
Transformers for Supervised Online Continual Learning

Jorg Bornschein¹ Yazhe Li¹ Amal Rannen-Triki¹

Abstract

Transformers have become the dominant architecture for sequence modeling tasks such as natural language processing or audio processing, and they are now even considered for tasks that are not naturally sequential such as image classification. Their ability to attend to and to process a set of tokens as context enables them to develop in-context few-shot learning abilities. However, their potential for online continual learning remains relatively unexplored. In online continual learning, a model must adapt to a non-stationary stream of data, minimizing the cumulative next-step prediction loss. We focus on the supervised online continual learning setting, where we learn a predictor $x_t \rightarrow y_t$ for a sequence of examples (x_t, y_t) . Inspired by the in-context learning capabilities of transformers and their connection to meta-learning, we propose a method that leverages these strengths for online continual learning. Our approach explicitly conditions a transformer on recent observations, while at the same time online training it with stochastic gradient descent, following the procedure introduced with Transformer-XL. We incorporate replay to maintain the benefits of multi-epoch training while adhering to the sequential protocol. We hypothesize that this combination enables fast adaptation through in-context learning and sustained long-term improvement via parametric learning. Our method demonstrates significant improvements over previous state-of-the-art results on CLOC, a challenging large-scale real-world benchmark for image geo-localization.

1. Introduction

We consider the problem of online supervised learning with modern deep neural networks, sometimes also referred to as online continual learning. Given a sequence

of observations x_t (e.g. images) and corresponding prediction targets y_t (e.g. labels), we aim to find a predictive model $p(y_t|x_t, \mathcal{D}_{<t})$ that minimizes the cumulative log-loss $\mathcal{L} = \sum_{t=1}^T -\log p(y_t|x_t, \mathcal{D}_{<t})$ ¹. We use $\mathcal{D}_{<t} = \{(x_1, y_1), \dots, (x_{t-1}, y_{t-1})\}$ to denote the sequence of observations and prediction targets up to (including) $t-1$. This formulation is quite general and allows for a wide variety of approaches to solve the problem. It is also theoretically well motivated and encourages us to find models with many desirable properties such as fast adaptation, sample-efficient learning, and maintaining plasticity. Furthermore, minimizing the cumulative next-step log-loss is a form of the Minimum Description Length (MDL) principle. Appendix A provides a deeper explanation of this connection and highlights some of its properties and benefits.

Traditionally, in online learning, the conditioning on previous examples comes from the learned model parameters θ and the model is otherwise independent across examples. Thus, the general objective is reformulated as $\mathcal{L} = \sum_{t=1}^T -\log p(y_t|x_t, \theta_{t-1})$, where the model parameters are a function of the previously observed data $\theta_{t-1} = \hat{\theta}(\mathcal{D}_{<t})$. Often, constraints are placed on how the parameter estimator can access previously observed data. In the streaming online-learning scenario, for instance, only the most recent example is used to update the model parameters: $\theta_t = \text{update-step}(\theta_{t-1}, x_t, y_t)$. For online gradient descent this simplifies to $\theta_t = \theta_{t-1} + \alpha_t \nabla_{\theta} \log p(y_t|x_t, \theta_{t-1})$. Within the deep-learning community, it is common to use experience replay or core-sets to perform multiple gradient steps on buffered examples to update the parameters (Chaudhry et al., 2019). Recent work argues that storing all previous data $\mathcal{D}_{<t}$ is often not a major technical challenge and that instead of limiting the size of the replay-buffer, limits on the compute requirements should be considered (Prabhu et al., 2023). To that end Bornschein et al. (2022) propose to replay the sequence in-order from long-term storage.

Another line of work considers memory- or retrieval inspired approaches. These methods often rely on pre-

¹Google Deepmind. Correspondence to: Jorg Bornschein <bornschein@google.com>.

¹Other loss-functions are feasible, for example the cumulative zero-one loss which would amount to minimizing the total number of prediction mistakes on the sequence.

trained feature extractors and similarity measures to retrieve relevant information from previously observed examples. Prabhu et al. (2023) for example apply approximate online k-nearest neighbor (kNN) to CLOC, a large-scale and real-world online geo-localization sequence. For each new observation x_t they retrieve a similar example (x_i, y_i) with $i < t$, and make a label prediction based on y_i .

In this work we propose a hybrid approach: We use a transformer model $p(y_{t+1}|x_{t+1}, \mathcal{D}_{(t-C)\dots t}, \theta)$ that is explicitly conditioned on the C most recent observations while at the same time training the model with online gradient descent on the chronologically ordered sequence $\{(x_1, y_1), \dots, (x_T, y_T)\}$. We use replay-streams (Bornschein et al., 2022) as a simple and effective way to implement chronological replay of previously observed examples – essentially re-introducing the benefits of multi-epoch training while strictly adhering to the online next-step evaluation protocol.

The motivation behind this approach is that gradient based learning has proven to be effective for stationary and slowly changing data. It promises steady learning progress even after seeing millions of examples. Explicitly conditioned transformer models on the other hand have demonstrated excellent performance when learning from short-term correlations within the window of the C most recent observations. Trained on data that exhibits meta-learning characteristics, transformers develop impressive in-context few-shot learning abilities. We conjecture that our approach can combine the complimentary strengths of these two approaches: In-context learning to enable rapid adaptation to re-occurring but potentially sudden changes in the data sequence; and parametric learning with SGD to enable sustained progress over long sequence lengths, far beyond the C most recent observations.

Contributions.

- We empirically investigate the behaviour and properties of online trained transformer models when applied to non-stationary supervised image classification problems.
- We propose the privileged information (pi) transformer, a variant of the standard transformer architecture that is tailored to sequential supervised prediction problems. We show that this architecture often achieves faster, more stable, and better prediction performance than the standard architecture while at the same time being more compute efficient.

Our primary evaluation is on CLOC, a large scale and real-world online continual learning benchmark. We show that both the standard transformer architecture, and espe-

cially the pi-transformer, achieve excellent prediction performance and far exceed previously reported prediction accuracies.

2. Architecture and Method

We experiment with two different model architectures:

- 1) 2-token approach: We use a plain decoder-only, causal transformer and present each example (x_t, y_t) as two consecutive tokens. The underlying transformer processes a sequence of length $2T$. We ignore the prediction loss on the x_t tokens: the model is trained to predict the y_t tokens only.
- 2) Privileged information (pi) transformer: For each example we present x_t as input token. But modify the basic transformer block such that each token x_t also receives additional privileged information y_t . The architecture is setup so that the prediction at time t can not access the information in y_t , however can access $y_{<t}$. Specifically, projections of y_i are added to the attention keys and values only. We also ensure that the attention mask has a zero diagonal, thus prohibiting the attention heads at time t to access values from time t . See B for the exact formulation.

For both architectures the input x_t is in our case the feature vector corresponding to an input image. We use ConvNets, ResNets and Vision Transformers (ViT) as feature extractors. We pick these because they are popular architectures in related works; we did not perform a systematic search for best performing feature extractors.

We perform Transformer-XL style (Dai et al., 2019) online training: The forward- and backward pass operates on a relatively small number of sequential tokens, typically $S \approx 100$. S is the chunk-length and with online SGD (without replay) we perform T/S gradient steps when fitting a sequence of length T . The self-attention heads, however, operate on a larger buffer of attendable keys and values. We store the keys and values for the C most recent tokens in a KV-cache ringbuffer. C is the (sliding) attention window size of the architecture. We use Multi Query Attention (MQA, Shazeer (2019)), and thus only store one key- and one value vector per transformer block and token.

With a 8-block deep transformer architecture and $dim_k = 128$ we store $8 \cdot 128 \cdot 1024 \approx 1\text{Mi}$ floats to keep $C=1024$ recent tokens accessible. We use MQA and MLP layers in a parallel configuration, following some recent models like CodeGen, PaLM, and GPT-NeoX.

Listing 1 Replay-streams training: `new_data_reader()` creates a deterministic data-reader that yields successive examples. `gradient_step(...)` performs a training step and returns the log-loss, state of the sequence model (KV-cache) and the updated parameters. In practice we perform the algorithm with chunks of S successive examples instead of individual ones. The reset probability to approximate uniform replay for chunk-size S is $\frac{S}{t}$ instead of $\frac{1}{t}$.

```

1 cumulative_nll = 0.
2 data_readers = [
3     new_data_reader() for _ in range(num_streams)]
4 kv_caches = [
5     empty_kv_cache() for _ in range(num_streams)]
6
7 # Iterate over all data
8 for pos in range(num_examples):
9     # Read new data, perform gradient step
10    # and accumulate NLL.
11    data = next(data_readers[0])
12    nll, kv_caches[0], model_params = gradient_step(
13        data, kv_caches[0], model_params)
14    cumulative_nll += nll
15
16 # Read data from remaining streams for replay.
17 for s in range(1, num_streams):
18     data = next(data_readers[s])
19     _, kv_caches[s], model_params = gradient_step(
20         data, kv_caches[s], model_params)
21
22 # Reset each replay stream with probability 1/pos
23 if bernoulli(1/pos):
24     data_readers[s] = new_data_reader()
25     kv_caches[s] = empty_kv_cache()

```

2.1. Replay-streams: In-order Replay

We use the approach from Bornschein et al. (2022) and introduce a hyper-parameter *num-streams* (E) that controls the replay factor, and therefore the number of effective epochs: We maintain E separate sequence states: for each its own KV-cache and its own deterministic, sequential data-reader that yields the examples in deterministic order. With each turn, all streams are advanced and perform a gradient step on chunk size S many examples. The first stream will proceed steadily through the data-sequence from 1 to T and suffer the (cumulative) log-loss and prediction error which we report. The remaining $E-1$ streams replay data that has been previously seen by the first one. By stochastically resetting some of the $E-1$ replay data-readers and their associated KV-caches back to position 0 we achieve, in expectation, uniform replay of the past without requiring large in-memory replay-buffers. See algorithm 1 for details.

Due to the approximate uniform replay, our method resembles a particular form of meta- or bi-level learning: At time t , the *outer-learner* optimizes the parameters θ_t to improve the prediction $p(y_{i+1}|x_{i+1}, \mathcal{D}_{i-C\dots i}, \theta_t)$ for all $i \leq t$ simultaneously; while the inner, fast learner (Schlag et al., 2021) is non-parametric and relies purely on in-context learning with C tokens in its attention window.

3. In-context and Meta-learning

Recent work has pointed out that the in-context learning abilities of transformer models can be understood and enhanced when taking a meta-learning perspective:

Different works have reported several relevant results while studying online learning and non-i.i.d. data and their relation to meta-learning and in-context learning and transformers (Lee et al., 2023). Ortega et al. (2019) shows that meta-learning of memory-based models leads to near Bayes-optimal solutions. The authors also highlight that meta-learning can occur spontaneously through online learning when the model capacity is bounded and the data is produced by a single generation process. They however warn against the challenges that this type of emergent meta-learning can lead to, due to the lack of control over it. Safety problems can arise, and need to be taken into consideration. Mikulik et al. (2020) and Genewein et al. (2023) reinforce the previous results both theoretically and empirically: The former shows that not only meta-learned models with memory converge to Bayes optimal agents, they are also indistinguishable from the perspective of predictions and computations. The latter shows that the meta-learning of memory-based models on non-stationary distributions leads to Bayes-optimal predictors. Chan et al. (2022) and Singh et al. (2023) take a difference stance and study the emergence of in-context learning behavior in transformers through the lens of data distribution. Chan et al. (2022) highlight the importance of both the architecture (i.e. the use of attention mechanisms) and the data distribution for this behavior to appear. The authors show that the few shot learning capabilities of transformers arise when data is ‘bursty’ and appear in clusters rather than uniformly across the training data, when data is highly skewed and the classes follow a Zipfian distribution, and when the meaning of labels is dynamic and context dependent. We note that all these characteristics make the transformer training divert from the classical i.i.d. setting, and gets closer to a non-stationary online setting, which connects to the results in Ortega et al. (2019). These papers also distinguishes between in-context learning and in-weights learning. Singh et al. (2023) highlight that the in-context learning behavior is transient, and that in-weight learning reemerges if the model is overtrained.

We note that none of these papers studies transformers in the online learning scenario, but they describe a perspective on in-context learning with transformers that inspired our approach here. In particular, we bridge the gap between in-context and in-weight learning, highlighting cases where they can work in synergy. We think that such synergy could be instrumental in solving the prospective learning problem discussed in De Silva et al. (2023).

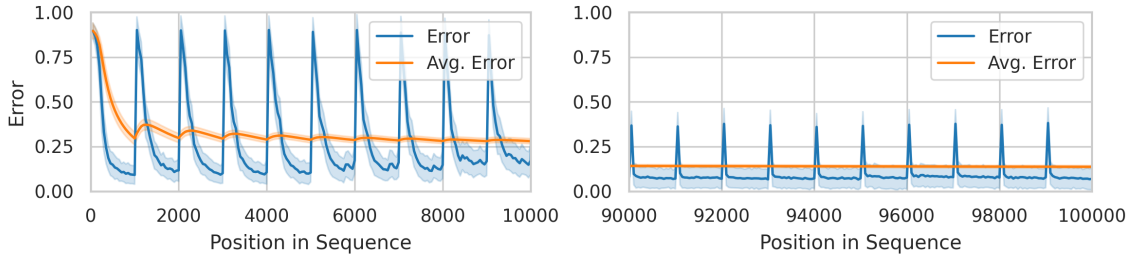


Figure 1: Instantaneous and averaged prediction performance for the first and last 10 tasks of Split-EMNIST: Image-to-label mappings are constant within each task, but randomly reassigned at task boundaries (averaged over 200 data generating random seeds).

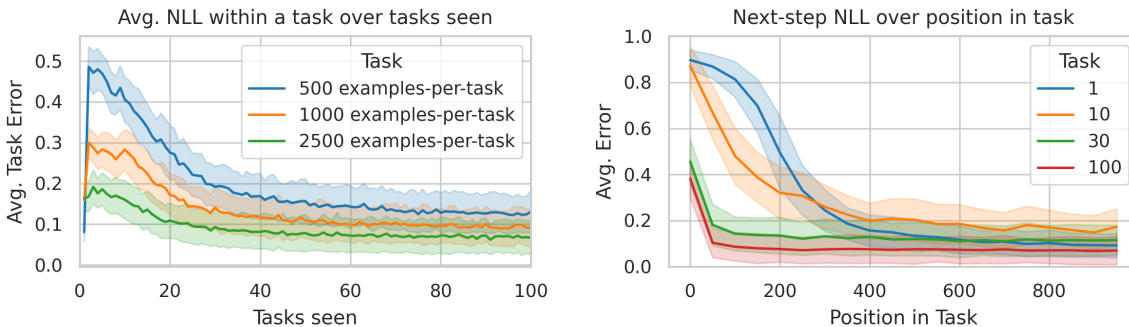


Figure 2: Alternative visualization of the Split-EMNIST experiments from Fig. 1: **Left:** Average performance per task shows strong forward-transfer after struggling during the first 10 to 20 tasks. **Right:** Detailed look at the within-task performance for the scenario with 1000 examples per task. The model is a strong few-shot learner after seeing 30 tasks and further improves until at least task 100

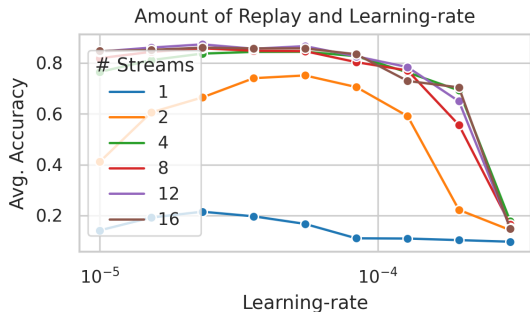


Figure 3: Average accuracy at the end of the sequence for different amount of replay (epochs) and learning-rates on Split-EMNIST (1000 examples per task, 100 tasks).

4. Experiments on Synthetic Toy Data

Piece-wise stationary Split-EMNIST. We create synthetic, piece-wise stationary sequences from underlying labeled data sets. The sequences are composed of 100 tasks. Each task is a 10-way classification problem and is generated by first randomly selecting 10 classes from the underlying dataset, assigning them to the 10 observed label-

classes and randomly selecting corresponding images until the desired number of examples for this task have been collected. From the perspective of the (task-agnostic) online-learner, a task boundary is a sudden event when the distributions of the observed data x and the associated target labels y change abruptly. Even if similar inputs x have been observed before, they are now with high probability assigned to different labels y . We create Split-EMNIST from the balanced EMNIST data set (Cohen et al., 2017), which provides 47 underlying classes.

We use the 6-layer convolutional neural network described by Blier & Ollivier (2018) as a feature extractor, however double the number of channels in all stages. On top of that we use 4 transformer blocks with a backbone width of 256 units. We first study the behaviour of the pi-transformer architecture. We choose a chunk size S of 50, a constant learning rate of $1e-4$ and 8 replay-streams. See Appendix C for details of the model architecture and hyperparameters. Figure 1 shows the next-example prediction performance averaged over 200 random seeds for the data generation. Averaging over many data generating seeds ensures we plot clear trends even though there is high across example- and across task variability. Note however that the *average* accu-

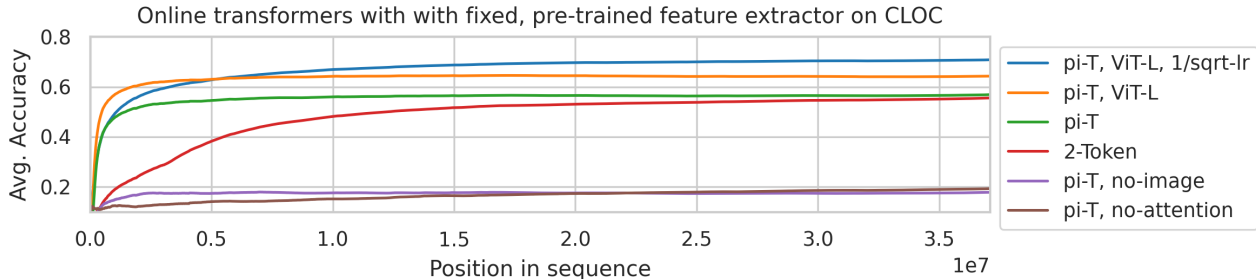


Figure 4: CLOC with pretrained and frozen feature extractors. We show the best performing models (in terms of final avg. accuracy) from the hyper-parameter cube in E.4 We also show pi-Transformer ablations: either without input features x_t , or without attention ($C = 0$).

racy has very low variability, even for a single run. Figure 2 shows a detailed view of the average NLL per task, or per position within a task. We observe that after initially struggling during the first ≈ 20 tasks, the model learns to be a highly efficient few-shot learner that only requires a few examples at the beginning of each task to make accurate predictions. Figure 3 shows the effect of the number of replay streams (epochs) and learning rate on the final result. We observe that without replay, the model usually does not learn to use in-context information effectively to make predictions. With replay, the model is forced to learn one set of parameters θ_t that work at $t + 1$, and at the replay positions $< t$. It is thus strongly encouraged to be an in-context learner.

We repeat the experiments with the 2-token architecture and observe generally the same behaviour. However, the 2-token approach is learning slower in the sense that it requires more samples and more tasks to achieve the same performance.

In Appendix D we presents analogous experiments with sequences based on CIFAR-100. We generally observe the same behaviour, however we use a deeper feature extractor and it requires more replay and more examples to obtain good few shot performance and is more sensitive to hyper-parameter choices. We note that these results confirm the observations previously highlighted in Lesort et al. (2022), but with a different replay strategy.

5. Large scale continuous geo-localization

Our primary evaluation with real-world data is on CLOC, the continual geo-localization sequence introduced by Cai et al. (2021). It consists of $\approx 39M$ chronologically ordered images with their geo-location as a categorical prediction target. Figure 2 in (Cai et al., 2021) shows that the data is strongly non-stationary, and that traditional i.i.d. training results in a held-out accuracy between 10 and 20%. We failed to download or decode around 5% of the images, which leaves us with a sequence of 37,093,769 im-

ages. Different protocols have been used to train and evaluate models. Most notable differences include a) whether pre-trained feature extractors are used; b) whether the feature extractor is fine-tuned.

We obtain substantial improvements over previous state-of-the-art for the two scenarios we consider here: Learning with frozen pre-trained feature extractors and learning the whole model from scratch. In both cases we almost double the archived average accuracy (Table 1). Note that initial results from Cai et al. (2021) considered a filtered versions of the task where some short term correlations are explicitly removed; recent works (Titsias et al., 2023; Bornschein et al., 2022; Prabhu et al., 2023) however focus on the unfiltered average online next-step accuracy. In the following we describe the experiments for the frozen pre-trained feature-extractor and online-learned feature extractor from scratch separately.

5.1. Pre-extracted image features.

Table 1: Results on CLOC: Experience Replay (ER, Cai et al. (2021)), ACM (Prabhu et al., 2023), Replay Streams (Bornschein et al., 2022), Kalman Filter (Titsias et al., 2023).

Method	Pretrained	Finetuning	Avg. Acc.
ER	✓	✓	20%
ACM	✓	-	26%
Replay Streams	✓	-	$\approx 38\%$
Replay Streams	-	✓	$\approx 37\%$
Kalman Filter	✓	-	30%
Kalman Filter	-	✓	37%
Ours (no-image)	-	-	18%
Ours (no-attention)	✓	-	19%
Ours (Superv. ResNet)	✓	-	59%
Ours (MAE ViT-L)	✓	-	70%
Ours	-	✓	67%

The choice of the pre-trained and frozen feature extractor has a significant impact on the achievable performance. We first present the results when using the top-level acti-

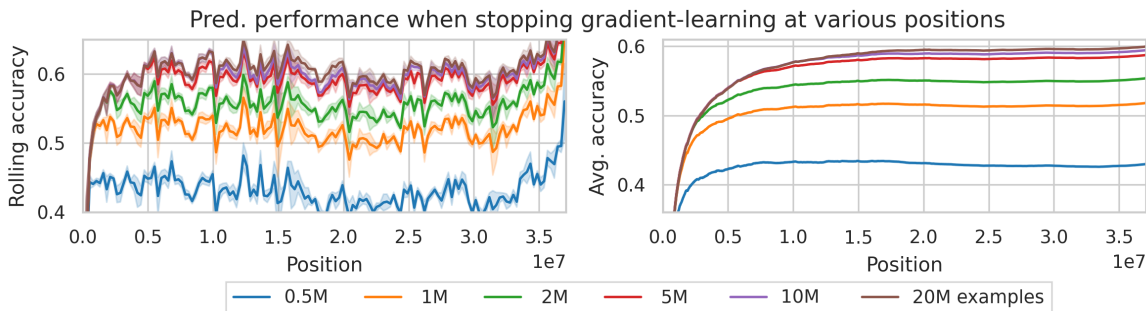


Figure 5: Stopping gradient updates at various positions to investigate the performance of *in-context* conditioning alone. Here for the pi-Transformer on CLOC with a fixed, pre-trained ResNet-50 feature extractor. Gradient based online-learning is most important at the beginning of the sequence, however keeps on contributing even after 10M examples.

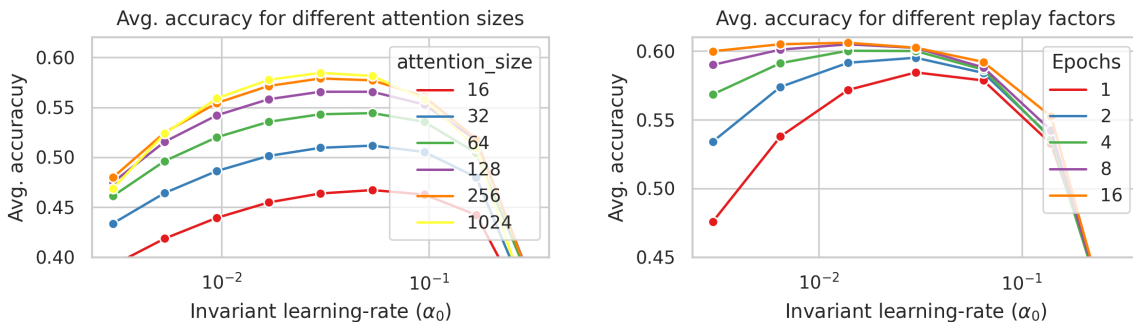


Figure 6: Pi-Transformer on CLOC with a fixed, pre-trained ResNet-50 feature extractor: **Left**) Increasing the attention window up to about 512 examples improves the results. **Right**) Performing more replay gradient steps improves performance and makes the model more robust to lower learning-rates.

variations of a supervised ImageNet trained ResNet-50, just like Prabhu et al. (2023). Consistent with their observation, we also obtain substantially better results when using self-supervised trained feature extractors such as a MAE ViT-L (He et al., 2022). With pre-extracted features and depending on the amount of replay and the size of the transformer it requires between 10^{14} and 10^{17} FLOPs to run CLOC. This corresponds to walltimes between 15 minutes and 12h hours on a single GPU. The base-architecture for all experiments is a 4-block deep transformer with a backbone width of 1024 units. We set the attention window size to 512 and choose a chunk size S of 256 examples with a constant learning rate for the AdamW optimizer. Figure 4 shows typical avg. accuracy curves. We include additional baselines: a) an online transformer that has no access to the image x_t and that makes prediction from the label sequence alone; b) an online transformer without attention, i.e., context size $C=0$ that relies on the current image features alone; and c) an oracle predictor, that predicts the correct location whenever an image with the same location was observed during the previous 100 examples.

We run extensive hyper-parameter sweeps to characterize the model’s performance. For these we first focus on the

pi-transformer architecture. Results for the 2-token approach can be found in the Appendix E.3. We observe that the optimal (constant) learning rate is typically $\alpha \approx \frac{\alpha_0}{D}$, where D is the width of the transformer architecture and $\alpha_0 \approx 3 \cdot 10^{-2}$. This is consistent with the results from Yang et al. (2021). Figure 6 (left) shows the average accuracy at the end of the sequence for various attention window sizes as a function of the learning rate. The plot on the right shows the influence of the number of replay streams (epochs). While more replay generally improves the predictive performance, we obtain already decent results completely without replay. Figure 12 in the Appendix shows the influence of the model size (width) and of the chunk size S on the final performance. The choice of the chunk-size is not critical to obtain good results, however it does correlate with the optimal learning rate, just like batch-sizes do for regular iid. training. Overall, we observe that learning a pi-transformer on top of pre-extracted features is very robust and well behaved under hyper-parameter variations.

Our approach combines in-weight and in-context learning into an online-learning algorithm where these two mechanism are synergistically intertwined. To get a sense of their

relative contributions we run two ablations: First, we train a model with $C=0$, effectively disabling the attention mechanism in the transformer. We obtain an average accuracy of only about 19% – comparable to previously reported results for experience replay (ER, see Figure 4). Secondly, in Figure 5, we train online transformers and disable gradient updates once the learner reaches positions 0.5M, 1M, 2M, 5M, 10M or 20M. From these points onward the weights are frozen and the predictor relies on in-context conditioning alone. As one would expect, parametric learning is most important at the beginning of the sequence, however still contributes to improved results after $> 10M$ datapoints.

We run the same suite of experiments with the 2-token approach. Learning curves often show distinct, non-smooth improvement steps (Figure 8). Olsson et al. (2022) associate such step-improvements with phase-changes in the learning dynamic and with the formation of induction-heads that are crucial for in-context learning. The exact occurrence of these steps is stochastic and sensitive to hyper-parameter choices. The pi-transformer learning curve shows in general more reliable improvements, presumably because the superimposed label information provides an adequate inductive bias where the label association does not have to be learned. Appendix E.3 presents more results for the 2-token approach.

In Figure 9 we directly compare the sensitivity of the different transformer based approaches to changes of the attention-window size C . Figure 10 shows the sensitivity to the number of replay-streams (epochs) E .

Pareto fronts. The total computational cost for fitting the model to the sequence varies substantially depending on the model size, depth, and number of replay streams. To compare approaches fairly, we try to characterize their pareto-front in terms of final avg. accuracy over required FLOPs. We define a hyper-parameter cube that spans a computational cost between around 10^{14} to 10^{17} FLOPs² (see Appendix E.4 for details). We show the results in Figure 7. Overall, with such a long sequence, the final average performance of the 2-token approach and the pi-transformer are very similar. Much larger gains can be obtained by switching to a different feature extractor, or adding learning-rate decay. We note that well performing pi-transformers were often relatively shallow (2 to 4 blocks) and wide, while well performing 2-token models preferred more depth (8 blocks), but became narrower.

5.2. From scratch.

Our second scenario for CLOC is to jointly online-learn the feature extractor from scratch. We consider ResNet style feature-extractors. However, off-the-shelf ResNet archi-

²We only count multiply-and-add operations (MACs)

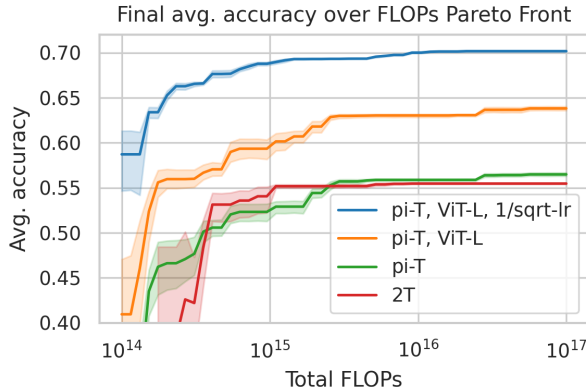


Figure 7: Pareto fronts for online learners with fixed, pre-trained feature extractors on CLOC: Final avg. accuracy at the end of the sequence over total number of FLOPs (MACs). Excluding the cost of the feature extractor. We run a broad hyper-parameter cube with 250 experiments for each method to obtain the pareto fronts.

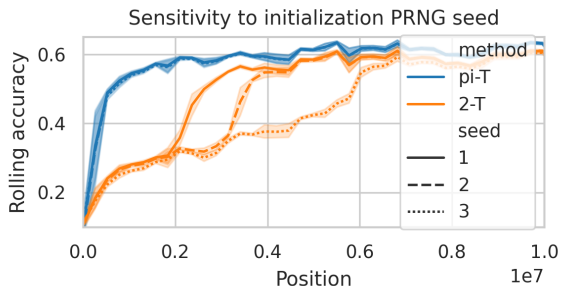


Figure 8: First 10M examples of CLOC with fixed, pre-extracted features: The pi-Transformer shows a smooth accuracy improvements. Accuracies for the 2-token approach exhibit discrete step improvements with their exact location depending on the random initialization seed.

tectures employ batch-norm layers, which introduce non-causal information-leakage between successive examples due to the shared activation-mean within a mini-batch. We therefore take inspiration from Brock et al. (2021) and Liu et al. (2022) and replace the batch-norm with layer-norm instead. The changes are consistent with the recommendations from Lyle et al. (2023). The details can be found in Appendix F. We observe that the the feature extractor and the temporal transformer prefer noticeable different learning rates in order to achieve optimal results: Reasonably sized ResNets prefer a learning rate between $3e-4$ and $1e-3$ while the transformer prefers again $\alpha = \frac{\alpha_0}{D}$ with $\alpha_0 \approx 1e-2$. Therefore, we tune the learning rates for these two model components independently and leave it to future work to find a more elegant solution. We set the feature extractor learning rate to $3e-4$ and choose a transformer backbone width of $D=1024$ with a transformer learning rate of

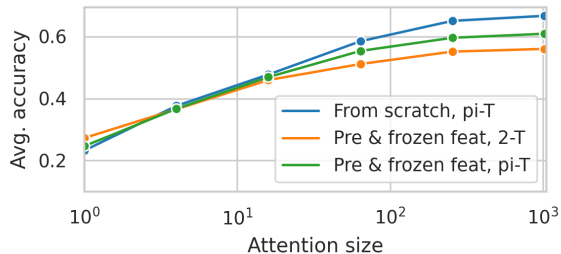


Figure 9: Sensitivity to the size of the attention window C for different models. For each experiment we sweep over 7 different transformer learning rates and report the best run.

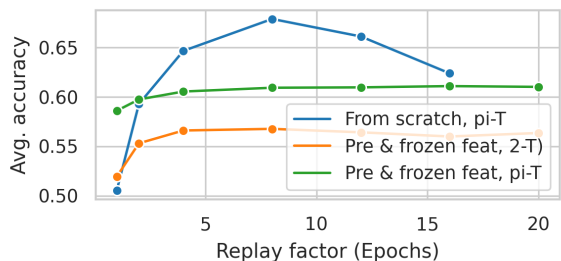


Figure 10: Sensitivity to the number of replay-streams (epochs) E for different models. For each experiment we sweep over 7 different transformer learning rates and report the best run.

$$\alpha = \frac{\alpha_0}{D} = 1e-5.$$

Figure 11 shows the average accuracy of our best performing pi-transformer and 2-token transformer architectures. We include the learning curve for a pi-transformer on pre-extracted ResNet-50 features for comparison. Figures 9 and 10 show the sensitivity to varying the attention size and number of replay streams compared to models with pre-trained and fixed feature extractors: Learned representations can take better advantage of large context windows C . However training is more sensitive to hyperparameters, as exemplified by the decline of predictive performance as a result of increasing E beyond ≈ 10 . We note that re-tuning the feature-extractor learning-rate allows for using more replay with marginal improvements for the final accuracy.

6. Conclusions

In this paper we study transformers in the challenging setting of supervised online continual learning: We combine the transformer (Vaswani et al., 2017) architecture, online-training with KV caching (Dai et al., 2019) and replay streams (Bornschein et al., 2022) into a practical algorithm for online learning. We propose two concrete methods: a 2-token variant where the input features and the label are fed as two separate tokens, and privileged information variant where the basic transformer block is modified to receive

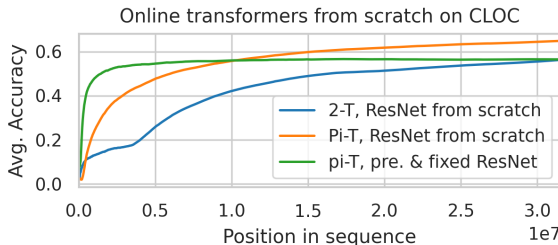


Figure 11: CLOC with ResNet feature extractor from scratch. For comparison we add the learning curve for a pi-Transformer with pretrained and frozen ResNet features from Figure 4.

additional information in the form of a projection of the labels.

We evaluate these variants in two online learning scenarios: a) learning on synthetic, task-agnostic and piecewise stationary sequences and b) on a naturally occurring large-scale sequence with continuous drift. The scenario with synthetic data allows us to study the emergent meta-learning-like behavior in detail: As the learner progresses in the sequence, we first observe a performance drop that is generally associated with *loss of plasticity* issues. However, after encountering more task boundaries, the model learns to be a highly effective few-shot learner that increasingly rapidly performs well on new tasks. In the second scenario, on CLOC real-world data, both the 2-token and the privileged-transformer approaches reach unprecedented levels of predictive performance; almost halving the state-of-the-art average error-rate with our best results. We obtain these results for both pretrained and frozen image features, and also with feature extractors trained online from scratch.

Overall this study sheds light on interesting interactions between in-context and in-weight learning in transformers. It shows scenarios where they work in tandem to improve the performance and the efficiency of online learning methods. With this paper, we take a step towards bridging the gap between these two mechanisms, and further studies are needed to continue in this direction and understand better the synergies between them.

6.1. Limitations and Future Work

There are several areas of improvements that we leave to future works: Training a deep feature extractor and the temporal transformer jointly sometimes requires two different learning rates for optimal results – which leads to expensive and inconvenient hyperparameter searches (Section 5). We also mostly stick to constant learning rates even though initial experiments suggest that better results can be obtained by introducing learning rate schedules. In this work uses

basic and potentially suboptimal pre-trained feature extractors for the images. We leave the study of the impact of more carefully designed embeddings to future works.

Impact Statement

The paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of this work, most of them are however not specific to our contribution and we think it is out of scope to discuss them here. Beyond these universal considerations, online-learning poses some additional opportunities and challenges. Most notably: The predictive performance and characteristics of online-learning systems changes over time. This is in contrast to static models, where only the data-distribution might change and as a result undesirable predictions might be produced. With online learning, it is the model itself that changes over time too. If the same online-model is shared among mutually non-trusting clients, every single one of them, maliciously or unintentionally, might cause severe and harmful mispredictions for the others. We believe that no such online system should be deployed without extensive additional safety assessments. Unfortunately, this precludes harnessing some of the benefits of online adapting systems: continuous tracking and steady improvements in face of a changing world, and improved efficiency of the learning systems.

References

- Blier, L. and Ollivier, Y. The description length of deep learning models. *Advances in Neural Information Processing Systems*, 31, 2018.
- Bornschein, J., Li, Y., and Hutter, M. Sequential learning of neural networks for prequential MDL. In *The Eleventh International Conference on Learning Representations*, 2022.
- Brock, A., De, S., Smith, S. L., and Simonyan, K. High-performance large-scale image recognition without normalization. In *International Conference on Machine Learning*, pp. 1059–1071. PMLR, 2021.
- Cai, Z., Sener, O., and Koltun, V. Online continual learning with natural distribution shifts: An empirical study with visual data. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8281–8290, 2021.
- Chaitin, G. J. On the intelligibility of the universe and the notions of simplicity, complexity and irreducibility. *arXiv preprint math/0210035*, 2002.
- Chan, S., Santoro, A., Lampinen, A., Wang, J., Singh, A., Richemond, P., McClelland, J., and Hill, F. Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891, 2022.
- Chaudhry, A., Rohrbach, M., Elhoseiny, M., Ajanthan, T., Dokania, P., Torr, P., and Ranzato, M. Continual learning with tiny episodic memories. In *Workshop on Multi-Task and Lifelong Reinforcement Learning*, 2019.
- Cohen, G., Afshar, S., Tapson, J., and Schaik, A. V. Emnist: Extending mnist to handwritten letters. *2017 International Joint Conference on Neural Networks (IJCNN)*, 2017. doi: 10.1109/ijcnn.2017.7966217.
- Dai, Z., Yang, Z., Yang, Y., Carbonell, J. G., Le, Q., and Salakhutdinov, R. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2978–2988, 2019.
- De Silva, A., Ramesh, R., Ungar, L., Shuler, M. H., Cowan, N. J., Platt, M., Li, C., Isik, L., Roh, S.-E., Charles, A., et al. Prospective learning: Principled extrapolation to the future. In *Conference on Lifelong Learning Agents*, pp. 347–357. PMLR, 2023.
- Genewein, T., Delétang, G., Ruoss, A., Wenliang, L. K., Catt, E., Dutordoir, V., Grau-Moya, J., Orseau, L., Hutter, M., and Veness, J. Memory-based meta-learning on non-stationary distributions. In *International Conference on Machine Learning*, 2023.
- Grunwald, P. and Vitanyi, P. Shannon information and kolmogorov complexity. October 2004.
- Grünwald, P. D. *The Minimum Description Length Principle*. The MIT Press, Cambridge, 2007.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
- Lee, S., Son, J., and Kim, G. Recasting continual learning as sequence modeling. October 2023.
- Lesort, T., Ostapenko, O., Rodriguez, P., Arefin, M. R., Misra, D., Charlin, L., and Rish, I. Challenging common assumptions about catastrophic forgetting. 2022.
- Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., and Xie, S. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11976–11986, 2022.
- Lyle, C., Zheng, Z., Nikishin, E., Pires, B. A., Pascanu, R., and Dabney, W. Understanding plasticity in neural networks. *International Conference on Machine Learning*, 2023.

- Mikulik, V., Delétang, G., McGrath, T., Genewein, T., Martic, M., Legg, S., and Ortega, P. Meta-trained agents implement bayes-optimal agents. *Advances in neural information processing systems*, 33:18691–18703, 2020.
- Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A., et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- Ortega, P. A., Wang, J. X., Rowland, M., Genewein, T., Kurth-Nelson, Z., Pascanu, R., Heess, N., Veness, J., Pritzel, A., Sprechmann, P., et al. Meta-learning of sequential strategies. *arXiv preprint arXiv:1905.03030*, 2019.
- Prabhu, A., Cai, Z., Dokania, P., Torr, P., Koltun, V., and Sener, O. Online continual learning without the storage constraint. *arXiv preprint arXiv:2305.09253*, 2023.
- Rathmanner, S. and Hutter, M. A philosophical treatise of universal induction. *Entropy*, 13(6):1076–1136, 2011. ISSN 1099-4300. doi: 10.3390/e13061076.
- Schlag, I., Irie, K., and Schmidhuber, J. Linear transformers are secretly fast weight programmers. In *International Conference on Machine Learning*, pp. 9355–9366. PMLR, 2021.
- Shazeer, N. Fast transformer decoding: One write-head is all you need. *arXiv preprint arXiv:1911.02150*, 2019.
- Singh, A. K., Chan, S. C. Y., Moskovitz, T., Grant, E., Saxe, A. M., and Hill, F. The transient nature of emergent In-Context learning in transformers. November 2023.
- Titsias, M. K., Galashov, A., Rannen-Triki, A., Pascanu, R., Teh, Y. W., and Bornschein, J. Kalman filter for online classification of non-stationary data. *arXiv preprint arXiv:2306.08448*, 2023.
- Vapnik, V. Principles of risk minimization for learning theory. In *Proceedings of the 4th International Conference on Neural Information Processing Systems, NIPS’91*, pp. 831–838, San Francisco, CA, USA, December 1991. Morgan Kaufmann Publishers Inc.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Wallace, C. S. *Statistical and Inductive Inference by Minimum Message Length*. Springer Science & Business Media, December 2005.
- Yang, G., Hu, E., Babuschkin, I., Sidor, S., Liu, X., Farhi, D., Ryder, N., Pachocki, J., Chen, W., and Gao, J. Tuning large neural networks via zero-shot hyperparameter transfer. *Advances in Neural Information Processing Systems*, 34:17084–17097, 2021.

Appendix

A. Learning on an Individual Sequence: Minimum Description Length

The majority of works in machine learning, especially in deep-learning, assume a distributional perspective and rely on the framework of *empirical risk minimization* (ERM, Vapnik (1991)) to reason about generalization and evaluation. In this view an unknown distribution Q is considered the source of the observed data $(x, y) \sim Q$ and we strive to learn a model $p(y|x, \theta)$ and parameters θ that minimize some loss, for example the log-loss $\mathcal{L} = -\mathbb{E}_{x, y \sim Q} \log p(y|x, \theta)$. We have access to a limited number of samples $(x, y) \sim Q$, split them into a training- and a test set, and the theory around ERM provides generalization bounds for future data from the same distribution. This framework relies crucially on a few assumptions; among others, that there is a stationary underlying distribution Q .

Non-stationary data. The distributional assumption can become a liability when considering non-stationary data. At best, there is some ambiguity about the objective a model is supposed to optimize, and how to choose a best performing model. In general however, the whole notion of a test *set* becomes nonsensical: Sets imply exchangeability, and a given problem might not have the invariances that justify exchangeability. Consider for example the problem of modeling an epidemic outbreak as a function of time and location. Selecting random temporal intervals as test set seems naive: interpolating past data between time steps observed in the training data might be much easier than predicting the future. Splitting-off test data by location raises a similar dilemma, because a model that generalizes well over spatial locations in the past is by no means expected to generalize well on any location into the future. It seems that choosing an evaluation protocol is a non-trivial and potentially subjective endeavour. These ambiguities arise because the concept of a test *set* implies stationary, exchangeable data; and in the absence of that, we might be better off not relying on test sets at all.

Compression based inference: Minimum Description Length. Independent of ERM, compression based inference and learning has been widely studied. It is based on the fundamental idea that learning and comprehension correspond to compression (Chaitin, 2002; Rathmanner & Hutter, 2011). Given data $\mathcal{D}=(y_t)_1^T$ and a hypothesis space $\mathcal{M} = \{M_1, M_2, \dots\}$, where each hypothesis M corresponds to a parametric probabilistic model $p(\mathcal{D}|\theta, M)$, we aim to identify the model that can compress the data \mathcal{D} best. Considering the close relationship between lossless coding and probability distributions, this can be achieved by associating a code-length function $L(\mathcal{D}|M)=-\log p(\mathcal{D}|M)$ with each hypothesis. Additionally, we would have to consider the code-length of the hypothesis M itself to obtain the total code length $L(\mathcal{D}) = L(M)+L(\mathcal{D}|M)$, which is often ignored if the hypothesis under consideration are very similar, or are small compared to the data part $L(\mathcal{D}|M)$. A vast body of literature argues that models with a shorter description length have a better chance of generalizing to future data (Wallace, 2005; Grünwald, 2007; Rathmanner & Hutter, 2011).

A crude way to obtain the description length of the data given a *parametric* model family is to consider $L(\mathcal{D}|M)=L_M(\theta)+L_M(\mathcal{D}|\theta)$, where $L_M(\theta)$ is the cost of encoding the parameters and $L_M(\mathcal{D}|\theta)=-\log p(\mathcal{D}|\theta, M)$ is the cost of compressing the data with the parameterized model. This *two-part code* approach might be intuitive but it is sub-optimal. It is also highly ambiguous because it does not specify how to encode the parameters. This crude two-part approach to MDL has been refined in three distinct but closely related ways:

- 1) The Bayesian marginal likelihood: $L_{\text{Bayes}}(\mathcal{D}|M) := -\log \int_{\theta} p(\mathcal{D}|\theta, M)p(\theta)d\theta$ and its variational upper bound $L_{\text{ELBO}}(\mathcal{D}|M) := \text{KL}(q(\theta)|p(\theta)) - \mathbb{E}_{\theta \sim q} \log p(\mathcal{D}|\theta, M)$. Both depend on the chosen prior $p(\theta)$ and require access to the posterior distribution over parameters given data; or an approximation $q(\theta)$ thereof.
- 2) Normalized Maximum Likelihood (NML): $L_{\text{NML}}(\mathcal{D}|M) := -\log [p(\mathcal{D}|\hat{\theta}(\mathcal{D}), M) / Z]$, where $\hat{\theta}(\cdot)$ denotes the maximum likelihood estimator $\hat{\theta}(\mathcal{D}) = \text{argmax}_{\theta} p(\mathcal{D}|\theta, M)$ and $Z = \int_z p(\mathcal{D}|\hat{\theta}(z), M)dz$ is a normalization constant. NML normalizes over all possible observable data $z \in \mathcal{Y}^T$, which is often intractable and even undefined for many model families of interest. The log-denominator Z is also called the complexity of M and measures how well the model could fit all possible data. Under NML, a model M achieves a short description length only when it fits the given data well (high $\log p(\mathcal{D}|\hat{\theta}(\mathcal{D}), M)$) as well as not fit well many different data (low denominator).
- 3) The sequential approach: $L_{\text{plugin}}(\mathcal{D}|M) := -\sum_{t=1}^T \log p(y_t|\hat{\theta}(\mathcal{D}_{<t}), M)$ decomposes the description length over datapoints and relies on the choice of a suitable plug-in parameter estimator $\hat{\theta}(\mathcal{D}_{<t})$. It emphasizes that in order to achieve a short codelength, a model M must not only be a good predictor given all training data \mathcal{D} , but already given only parts of the data $\mathcal{D}_{<t}$ for all t , i.e. be a sample efficient predictor. Model complexity and sample efficiency can thus be considered two sides of the same coin.

Approaches for computing description lengths for some data \mathcal{D} under a model family $p(\mathcal{D}|\theta, M)$ are called *universal codes* when they result in description lengths that stay close to the (generally unachievable) maximum likelihood codelength for that family: $L_{\text{universal}}(\mathcal{D}|M) \leq -\log p(\mathcal{D}|\hat{\theta}(\mathcal{D}), M) + \mathcal{O}(\log T)$. $L_{\text{NML}}(\mathcal{D}|M)$ is by definition universal because the log-denominator contributes a model-architecture dependent gap that is constant in respect to the sequence length. $L_{\text{Bayes}}(\mathcal{D}|M)$ and $L_{\text{plugin}}(\mathcal{D}|M)$ are universal for many reasonable choices of the prior $p(\theta)$ and the estimator $\hat{\theta}(\mathcal{D}_{<t})$ respectively (Grünwald, 2007).

Note that none of these make an assumption about the i.i.d-ness of the data \mathcal{D} or the model $p(\mathcal{D}|\cdot)$. They can all be used to fit a model to an individual sequence $(y_t)_{t=1}^T$ and do not rely on held-out sets to evaluate generalization. Instead, all three versions of $L(\mathcal{D}|M)$ contain a build-in form of Occams Razor that ensures that models with shorter description length have a better chance to generalize to future data.

Compression based inference is also at the heart of *Algorithmic Complexity Theory* (ACT, Grunwald & Vitanyi (2004)), which provides an alternative foundation to reason about information: Where Shannon information is based on distributions, ACT is based on Turing machines and code-lengths. Most concepts from Shannons information theory have counterparts in ACT; e.g: Entropy vs. Kolmogorov complexity, cross-entropy vs. description length, Shannon mutual information vs. algorithmic (Kolmogorov) mutual information etc. (Grunwald & Vitanyi, 2004). In ACT, Solomonoff-Induction is the optimal but uncomputable learning and inference machine. It is closely related to the prequential MDL approach, where we chose the hypothesis space \mathcal{M} to be the set of *all* computable models (that can be expressed with Turing machines).

The prequential approach

With prequential MDL, learning and model-selection reduce to a specific form of online- or sequential learning. The goal is to find a model M that minimizes the description length $L(\mathcal{D})$:

$$\begin{aligned} L(\mathcal{D}) &= L(M) + L(\mathcal{D}|M) \\ &= L(M) - \log p(\mathcal{D}|M) \\ &= L(M) - \sum_{t=1}^T \log p(y_t|\mathcal{D}_{<t}, M) \end{aligned}$$

M denotes all modeling choices: The type of model, model-family, model assumptions, as well as inference and optimization choices – everything necessary to make a prediction for y_t given $\mathcal{D}_{<t}$. For neural networks, it includes the model architecture, parameter initialization, the optimizer and its hyperparameters, data augmentation strategies, etc. Formally, M is the program that implements the model and its training procedure on a universal computer such as a Turing machine. The search space for M is vast, and we can expect that some choices for M lead to significantly better models in terms of $\log p(\mathcal{D}|M) = \sum_{t=1}^T \log p(y_t|\mathcal{D}_{<t}, M)$ than others. However, the fact that we also consider the length $L(M)$ provides a form of Occam’s razor and ensures a good chance of generalizing to future data: $L(M)$ grows large if we encode specific information about \mathcal{D} into M that does not generalize.

An upper bound for $L(M)$ can be the length of the Python source code of our algorithm, plus the framework libraries and all the support code necessary to execute on a general Turing machine – potentially as a well compressed and self-extracting executable. In practice, we here compare methods M_i that differ by only a few lines of Python code. For our purpose the differences in $L(M_i)$ are thus very small compared to $L(\mathcal{D}|M)$ and we ignore $L(M_i)$. $L(M)$ should however be taken into account when comparing very different kind of learning approaches.

B. Privileged Information Transformer

We maintain the standard architecture for self-attention transformer blocks (Vaswani et al., 2017) with parallel feed-forward and attention pathways, however add projections of some privileged information y_t to the keys and values:

$$\begin{aligned}
 h_t^{l+1} &:= h_t^l + \text{FFW}(\bar{h}_t^l) + \text{Attn}(\bar{h}_t^l) \\
 \bar{h}_t^l &:= \text{LayerNorm}(h_t^l) \\
 \text{FFW}(\bar{h}_t^l) &:= W_{\text{down}} \sigma(W_{\text{up}} \bar{h}_t^l) \\
 \text{Attn}(\bar{h}_t^l) &:= \sum_{t'} a_{t,t'} v_{t'} \\
 a_{t,t'} &:= M \odot \text{softmax}_{t'}(q_t \times k_{t'}^l) \\
 q_t &:= W^q \bar{h}_t^l \\
 k_{t'} &:= W^k \bar{h}_{t'}^l + W^{\bar{k}} y_{t'} \\
 v_{t'} &:= W^v \bar{h}_{t'}^l + W^{\bar{v}} y_{t'}
 \end{aligned}$$

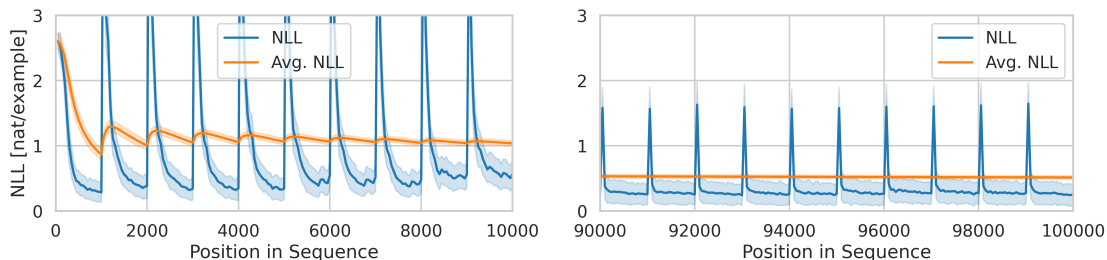
M is a causal attention mask with the diagonal set to 0. Setting the diagonal to 0 ensures that a token at time step t can not access the privileged information in y_t when generating predictions for the current time step.

C. Piecewise stationary Split-EMNIST experiments

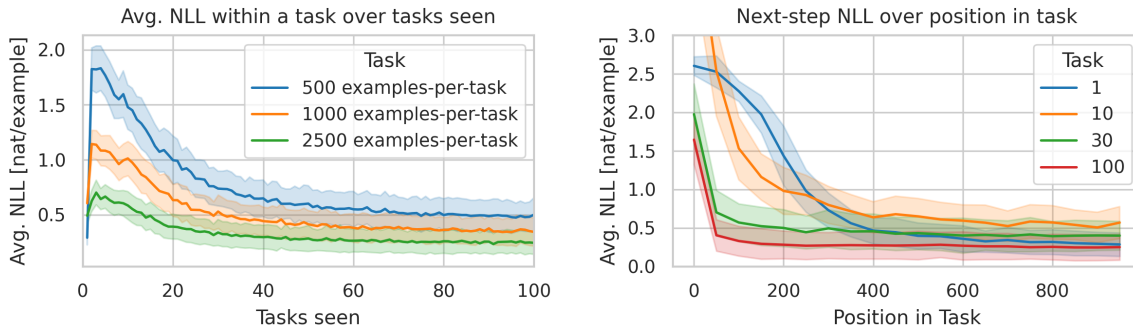
We use the convolutional neural network described by Blier & Ollivier (2018) as a feature extractor, however double the number of channels in all stages: The convolutions have 32, 32, 64, 64, 128, 128, 256 and 256 channels per layer respectively and max pooling operators after every second layer. The convolutions are followed by two fully connected layers of size 256. This architecture was used by Blier & Ollivier to evaluate neural network description lengths of MNIST. On top, we use 4 transformer blocks and choose a backbone width of 256 units, with 4 query-heads with a key-, value- and query- width of 64.

C.1. pi-Transformer.

For comparison we here show the plots from Figures 1 and 2 in terms of negative log-loss, instead of error-rates:



Alternative visualization of the Split-EMNIST experiments from Fig. 1: Instantaneous and averaged next-example log-losses for the first and last 10 tasks of Split-EMNIST: Image-to-label mappings are constant within each task, but randomly reassigned at task boundaries (averaged over 200 data generating random seeds).

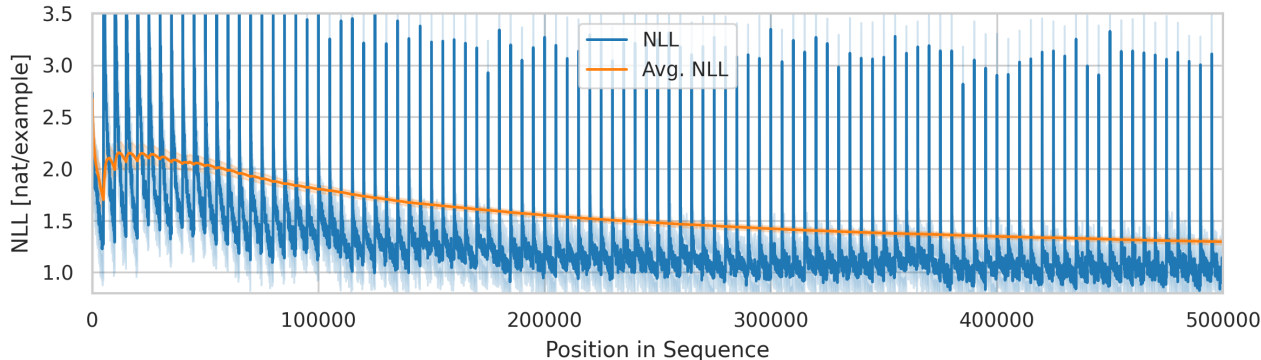


Alternative visualization of the same experiments: **Left:** Average performance per task shows strong forward-transfer after struggling during the first 10 to 20 tasks. **Right:** Detailed look at the within-task performance for the scenario with 1000 examples per task. The model is a strong few-shot learner after seeing 30 tasks and further improves until at least task 100

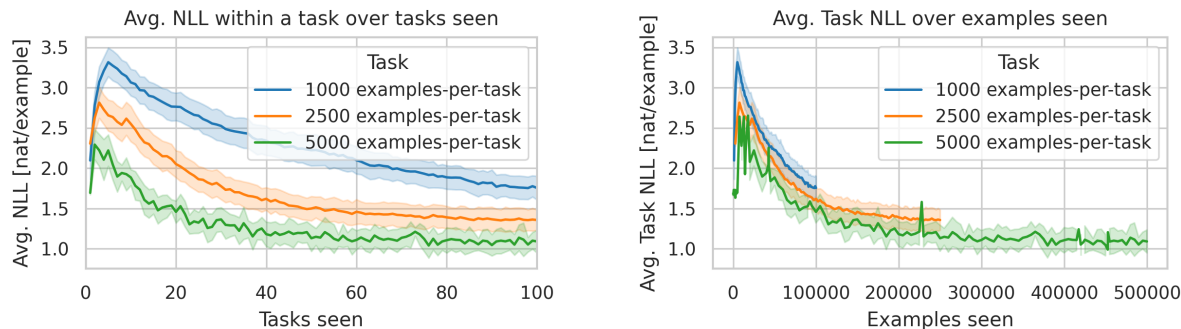
D. Piecewise stationary Split-CIFAR experiments

For CIFAR we use the VGG++ architecture from (Bornschein et al., 2022) as a feature extractor, a VGG-inspired ConvNet with normalization and 10% dropout added. On top of that feature extractor, as the temporal model, we use the same 4-block transformer architecture as in the Split-EMNIST experiments (Section 4.)

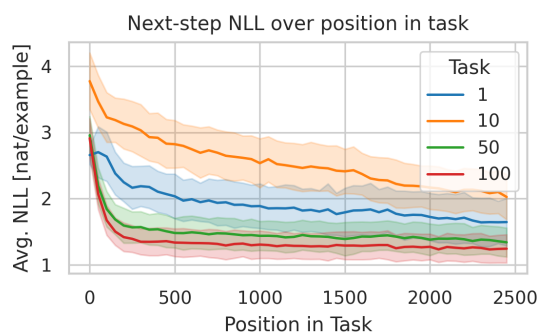
We train with 20 replay streams but otherwise use the same hyper-parameters from the previous section (see C for details). We observe that the results generally resemble those of the Split-EMNIST experiments. In particular, the task-averaged predictive performance degrades during the first couple of tasks; however then rapidly improves and the model soon learns to be an effective few shot learner that makes accurate predictions within a few examples after each task switch.



Predictive performance on a piecewise-iid CIFAR-100 sequence. The sequence consists of 100 tasks with 5000 examples. Each task is a 10-way classification problem, with the 10 classes randomly sourced from the 100 classes of the CIFAR-100 data set. We plot the the instantaneous and averaged next-example prediction performance as a function of examples seen.



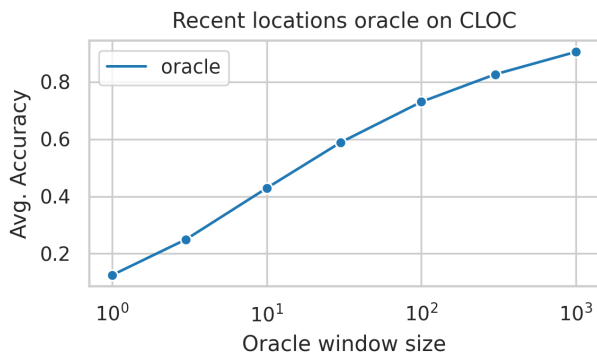
Detailed look at predictive performance for piecewise-i.i.d. CIFAR-100: We plot the the average per task next-step performance either as a function of tasks seen (left), or as a function of total examples seen (right).



Detailed look at the within-task performance for piecewise-iid CIFAR-100 with 2500 examples per task. The model is a strong few-shot learner after seeing 50 tasks and further improves until task 100

E. CLOC with fixed, pre-trained feature extractor

E.1. Oracle model



CLOC sequence: Predictive performance of an oracle model that predicts the correct location if and only if that location was observed within a window of the W most recent examples. The left most data point ($W=1 \Leftrightarrow 12\%$ avg. acc.) for example corresponds to a model that predicts the correct location if it is the same location as the previous image. It corresponds to the prediction performance of a model that simply predicts $\hat{y}_t = y_{t-1}$.

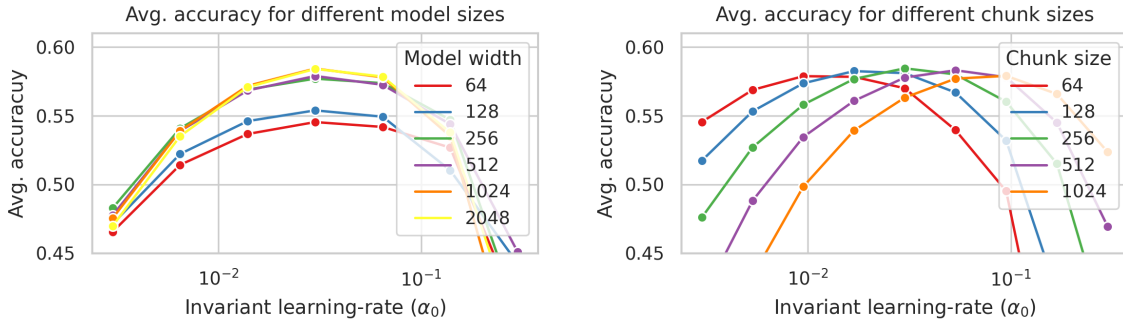
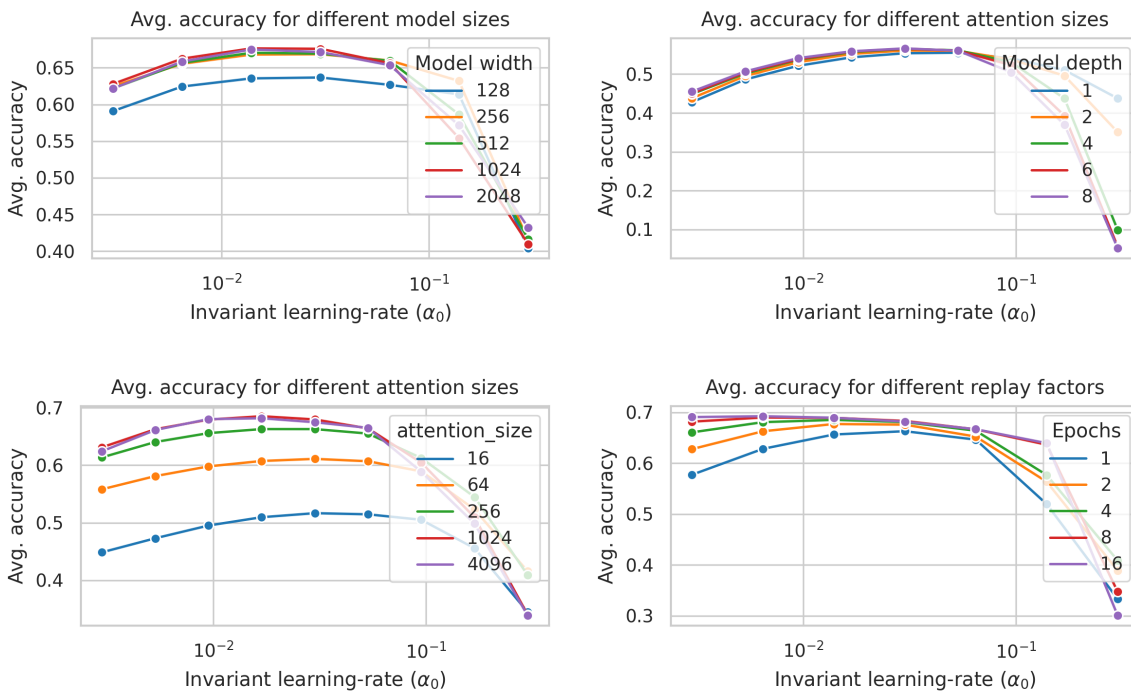


Figure 12: Pi-Transformer on CLOC with a fixed, pre-trained ResNet-50 feature extractor: **Left)** Larger models are generally better. **Right)** Average accuracy is comparably insensitive to the chunk sizes considered here.

E.2. Pi-Transformer with pretrained and fixed Supervised Resnet features

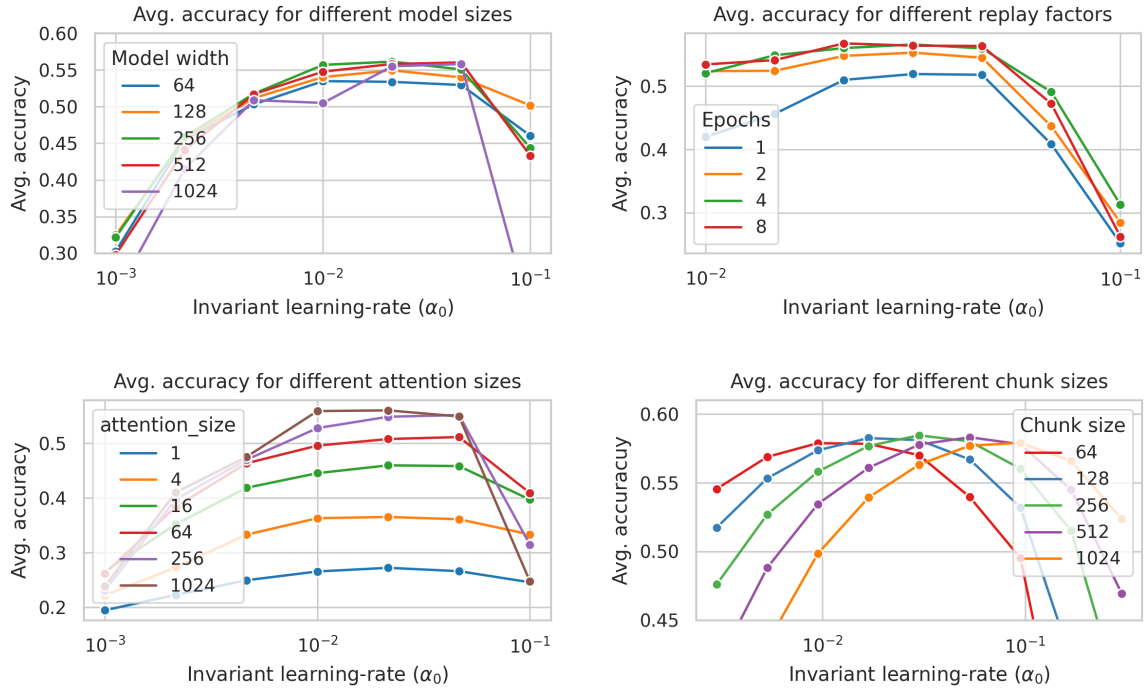
E.2.1. PI-TRANSFORMER WITH PRETRAINED AND FIXED MAE VIT-B FEATURES

We run the pi-Transformer with features extracted with the MAE pretrained VIT-B network from He et al. (2022). We show the sensitivity to various hyperparameter changes:



E.3. 2-token approach with pretrained and fixed Resnet features

We online train a Transformer with the 2-token approach on a with features extracted with a Resnet-50 trained on supervised ImageNet 1k. The transformer has a backbone-width of $D = 1024$ and is 8 blocks deep, which is twice as deep as the 4 block deep architecture we usually use for pi-Transformers on pre-extracted CLOC features. We show the sensitivity to various hyperparameter changes:



E.4. Hyper-parameter cube for Pareto fronts with pretrained and fixed feature extractors

We use the same hyper parameter search space for all experiments:

Hyper Parameter	Range
Model Width (D)	[16, \dots , 4096)
Model Depth	[1, \dots , 12]
Num Streams	[1, \dots , 16]
Attention window (C)	[128, \dots , 2048]
Learning Rate (α_0)	[1e-2, \dots , 1e-2]
AdamW Weight Decay	[1e-2, \dots , 3e-2]

F. CLOC with ConvNet feature extractor from scratch

F.1. Description of the feature extractor

We use the ConvNext tiny architecture (Liu et al., 2022) as a feature extractor and feed the top-level activations as inputs to the temporal transformer. We scale input images to a resolution of $160 \times 160 \times 3$ channels and run them initially through a convolution with stride 4 and 96 channels. These features are then processed by 4 successive stages with 96, 192, 384, 768 channels respectively. Between each stage we downscale the spatial resolution by a factor of 2 and apply LayerNorm.. The stages have 3, 3, 9, 3 blocks each. The pseudo-code for each block reads:

```
def forward(input)
    x = DepthwiseConv2D(1, kernel_shape=7)(input)
    x = LayerNorm(axis=-1, create_scale=True, create_offset=True)(out)
    x = self.activation(
    x = Linear(4 * channels)(out)
    x = Linear(channels)(out)
    return input + x * get_parameter('skip_gain', x.shape[-1])
```

F.2. Hyper-parameters for Pi-Transformer experiments

Hyper Parameter	Value
Image Size	160×160
Feature Extractor Learning Rate	$3e-4$
Transformer Width (D)	1024
Model Depth	2
Num Streams	8
Chunk Size (S)	256
Attention window (C)	2048
Transformer Learning Rate (α_0)	$1e-2$

F.3. Hyper-parameters for online Transformers with 2-token approach

Hyper Parameter	Value
Image Size	160×160
Feature Extractor Learning Rate	$3e-4$
Transformer Width (D)	1024
Model Depth	10
Num Streams	8
Chunk Size (S)	256
Attention window (C)	2048
Transformer Learning Rate (α_0)	$1e-2$

$$p(y_t | x_t, \mathcal{D}_{<t}, \theta_{t-1})$$