

IN-CONTEXT LEARNING AND OCCAM’S RAZOR

Anonymous authors

Paper under double-blind review

ABSTRACT

A central goal of machine learning is generalization. While the No Free Lunch Theorem states that we cannot obtain theoretical guarantees for generalization without further assumptions, in practice we observe that *simple* models which explain the training data generalize best—a principle called *Occam’s razor*. Despite the need for simple models, most current approaches in machine learning only minimize the training error, and at best indirectly promote simplicity through regularization or architecture design. Here, we draw a connection between Occam’s razor and in-context learning—an emergent ability of certain sequence models like Transformers to learn at inference time from past observations in a sequence. In particular, we show that the next-token prediction loss used to train in-context learners is directly equivalent to a data compression technique called prequential coding, and that minimizing this loss amounts to jointly minimizing both the training error *and* the complexity of the model that was implicitly learned from context. Our theory and the empirical experiments we use to support it not only provide a normative account of in-context learning, but also elucidate the shortcomings of current in-context learning methods, suggesting ways in which they can be improved.

1 INTRODUCTION

The goal of machine learning (ML) is to learn models that generalize to unseen data. Longstanding theory shows that minimizing training error alone can lead to overfitting and poor generalization (Bishop & Nasrabadi, 2006). To enable better generalization, ML follows the principle of *Occam’s razor*—the best explanation is the simplest one that explains the observations (Rathmanner & Hutter, 2011; Sunehag & Hutter, 2014; Hutter, 2010). The intuition is that simple rules that explain the data cannot simply memorize observations, and must instead capture more general patterns. Consequently, learning algorithms usually trade off low training error and low model complexity with *ad hoc* approaches (e.g., via regularization and inductive biases), motivating the need for notions of complexity that can be tractably minimized directly.

Although there exist mathematical notions of model complexity such as VC dimension or Kolmogorov complexity, these quantities cannot be directly minimized, or even tractably computed for the latter. In practice, we instead learn predictors that minimize training error as well as *proxies* of the model’s complexity, such as the L_1 norm of the parameters, or rely on inductive biases for low-complexity solutions that are implicit in the model class and learning algorithm. Defying this trend, however, pretrained large language models (LLMs) have a surprising ability to rapidly learn and generalize from small amounts of data presented in their context (or *prompt*) (Radford et al., 2019). This ability called *in-context learning* (ICL) is typically explained through the lens of *memory-based meta-learning* (e.g., Xie et al., 2022; Chan et al., 2022), a theoretical framework where sequence models are explicitly trained to learn statistical models from sequences of observations.

The main contribution of this paper is to provide theoretical arguments linking ICL to Occam’s razor and a preference for simple models. Briefly, our theory frames ICL as a meta-learning algorithm whose next-token prediction objective is directly equivalent to a powerful compression method called prequential coding (Blier & Ollivier, 2018). Given the relationship between optimal compression and Kolmogorov complexity, we show that the meta-objective in ICL is to find a learner capable of jointly minimizing both training error *and* model complexity across a diverse range of tasks. Our theory, along with the empirical experiments that we use to support it, explain why ICL has proven so effective in meta-learning settings, and also explain the shortcomings of current ICL methods.

Namely, we find that current methods produce learning algorithms which are susceptible to underfitting and can fail to generalize to novel tasks, suggesting principled avenues for future research.

2 OCCAM’S RAZOR AND IN-CONTEXT LEARNING

In this section, we introduce a meta-learning objective that directly targets simple models, and then show that it is equivalent to the next-token prediction objective underlying ICL. We reach this result via four key steps:

1. We begin by formalizing both training error and model simplicity through the lens of Kolmogorov complexity, which deals with optimal data and model compression.
2. We then show how learning algorithms can be used to compress data through a technique called prequential coding (Blier & Ollivier, 2018), and that minimizing the resulting “prequential code length” achieved by a learning algorithm is equivalent to jointly minimizing the training error and complexity of the model it fits.
3. We then introduce the idea of finding a learning algorithm that minimizes prequential code length by formalizing a meta-learning problem that appears difficult to optimize.
4. Finally, we show that the next-token prediction objective underlying ICL *already* solves this meta-learning problem in an efficient and scalable way.

2.1 KOLMOGOROV COMPLEXITY AND DATA COMPRESSION

Kolmogorov complexity (Kolmogorov, 1965; Li et al., 2008) is a notion of information quantity. Intuitively, the Kolmogorov complexity $K(x)$ of an object x is the length of the shortest program (in some programming language) that outputs x . A related notion is the conditional Kolmogorov complexity $K(x|y)$ of the object x given another object y , which is the length of the shortest program that takes y as input and outputs x . Finally, the Kolmogorov complexity of encoding two objects jointly is denoted $K(x, y)$. While quite abstract, this notion of complexity has deep ties to *compression*, making it intuitive as a measure of information quantity. The smaller and more “structured” an object is—regularity, patterns, rules, etc.—the more easily it can be described by a short program, correspondingly having lower Kolmogorov complexity. Although Kolmogorov complexity is very general—objects x, y can be datasets, programs, models—it is intractable to compute. However, it can often be tractably estimated or bounded, as we will show below.

A quantity relevant to ML is the joint Kolmogorov complexity of a dataset $D = (d_1, \dots, d_n)$ and of a generative model $p(d)$, where each sample $d_i \in \mathcal{D}$ is drawn *iid*:

$$K(D, p) = K(D|p) + K(p), \quad (1)$$

where $K(p)$ refers to the complexity of the model (i.e., the length of the shortest program that outputs function $p : \mathcal{D} \rightarrow \mathbb{R}^+$). This term is intractable to compute as it requires an enumeration over all programs that output p , but the conditional complexity $K(D|p)$ can be easily computed. According to (Grünwald, 2007), if the dataset is sufficiently large, the optimal method for compressing a data point d_i uses only $-\log_2 p(d_i)$ bits (e.g., using an arithmetic coding scheme, Witten et al., 1987), as in the case of Shannon information (Shannon, 2001). As such, we have $K(D|p) = -\sum_D \log_2 p(d_i)$ which is the negative log-likelihood of the data under model $p(d)$, a commonly used objective function in ML. It follows that models which achieve lower error under this objective better compress data. We provide further background on Kolmogorov complexity in Appendix A.

As we are interested in model optimization, we henceforth consider parameterized models p_θ with parameters θ . We denote a learning algorithm by a function $T : \mathcal{P}(\mathcal{D}) \rightarrow \Theta$, where $\mathcal{P}(\mathcal{D})$ denotes the power-set, which maps a dataset D to a model $p_{T(D)}$. Maximum likelihood training, which is the norm in ML, is a learning algorithm T^{ml} which fits a model that best compresses the training data:

$$T^{ml}(D) = \arg \min_{\theta'} - \sum_{d \in D} \log_2 p_{\theta'}(d) = \arg \min_{\theta'} K(D|p_{\theta'}). \quad (2)$$

However, Occam’s razor says that we also need simple models. Thus, we consider the learning algorithm T^{oc} , which defines “simple” via complexity:

$$T^{oc}(D) = \arg \min_{\theta'} [K(D|p_{\theta'}) + K(p_{\theta'})]. \quad (3)$$

In reality, T^{oc} is intractable since $K(p_{\theta'})$ cannot be computed. In practice, maximum log-likelihood training T^{ml} is often enhanced with regularizers (e.g., L_2) and inductive biases (e.g., restricting the model class) to implicitly favor low-complexity models that combat overfitting and improve generalization. For instance, deep neural networks (DNNs) trained through SGD tend to be biased towards simple solutions (Blier & Ollivier, 2018; Goldblum et al., 2023). However, existing regularizers at most amount to *indirect* methods that roughly penalize model complexity $K(p_{\theta})$ along with training error. No known learning algorithm (which we will often call a “learner” for brevity) directly attempts to minimize Equation (3) as T^{oc} would. In what follows, we introduce learners T_{ϕ} that have learnable parameters ϕ , estimated via meta-optimization, to approximate the ideal learner T^{oc} .

2.2 PREQUENTIAL CODING

While a learner T that adheres to Occam’s razor and solves Equation (3) would improve generalization, it is difficult to design one in practice. Even if $K(p_{\theta})$ could be computed efficiently, there is the further challenge of minimizing it. We will first describe an approach to the problem of estimating $K(p_{\theta})$, and then consider the optimization problem in the next section.

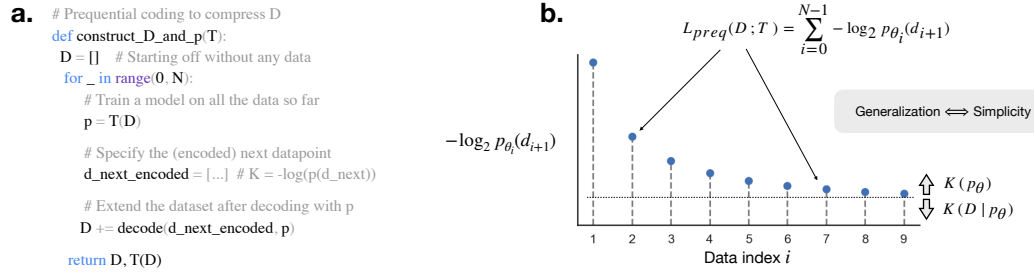


Figure 1: **Illustration of prequential coding, a method for estimating $K(D, \theta) = K(D|p_{\theta}) + K(p_{\theta})$ using p_{θ} ’s learning algorithm T .** **a.** Pseudocode of the prequential coding program, which jointly compresses D and p_{θ} by incrementally training a model using T on increasingly more data. The primary contribution to total program length comes from specifying each next datapoint d_{i+1} using the current model p_{θ_i} , which takes $-\log_2 p_{\theta_i}(d_{i+1})$ bits. `decode()` is a short function that decodes an compressed object using arithmetic coding (Witten et al., 1987) **b.** A visual illustration of prequential coding. As the learner T sees more data, it outputs models that assign a higher likelihood to new observations, and can thus better compress them. The total prequential code length $L_{preseq}(D; T)$ is given by the area under the curve. The area underneath the curve’s last point is equal to the complexity of the dataset given the final model, $K(D|p_{\theta})$. Since $L_{preseq}(D; T) = K(D|p_{\theta}) + K(p_{\theta})$, the area above the curve’s last point is equal to $K(p_{\theta})$. Prequential coding formalizes the intuition that simple models generalize better from less data.

While $K(p_{\theta})$ is difficult to measure directly, it turns out that we can estimate the joint complexity $K(D, p_{\theta}) = K(D|p_{\theta}) + K(p_{\theta})$ using a compression algorithm called *prequential coding* (illustrated in Figure 1) that leverages the learner T which gave p_{θ} (i.e., $p_{\theta} = T(D)$). Consider an ordering of *iid* datapoints $D = \{d_1, \dots, d_N\}$, and denote $D_{1:i} = \{d_1, \dots, d_i\}$. Prequential coding uses the learner T to train models on increasing amounts of data. First, we train a model on just the first data point to get $p_{\theta_1} = T(d_1)$. Because the model is trained on a single datapoint, it will not be very accurate; however, it should be better than a random model that has seen no data at all. We can then use this model p_{θ_1} to compress the next (unseen) datapoint d_2 , which takes $-\log_2 p_{\theta_1}(d_2)$ bits. At this point, we can train a new model $p_{\theta_2} = T(D_{1:2})$. Having seen more data, this model should assign a higher likelihood to a new datapoint d_3 , which we can compress using $-\log_2 p_{\theta_2}(d_3)$ bits. This process repeats until the entire dataset has been covered. At this point, the model p_{θ} can be obtained simply by applying the learning algorithm to the complete dataset $p_{\theta} = T(D)$.

The total number of bits that it takes to jointly compress D and p_{θ} using prequential coding is the sum of how many bits it takes to compress each datapoint using a model that was trained on all previous ones. Visually, it is the area under the *prequential coding curve* shown in Figure 1b. The length of this program is called the *prequential code length* $L_{preseq}(D; T)$ (Blier & Ollivier, 2018):

$$L_{preseq}(D; T) = \sum_{i=0}^{N-1} -\log_2 p_{\theta_i}(d_{i+1}) \geq K(D, p_{\theta}) = K(D|p_{\theta}) + K(p_{\theta}). \quad (4)$$

$L_{preq}(D; T)$ is an upper-bound on $K(D, p_\theta)$: prequential coding is *one* way to jointly compress the data and model, but it is not necessarily the optimal way. However, in [Section 2.3](#) we will minimize this upper-bound with respect to the learner T , and thus minimize the joint data and model complexity $K(D, p_\theta)$.

Prequential coding relates Kolmogorov complexity to intuitions about generalization in ML: the simpler a model is, the quicker it generalizes from limited amounts of training data. Although the relationship in [Equation \(4\)](#) offers a promising way forward to operationalize the idealized learner T^{oc} , there is a problem. The prequential code length given by [Equation \(4\)](#) conditions on the choice of a learner T . However, prequential coding also requires us to encode the learning algorithm itself. When we take the description length of T into account, the quantity $L_{preq}(D; T) + K(T)$ is an upper-bound on $K(D|p_\theta) + K(p_\theta)$ (see [Appendix B](#)). Since we will optimize for learners T_ϕ that minimize $L_{preq}(D; T)$, we will need to ensure that T_ϕ has low complexity.

2.3 MINIMIZING PREQUENTIAL CODE LENGTH THROUGH META-LEARNING

Consider a parameterized learner T_ϕ that minimizes the prequential code length $L_{preq}(D; T_\phi)$ of a dataset D . This objective upper-bounds the objective that the idealized learner T^{oc} minimizes, but only when $K(T_\phi)$ is low. This second criteria is violated if T_ϕ overfits to a single dataset D . To forbid T_ϕ from memorizing a single dataset, we consider a meta-dataset $\mathcal{D} = \{D^1, \dots, D^M\}$ coming from M different tasks and meta-learn T_ϕ to minimize prequential code length on average across the meta-dataset \mathcal{D} . This allows us to write the following objective for the learner T_ϕ :

$$\mathcal{L}(\mathcal{D}; \phi) = \sum_{i=1}^M L_{preq}(D^i; T_\phi) \geq \sum_{i=1}^M K(D^i, p_\theta | T_\phi) \quad (5)$$

$$= \left[\sum_{i=1}^M K(D^i | p_\theta, T_\phi) + K(p_\theta | T_\phi) \right] \quad (6)$$

$$= \left[\sum_{i=1}^M K(D^i | p_\theta) + K(p_\theta | T_\phi) \right], \quad (7)$$

where $p_\theta = T_\phi(D^i)$, and the last line is obtained from noticing that all the relevant information about D^i contained in T_ϕ is already encoded in the model $p_\theta = T_\phi(D^i)$.

By minimizing $\mathcal{L}(\mathcal{D}; \phi) = \sum_{i=1}^M L_{preq}(D^i; T_\phi)$, we thus minimize an upper-bound on the training error $K(D^i | p_\theta)$ and model complexity given the learner $K(p_\theta | T_\phi)$ in expectation over datasets. This approach of minimizing an upper-bound on an objective is a common practice when dealing with intractable objectives, as in the case of the evidence-lower-bound (ELBO) in variational inference ([Kingma, 2013](#)). As a result, minimizing expected prequential code length in [Equation \(5\)](#) meta-trains a learner $T_{\phi^*} = \arg \min_{\phi} \mathcal{L}(\mathcal{D}; \phi)$ which fits simple models that explain their training data. After obtaining T_{ϕ^*} through meta-training, the prequential code length of a new dataset of interest D is then:

$$L_{preq}(D; T_{\phi^*}) \geq K(D, p_{\theta^*} | T_{\phi^*}) \quad (8)$$

$$= K(D | p_{\theta^*}, T_{\phi^*}) + K(p_{\theta^*} | T_{\phi^*}) \quad (9)$$

$$= K(D | p_{\theta^*}) + K(p_{\theta^*} | T_{\phi^*}). \quad (10)$$

Note that the learners T_{ϕ^*} and T^{oc} ($= \arg \min_{\theta'} [K(D | p_{\theta'}) + K(p_{\theta'})]$) are not equivalent: T^{oc} aims to minimize $K(p_{\theta'})$ directly whereas T_{ϕ^*} fits models that are simple *given* T_{ϕ^*} (i.e. low $K(p_\theta | T_\phi)$). Despite these differences, the two learners are deeply related. As a result of its meta-objective in [Equation \(5\)](#), the learner T_{ϕ^*} attempts to minimize training error across many datasets while fitting compressible models. The learner T_{ϕ^*} will succeed in doing this on a *novel* dataset D when it *generalizes* to that novel dataset.

2.4 TRAINING FOR ICL META-LEARNS A PREQUENTIAL CODE LENGTH MINIMIZER

In practice, solving the meta-learning problem in [Equation \(5\)](#) involves several constraints:

1. The performance of $T_\phi(\cdot)$ must be evaluated w.r.t. a dataset’s prequential code length.
2. $T_\phi(\cdot)$ must be fast to evaluate because it is iteratively called on multiple datasets.
3. To meta-optimize ϕ , it must be easy to take gradients of $L_{preq}(\cdot; T_\phi)$ w.r.t. ϕ .
4. ϕ must parameterize an expressive class of learning algorithms, capable of minimizing prequential code length on a broad distribution of tasks and generalizing to unseen ones.

While this may appear daunting, it turns out that these desiderata are readily addressed by ICL in probabilistic sequence models. Such models are trained to predict the distribution over the next element in a sequence given its past context: $F(d_t|D_{1:t-1})$. Crucially, the sequence model F is both the learner T_ϕ and the inner model p_θ . Indeed, ϕ corresponds to the parameters of the sequence model F (e.g. weights in a Transformer), and $\theta = T_\phi(D_{1:t-1})$ is encoded by the activations of hidden units in the model when presented with the context $D_{1:t-1}$. Thus, the predicted distribution over the next token is given by: $F(d_t|D_{1:t-1}) = p_{T_\phi(D_{1:t-1})}(d_t)$. The model is trained to minimize the cumulative next-token prediction error: $\mathcal{L}(D; \phi) = \sum_{t=1}^N -\log p_{T_\phi(D_{1:t-1})}(d_t)$, which corresponds exactly to the prequential code length in Equation (4).

The dual nature of the sequence model as both the learner and the learned model offers a natural solution to the constraints above, enabling fast and differentiable evaluation of $T_\phi(\cdot)$ (2 & 3 above) with respect to cumulative next-token prediction loss (1 above). Moreover, modern sequence models can parameterize a rich class of learning algorithms, which is crucial to minimizing Equation (5) (4 above). Notably, architectures such as Transformers are known to have components which make them especially good meta-learners, such as multi-head attention (Olsson et al., 2022). It is thus no surprise that sequence models are leveraged in settings outside of the language domain (Von Oswald et al., 2023a; Bauer et al., 2023; Kirsch et al., 2022), making them general-purpose meta-learners.

This predictive formulation is quite flexible as it can be used to model data which contains sequential correlations, such as language, but can also be used to process any *iid* dataset. Indeed, consider $D = \{(x_1, y_1), \dots, (x_T, y_T)\}$ and the supervised task of learning a function $y = f(x)$. In this setting, a data point is given by the pair $d_t = (x_t, y_t)$, and straightforward tokenization schemes can be used to append a novel query x^* to the context D such that the predicted output \hat{y}^* is given by the next token in the sequence. This ICL setup is well-suited for regression-type tasks (see e.g. (see e.g., Von Oswald et al., 2023a;b)) but can be used for most supervised tasks. ICL thus turns the training of a sequence model into a meta-optimization problem over datasets—an approach also called *memory-based* meta-learning (Hochreiter et al., 2001; Santoro et al., 2016; Ortega et al., 2019). It is assumed here that (x_t, y_t) are *iid*. Although pretrained LLMs that can execute tasks with instructions given via context (or prompt) (Radford et al., 2019) break this *iid* data assumption, prequential code length is well-defined over arbitrary sequences, and our theory can possibly be adapted to settings with non-stationary data. Further exploration of this topic is left for future work.

Summary. We showed that sequence models trained on cumulative next-token prediction losses explicitly optimize a meta-learning objective that jointly minimizes training error and model complexity. This provides a normative account of ICL in terms of Occam’s razor, and explains recent experimental findings showing that LLMs are good universal compressors (Delétang et al., 2023).

3 EXPERIMENTS

Our experiments are designed to illustrate the benefits of ICL in terms of fitting simple models that generalize on *iid* examples. In Section 3.1, we compare ICL’s standard next-token prediction objective to an alternative that minimizes training error alone, rather than prequential code length. Section 3.2 then compares ICL to standard gradient-based learners that minimize training error, such as SGD. Section E.2 shows the impact of regularization on gradient-based learners from a compression perspective. In Section 3.3, we explore the impact of learner T_ϕ ’s architecture on prequential code length minimization. Section 3.4 explores the ability of T_ϕ to generalize to novel tasks. Finally, in Section 3.5 we use insights from our theory to control the data distribution seen by T_ϕ in order to better minimize prequential code length. Experimental details not described in the main paper (e.g., precise architectures, hyperparameters, etc.) can be found in Appendix E.

Tasks. In line with similar work studying ICL in a controlled setting (Mahankali et al., 2023; Garg et al., 2023; Akyürek et al., 2023), we use synthetically-generated tasks. Each task consists of a supervised learning dataset $D^i = \{(x_1, y_1), \dots, (x_k, y_k)\}$, where the labels are a (potentially

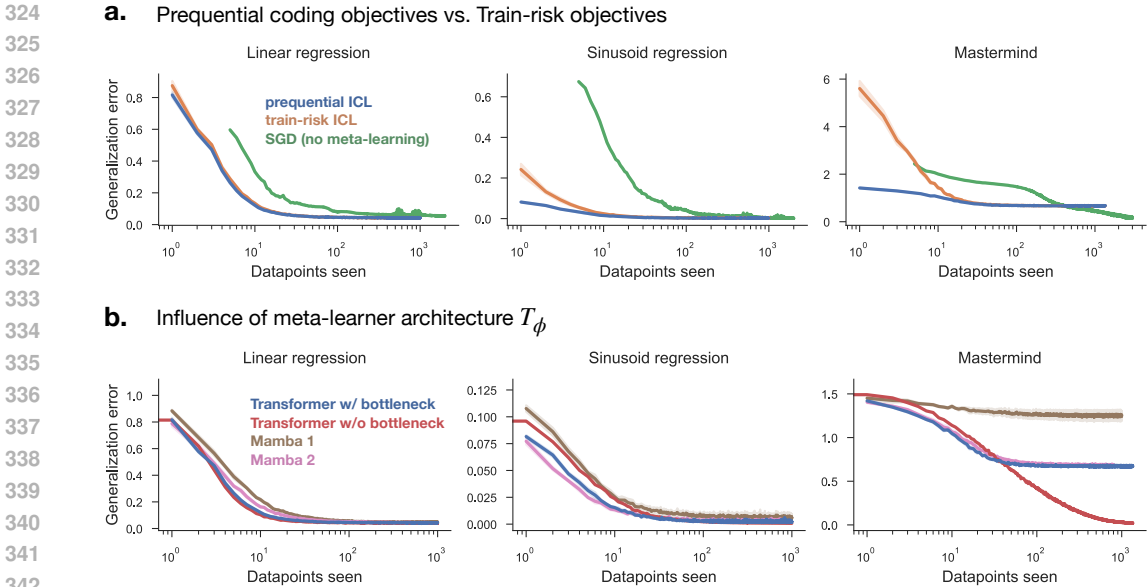
stochastic) function of the input $y_j = f^i(x_j, \epsilon_j)$. ICL learners T_ϕ are trained on a meta-dataset $\mathcal{D} = \{D^1, \dots, D^N\}$, where each D^i is associated with a different ground-truth data-generating function f^i . We primarily study three meta-datasets: **(1) Linear regression** problems where $x \in \mathbb{R}^3$ and $y \in \mathbb{R}$. The ground-truth functions f^i are noisy linear mappings $y_j = W^i x_j + b^i + \epsilon_j$, where each $\{W^i, b^i\}$ is sampled from a standard Normal distribution and ϵ_j is Gaussian noise with $\sigma^2 = 0.04$. **(2) Sinusoidal regression** problems where $x_j \in \mathbb{R}$ and functions f^i are linear combinations $y_j = \sum_{l=1}^L \alpha^{i,l} \sin(\omega^l x_j)$. We use $L = 3$ with frequencies $\omega^l \sim U(0, 5)$ that are shared across tasks, varying only the amplitudes $\alpha_{i,l} \sim \mathcal{N}(0, 1)$. **(3) Mastermind**: a multi-label classification problem inspired by the code-breaking game *Mastermind*. Each f^i is associated with an underlying discrete code (a fixed-size sequence of digits) that needs to be inferred from random guesses that return partial information. The inputs x_j are random guesses for the code, and y_j is a tuple of two class labels where the first specifies the number of digits in x_j that are correct in terms of both position and value, and the second label specifies the number of digits that are correct in value but not necessarily position. We use randomly sampled codes of length 8 with digits varying from 1..6.

3.1 COMPARISONS TO IN-CONTEXT LEARNING WITH A TRAIN-RISK OBJECTIVE

We have argued that standard ICL can be seen as a meta-learning method whose meta-objective is to minimize training error and model complexity through cumulative next-token prediction (prequential code length). However, this is not the only meta-objective that one could design for ICL. In particular, we can design an alternative meta-objective that minimizes *only* training error simply by training T_ϕ to predict *past* datapoints in the context rather than future unseen ones. In both cases, the learner T_ϕ is some function that takes a context (i.e., a partial dataset) as input, and outputs a model p_θ capable of making predictions for arbitrary datapoints. For supervised learning, this can be represented as $\hat{y}_q = T_\phi((x, y)_{1:j}, x_q)$ where $(x, y)_{1:j}$ corresponds to an observed context, x_q is the queried input, and the model p_θ is implicitly encoded in T_ϕ 's weights and latent activations given the context. In standard ICL (which we will refer to as *prequential ICL*), the query x_q is a novel input that does not appear in the context. In the alternative form of ICL (which we will call *train-risk ICL*), the query x_q is a randomly-selected input that appeared previously in the context $x_{1:j}$. Note the similarities of train-risk ICL to standard objectives of learners that minimize training error: it processes some fixed-sized training set (here a context) and attempts to minimize the empirical risk on a subset of that very same data (here a single query that appeared in the context). While nobody uses train-risk ICL in practice, it serves as an ideal control to illustrate our theory of ICL and the generalization benefits of minimizing prequential code length as opposed to only training error. One can use an identical architecture for T_ϕ in both cases and train using precisely the same methodology and loss function; the only difference is which query the loss function is evaluated on.

In our experiments, we parameterize T_ϕ using a Transformer. For the train-risk case, a standard Transformer could simply attend to the context position that matches x_q and retrieve the corresponding label. To prevent this trivial solution, we instead use a bottlenecked architecture for T_ϕ described in [Mittal et al. \(2024\)](#). In this architecture, a Transformer first summarizes the context into a low-dimensional vector $z = \text{Transformer}_\phi((x, y)_{1:j})$, and a separate prediction head—here a multi-layer perceptron (MLP)—subsequently outputs a prediction for the query $\hat{y}_q = \text{MLP}_\phi(x_q, z)$. For fair comparison, we use the same bottleneck architecture for train-risk ICL and prequential ICL in all experiments, unless otherwise stated. [Figure 2a](#) shows our comparisons between prequential ICL to train-risk ICL, where we plot the prequential coding curves for each ICL method after loss convergence on a meta-dataset. The curves are constructed at inference time by evaluating the average *iid* generalization error (i.e., unseen next-token prediction loss) on *unseen* tasks from the meta-dataset, for varying context lengths.

Findings. Two findings follow directly from our theory. The first is that for large context lengths, generalization error is identical for both prequential ICL and train-risk ICL. This is because with significant data, overfitting is less likely to occur, even when minimizing training error alone. The benefits of simple models are instead expected to be most prominent in *low-data* regimes where generalization is difficult, and this is precisely what we observe. Across all tasks, prequential ICL consistently outperforms train-risk ICL in terms of generalization for short context lengths, and this performance gap extends further the more difficult the task (e.g., it is small for linear regression, and larger for sinusoid regression and mastermind). We confirm that the performance gap widens with increasing task difficulty by fixing the function class and increasing the dimensionality of the inputs x in [Appendix C](#), which is expected given that harder tasks require more data for generalization.



343 **Figure 2: Experimental results comparing different learners.** Figures show average prequential

344 coding curves for a meta-dataset, which is the mean prediction error on unseen data (generalization

345 error, y-axis) given observed contexts of increasing length (datapoints seen, x-axis). The area under-

346 neath these curves corresponds to prequential code length. Error is measured using MSE for linear

347 and sinusoid regression and cross-entropy for Mastermind. Error bars show standard error across

348 seeds (5 for ICL, 15 for SGD). **a.** ICL from next-token prediction objectives (prequential ICL, blue)

349 yields lower prequential code lengths than ICL from past-token prediction objectives (train-risk ICL,

350 orange), with greater effects in low-data regimes. An SGD-based learner (green) fits more complex

351 models than prequential ICL and performs poorly in low-data regimes, but can generalize better in

352 large-data regimes on a difficult Mastermind task due to underfitting in ICL. **b.** The architecture used

353 to parameterize T_ϕ has substantial influence on ICL’s ability to minimize prequential code length.

354

355 **3.2 COMPARISONS TO TRADITIONAL GRADIENT-BASED LEARNERS**

356 We next consider whether there are empirical advantages of meta-learning a learner T_ϕ to min-

357 imize prequential code length through ICL, compared to using standard out-of-the-box learning

358 algorithms. In particular, we know that traditional SGD-based learners can optimize DNN models

359 that generalize well across a wide range of tasks, despite only explicitly minimizing training error.

360 We consider a standard SGD-based learner that fits a randomly-initialized MLP to the training set

361 until validation loss converges. We repeatedly sample a dataset from our meta-dataset, truncate it

362 to a specified number of observed datapoints, apply the SGD-based learner to the truncated dataset,

363 and evaluate the resulting model’s generalization error on new datapoints.

364 **Findings.** Figure 2a compares this SGD-based learner to prequential (and train-risk) ICL learn-

365 ers. Across all tasks, the models obtained through ICL generalize better in low-data regimes as a

366 result of directly minimizing model complexity. With enough training data, however, models ob-

367 tained through the SGD-based learner generalize just as well. In fact, on the Mastermind task, SGD

368 performs *better* in large-data regimes. This result demonstrates that even though the next-token pre-

369 diction objective in ICL is well-motivated from a theoretical perspective, the degree to which that

370 objective can successfully be minimized strongly depends on the architecture of T_ϕ and the methods

371 used to train it. For instance, when T_ϕ is a Transformer, the expressivity of the model it implicitly

372 fits to the context scales with the number of activations in the network (N), whereas the expressivity

373 of a DNN trained through SGD scales with the number of weights (N^2). Furthermore, the amount

374 of compute that T_ϕ uses to fit the context amounts to one forward pass of a network, whereas the

375 amount of compute that goes into fitting a dataset using SGD can be arbitrarily large.

376 **3.3 INFLUENCE OF THE IN-CONTEXT LEARNING ARCHITECTURE**

377 The previous section argued that the structure of T_ϕ can influence its ability to minimize prequential

code length. In this section, we further illustrate this point by considering a wider breadth of neural

architectures for T_ϕ . Since state-space models (SSMs) have recently been shown to exhibit ICL (Lu et al., 2024), we test Mamba 1 (Gu & Dao, 2023) and Mamba 2 (Dao & Gu, 2024). We also test a standard causal Transformer in addition to the bottlenecked Transformer from previous sections. We refer to Appendix E for additional information about the specificity of each architecture. Prequential code length comparisons in Figure 2b show that the architecture for T_ϕ indeed plays a substantial role, with the Transformers and Mamba 2 performing best across our tasks, and only the Transformer without bottleneck doing well on Mastermind in large-data regimes. Analyzing why this is the case is out of scope for this work; we only intend to show that having a next-token prediction objective alone does not guarantee that prequential code length can successfully be minimized in practice through ICL.

3.4 LARGE PRETRAINED MODELS

A core element of our theory of ICL is that T_ϕ is trained to minimize average prequential code length on a meta-dataset \mathcal{D} . There is no guarantee, however, that prequential code length will be small on a novel dataset D that was unseen at training time: this depends on the generalization abilities of the learner T_ϕ . In this section, we look at the task-generalization abilities of a large pretrained LLM (GPT-4 Achiam et al., 2023) on the Mastermind task. We do this by prompting the LLM with a description of the task and a number of in-context examples, then obtaining the logits and prediction error for a novel example. In Figure 3a, we find that despite its massive pretraining across a breadth of tasks, the LLM is unable to meaningfully minimize prequential code length on Mastermind. Not only is its prequential code length substantially higher than for a much smaller model trained on a distribution of Mastermind tasks, but it is also higher than for a naive baseline that just predicts the empirical marginal distribution over class labels in the context. These results demonstrate that even when the size of the model and meta-dataset used to train T_ϕ are scaled significantly, current methods for ICL can still struggle to minimize prequential code length on a novel task.

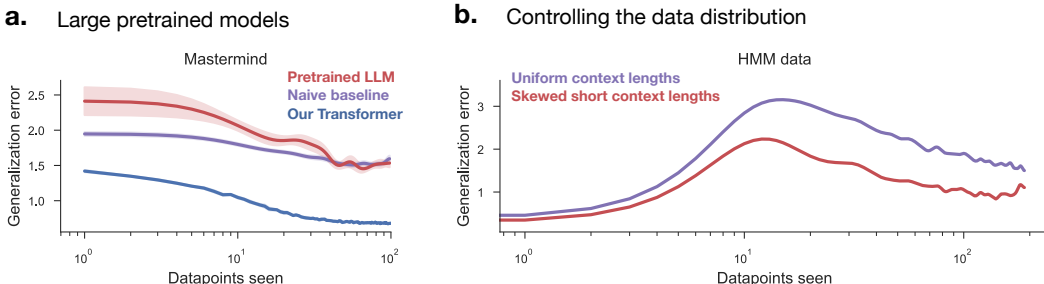


Figure 3: **Experimental results for LLM and data manipulation strategies.** Figures show average prequential coding curves for a meta-dataset, which is the mean prediction error on unseen data (generalization error, y-axis) given observed contexts of increasing length (datapoints seen, x-axis). The area underneath these curves corresponds to prequential code length. Error bars show standard error across 5 seeds. **a.** An LLM (GPT-4, red) fails to meaningfully minimize prequential code length on a novel Mastermind task, performing far worse than small ICL models trained on a distribution of Mastermind tasks (blue) and a naive baseline that predicts the marginal class distribution over the context (purple). Error is measured using cross-entropy. **b.** On a synthetic HMM dataset designed to mimic natural language, preferentially training on shorter contexts (red) yields lower prequential code lengths than training uniformly over context lengths (purple). Error is measured using reverse KL divergence between model and oracle conditioned on seen context.

3.5 IMPROVING ICL BY CONTROLLING THE DATA DISTRIBUTION

In addition to improving architectures used for T_ϕ or scaling the diversity of tasks on which it is trained, a complementary approach is to manipulate the distribution of data presented in-context at training time. This approach can be especially useful in non-*iid* settings; for instance Chan et al. (2022) found that in order for ICL to emerge in an image classification setting, the distribution over classes needed to be “bursty”, or Zipfian. In this section, we consider a simple manipulation of the data distribution that is inspired by our theory, with a particular focus on improving ICL

in language-like data modalities relevant to LLMs. In prequential coding, model complexity is related to the speed of convergence in generalization error as context length increases. We might therefore be able to further bias ICL towards simple models by sampling *short* contexts, such that downstream prediction errors on larger context lengths (after which the prequential coding curve has already converged) do not disproportionately dominate the loss.

We attempt this on synthetically-generated data from Hidden Markov Models (HMMs) that were designed to mimic the statistical properties of natural language in a simplified and controlled setting (see [Appendix E](#) for details). Briefly, we generate a family of HMMs parameterized by compositional latent attributes and train a Transformer to predict the next observation in a sequence. The model is evaluated on unseen HMMs with novel compositions of latents. Our results, presented in [Figure 3b](#), show that this data-manipulation strategy is effective. Generalization error is lower when preferentially training on short context lengths, with the gap narrowing the more tokens are seen during training as shown in [Figure E.3](#). Surprisingly, biasing the data distribution in this way not only decreases generalization error for short context lengths, but also for long ones. In general, these results show how our theory can lead to practical improvements for ICL, where we look at prequential coding curves and compression ability to guide method design.

4 RELATED WORK

Sequence modeling and compression. The idea that probabilistic models can be used to efficiently compress data is a topic widely studied in machine learning across different modalities and settings (Ollivier, 2015; Delétang et al., 2023; Blier & Ollivier, 2018; Veness et al., 2014), specifically in sequence modeling (Goyal et al., 2018; Valmeekam et al., 2023; Delétang et al., 2023) due to its close similarities to prequential coding (Blier & Ollivier, 2018). In this area, the generic sequence modeling capabilities of certain foundation models are crucial for defining effective “universal” compressors. While Goyal et al. (2018) and Valmeekam et al. (2023) claim that learned sequence models can outperform simple compressors like JPEG or gzip, they overlook model complexity in their analysis, adhering strictly to Shannon’s notion of compression. In contrast, more recent studies from Delétang et al. (2023) and Bornschein et al. (2022) opted for the Kolmogorov approach, incorporating model size to account for model complexity. Delétang et al. (2023), in particular, add nuance to the claimed advantages of foundation models due to the substantial memory allocation required to store their weights. Our theory builds on these works by relating compression and sequence modeling to the approach of meta-learning across tasks using ICL, which we show yields simple models that adhere to Occam’s razor.

In-context learning as Bayes-optimal prediction. One of the dominant perspectives of ICL and related meta-learning approaches is that they yield Bayes-optimal learners (Ortega et al., 2019; Mikulik et al., 2020; Müller et al., 2021; Hollmann et al., 2022; Binz et al., 2023; Wang et al., 2024), in the sense that they learn a prior distribution over tasks during training, and then compute a posterior given data presented in-context at inference time. This posterior can then be used to make predictions with minimum Bayes’ risk. Various studies have tested this in controlled settings with tractable posteriors (Xie et al., 2022; Panwar et al., 2024; Genewein et al., 2023; Mittal et al., 2023). Xie et al. (2022) assume a *concept* latent that parameterizes the generation of dependent samples through a Hidden Markov Model (HMM) and provide formal conditions for ICL to effectively approximate the Bayes-optimal predictor on the prompt, specifically, requiring the pretraining distribution to be structured similarly to a HMM. In a supervised fashion, Akyürek et al. (2023) construct sequence of labeled examples $(x, f(x))$ and shows that under uncertainty, ICL behaves as the Bayes-optimal predictor on noisy linear regression. Additionally, they argue that with limited capacity, ICL does not necessarily match the Bayes predictor but can meta-learn other learning algorithms, such as gradient-based algorithms on linear models and closed-form ridge regressors (Panwar et al., 2024). Grau-Moya et al. (2024) induce a prior for model simplicity in ICL by generating tasks from short programs run on Universal Turing Machines. Finally, (Raventós et al., 2024) find that under a sufficiently diverse set of pretraining tasks, ICL does *not* yield Bayes-optimal predictors, but instead infers a more uniform prior. While the Bayesian perspective of ICL is very useful and complementary to the Kolmogorov one that we have proposed, we argue in [Appendix D](#) that the Kolmogorov perspective generalizes the Bayesian one and more easily accounts for diverse findings in ICL (e.g., cases where ICL does not yield Bayes-optimal predictors).

In-context learning as a direct meta-learned optimizer. Elaborating on the possibility that ICL emulates non-Bayesian learning algorithms, Von Oswald et al. (2023a) show that k -layer linear Transformers with a specific weight parameterization can mimic k steps of gradient descent for a least squares loss. Ahn et al. (2023) provide a theoretical foundation for these observations, provably showing that the optimization of the parameters of a linear Transformer under certain assumptions about the data distribution effectively implements this learning algorithm. Concurrent studies by Zhang et al. (2023) and Mahankali et al. (2023) report similar findings, albeit under slightly different assumptions regarding weight initialization or data generation processes. Beyond the scope of linear regression, Kirsch et al. (2022) explore this phenomenon on augmented natural data (MNIST, CIFAR10) and provide insightful empirical conditions for the emergence of ICL as a general-purpose learning algorithm. Other works empirically show that Transformers can learn more complex function classes in-context, such as sinusoidal regression (Von Oswald et al., 2023a), decision trees (Garg et al., 2023), and RASP-programmable functions (Zhou et al., 2023). While prior works such as these attest to the powerful meta-learning capabilities of ICL, our work differs in that it identifies the precise meta-*objective* as an implementation of Occam’s razor.

5 DISCUSSION AND FUTURE WORK

In this work, we introduced novel theoretical arguments linking ICL and the next-token prediction objective to Occam’s razor. Our theory provides a normative account of the strong generalization abilities of in-context learners at inference time, especially in low-data regimes when compared to traditional optimizers. These theoretical insights were supported by a number of empirical experiments, some of which also identified shortcomings of current methods for ICL that should be addressed in future work.

One such shortcoming is that models learned through current ICL methods can underfit data presented in-context, and that this can hamper generalization in large-data regimes on difficult tasks. We also found that the degree of underfitting was highly dependent on the architecture used to parameterize the in-context learner (i.e., the sequence model)—a finding corroborated by Ding et al. (2024). In light of this, we hypothesize that ICL can be improved through the design of novel sequence model architectures that explicitly target prequential code length. For example, current methods learn in-context through a single forward pass of a sequence model with fixed layer depth. In contrast, DNNs can be trained using gradient-based methods until training loss converges, which can take weeks and substantial compute. One improvement to ICL might therefore be to augment current sequence model architectures with “layers” that use built-in optimization primitives with variable compute budgets, as was done in Von Oswald et al. (2023b). Another promising approach is to combine ICL and SGD through a “mixture of learners” that reaps their complementary benefits. ICL is sample-efficient and generalizes well in low-data regimes, while SGD-based methods that optimize the weights of a DNN excel on difficult tasks when significant training data is available. Recent work by Bornschein et al. (2024) explored a simple method for combining both learners by presenting a smaller number of *recent* tokens in-context to a sequence model for ICL, while at the same time using a large number of earlier tokens to fine-tune the weights of the sequence model using gradient methods, finding significant performance gains.

Another challenge of ICL that follows directly from our theory is that the in-context learner must generalize to novel tasks and datasets. While we found that task generalization was successful over narrow task distributions (e.g. a distribution of linear regression tasks), we also found that task generalization was more difficult in open-ended cases, in which even a large pretrained LLM was unable to learn in-context on a novel task that was easily solved by a small MLP trained using SGD. One possible path forward is to have many domain-specific in-context learners that each specialize in compressing data from a given task distribution. Another option is to learn *simple learners* that are more likely to generalize to novel tasks, which could be achieved through inductive biases, regularization, or, intriguingly, through an additional meta-layer of ICL at the task level that would minimize the Kolmogorov complexity of the learner itself (and not only the model it fits).

Finally, our work only provides a theoretical framework for ICL on *iid* data. Relaxing these *iid* assumptions opens up two avenues for future work: connecting ICL to generalization on out-of-distribution samples, and studying the effect of nonstationary data presented in context, as is the case in language and the HMM experiment presented here.

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement
546 preconditioned gradient descent for in-context learning, 2023. URL [https://arxiv.org/
547 abs/2306.00297](https://arxiv.org/abs/2306.00297).
- 548 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learn-
549 ing algorithm is in-context learning? investigations with linear models. In *Proceedings of the
550 Eleventh International Conference on Learning Representations (ICLR)*, 2023. URL [https:
551 //openreview.net/forum?id=0g0X4H8yN4I](https://openreview.net/forum?id=0g0X4H8yN4I).
- 552 Jakob Bauer, Kate Baumli, Feryal Behbahani, Avishkar Bhoopchand, Nathalie Bradley-Schmieg,
553 Michael Chang, Natalie Clay, Adrian Collister, Vibhavari Dasagi, Lucy Gonzalez, et al. Human-
554 timescale adaptation in an open-ended task space. In *International Conference on Machine Learn-
555 ing*, pp. 1887–1935. PMLR, 2023.
- 556 Marcel Binz, Ishita Dasgupta, Akshay K Jagadish, Matthew Botvinick, Jane X Wang, and Eric
557 Schulz. Meta-learned models of cognition. *Behavioral and Brain Sciences*, pp. 1–38, 2023.
- 558 Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, vol-
559 ume 4. Springer, 2006.
- 560 Léonard Blier and Yann Ollivier. The description length of deep learning models. *Advances in
561 Neural Information Processing Systems*, 31, 2018.
- 562 Jorg Bornschein, Yazhe Li, and Marcus Hutter. Sequential learning of neural networks for prequen-
563 tial mdl. *arXiv preprint arXiv:2210.07931*, 2022.
- 564 Jorg Bornschein, Yazhe Li, and Amal Rannen-Triki. Transformers for supervised online continual
565 learning. *arXiv preprint arXiv:2403.01554*, 2024.
- 566 Gregory J Chaitin. On the length of programs for computing finite binary sequences. *Journal of the
567 ACM (JACM)*, 13(4):547–569, 1966.
- 568 Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond,
569 James McClelland, and Felix Hill. Data distributional properties drive emergent in-context learn-
570 ing in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891, 2022.
- 571 Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through
572 structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- 573 Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christo-
574 pher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, et al.
575 Language modeling is compression. *arXiv preprint arXiv:2309.10668*, 2023.
- 576 D. Deutsch. *The Beginning of Infinity: Explanations that Transform the World*. Always learning.
577 Penguin Books, 2012. ISBN 9780140278163. URL [https://books.google.ca/books?
578 id=WFZl7YvsiuIC](https://books.google.ca/books?id=WFZl7YvsiuIC).
- 579 Nan Ding, Tomer Levinboim, Jialin Wu, Sebastian Goodman, and Radu Soricut. CausalLM is not
580 optimal for in-context learning. In *The Twelfth International Conference on Learning Represen-
581 tations*, 2024. URL <https://openreview.net/forum?id=guRNebwZBb>.
- 582 Lance Fortnow. Kolmogorov complexity. In *Aspects of Complexity, Minicourses in Algorithmics,
583 Complexity, and Computational Algebra, NZMRI Mathematics Summer Meeting, Kaikoura, New
584 Zealand*, pp. 73–86, 2000.
- 585 Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn
586 in-context? a case study of simple function classes, 2023. URL [https://arxiv.org/abs/
587 2208.01066](https://arxiv.org/abs/2208.01066).

- 594 Tim Genewein, Grégoire Delétang, Anian Ruoss, Li Kevin Wenliang, Elliot Catt, Vincent Dutoirdoir,
595 Jordi Grau-Moya, Laurent Orseau, Marcus Hutter, and Joel Veness. Memory-based meta-learning
596 on non-stationary distributions. In *International conference on machine learning*, pp. 11173–
597 11195. PMLR, 2023.
- 598
599 Micah Goldblum, Marc Finzi, Keefer Rowan, and Andrew Gordon Wilson. The no free lunch
600 theorem, kolmogorov complexity, and the role of inductive biases in machine learning. *arXiv*
601 *preprint arXiv:2304.05366*, 2023.
- 602 Mohit Goyal, Kedar Tatwawadi, Shubham Chandak, and Idoia Ochoa. Deepzip: Lossless data com-
603 pression using recurrent neural networks, 2018. URL [https://arxiv.org/abs/1811.](https://arxiv.org/abs/1811.08162)
604 [08162](https://arxiv.org/abs/1811.08162).
- 605
606 Jordi Grau-Moya, Tim Genewein, Marcus Hutter, Laurent Orseau, Grégoire Delétang, Elliot Catt,
607 Anian Ruoss, Li Kevin Wenliang, Christopher Mattern, Matthew Aitchison, et al. Learning uni-
608 versal predictors. *arXiv preprint arXiv:2401.14953*, 2024.
- 609 Peter D Grünwald. *The minimum description length principle*. MIT press, 2007.
- 610
611 Peter D Grünwald and Paul MB Vitányi. Kolmogorov complexity and information theory. with an
612 interpretation in terms of questions and answers. *Journal of Logic, Language and Information*,
613 12:497–529, 2003.
- 614
615 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
616 *preprint arXiv:2312.00752*, 2023.
- 617
618 Sepp Hochreiter, A Steven Younger, and Peter R Conwell. Learning to learn using gradient descent.
619 In *Artificial Neural Networks—ICANN 2001: International Conference Vienna, Austria, August*
620 *21–25, 2001 Proceedings 11*, pp. 87–94. Springer, 2001.
- 621
622 Noah Hollmann, Samuel Müller, Katharina Eggenberger, and Frank Hutter. TabPFN: A transformer
623 that solves small tabular classification problems in a second. *arXiv preprint arXiv:2207.01848*,
624 2022.
- 625
626 Marcus Hutter. A complete theory of everything (will be subjective). *Algorithms*, 3(4):329–350,
627 2010. ISSN 1999-4893. doi: 10.3390/a3040329. URL [http://arxiv.org/abs/0912.](http://arxiv.org/abs/0912.5434)
628 [5434](http://arxiv.org/abs/0912.5434).
- 629
630 Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 631
632 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL
633 <https://arxiv.org/abs/1412.6980>.
- 634
635 Louis Kirsch, James Harrison, Jascha Sohl-Dickstein, and Luke Metz. General-purpose in-context
636 learning by meta-learning transformers. *arXiv preprint arXiv:2212.04458*, 2022.
- 637
638 Andrei N Kolmogorov. Three approaches to the quantitative definition of information’. *Problems of*
639 *information transmission*, 1(1):1–7, 1965.
- 640
641 Ming Li, Paul Vitányi, et al. *An introduction to Kolmogorov complexity and its applications*, vol-
642 *ume 3*. Springer, 2008.
- 643
644 Chris Lu, Yannick Schroecker, Albert Gu, Emilio Parisotto, Jakob Foerster, Satinder Singh, and
645 Feryal Behbahani. Structured state space models for in-context reinforcement learning. *Advances*
646 *in Neural Information Processing Systems*, 36, 2024.
- 647
648 Arvind Mahankali, Tatsunori B. Hashimoto, and Tengyu Ma. One step of gradient descent is
649 provably the optimal in-context learner with one layer of linear self-attention, 2023. URL
650 <https://arxiv.org/abs/2307.03576>.
- 651
652 Vladimir Mikulik, Grégoire Delétang, Tom McGrath, Tim Genewein, Miljan Martić, Shane Legg,
653 and Pedro Ortega. Meta-trained agents implement bayes-optimal agents. *Advances in neural*
654 *information processing systems*, 33:18691–18703, 2020.

- 648 Sarthak Mittal, Niels Leif Bracher, Guillaume Lajoie, Priyank Jaini, and Marcus A Brubaker. Ex-
649 ploring exchangeable dataset amortization for bayesian posterior inference. In *ICML 2023 Work-*
650 *shop on Structured Probabilistic Inference* {\&} *Generative Modeling*, 2023.
- 651 Sarthak Mittal, Eric Elmoznino, Leo Gagnon, Sangnie Bhardwaj, Dhanya Sridhar, and Guillaume
652 Lajoie. Does learning the right latent variables necessarily improve in-context learning?, May
653 2024. URL <http://arxiv.org/abs/2405.19162>. arXiv:2405.19162 [cs].
- 654 Samuel Müller, Noah Hollmann, Sebastian Pineda Arango, Josif Grabocka, and Frank Hutter. Trans-
655 formers can do bayesian inference. *arXiv preprint arXiv:2112.10510*, 2021.
- 656 Yann Ollivier. Auto-encoders: reconstruction versus compression, 2015. URL <https://arxiv.org/abs/1403.7752>.
- 657 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
660 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli,
661 Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane
662 Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish,
663 and Chris Olah. In-context learning and induction heads, 2022. URL <https://arxiv.org/abs/2209.11895>.
- 664 Pedro A Ortega, Jane X Wang, Mark Rowland, Tim Genewein, Zeb Kurth-Nelson, Razvan Pascanu,
665 Nicolas Heess, Joel Veness, Alex Pritzel, Pablo Sprechmann, et al. Meta-learning of sequential
666 strategies. *arXiv preprint arXiv:1905.03030*, 2019.
- 667 Madhur Panwar, Kabir Ahuja, and Navin Goyal. In-context learning through the bayesian prism,
668 2024. URL <https://arxiv.org/abs/2306.04891>.
- 669 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
670 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 671 Samuel Rathmanner and Marcus Hutter. A philosophical treatise of universal induction. *Entropy*,
672 13(6):1076–1136, 2011. ISSN 1099-4300. doi: 10.3390/e13061076. URL <http://arxiv.org/abs/1105.5721>.
- 673 Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the
674 emergence of non-bayesian in-context learning for regression. *Advances in Neural Information*
675 *Processing Systems*, 36, 2024.
- 676 Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-
677 learning with memory-augmented neural networks. In *International conference on machine learn-*
678 *ing*, pp. 1842–1850. PMLR, 2016.
- 679 Claude Elwood Shannon. A mathematical theory of communication. *ACM SIGMOBILE mobile*
680 *computing and communications review*, 5(1):3–55, 2001.
- 681 Ray J Solomonoff. A formal theory of inductive inference. part i. *Information and control*, 7(1):
682 1–22, 1964.
- 683 Peter Sunehag and Marcus Hutter. Intelligence as inference or forcing Occam on the world. In *Proc.*
684 *7th Conf. on Artificial General Intelligence (AGI'14)*, volume 8598 of *LNAI*, pp. 186–195, Quebec
685 City, Canada, 2014. Springer. ISBN 978-3-319-09273-7. doi: 10.1007/978-3-319-09274-4_18.
- 686 Chandra Shekhara Kaushik Valmeekam, Krishna Narayanan, Dileep Kalathil, Jean-Francois Cham-
687 berland, and Srinivas Shakkottai. Llmzip: Lossless text compression using large language models,
688 2023. URL <https://arxiv.org/abs/2306.04050>.
- 689 Joel Veness, Marc G. Bellemare, Marcus Hutter, Alvin Chua, and Guillaume Desjardins. Compress
690 and control, 2014. URL <https://arxiv.org/abs/1411.5326>.
- 691 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordv-
692 intsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient
693 descent. In *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023a.

702 Johannes Von Oswald, Eyvind Niklasson, Maximilian Schlegel, Seijin Kobayashi, Nicolas Zucchet,
703 Nino Scherrer, Nolan Miller, Mark Sandler, Max Vladymyrov, Razvan Pascanu, et al. Uncovering
704 mesa-optimization algorithms in transformers. *arXiv preprint arXiv:2309.05858*, 2023b.
705

706 Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large lan-
707 guage models are latent variable models: Explaining and finding good demonstrations for in-
708 context learning, 2024. URL <https://arxiv.org/abs/2301.11916>.

709 Ian H Witten, Radford M Neal, and John G Cleary. Arithmetic coding for data compression. *Com-*
710 *munications of the ACM*, 30(6):520–540, 1987.
711

712 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context
713 learning as implicit bayesian inference, 2022. URL <https://arxiv.org/abs/2111.02080>.
714

715 Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models in-
716 context, 2023. URL <https://arxiv.org/abs/2306.09927>.
717

718 Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Josh Susskind, Samy Bengio,
719 and Preetum Nakkiran. What Algorithms can Transformers Learn? A Study in Length Gener-
720 alization, October 2023. URL <http://arxiv.org/abs/2310.16028>. arXiv:2310.16028
721 [cs, stat].
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
755

APPENDIX A BACKGROUND ON KOLMOGOROV COMPLEXITY

Kolmogorov complexity was independently developed in the 1960s by Kolmogorov (1965), Solomonoff (1964), and Chaitin (1966), and defines a notion of “information quantity”.

Intuitively, the Kolmogorov complexity of an object is the length of the shortest program (in some programming language) that outputs that object. Specifically, given some finite string x , $K(x)$ is the length $l(r)$ (in bits) of the shortest binary program r that prints x and halts. Let U be a universal Turing machine that executes these programs. The Kolmogorov complexity of x is then:

$$K(x) = \min_r \{l(r) : U(r) = x, r \in \{0, 1\}^*\}, \quad (11)$$

where $\{0, 1\}^*$ denotes the space of finite binary strings. A related notion is the conditional Kolmogorov complexity of a string x given another string y , which is the length of the shortest program that takes y as input and outputs x :

$$K(x|y) = \min_r \{l(r) : U(r(y)) = x, r \in \{0, 1\}^*\}, \quad (12)$$

where $r(y)$ denotes a program taking y as input. Finally, we can also define a “joint” Kolmogorov complexity $K(x, y)$, which denotes the length of the shortest program that jointly outputs both x and y . Surprisingly, joint Kolmogorov complexity is related to conditional Kolmogorov complexity (up to an additive logarithmic term, which we will ignore) by the Symmetry of Information theorem (Li et al., 2008):

$$K(x, y) = K(y|x) + K(x) = K(x|y) + K(y). \quad (13)$$

Kolmogorov complexity has many intuitive properties that make it attractive as a measure of information quantity, and although it is less common than notions from Shannon information theory (Shannon, 2001), it is strictly more general (as we will show later below). The smaller and the more “structure” an object has—regularity, patterns, rules, etc.—the more easily it can be described by a short program and the lower its Kolmogorov complexity. Kolmogorov complexity therefore is deeply rooted in the idea of compression. For instance, a sequence with repeating patterns or a dataset that spans a low-dimensional subspace can be significantly compressed relative to its original size, and this results in low Kolmogorov complexity. In contrast, a random string devoid of any structure cannot be compressed at all and must in effect be “hard-coded”, making its Kolmogorov complexity equal to its original size in bits.

While powerful, Kolmogorov complexity has certain limitations. First and foremost, Kolmogorov is intractable to compute exactly because it requires a brute force search over an exponentially large space of possible programs. It is therefore often of conceptual rather than practical value, although it can nevertheless be upper-bounded using more efficient compression strategies. Second, Kolmogorov complexity depends on the programming language of choice. For instance, if a programming language has a built-in primitive for the object being encoded, Kolmogorov complexity is trivially small. This concern, however, is often overblown: given any two Turing-complete programming languages, the difference in Kolmogorov complexity that they assign to an object is upper-bounded by a constant that is independent of the object itself, because any Turing-complete programming language can simulate another (Grünwald & Vitányi, 2003; Fortnow, 2000). In practice, we can simply consider “reasonable” Turing-complete programming languages that don’t contain arbitrary object-specific primitives, in which case this simulation constant will be relatively small and the particular programming language of choice will have little effect. Finally, Kolmogorov complexity is only defined for discrete objects because no terminating program can output a continuous number with infinite precision. This concern is also less consequential in practice, because we can always represent continuous objects using finite (e.g., floating-point) precision.

Important properties for machine learning. In ML, we are often concerned with datasets and probabilistic models. Kolmogorov complexity relates to these two concepts in several interesting

ways. First, we can ask about the Kolmogorov complexity of a finite dataset $X = (x_1, \dots, x_n)$ where each sample is drawn *iid* from a distribution $p(x)$. It turns out that if we have access to the true distribution $p(x)$, optimal algorithms such as arithmetic coding (Witten et al., 1987) can encode each sample using only $\log_2 p(x_i)$ bits. Intuitively, this is because samples that occur more frequently can be encoded using shorter codes in order to achieve an overall better compression. We thus have that:

$$K(X|p) = - \sum_{i=1}^n \log_2 p(x_i). \quad (14)$$

If instead of access to the true distribution $p(x)$ we only have a probabilistic model of the data $p_\theta(x)$, we have that:

$$K(X|p) \leq K(X|p_\theta) \leq - \sum_{i=1}^n \log_2 p_\theta(x_i), \quad (15)$$

where we have equality on the LHS when $p_\theta = p$ and equality on the RHS when the cost of improving p_θ (in bits of written code) would be greater than the benefits from more accurate modeling. In practice, if p_θ is close to p , we can say that $K(X|p_\theta) \approx - \sum_{i=1}^n \log_2 p_\theta(x_i)$.

This insight is significant. Notice that $-\sum_{i=1}^n \log_2 p_\theta(x_i)$ is the negative log-likelihood of the data under the model, which is a common loss function used in ML. This tells us that models with lower error better compress their data, and directly relates Kolmogorov complexity to optimization in ML. However, what if we do not have a model? What is the Kolmogorov complexity of the data itself? Intuitively, if the dataset is sufficiently large, the optimal method for encoding it should be to first specify a model and then encode the data using that model as in Equation (15). Specifically, using identities in Fortnow (2000), we have:

$$K(X) \leq K(X|p_\theta) + K(p_\theta). \quad (16)$$

This encoding scheme on the RHS is referred to as a 2-part code (Grünwald, 2007). For large datasets, we have equality when the model’s description length and error are jointly minimized, which occurs when the model $p_\theta(x)$ is equivalent to the true distribution $p(x)$:

$$K(X) = \arg \min_{p_\theta} [K(X|p_\theta) + K(p_\theta)] = \arg \min_{p_\theta} \left[- \sum_{i=1}^n \log_2 p_\theta(x_i) + K(p_\theta) \right] \quad (17)$$

$$= K(X|p) + K(p) = - \sum_{i=1}^n \log_2 p(x_i) + K(p). \quad (18)$$

Again, we can draw important connections to ML. Equation (16) says that the Kolmogorov complexity of a dataset is upper-bounded by the a model’s error and complexity. In addition, Equations (17) and (18) tell us that the simplest model that explains the data is most likely to be the true one, which draws a theoretical link between compression, maximum likelihood training, model complexity, and generalization (Goldblum et al., 2023).

Relation to Shannon information. In Shannon information theory (Shannon, 2001), the notion of information quantity is entropy. Given a random variable $X \sim p(x)$, entropy is defined as: $H(X) = \mathbb{E}_{x \sim p(x)} - \log_2(p(x))$. Notice that the $-\log_2(p(x))$ inside the expectation is equal the quantity inside the sum of Equation (14), which specified the minimum number of bits needed to encode a sample from a dataset given the distribution that sample was drawn from. This is no accident: entropy can be seen as the average number of bits needed to compress events from a distribution using an optimal encoding scheme when the distribution $p(x)$ is known. If we simply sum these bits

864 for a finite number of samples instead of taking an expectation, we get exactly $K(X|p)$ as defined
865 in Equation (14).
866

867 As we have seen, though, the assumption about a known distribution $p(x)$, need not be made in the
868 Kolmogorov complexity framework. In this sense, Kolmogorov complexity is a strict generalization
869 of Shannon information theory: $K(X)$ as defined in Equation (18) is equivalent to summed entropy
870 plus the complexity of the distribution $p(x)$, which is unknown and needs to be encoded. In the
871 Shannon framework, it is difficult to derive a meaningful notion for the information quantity in
872 the distribution $p(x)$ because it is an individual object—a function, in particular—and Shannon
873 information is only defined for random variables (Grünwald & Vitányi, 2003). A second drawback of
874 Shannon information is that entropy is a measure of statistical determinability of states; information
875 is fully determined by the probability distribution on states and unrelated to the representation,
876 structure, or content of the individual states themselves (Grünwald & Vitányi, 2003). For this current
877 work, we require a notion of complexity that can account for representations and functions, making
878 Kolmogorov complexity better suited to the task.

879 APPENDIX B PREQUENTIAL CODING AND COMPRESSION WITHOUT A 880 KNOWN LEARNING ALGORITHM 881

882 When introducing the relationship between prequential coding and optimal compression in Equa-
883 tion (4), we mentioned that a key assumption is that the learning algorithm T is known. In reality,
884 then, we have that:
885

$$886 K(D|p_\theta) + K(p_\theta) = K(D, p_\theta) \tag{19}$$

$$887 \leq K(D, p_\theta, T) \tag{20}$$

$$888 = K(D, p_\theta|T) + K(T) \tag{21}$$

$$889 \leq L_{preq}(D; T) + K(T) \tag{22}$$

$$890 \implies L_{preq}(D; T) + K(T) \geq K(D|p_\theta) + K(p_\theta), \tag{23}$$

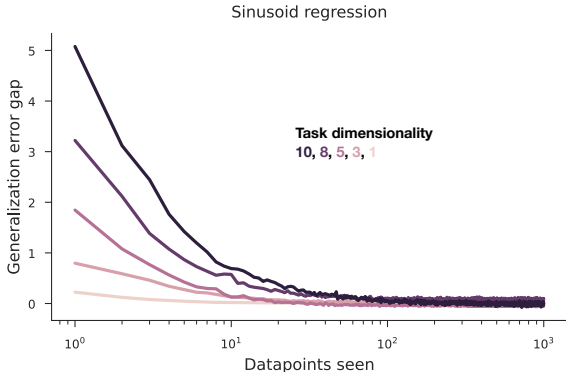
891 where the first inequality on line Equation (20) appears because compressing additional objects
892 can only take more bits, and the second inequality on line Equation (22) comes from the fact that
893 prequential coding is not necessarily the optimal way to compress a dataset and model given a
894 learning algorithm. If the learning algorithm is a short program like SGD, however, then $K(T) \approx 0$
895 and $L_{preq}(D; T)$ is an upper-bound of $K(D|p_\theta) + K(p_\theta)$. For simple learning algorithms, then,
896 Equation (4) holds.
897
898
899

900 APPENDIX C EFFECT OF TASK DIFFICULTY ON PREQUENTIAL CODE LENGTH 901

902 In Section 3.1 Figure 2a, we found that a meta-learned in-context learner trained to minimize pre-
903 quential code length (prequential ICL) was better able to generalize than one that only minimized
904 training error (train-risk ICL). We further noted that the gap in generalization error between these
905 two learners was greater in low-data regimes, and that the gap extended further as a function of
906 task difficulty (i.e., more in-context data was required to close the gap going from linear regres-
907 sion, to sinusoid regression, to Mastermind). This result is predicted by our theory relating ICL to
908 Occam’s razor. A complex task requires the algorithm to learn more complex functions to success-
909 fully minimize train risk. However, learning more complex functions with very limited data leads to
910 overfitting, which is the basis for our hypothesis that as task complexity increases, simple predictors
911 learned by minimizing prequential code length enjoy a bigger advantage over predictors learned by
912 minimizing train risk.

913 To investigate the effect of task difficulty more systematically in this section, we fix the underlying
914 meta-dataset (sinusoid regression tasks) and vary the dimensionality of the input data $dim(x)$. We
915 plot our results in Figure C.1, showing the difference in generalization error between train-risk ICL
916 learners and prequential ICL learners. As expected, as task difficulty increases, this generalization
917 gap extends further, and the train-risk learners must observe more data in-context in order to close
it.

918
919
920
921
922
923
924
925
926
927
928
929
930



931
932
933
934
935
936
937
938
939
940
941
942

Figure C.1: **Comparison of gap between prequential ICL and train-risk ICL as a function of task difficulty.** Figure shows the difference in average prequential coding curves (i.e., generalization error for train-risk ICL – generalization error for prequential ICL) for sinusoid regression tasks of increasing input dimensionality. Error is measured using MSE. Error bars show standard error across 5 seeds. For all task dimensionalities, the performance gap is positive: ICL from next-token prediction objectives (prequential ICL) yields lower prequential code lengths than ICL from past-token prediction objectives (train-risk ICL), with greater effects in low-data regimes. This gap in generalization error increases with task dimensionality, demonstrating that learners which minimize prequential code length generalize better in virtue of fitting simpler models, and that these simpler models are most important when generalization is difficult (i.e., when the task difficulty is too great for the amount of training data observed).

943
944

APPENDIX D ADVANTAGES OVER THE BAYESIAN PERSPECTIVE

945
946
947
948
949
950
951

The Bayes-optimal prediction perspective of ICL and meta-learning says that by meta-training on some set of tasks \mathcal{D} , the learner infers some prior over latent task variables—or, equivalently, a prior over models— $p(p_\theta|\mathcal{D})$. On some novel task D , the learner then infers a posterior over models that both explain the training data (i.e., assign it a high likelihood) and are consistent with the prior: $p_{\mathcal{D}}(p_\theta|D) = p(D|p_\theta)p(p_\theta|\mathcal{D})/Z$, where Z is a normalizing constant. According to the theory, subsequent predictions are then done through implicit Bayesian averaging under this posterior model distribution.

952
953
954
955
956

Crucial differences in our theory are that \mathcal{D} does not need to be drawn from a well-defined distribution over tasks for us to reason about the meta-learning problem—the Kolmogorov framework does not require this—and $K(p_\theta|T_\phi)$ is not *literally* a prior probability distribution over models given \mathcal{D} —it only implicitly defines a prior based on the meta-learned T_ϕ . As a result, our theory generalizes the Bayesian perspective.

957
958
959
960
961
962
963
964
965
966

To see why these generalizations provide value, consider where the prior in the Bayesian framework $p(p_\theta|\mathcal{D})$ comes from. This prior is not defined explicitly in the ICL framework; instead, it is implicitly defined based on \mathcal{D} , the implicit *initial* prior $p(p_\theta)$, and the implicit inference machinery that approximates $p(p_\theta|\mathcal{D}) = p(\mathcal{D}|p_\theta)p(p_\theta)/Z$. All of these implicit components make any meaningful analysis difficult, since it is difficult to characterize them. However, these implicit components are all intrinsic properties of the meta-learning algorithm (the meta-learner’s architecture, the meta-objective, etc.), which we *do* have *explicit* control over. Our theory only makes reference to this meta-learner T_ϕ and the description length of models under it $K(p_\theta|T_\phi)$, rather than to objects that are only implicitly defined (and never known). As such, we argue that our theory is more amenable to analysis and provides more explanatory value.

967
968
969
970
971

For example, in the Kolmogorov framework that we have proposed, it is easy to see how ICL might in some cases generalize to a novel dataset D that is entirely out-of-domain with respect to \mathcal{D} . Perhaps, for instance, the tasks have compositional structure and T_ϕ has some inductive biases for compositional generalization. In contrast, it is far more difficult to find a good explanation for such a phenomenon in the Bayesian framework. The explanation would have to be in terms of some implicit initial prior $p(p_\theta)$ (which we never defined) and the subsequent prior $p(p_\theta|\mathcal{D})$ that

972 it induced. Proponents of the Bayesian framework would thus have to say “ahh, generalization
 973 here must have been possible because $p(p_\theta)$ had the right kind of structure”. However, this same
 974 rationale could be used to explain *any* outcome (positive or negative), and therefore is a bad scientific
 975 explanation (Deutsch, 2012).

976 Another problem with the Bayesian perspective is that its predictions do not always hold in practice.
 977 Notably, Raventós et al. (2024) found that when the diversity in pretraining tasks is sufficiently large,
 978 solutions emerge that are *not* consistent with a Bayes-optimal predictor that uses the pretraining task
 979 distribution as its prior. Instead, the solution is consistent with a much broader prior, which allows
 980 the learner to adapt to novel tasks that are outside of the pretraining task distribution. Our theory,
 981 in contrast, permits explanations for this phenomenon. For instance, perhaps that model used to
 982 parameterize T_ϕ had insufficient capacity to encode a diverse (and potentially complex) prior over
 983 tasks, and instead learned a simpler approximation with more broad coverage over a larger space of
 984 tasks.

986 APPENDIX E EXPERIMENT DETAILS

988 In this section, we provide additional experimental details, including a comprehensive overview of
 989 the model architectures and hyperparameters used during training.

991 E.1 META-LEARNER ARCHITECTURES

993 We considered different architectures which exhibit ICL to study and compare their ability to min-
 994 imize prequential code length (Section 3.3). Each architecture described here parameterizes the
 995 meta-learner T_ϕ .

997 **Transformer with bottleneck.** We use a standard causal decoder-only Transformer with 4 layers,
 998 4 attention heads, 256 latent dimensions and a feed-forward network with 512 dimensions. Addi-
 999 tionally, it has linear projection that bottlenecks the Transformer to 128 dimension. A 5-layer MLP
 1000 with RELU activations and 256 latent dimensions is used as a separate prediction head.

1001 The Transformer takes a dataset D as input in the format $[x_1, y_1], [x_2, y_2], \dots, [x_n, y_n]$ (where x_i
 1002 and y_i are concatenated and each $[\cdot]$ is a token) and computes $T_\phi(D_{1:t-1})$ for each context size
 1003 starting from 1 to $n - 1$. The computation of $T_\phi(D_{1:t-1})$ is based on the encoding of the t -th token,
 1004 which attends only to tokens that appear to the left of $[x_t, y_t]$ and itself. Information leakage from
 1005 future tokens is prevented using a causal mask. After computing $T_\phi(D_{1:t-1})$, we concatenate it
 1006 with x_t (i.e., $[T_\phi(D_{1:t-1}), x_t]$) and pass this combined input to an MLP prediction head to predict
 1007 the next y -token.

1008 **Transformer without bottleneck.** We use a custom encoder-decoder Transformer with 4 layers,
 1009 4 attention heads, 256 latent dimensions and a feed-forward network with 512 dimensions. Also, in
 1010 contrast to the previous architecture we don’t use a separate prediction head.

1012 To allow for parallel processing at each position x without leaking information about the cor-
 1013 responding y in a model without bottleneck, we augment a standard Transformer architec-
 1014 ture in the following manner. It considers two sets of tokens, namely (a) D in the format
 1015 $[0, 0], [x_1, y_1], [x_2, y_2], \dots, [x_n, y_n]$ (where x_i and y_i are concatenated for each token), and (b) X
 1016 in the format $[x_1], [x_2], \dots, [x_n]$ (where each token only has x information). Note that $[\cdot]$ describes a
 1017 token, and the first token in D represents an empty context.

1018 Each layer of this Transformer performs the following attention procedures:

$$1019 X^{(l)} = \text{Attention} \left(\text{Query} = X^{(l-1)}, \text{Key} = D^{(l-1)}, \text{Value} = D^{(l-1)}, \text{Mask} = \mathcal{M}^X \right) \quad (24)$$

$$1021 D^{(l)} = \text{Attention} \left(\text{Query} = D^{(l-1)}, \text{Key} = D^{(l-1)}, \text{Value} = D^{(l-1)}, \text{Mask} = \mathcal{M}^D \right) \quad (25)$$

1023 where \mathcal{M}^X ensures that $X_t^{(l-1)}$ can only attend to $D_{1:t-1}^{(l-1)}$ and \mathcal{M}^D ensures that $D_t^{(l-1)}$ can only
 1024 attend to $D_{1:t}^{(l-1)}$. Both $X^{(l)}$ and $D^{(l)}$ go through a residual feed-forward network after the attention
 1025 operations.

Note that the above operation achieves two distinct properties: (a) it prevents the token $[x_t]$ from accessing information about y_t while allowing access to all $x_{1:t-1}$ and $y_{1:t-1}$ in making the corresponding prediction, and (b) akin to standard Transformers the $[x_t, y_t]$ token can attend to $x_{1:t}$ and $y_{1:t}$.

Mamba. We experiment with two state-space model (SSM) architectures, Mamba 1 and Mamba 2, both composed of 4 layers, 256 latent dimensions, state dimensions 8, and local convolution dimension of 4. Additionally, each layer includes a gated MLP with 256 latent dimensions. Similar to the Transformer with bottleneck, the prediction model is a 5-layer MLP with RELU activations and 256 latent dimensions is used as a separate prediction head.

The SSM takes a dataset D as input in the format $[x_1, y_1], [x_2, y_2], \dots, [x_n, y_n]$ (where x_i and y_i are concatenated and each $[\cdot]$ is a token). For each context of size $t - 1$, we compute the $T_\phi(D_{1:t-1})$ which is a vector that represents the parameters of the output model obtained after processing the first $t - 1$ data points. After computing $T_\phi(D_{1:t-1})$, we concatenate it with x_t (i.e., $[T_\phi(D_{1:t-1}), x_t]$) and pass this combined input to an MLP prediction head to predict the next y -token.

E.2 META-TRAINING AND EVALUATION SETUP

In this section, we outline the complete set of hyperparameters and configurations used across different training objectives and model architectures in our experiments.

In-context learner (prequential and train-risk). We trained both the Transformer-based meta-learners (with and without bottleneck) for 50 epochs and the Mamba-based meta-learners for 120 epochs. All results were averaged across 5 different random seeds to mitigate the effect of randomness in the pipeline. The training was conducted on a meta-dataset consisting of 10,000 tasks, each with 1,000 data points that serve as context. We used the Adam optimizer (Kingma & Ba, 2017) with a learning rate of $\eta = 0.0001$ and a batch size of 256, without any early stopping. After meta-training, we evaluated the learners on a distinct meta-dataset of 100 tasks, each with 1,000 data points.

Gradient based learner. Since gradient-based learner are off-the-shelf learning algorithms which don't require meta-training. The prediction model used is a 5-layers MLP with RELU activations and latent dimensions of 64 or 256 depending on the complexity of the task. We used a meta-dataset of 10000 tasks (with 2000 data points each) split into training (80%) and validation (20%). At each step of prequential coding, we train and evaluate a model by randomly sampling a dataset of fixed size across each of the tasks, starting from 20 to 2000 datapoints. We used an early stopping criteria with minimum loss delta of 0.001 and patience of 10 epochs to avoid overfitting. On each of them, the prediction model was fit using the Adam optimizer (Kingma & Ba, 2017) with a learning rate of $\eta = 0.0001$ and a batch size of 64. All results were averaged across 15 different random seeds.

Regularization techniques. Regularization techniques are widely used for gradient-based learners to prevent over-fitted solutions. In this experiment we fit prediction models considering different regularization techniques, namely early-stopping combined with validation data, and weight-decay (L2 regularization). The results are presented in Figure E.1. Experiments with early-stopping halt training when the validation loss does not decrease by more than $1e - 4$ over 10 consecutive steps. Experiments with weight-decay consider a regularization parameter $\lambda \in \{0.05, 0.005\}$ and were trained for 1000 epochs. The prediction models used are 5-layers MLPs with RELU activations and latent dimensions of 64. The different prediction models were fit using an Adam optimizer (Kingma & Ba, 2017) with a learning rate of $\eta = 0.0001$ and a batch size of 64. All results were averaged across 15 different random seeds.

E.3 PRETRAINED LLM ON MASTERMIND

As described in Section 3.4, we evaluate the performance of a pretrained LLM on the Mastermind task using one of the latest OpenAI models GPT-4 (i.e., gpt-4o). To query the model, we used the OpenAI API with a temperature of 0, ensuring that the outputs are deterministic. Along with the responses, we also obtained the log probabilities using the API for calculating the prediction error

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

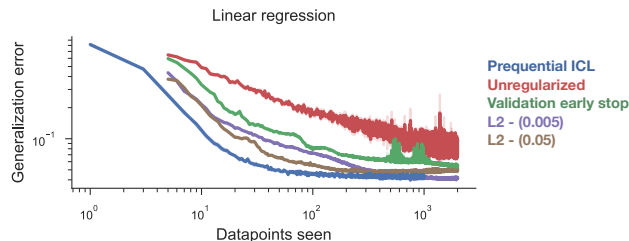


Figure E.1: **Experimental results comparing different regularization techniques.** Figure show average prequential coding curves obtained using both unregularized and regularized Adam optimizers on a linear regression task. Regularized learners exhibit better compression rate (i.e. lower PCL), which implies a stronger incentive toward simple models according to our theory. This experiment confirms the claim that regularization techniques serve as indirect Occam’s aligned methods to learn simple models. Analogous to the meta-learning setting, PCL could be minimized with respect to the hyperparameters of the regularization technique.

with respect to each query. This was possible using `logprobs` (boolean) and `top_k_logprobs` (integer) attributes in the API that returns log probabilities for each token in the response and the k tokens with the top log probabilities corresponding to each token in response. By using a structured prompting technique and a retry mechanism (up to 10 retries in case of failure to adhere to the required output format), we were able to consistently obtain appropriate responses to our queries. An example prompt, which includes the task description, context, and the query, is provided below. To calculate the prequential code length, we iteratively query novel examples with an increasing number of in-context examples and obtain the prediction errors. This process emulates the prequential ICL objective.

Example Prompt

I have a secret code in mind. It’s a 8-digit code with each digit ranging between 0 and 5. I’ll give you a couple example guesses, and for each guess I’ll tell you two numbers:

- First number: the number of correct correct digits at their correct position. - Second number: the number of correct digits, which aren’t necessarily in the correct position.

Here’s a demo to show you what a guess and response would look like. Imagine my secret code was:

0 5 2 1 3 4 2 4

And imagine the guess I presented you was:

0 2 1 1 0 2 0 4

Then, the response would be:

3 5

The response is the way it is because the first, forth and last digit were in the correct place (first response number is therefore 3) and additionally the second and sixth digit were in the guess but at the wrong position (second response number is therefore 5).

The game is about to start. I’ll present you with a series of guesses and their responses. Finally, I will present you with a new guess, and you’ll have to predict the correct response. Make sure your response is formatted the same

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151

way as in the examples (i.e., with 2 digits between 0-8, separated by a space). Let's begin.

Guess: 4 2 1 3 4 0 0 5
Response: 3 7

Guess: 1 1 4 3 5 5 0 1
Response: 2 5

Guess: 3 0 2 2 0 5 3 4
Response: 2 6

Guess: 0 2 5 0 4 2 0 1
Response: 1 5

Guess: 4 1 3 2 5 4 2 3
Response: ? ?

What do you think the response is for this final guess? Make sure to reply with just 2 digits between 0-8, separated by a single space character.

1152
1153
1154
1155
1156
1157

E.4 HIDDEN MARKOV MODEL EXPERIMENT

1158
1159
1160
1161
1162
1163
1164
1165
1166
1167

A prominent theory for why ICL emerges from the next-token prediction objective of LLMs is that sequences $x_{1:n}$ in the pre-training dataset (e.g. large corpuses of text) can be interpreted as implicitly being sampled from a latent variable generative model $Q(x_{1:n} | \tau)$ where τ are some abstract *concepts* underlying samples (Chan et al., 2022; Xie et al., 2022). τ can range from abstract *style* attributes in natural language (Xie et al., 2022) to *task parameters* such as the teacher weight matrix in linear regression ICL task (Von Oswald et al., 2023a); the important part is that some latent variables can be inferred from the context and subsequently aid prediction. ICL would then emerge as the ability of performing implicit Bayesian inference (i.e. learn from the context) in order to predict x_t :

1168
1169
1170

$$Q(x_t | x_{<t}) = \sum_{\tau} \underbrace{Q(x_t | x_{<t}, \tau)}_{\text{Condition on the latent}} \underbrace{Q(\tau | x_{<t})}_{\text{Infer latent}} \tag{26}$$

1171
1172
1173

We propose to leverage this conceptual framework to devise a novel generation procedure for synthetic LLM pre-training dataset. The general idea is to design a family of sequence models $Q_{\tau}(x_{1:n})$ parameterized by task latents τ , leading to the latent variable generative distribution

1174
1175

$$Q(x_{1:n} | \tau) = Q_{\tau}(x_{1:n}).$$

1176
1177
1178
1179
1180

Specifically, we use hidden markov models (HMMs) as the sequences models, and we parameterize the HMMs $Q_{\tau}(x_{1:n})$ with parameters $f_{\xi}(\tau) = \psi_{\tau}$. We use this function f to introduce hyperparameters ξ which define the whole family of sequence models; i.e. the dataset. Below, we define in details a specific *ad-hoc* function $f_{\xi}(\tau)$ which generates a family of HMM where each member share non-trivial structure.

1181
1182

E.4.1 DETAILED DESCRIPTION OF THE GENERATIVE PROCESS

1183
1184
1185
1186
1187

A HMM defines a probability distribution over sequences of *observations* $x_i \in \mathcal{X}$ with a discrete-time probabilistic process over *hidden states* $z_i \in \mathcal{Z}$ paired with a mapping $\mathcal{Z} \rightarrow \mathcal{X}$. Both \mathcal{X} and \mathcal{Z} are discrete sets. The hidden process is defined by an initial state distribution $\pi(z)$ and a transition matrix $A \in \mathbb{R}^{|\mathcal{Z}| \times |\mathcal{Z}|}$ such that

$$Q(z_i | z_j) = A_{ji}$$

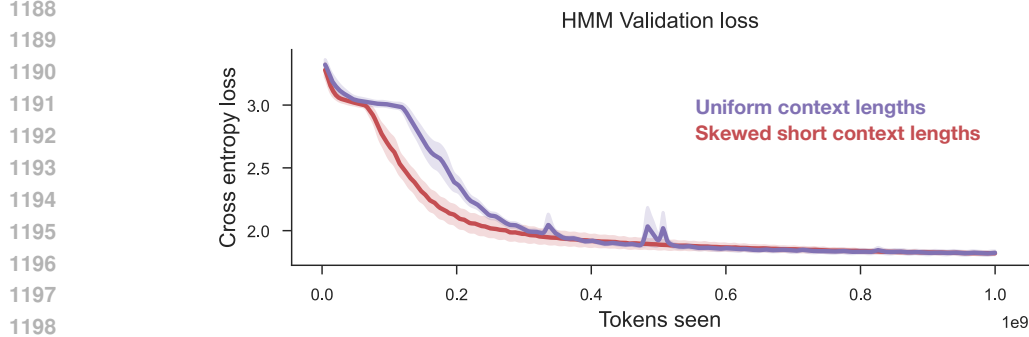


Figure E.2: **Validation loss as a function of the number of tokens seen during training.** The curve is averaged over 5 different datasets (seeds). We can see that the models trained on sequences with shorter length converge faster.

Lastly, the mapping between states and observations is governed by the emission matrix $B \in \mathbb{R}^{|\mathcal{Z}| \times |\mathcal{X}|}$ such that

$$Q(x_j | z_i) = B_{ji}$$

In the rest of the section, we will explicitly define how $f_{\xi}(\tau)$ generates $\psi_{\tau} = (\pi^{\tau}, A^{\tau}, b^{\tau})$. We first give a high level description.

The *hyper-parameters* ξ will define a number of building blocks which will be used to create the transition and emission matrix of all HMMs. Then τ will specify a specific way to combine and manipulate these building blocks to instantiate a specific HMM Q_{τ} . For the transition matrix A^{τ} , the building blocks are pre-defined cycles; which are combined, flipped and accelerated based on τ . For the emission matrix B^{τ} , the building blocks are groups of sub-emission matrices which each only affect a subset of $|\mathcal{X}|$; which are combined and possibly internal shifted based on τ . Overall, we will have

$$\xi = (\text{N_BASE_CYCLES}, \text{N_BASE_SPEEDS}, \text{N_CYCLE_FAMILIES}, \\ \text{N_GROUP_PER_FAMILY}, \text{N_FAMILY_SPEEDS}, \text{N_EMISSION_GROUPS}, \\ \text{N_EMISSION_PER_GROUP}, \text{N_EMISSION_SHIFT})$$

and

$$\tau = (\text{BASE_ID}, \text{BASE_SPEED}, \text{FAMILIES_IDS}, \\ \text{FAMILIES_SPEED}, \text{EMISSION_IDS}, \text{EMISSION_SHIFT})$$

We will refer to the dimensions of ξ , τ as ξ_i , τ_i to avoid clutter and discuss further details below.

Transition matrix A^{τ} . We define a cycle as sequence of hidden states $\mathbf{c} = (c_0, \dots, c_{|c|-1})$, $c_i \in \mathcal{Z}$, and the following manipulation functions

$$\text{DIR}(\mathbf{c}, k) = \begin{cases} (c_0, c_{|c|-1}, \dots, c_1) & \text{if } k = 1 \\ \mathbf{c} & \text{otherwise.} \end{cases}$$

$$\text{SPEED}(\mathbf{c}, k) = (c_0, c_{k \pmod{|c|}}, c_{2k \pmod{|c|}}, \dots)$$

In words, $\text{SPEED}(\mathbf{c}, k)$ changes the speed at which the cycle is traversed and $\text{DIR}(\mathbf{c}, k)$ change its direction. We finally define the transition matrix $\mathcal{T}(\mathbf{c})$ associated with cycle \mathbf{c} such that

$$\mathcal{T}(\mathbf{c})_{ij} = \begin{cases} 1 & \text{if } \exists k < n \text{ s.t. } (i, j) = (c_k, c_{k+1 \pmod{n}}) \\ 0 & \text{otherwise.} \end{cases}$$

Initially, we randomly generate ξ_0 base cycles \mathbf{b}_i which go through all states z_i . Further, we initialize ξ_2 families of ξ_3 groups of cycles \mathbf{g}_j^i , $i \in [\xi_1]$, $j \in [\xi_2]$. Each HMM's transition matrix is then built from these "building blocks" cycles. Specifically,

$$A^{\tau} = \mathcal{T}(\text{SPEED}(\text{DIR}(\mathbf{b}_{\tau_0}, \tau_1), \tau_2)) + \sum_{i=1}^{\xi_2} \tau_{4,i} \sum_{j=1}^{\xi_3} \mathcal{T}(\text{SPEED}(\text{DIR}(\mathbf{g}_j^i, \tau_5), \tau_6))$$

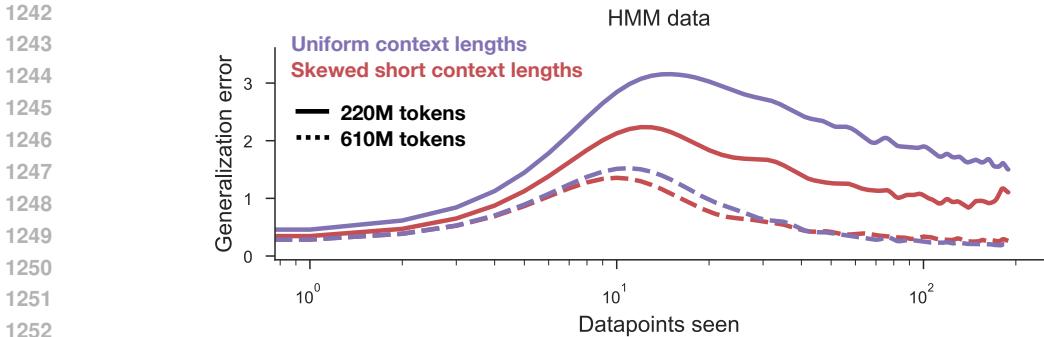


Figure E.3: **Prequential code curves at different stages of training** Reproduction of Figure 3b but with the prequential curve at 610M tokens also. At this point, the models trained with uniform context length have essentially the same performance as the ones trained with smaller context lengths.

In words, each transition matrix is made of a) one of ξ_0 base cycle, possibly sped up and flipped and b) ξ_2 groups of smaller cycles (each from a pool of ξ_3 groups), possibly sped up and flipped. The number of possible speeds for the base cycle is defined by ξ_1 . For the cycle families, it is defined by ξ_4

Emission matrix B^τ . We separate the states $z \in \mathcal{Z}$ in ξ_5 groups $h_i \subset \mathcal{Z}$ and for each group we initialize ξ_6 sub-emission matrices $H_j^i \in \mathbb{R}^{|h_i| \times |\mathcal{Z}|}$. Then, we define the manipulation function $\text{SHIFT}(H, k)$ which applies a circular shift of k to the indices of the matrix. Finally, we have

$$B^\tau = \sum_{i=1}^{\xi_5} \text{SHIFT}(B_{\tau,i}^i, \tau_8)$$

In words, each emission matrix is made of ξ_5 possibly overlapping sub-emission matrix, each picked from a pool of ξ_6 unique ones. The number of possible shifts is ξ_7 .

Initial distribution. We always use the uniform distribution.

E.4.2 HMM HYPER-PARAMETERS

For experiments in this paper, we use $|\mathcal{X}| = 50$ and $|\mathcal{Z}| = 20$. The hyper-parameters of f, ξ , are given in Table E.1. This results in a total of 512 different transition matrices and 24 different emission matrices, for a total of 12,228 different HMMs. We show results averaged from 5 different seed.

N_BASE_CYCLES (ξ_0)	4
N_BASE_SPEEDS (ξ_1)	2
N_CYCLE_FAMILIES (ξ_2)	3
N_GROUP_PER_FAMILY (ξ_3)	2
N_FAMILY_SPEEDS (ξ_4)	2
N_EMISSION_GROUPS (ξ_5)	3
N_EMISSION_PER_GROUP (ξ_6)	2
N_EMISSION_SHIFT (ξ_7)	3

Table E.1: **HMM dataset hyper-parameters**

E.4.3 TRAINING

We hold out 1,000 HMMs for validation and train on the 11,228 others. Training consists on next-token prediction with a cross-entropy loss, using sequences coming from the training HMMs. Specifically, each epochs consists of one sequence sampled from each training HMM. Every epochs, the sequence sampled from a given HMM is different (using a different seed). As such, the model most likely never sees the same sequence twice. We evaluate on sequences from the 1,000 held-out

HMMs. Finally, we use Transformers with 6 layers, 8 heads and embedding dimension of 512. We use a batch size of 512 and a learning rate of 0.001 with Adam.

E.4.4 EVALUATION

To obtain the curve in Figure 3b, we compute the KL divergence between the next-token distribution of trained models to the ground truth which we can compute explicitly with Equation (26):

$$KL[p_{\text{model}}(x_t | x_{<t}), p_{\text{true}}(x_t | x_{<t})] \quad (27)$$

We can compute Equation (26) explicitly because HMMs afford very efficient and parallelizable inference through the forward algorithm. Also, we observe that this "backward" KL divergence is simply a better version of the cross-entropy loss, used to train the model. Indeed, in the cross-entropy loss, $p_{\text{true}}(x_t | x_{<t})$ is replaced by a delta-dirac distribution on the observed x . While training on it also ends up minimizing Equation (27), it is not the best evaluation metric. Indeed, cross-entropy doesn't take into account the stochasticity of the ground-truth, while Equation (27) does.

Note that the non-monotonicity of the KL prequential coding curve is a consequence of using the above KL . Indeed, when very few datapoints have been seen, the model can learn memorise the marginal probability $p_{\text{true}}(x_t | x_{<t})$ quite easily, bypassing the to perform ICL. This doesn't show when displaying cross-entropy because $p_{\text{true}}(x_t | x_{<t})$ has often very high entropy for small t .

E.4.5 TRAINING WITH SHORTER SEQUENCES

When training sequence models like LLMs, the typical approach is to fill the maximum context window of the model with sequences, possibly concatenating multiple ones. This ensures that every batch contains as much tokens—i.e. training signal—as possible. However, because of this, most tokens seen during training are preceded by a lot of tokens: putting more pressure on correctly predicting late tokens than rapidly adapting with small amount of context. According to our theory, this leads to more complex models, possibly worse at generalizing.

Based on this reasoning, we propose a simple way to bias the meta-learner towards simpler models: training on sequences with random context length, typically much shorter than the maximal one. We show the efficacy of our method using our HMM dataset: models trained with **uniform context length** (i.e. all sequences have maximal length) need less tokens to arrive at simple models than the ones trained with **skewed short context lengths** (i.e. sequences of random lengths), as shown in Figure E.2 and Figure E.3. However, there are diminishing returns: with enough data training on long context catches up. Exploring this approach on large-scale language modeling is an interested future work.