:N}) \geq \epsilon) < \delta \tag{35}$$

802 Firstly, note that when  $\epsilon > 1$  or  $\delta \geq 1$ , the above  
803 formula holds trivially. When  $0 < \epsilon \leq 1$ , define:

$$\hat{N}(\epsilon) = \log_{c_1} \frac{|A|(1-\epsilon)p_{\max}(\hat{\alpha})}{tp(\hat{\alpha})} \tag{36}$$

805 If  $N \geq \hat{N}(\epsilon)$ , then

$$\epsilon \geq \frac{|A|p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}} \tag{37}$$

806 Therefore, according to proposition 1, we have:

$$\text{P}(1 - p(\hat{\alpha}|x_{1:N}) \geq \epsilon) \leq f(N) \tag{38}$$

807 where

$$f(N) = |A|e^{-\frac{(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(1-\epsilon)p_{\max}(\hat{\alpha})})^2}{8N \log^2(c_4/c_3)}} \tag{39}$$

810 Since  $f(N) \in (0, |A|^2]$  is a monotonic decreasing  
811 function in the domain of  $[\hat{N}(\epsilon), \infty]$ ,  $\forall \delta \in (0, 1)$   
812 there must exists  $N' \geq \hat{N}(\epsilon)$  such that  $\delta = f(N')$ ,  
813 or equivalently,  $N' = f^{-1}(\delta)$ . Let's set  $N_0 =$   
814  $\lceil f^{-1}(\delta) \rceil + 1$ . If  $N \geq N_0$ ,  
815

$$\text{P}(1 - p(\hat{\alpha}|x_{1:N}) \geq \epsilon) \leq f(\lceil f^{-1}(\delta) \rceil + 1) < \delta \tag{40}$$

816 Therefore, the result is proven.  $\square$   
817

818 **Corollary 2.** If  $t$  satisfies:

$$\frac{|A|^{5/2}p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}} \leq t < 1 \tag{41}$$

819 and  $\hat{\beta} \in R$ , we have:

$$\begin{aligned}
\text{P}(|p(\beta|x_{1:N}) - p(\beta|x_{1:N}, \hat{\alpha})| \geq t) \\
\leq |A|^2 e^{-\frac{(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(|A|^{3/2} - t)p_{\max}(\hat{\alpha})})^2}{8N \log^2(c_4/c_3)}}
\end{aligned} \tag{42}$$

820 *Proof.* Let  $\mathbf{p}_N^\alpha \in \Delta^{|A|}$  be the topic posterior vec-  
821 tor:

$$\mathbf{p}_N^\alpha = \begin{bmatrix} \dots \\ p(\alpha|x_{1:N}) \\ \dots \end{bmatrix} \in \Delta^{|A|} \tag{43}$$

822 and  $\delta^{\hat{\alpha}}$  be the one-hot vector peaking at  $\hat{\alpha}$ .  $\forall 0 <$   
823  $t < 1$ , Obviously:  
824

$$\begin{aligned}
&\text{P}\left(\|\mathbf{p}_N^\alpha - \delta^{\hat{\alpha}}\|_2 \geq t\right) \\
&\leq \text{P}\left(\sum_{\alpha \in A/\hat{\alpha}} p(\alpha|x_{1:N}) + 1 - p(\hat{\alpha}|x_{1:N}) \geq t\right) \\
&\leq \sum_{\alpha \in A/\hat{\alpha}} \text{P}\left(p(\alpha|x_{1:N}) \geq \frac{t}{|A|}\right) \\
&\quad + \text{P}\left(1 - p(\hat{\alpha}|x_{1:N}) \geq \frac{t}{|A|}\right)
\end{aligned} \tag{44}$$

If

$$\begin{aligned} \frac{t}{|A|} &\geq \frac{|A|p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}} \\ \Rightarrow t &\geq \frac{|A|^2p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}} \end{aligned} \quad (45)$$

then we can apply formula (25) and (26) to get the following:

$$\begin{aligned} &\mathbb{P}\left(\|\mathbf{p}_N^\alpha - \delta^{\hat{\alpha}}\|_2 \geq t\right) \\ &\leq |A|^2 e^{-\frac{(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(|A|-t)p_{\max}(\hat{\alpha})})^2}{8N \log^2(c_4/c_3)}} \end{aligned} \quad (46)$$

Now, denote:

$$\mathbf{p}_{\cdot|N,\alpha}^\beta = \begin{bmatrix} \dots \\ p(\beta|x_{1:N}, \alpha) \\ \dots \end{bmatrix} \in [0, 1]^{|A|} \quad (47)$$

Then,  $\forall 0 < t < 1$ , we have:

$$\begin{aligned} &\mathbb{P}(|p(\beta|x_{1:N}) - p(\beta|x_{1:N}, \hat{\alpha})| \geq t) \\ &= \mathbb{P}\left(\left| \left(\mathbf{p}_N^\alpha - \delta^{\hat{\alpha}}\right)^T \mathbf{p}_{\cdot|N,\alpha}^\beta \right| \geq t\right) \\ &\leq \mathbb{P}\left(\left\|\mathbf{p}_N^\alpha - \delta^{\hat{\alpha}}\right\|_2 \left\|\mathbf{p}_{\cdot|N,\alpha}^\beta\right\|_2 \geq t\right) \\ &\leq \mathbb{P}\left(\left|\mathbf{p}_N^\alpha - \delta^{\hat{\alpha}}\right| \geq \frac{t}{\sqrt{|A|}}\right) \end{aligned} \quad (48)$$

If  $t \geq \frac{|A|^{5/2}p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}}$ , we can then apply formula (46) to obtain the result.  $\square$

## B.2 Property of topic posterior under repetition mode

Secondly, we study the property of the topic posterior under the repetition mode  $p(\beta|x_{1:N}, \hat{\alpha})$ .

**Proposition 2.** *Let*

$$p_{\max}(\hat{\beta}) = \max_{\beta \in B/\hat{\beta}} p(\beta|\hat{\alpha}) \quad (49)$$

If  $t$  satisfies:

$$\frac{|B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}}{p(\hat{\beta}|\hat{\alpha}) + |B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}} \leq t < 1 \quad (50)$$

Then, for the ground-truth topic  $\hat{\beta}$ , if  $\hat{\beta} \in R$ , we have:

$$\begin{aligned} &\mathbb{P}(1 - p(\hat{\beta}|x_{1:N}, \hat{\alpha}) \geq t) \leq \sum_{\beta \in B/\hat{\beta}} \\ &\leq |B|e^{-\frac{2(N \log c_2 + \log \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(1-t)p_{\max}(\hat{\beta}|\hat{\alpha})})^2}{N \log^2(c_4/c_3)}} \end{aligned} \quad (51)$$

For any other topic  $\beta \in R/\hat{\beta}$ , we also have:

$$\begin{aligned} &\mathbb{P}(p(\beta|x_{1:N}, \hat{\alpha}) \geq t) \\ &\leq |B|e^{-\frac{2(N \log c_2 + \log \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(1-t)p_{\max}(\hat{\beta}|\hat{\alpha})})^2}{N \log^2(c_4/c_3)}} \end{aligned} \quad (52)$$

*Proof.* For any topic  $\beta \in B/\hat{\beta}$ , let

$$s_n = \sum_{i=1}^n \log \frac{p(x_i|\beta)}{p(x_i|\hat{\beta})} \quad (53)$$

Since each demonstration  $x_n$  is independently sampled from  $p(x|\hat{\beta})$ , all the addends in the above formula are independent. Also, note that:

$$\begin{aligned} \mathbb{E}[s_n] &= \sum_{i=1}^n \mathbb{E}\left[\log \frac{p(x_i|\beta)}{p(x_i|\hat{\beta})}\right] = n\text{KL}(\hat{\beta}||\beta) \\ &\geq n \log c_2 \end{aligned}$$

$$\left|\log \frac{p(x_i|\beta)}{p(x_i|\hat{\beta})}\right| \leq \log \frac{c_4}{c_3} \quad (54)$$

Then, according to Hoeffding's inequity (Hoeffding, 1994),  $\forall \epsilon > 0$ ,

$$\begin{aligned} &\mathbb{P}\left(\sum_{i=1}^N \log \frac{p(x_i|\beta)}{p(x_i|\hat{\beta})} \geq \epsilon - N \log c_2\right) \\ &\leq \mathbb{P}\left(\sum_{i=1}^N \log \frac{p(x_i|\beta)}{p(x_i|\hat{\beta})} \geq \epsilon - N\text{KL}(\hat{\beta}||\beta)\right) \\ &= \mathbb{P}\left(\prod_{i=1}^N \frac{p(x_i|\beta)}{p(x_i|\hat{\beta})} \geq e^{\epsilon - N\text{KL}(\hat{\beta}||\beta)}\right) \\ &\leq e^{-\frac{2\epsilon^2}{N \log^2(c_4/c_3)}} \end{aligned} \quad (55)$$

Let  $t = e^{\epsilon - N \log c_2} \geq c_2^{-N}$ , we have:

$$\mathbb{P}\left(\prod_{n=1}^N \frac{p(x_n|\beta)}{p(x_n|\hat{\beta})} \geq t\right) \leq e^{-\frac{2(N \log c_2 + \log t)^2}{N \log^2(c_4/c_3)}} \quad (56)$$

The rest of proof of is very similar to that of proposition 1,  $\forall t \geq \frac{|B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}}{p(\hat{\beta}|\hat{\alpha}) + |B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}}$ ,

$$\begin{aligned} &\mathbb{P}(1 - p(\hat{\beta}|x_{1:N}, \hat{\alpha}) \geq t) \leq \sum_{\beta \in B/\hat{\beta}} \\ &\mathbb{P}\left(\prod_{n=1}^N \frac{p(x_n|\beta)}{p(x_n|\hat{\beta})} \geq \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(1-t)p_{\max}(\hat{\beta}|\hat{\alpha})}\right) \\ &\leq |B|e^{-\frac{2(N \log c_2 + \log \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(1-t)p_{\max}(\hat{\beta}|\hat{\alpha})})^2}{N \log^2(c_4/c_3)}} \end{aligned} \quad (57)$$

And  $\forall \beta \in R/\hat{\beta}$ ,

$$\begin{aligned} & \mathbb{P}(p(\beta|x_{1:N}, \hat{\alpha}) \geq t) \\ & \leq |B|e^{-\frac{2\left(N \log c_2 + \log \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(1-t)p_{\max}(\hat{\beta}|\hat{\alpha})}\right)^2}{N \log^2(c_4/c_3)}} \end{aligned} \quad (58)$$

□

Likewise, we can also obtain the following corollary:

**Corollary 3.** *If  $\hat{\beta} \in R$ ,  $\text{plim}_{N \rightarrow \infty} p(\hat{\beta}|x_{1:N}, \hat{\alpha}) = 1$ .*

*Proof.* The proof is identical to the proof of corollary 4 so we omit it. □

### B.3 Property of topic predictive distribution

Based on the above results, we are able to investigate the property of the topic predictive distribution  $p(\beta|x_{1:N})$ .

**Proposition 3.** *If  $t$  satisfies:*

$$1 > t \geq \max \left\{ \frac{2|A|^{5/2}p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha})+|A|p_{\max}(\hat{\alpha})c_1^{-N}}, \frac{2|B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}}{p(\hat{\beta}|\hat{\alpha})+|B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}} \right\} \quad (59)$$

*Then, for the ground-truth topic  $\hat{\beta}$ , if  $\hat{\beta} \in R$ , we have:*

$$\begin{aligned} & \mathbb{P}(1 - p(\hat{\beta}|x_{1:N}) \geq t) \\ & \leq |A|^2 e^{-\frac{\left(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(2|A|^{3/2}-t)p_{\max}(\hat{\alpha})}\right)^2}{8N \log^2(c_4/c_3)}} \\ & \quad + |B|e^{-\frac{2\left(N \log c_2 + \log \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(2-t)p_{\max}(\hat{\beta}|\hat{\alpha})}\right)^2}{N \log^2(c_4/c_3)}} \end{aligned} \quad (60)$$

*For other topics  $\beta \in B/\hat{\beta}$ , we also have:*

$$\begin{aligned} & \mathbb{P}(p(\beta|x_{1:N}) \geq t) \\ & \leq |A|^2 e^{-\frac{\left(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(2|A|^{3/2}-t)p_{\max}(\hat{\alpha})}\right)^2}{8N \log^2(c_4/c_3)}} \\ & \quad + |B|e^{-\frac{2\left(N \log c_2 + \log \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(2-t)p_{\max}(\hat{\beta}|\hat{\alpha})}\right)^2}{N \log^2(c_4/c_3)}} \end{aligned} \quad (61)$$

*Proof.* For the ground-truth topic  $\hat{\beta}$  and any  $0 <$

$t < 1$ , we have:

$$\begin{aligned} & \mathbb{P}(1 - p(\hat{\beta}|x_{1:N}) \geq t) \\ & = \mathbb{P}(p(\hat{\beta}|x_{1:N}, \hat{\alpha}) - p(\hat{\beta}|x_{1:N}) + \\ & \quad 1 - p(\hat{\beta}|x_{1:N}, \hat{\alpha}) \geq t) \\ & \leq \mathbb{P}(|p(\hat{\beta}|x_{1:N}, \hat{\alpha}) - p(\hat{\beta}|x_{1:N})| + \\ & \quad 1 - p(\hat{\beta}|x_{1:N}, \hat{\alpha}) \geq t) \\ & \leq \mathbb{P}\left(|p(\hat{\beta}|x_{1:N}, \hat{\alpha}) - p(\hat{\beta}|x_{1:N})| \geq \frac{t}{2}\right) \\ & \quad \mathbb{P}\left(1 - p(\hat{\beta}|x_{1:N}, \hat{\alpha}) \geq \frac{t}{2}\right) \end{aligned} \quad (62)$$

Therefore, if

$$1 > t \geq \max \left\{ \frac{2|A|^{5/2}p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha})+|A|p_{\max}(\hat{\alpha})c_1^{-N}}, \frac{2|B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}}{p(\hat{\beta}|\hat{\alpha})+|B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}} \right\} \quad (63)$$

we can then apply corollary 2 and proposition 2 to prove formula (60). Meanwhile, for other topics  $\beta \in B/\hat{\beta}$ , we have:

$$\begin{aligned} & \mathbb{P}(p(\beta|x_{1:N}) \geq t) \leq \mathbb{P}\left(\sum_{\beta \in B/\hat{\beta}} p(\beta|x_{1:N}) \geq t\right) \\ & = \mathbb{P}(1 - p(\hat{\beta}|x_{1:N}) \geq t) \end{aligned} \quad (64)$$

Then, if  $t$  satisfies formula (63), we can obtain formula (61). □

The property of the topic predictive distribution can be summarized more compactly via the following corollary:

**Corollary 4.** *If  $\hat{\beta} \in R$ ,  $\text{plim}_{N \rightarrow \infty} p(\hat{\beta}|x_{1:N}) = 1$ .*

*Proof.* The proof is identical to the proof of corollary 4 so we omit it. □

### B.4 Property of in-context generative distribution

According the property of the topic predictive distribution, we can finally study the property of the in-context generative distribution.

**Proposition 4.** *If  $t$  satisfies:*

$$1 > t \geq \max \left\{ \frac{2c_4|A|^{5/2}|B|^{3/2}p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha})+|A|p_{\max}(\hat{\alpha})c_1^{-N}}, \frac{2c_4|B|^{3/2}p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}}{p(\hat{\beta}|\hat{\alpha})+|B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}} \right\} \quad (65)$$

and  $\hat{\beta} \in R$ , for any candidate paragraph  $x \in \Sigma^*$ , we have:

$$\begin{aligned} & \mathbb{P}(|p(x|x_{1:N}) - p(x|\hat{\beta})| \geq t) \\ & \leq |A|^2 |B| e^{-\frac{\left(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(2|A|^{\frac{3}{2}}|B|^{\frac{3}{2}}c_4 - t)p_{\max}(\hat{\alpha})}\right)^2}{8N \log^2(c_4/c_3)}}} \\ & \quad + |B|^2 e^{-\frac{2\left(N \log c_2 + \log \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(2|B|^{\frac{3}{2}}c_4 - t)p_{\max}(\hat{\beta}|\hat{\alpha})}\right)^2}{N \log^2(c_4/c_3)}}} \end{aligned} \quad (66)$$

*Proof.* Let  $\mathbf{p}_N^\beta \in \Delta^{|B|}$  be the topic predictive vector:

$$\mathbf{p}_N^\beta = \begin{bmatrix} \dots \\ p(\beta|x_{1:N}) \\ \dots \end{bmatrix} \in \Delta^{|B|} \quad (67)$$

and  $\delta^{\hat{\beta}}$  be the one-hot vector peaking at  $\hat{\beta}$ . For all  $0 < t < 1$ , we have:

$$\begin{aligned} & \mathbb{P}\left(\|\mathbf{p}_N^\beta - \delta^{\hat{\beta}}\|_2 \geq t\right) \\ & \leq \mathbb{P}\left(\sum_{\beta \in B/\hat{\beta}} p(\beta|x_{1:N}) + 1 - p(\hat{\beta}|x_{1:N}) \geq t\right) \\ & \leq \sum_{\beta \in B/\hat{\beta}} \mathbb{P}\left(p(\beta|x_{1:N}) \geq \frac{t}{|B|}\right) \\ & \quad + \mathbb{P}\left(1 - p(\hat{\beta}|x_{1:N}) \geq \frac{t}{|B|}\right) \end{aligned} \quad (68)$$

If

$$\begin{aligned} \frac{t}{|B|} & \geq \max \left\{ \frac{2|A|^{5/2}p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}}, \frac{2|B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}}{p(\hat{\beta}|\hat{\alpha}) + |B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}} \right\} \\ \Rightarrow t & \geq \max \left\{ \frac{2|A|^{5/2}|B|p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}}, \frac{2|B|^2p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}}{p(\hat{\beta}|\hat{\alpha}) + |B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}} \right\} \end{aligned} \quad (69)$$

Then we can apply results from proposition 3 to get the following:

$$\begin{aligned} & \mathbb{P}\left(\|\mathbf{p}_N^\beta - \delta^{\hat{\beta}}\|_2 \geq t\right) \\ & \leq |A|^2 |B| e^{-\frac{\left(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(2|A|^{\frac{3}{2}}|B|^{\frac{3}{2}}|B| - t)p_{\max}(\hat{\alpha})}\right)^2}{8N \log^2(c_4/c_3)}}} \\ & \quad + |B|^2 e^{-\frac{2\left(N \log c_2 + \log \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(2|B|^{\frac{3}{2}}|B| - t)p_{\max}(\hat{\beta}|\hat{\alpha})}\right)^2}{N \log^2(c_4/c_3)}}} \end{aligned} \quad (70)$$

Now, denote:

$$\mathbf{p}_{\cdot|\beta}^x = \begin{bmatrix} \dots \\ p(x|\beta) \\ \dots \end{bmatrix} \in [c_3, c_4]^{|B|} \quad (71)$$

Therefore, For all  $0 < t < 1$ ,

$$\begin{aligned} & \mathbb{P}(|p(x|x_{1:N}) - p(x|\hat{\beta})| \geq t) \\ & = \mathbb{P}\left(\left| \left(\mathbf{p}_N^\beta - \delta^{\hat{\beta}}\right)^T \mathbf{p}_{\cdot|\beta}^x \right| \geq t\right) \\ & \leq \mathbb{P}\left(\left\|\mathbf{p}_N^\beta - \delta^{\hat{\beta}}\right\|_2 \left\|\mathbf{p}_{\cdot|\beta}^x\right\|_2 \geq t\right) \\ & \leq \mathbb{P}\left(\left\|\mathbf{p}_N^\beta - \delta^{\hat{\beta}}\right\|_2 \geq \frac{t}{\sqrt{|B|c_4}}\right) \end{aligned} \quad (72)$$

Therefore, if  $t$  satisfies formula (65), we can then apply formula (66) to prove the result.  $\square$

Proposition 4 directly supports the following corollary:

**Corollary 5.** If  $\hat{\beta} \in R$ ,  $\text{plim}_{N \rightarrow \infty} p(x|x_{1:N}) = p(x|\hat{\beta})$ .

*Proof.* The proof is identical to the proof of corollary 4 so we omit it.  $\square$

## B.5 Property of in-context predictive distribution

We can generalize the property of ICG distribution to the in-context predictive distribution as well, which forms the theoretical foundation of ICL.

**Proposition 5.** If  $t$  satisfies:

$$1 > t \geq \max \left\{ \frac{4c_3^2c_4^2|A|^{5/2}|B|^{3/2}p_{\max}(\hat{\alpha})c_1^{-N}}{p(\hat{\alpha}) + |A|p_{\max}(\hat{\alpha})c_1^{-N}}, \frac{4c_3^2c_4^2|B|^{3/2}p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}}{p(\hat{\beta}|\hat{\alpha}) + |B|p_{\max}(\hat{\beta}|\hat{\alpha})c_2^{-N}} \right\} \quad (73)$$

and  $\hat{\beta} \in R$ , we have

$$\begin{aligned} & \mathbb{P}\left(\left|p(y|(x, y)_{1:N}, x) - p(y|x, \hat{\beta})\right| \geq t\right) \\ & \leq |A|^2 |B| e^{-\frac{\left(N \log c_1 + \log \frac{tp(\hat{\alpha})}{|A|(4|A|^{\frac{3}{2}}|B|^{\frac{3}{2}}c_3^2c_4^2 - t)p_{\max}(\hat{\alpha})}\right)^2}{8N \log^2(c_4/c_3)}}} \\ & \quad + |B|^2 e^{-\frac{2\left(N \log c_2 + \log \frac{tp(\hat{\beta}|\hat{\alpha})}{|B|(4|B|^{\frac{3}{2}}c_3^2c_4^2 - t)p_{\max}(\hat{\beta}|\hat{\alpha})}\right)^2}{N \log^2(c_4/c_3)}}} \end{aligned} \quad (74)$$

*Proof.*  $\forall 0 < t < 1$ , we have

$$\begin{aligned}
& \mathbb{P} \left( \left| p(y|(x, y)_{1:N}, x) - p(y|x, \hat{\beta}) \right| \geq t \right) \\
&= \mathbb{P} \left( \left| \frac{p(x, y|(x, y)_{1:N})}{p(x|(x, y)_{1:N})} - \frac{p(x, y|\hat{\beta})}{p(x|\hat{\beta})} \right| \geq t \right) \\
&= \mathbb{P} \left( \left| \frac{p(x|\hat{\beta})p(x, y|(x, y)_{1:N})}{p(x|(x, y)_{1:N})p(x|\hat{\beta})} \right. \right. \\
&\quad \left. \left. - \frac{p(x, y|\hat{\beta})p(x|(x, y)_{1:N})}{p(x|\hat{\beta})p(x|(x, y)_{1:N})} \right| \geq t \right) \\
&\leq \mathbb{P} \left( \left| p(x|\hat{\beta})p(x, y|(x, y)_{1:N}) \right. \right. \\
&\quad \left. \left. - p(x, y|\hat{\beta})p(x|(x, y)_{1:N}) \right| \geq \frac{t}{c_3^2} \right) \\
&= \mathbb{P} \left( \left| p(x|\hat{\beta}) \left( p(x, y|(x, y)_{1:N}) - p(x, y|\hat{\beta}) \right) \right. \right. \\
&\quad \left. \left. + p(x, y|\hat{\beta}) \left( p(x|\hat{\beta}) - p(x|(x, y)_{1:N}) \right) \right| \geq \frac{t}{c_3^2} \right) \\
&\leq \mathbb{P} \left( \left| p(x|(x, y)_{1:N}) - p(x|\hat{\beta}) \right| \geq \frac{t}{2c_3^2 c_4} \right) \\
&\quad + \mathbb{P} \left( \left| p(x, y|(x, y)_{1:N}) - p(x, y|\hat{\beta}) \right| \geq \frac{t}{2c_3^2 c_4} \right) \tag{75}
\end{aligned}$$

946

947

Therefore, if  $t$  satisfies:

$$1 > t \geq \max \left\{ \frac{4c_3^2 c_4^2 |A|^{5/2} |B|^{3/2} p_{\max}(\hat{\alpha}) c_1^{-N}}{p(\hat{\alpha}) + |A| p_{\max}(\hat{\alpha}) c_1^{-N}}, \frac{4c_3^2 c_4^2 |B|^{3/2} p_{\max}(\hat{\beta}) c_2^{-N}}{p(\hat{\beta}) + |B| p_{\max}(\hat{\beta}) c_2^{-N}} \right\} \tag{76}$$

948

949

950

we can use the results of proposition 4 to obtain the results.  $\square$

951

952

We can also obtain the following convergence corollary from proposition 5:

953

954

**Corollary 6.** *If  $\hat{\beta} \in R$ ,  $\text{plim}_{N \rightarrow \infty} p(y|x_{1:N}, x) = p(y|x, \hat{\beta})$ .*

955

956

*Proof.* The proof is identical to the proof of corollary 4 so we omit it.  $\square$

957

## C Convergence Speed

958

959

960

961

962

We can also observe the convergence speed from  $p(\hat{\beta}|x_{1:N})$  to 1 from proposition 3. Specifically, take the derivative of the upper-bound to  $N$  in formula (60), we can see that the convergence speed is around

$$O \left( - \left( e^{\frac{\log^2 c_1}{8 \log^2 (c_4/c_3)}} \right)^{-N} - \left( e^{\frac{2 \log^2 c_2}{\log^2 (c_4/c_3)}} \right)^{-N} \right) \tag{77}$$

963

Therefore, the higher the distinction between different modes and topics, i.e, the higher  $\log c_1$  and  $\log c_2$ , the faster the convergence of the data ICTR.

964

965

966

## D Expectation of $\text{KL}(\hat{\beta}||\beta)$

967

According to the settings, each topic  $\beta \in B$  contains a few sub-topics, then the expectation of  $\text{KL}(\hat{\beta}||\beta)$  depends on KL divergences of those sub-topics:

968

969

970

971

$$\begin{aligned}
\mathbb{E} [\text{KL}(\hat{\beta}||\beta)] &= \sum_{m=1}^M \mathbb{E}_{\hat{\rho}_m, \rho_m} [\text{KL}(\hat{\rho}_m||\rho_m)] \\
&= \sum_{m=1}^M \mathbb{E}_{\hat{\rho}_m, \rho_m} \left[ \sum_s p(s|\hat{\rho}_m) \log \frac{p(s|\hat{\rho}_m)}{p(s|\rho_m)} \right] \tag{78}
\end{aligned}$$

972

973

974

Given that  $\hat{\beta}$  and  $\beta$  are different, there at least exists one subtopic is different between them, so:

$$\mathbb{E} [\text{KL}(\hat{\beta}||\beta)] \geq \mathbb{E}_{\hat{\rho}, \rho} [\text{KL}(\hat{\rho}||\rho)] \tag{79}$$

975

Note that for each  $\rho \in B_*$ , the sub-paragraph distribution  $p(s|\rho) = p(s|\tilde{\mathbf{A}}_\rho)$  is Markovian, where  $\tilde{\mathbf{A}}_\rho = [\boldsymbol{\pi}_\rho, \mathbf{A}_\rho]$  is a row concatenation of the initial probability vector  $\boldsymbol{\pi}_\rho$  and transition matrix  $\mathbf{A}_\rho$  sampled from  $\text{Dir}([\gamma]^{|\Sigma|})$ . Let  $T$  be the length of  $s$ . Expand the KL divergence, we have

976

977

978

979

980

981

$$\begin{aligned}
\mathbb{E}_{\hat{\rho}, \rho} [\text{KL}(\hat{\rho}||\rho)] &= \mathbb{E}_{\hat{\rho}, \rho}^T [\text{KL}(\hat{\rho}||\rho)] \\
&= \mathbb{E}_{\tilde{\mathbf{A}}_{\hat{\rho}}, \tilde{\mathbf{A}}_\rho} \left[ \text{KL} \left( p(\cdot|\tilde{\mathbf{A}}_{\hat{\rho}}) || p(\cdot|\tilde{\mathbf{A}}_\rho) \right) \right] \\
&= \mathbb{E}_{\tilde{\mathbf{A}}_{\hat{\rho}}, \tilde{\mathbf{A}}_\rho} \left[ \sum_{s_{1:T-1}} \sum_{s_T} \right. \\
&\quad \left. p(s_{1:T-1}|\tilde{\mathbf{A}}_{\hat{\rho}}) \tilde{\mathbf{A}}_{\hat{\rho}}^{s_{1:T-1}, s_T} \log \frac{p(s_{1:T-1}|\tilde{\mathbf{A}}_{\hat{\rho}}) \tilde{\mathbf{A}}_{\hat{\rho}}^{s_{1:T-1}, s_T}}{p(s_{1:T-1}|\tilde{\mathbf{A}}_\rho) \tilde{\mathbf{A}}_\rho^{s_{1:T-1}, s_T}} \right] \\
&= \mathbb{E}_{\hat{\rho}, \rho}^{T-1} [\text{KL}(\hat{\rho}||\rho)] + \mathbb{E}_{\tilde{\mathbf{A}}_{\hat{\rho}}, \tilde{\mathbf{A}}_\rho} \left[ \sum_{s_{T-1}, s_T} \right. \\
&\quad \left. p(s_{T-1}|\tilde{\mathbf{A}}_{\hat{\rho}}) \tilde{\mathbf{A}}_{\hat{\rho}}^{s_{T-1}, s_T} \log \frac{\tilde{\mathbf{A}}_{\hat{\rho}}^{s_{T-1}, s_T}}{\tilde{\mathbf{A}}_\rho^{s_{T-1}, s_T}} \right] \tag{80}
\end{aligned}$$

982

983

984

985

986

Note that Assumption 3 actually implicit that  $p(s_T|\tilde{\mathbf{A}}_\rho)$  is bounded for all  $T$  and  $\rho \in B_*$ . We assume the lower bound is  $c_5$ . Then, the second term of the above formula has the following lower



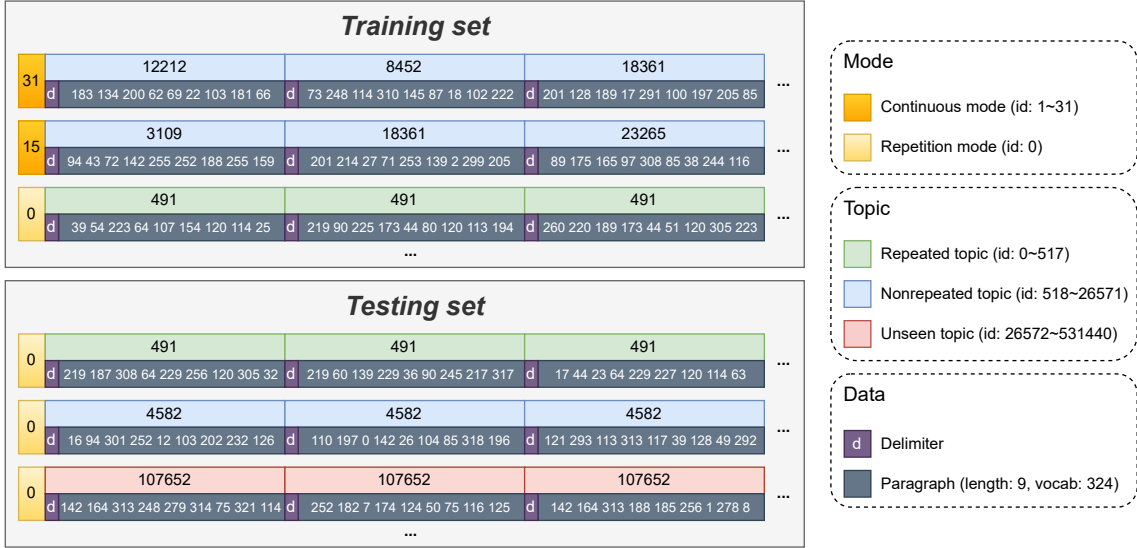


Figure 6: Examples in the synthetic dataset, where we set  $M = 3$ ,  $r_R = 1/1024$  and  $\gamma = 0.01$ .

bound:

$$\begin{aligned}
& \mathbb{E}_{\tilde{\mathbf{A}}_{\hat{\rho}}, \tilde{\mathbf{A}}_{\rho}} \left[ \sum_{s_{T-1}, s_T} \right. \\
& \quad \left. p(s_{T-1} | \tilde{\mathbf{A}}_{\hat{\rho}}) \tilde{\mathbf{A}}_{\hat{\rho}}^{s_{T-1}, s_T} \log \frac{\tilde{\mathbf{A}}_{\hat{\rho}}^{s_{T-1}, s_T}}{\tilde{\mathbf{A}}_{\rho}^{s_{T-1}, s_T}} \right] \\
& \geq c_5 \mathbb{E}_{\tilde{\mathbf{A}}_{\hat{\rho}}, \tilde{\mathbf{A}}_{\rho}} \left[ \sum_{s_{T-1}, s_T} \tilde{\mathbf{A}}_{\hat{\rho}}^{s_{T-1}, s_T} \log \frac{\tilde{\mathbf{A}}_{\hat{\rho}}^{s_{T-1}, s_T}}{\tilde{\mathbf{A}}_{\rho}^{s_{T-1}, s_T}} \right] \\
& = c_5 \mathbb{E}_{\tilde{\mathbf{A}}_{\hat{\rho}}} \left[ \sum_{s_{T-1}, s_T} \tilde{\mathbf{A}}_{\hat{\rho}}^{x_{T-1}, x_T} \log \tilde{\mathbf{A}}_{\hat{\rho}}^{x_{T-1}, x_T} \right] \\
& - c_5 \mathbb{E}_{\tilde{\mathbf{A}}_{\hat{\rho}}, \tilde{\mathbf{A}}_{\rho}} \left[ \sum_{s_{T-1}, s_T} \tilde{\mathbf{A}}_{\hat{\rho}}^{x_{T-1}, x_T} \log \tilde{\mathbf{A}}_{\rho}^{x_{T-1}, x_T} \right] \\
& = c_5 |\Sigma| [\psi(\gamma + 1) - \psi(|\Sigma|\gamma + 1)] \\
& \quad - c_5 |\Sigma| [\psi(\gamma) - \psi(|\Sigma|\gamma)] \\
& = \frac{c_5 (|\Sigma| - 1)}{\gamma}
\end{aligned} \tag{81}$$

where  $\psi(x)$  is the digamma function, and we use the property  $\psi(x + 1) = \psi(x) + 1/x$  to simplify

the above formula. Therefore, we have:

$$\begin{aligned}
\mathbb{E}_{\hat{\rho}, \rho}^T [\text{KL}(\hat{\rho} \| \rho)] & \geq \mathbb{E}_{\hat{\rho}, \rho}^{T-1} [\text{KL}(\hat{\rho} \| \rho)] + \frac{c_5 (|\Sigma| - 1)}{\gamma} \\
& \geq \mathbb{E}_{\hat{\rho}, \rho}^{T-2} [\text{KL}(\hat{\rho} \| \rho)] + \frac{2c_5 (|\Sigma| - 1)}{\gamma} \\
& \quad \dots \\
& \geq \frac{Tc_5 (|\Sigma| - 1)}{\gamma}
\end{aligned} \tag{82}$$

Therefore, the expectation of  $\text{KL}(\hat{\beta} \| \beta)$  is bounded:

$$\mathbb{E} [\text{KL}(\hat{\beta} \| \beta)] \geq \frac{Tc_5 (|\Sigma| - 1)}{\gamma} \tag{83}$$

We can see that the lower the value of  $\gamma$ , the larger the expected topic-wise KL divergence, and the more significant the topic distinction is.

## E Synthetic Dataset Illustration

Figure 6 shows examples in the synthetic dataset, where we also visualize the latent variables mode  $\alpha$  and outline  $\beta_{1:N}$  for a better understanding.

## F Computation of Prompt and Topic-wise ICTR

According to the definition, given an in-context prompt  $x_{1:N}$ , where each sample  $x_n \sim p(x | \hat{\beta})$ , ICTR is the probability that the language model generates a paragraph also belongs to topic  $\hat{\beta}$ . Thus, to measure the belongness of the generated paragraph, we use the mixture of topic-paragraph models  $\sum_{\beta \in B} \pi_{x_{1:N}}^{\beta} p(x | \beta)$  to fit the ICG distribution of the target language model  $p_{\text{LM}}(x | x_{1:N})$ . Here,

1012  $p(x|\beta)$  is fixed, and we sample  $L_1$  paragraphs  
 1013 from  $p_{\text{LM}}(x|x_{1:N})$  to fit  $\pi_{x_{1:N}}^\beta$  using EM algorithm  
 1014 (Bishop and Nasrabadi, 2006) as shown in Algo-  
 1015 rithm 1. As a result, the estimated  $\pi_{x_{1:N}}^{\hat{\beta}}$  can repre-  
 1016 sent the ICTR given the in-context prompt  $x_{1:N}$ .

1017 We further compute the topic-wise ICTR to sum-  
 1018 marize the ICG ability of a specific topic. Topic-  
 1019 wise ICTR is the expectation of prompt-wise ICTR:

$$1020 \quad \pi_N^\beta = \mathbb{E}_{p(x_{1:N}|\beta^N)} \left[ \pi_{x_{1:N}}^\beta \right] \simeq \frac{1}{L_2} \sum_{l=1}^{L_2} \pi_{x_{1:N}}^{\beta^l} \quad (84)$$

1021 Here, we use Monte-Carlo sampling to estimate  
 1022 the expectation, where  $x_{1:N}^l$  is the  $l$ -th sample of  
 1023  $\prod_{n=1}^N p(x_n|\hat{\beta})$ . Due to the large number of the  
 1024 topics (531441) in the pretrained distribution, for  
 1025 simplicity,  $L_1$  and  $L_2$  are both set to 1. Thus, the  
 1026 evaluation of a model just requires 531441 forward  
 1027 passes, where the time consumption is acceptable.  
 1028 In-context prompts for evaluation is shown in Fig-  
 1029 ure 6.

---

**Algorithm 1** Prompt-wise ICTR computation

---

Randomly initialize  $\pi_{x_{1:N}}^\beta$ .  
**for**  $l = 1, \dots, L_1$  **do**  
 $x^l \sim p_{\text{LM}}(x|x_{1:N})$   
**end for**  
**while** not convergence **do**  
**for**  $l = 1, \dots, L_1$  **do**  
 $\omega_{x_{1:N}}^{\beta,l} = \frac{\pi_{x_{1:N}}^\beta p(x^l|\beta)}{\sum_{\beta' \in B} \pi_{x_{1:N}}^{\beta'} p(x^l|\beta)}$   
**end for**  
 $\pi_{x_{1:N}}^\beta = \frac{\sum_{l=1}^{L_1} \omega_{x_{1:N}}^{\beta,l}}{L_1}$   
**end while**  
 $p_{\text{LM}}(\beta|x_{1:N}) \leftarrow \pi_{x_{1:N}}^\beta$   
**return**  $p_{\text{LM}}(\beta|x_{1:N})$

---