
Not All Memories are Created Equal: Learning to Forget by Expiring

Sainbayar Sukhbaatar¹ Da Ju¹ Spencer Poff¹ Stephen Roller¹ Arthur Szlam¹ Jason Weston¹ Angela Fan^{1,2}

Abstract

Attention mechanisms have shown promising results in sequence modeling tasks that require long-term memory. Recent work investigated mechanisms to reduce the computational cost of preserving and storing memories (Rae et al., 2020). However, not all content in the past is equally important to remember. We propose *Expire-Span*, a method that learns to retain the most important information and *expire* the irrelevant information. This forgetting of memories enables Transformers to scale to attend over tens of thousands of previous timesteps efficiently, as not all states from previous timesteps are preserved. We demonstrate that Expire-Span can help models identify and retain critical information and show it can achieve strong performance on reinforcement learning tasks specifically designed to challenge this functionality. Next, we show that Expire-Span can scale to memories that are tens of thousands in size, setting a new state of the art on incredibly long context tasks such as character-level language modeling and a frame-by-frame moving objects task. Finally, we analyze the efficiency of Expire-Span compared to existing approaches and demonstrate that it trains faster and uses less memory.

1. Introduction

Transformer architectures (Vaswani et al., 2017) have demonstrated strong performance across a variety of tasks (Devlin et al., 2019; Roller et al., 2020; Brown et al., 2020), including those that require learning long term relationships (Zhang et al., 2018; Fan et al., 2019a; Izacard & Grave, 2020). Recent work has focused on scaling attention mechanisms efficiently to longer memory sizes, enabling large improvements on long context tasks (Dai et al., 2019;

Sukhbaatar et al., 2019a). However, a critical component of human memory is not just the ability to remember, but also *forgetting* irrelevant information to focus on the salient, relevant bits. Most studies of long-term memory in humans indicate that not everything is remembered (Murre & Dros, 2015; Bahrick et al., 2008) — instead, only vivid, remarkable memories are retained from the far past (Wixted, 2004).

Standard Transformer architectures lack the ability to search over extremely large memories, as the self-attention mechanism is computationally intensive and the storage cost of preserving the large memory grows quickly. Recent work (Child et al., 2019; Rae et al., 2020) has proposed learning how to extend to greater context through sparse mechanisms or through compression, to more compactly represent the past. However, there exists a fundamental problem with large memories beyond strict computational concerns: as the amount of information stored increases, deciding which information is relevant becomes more challenging. Other work (Lample et al., 2019) approaches this by considering how to efficiently search large memories. We focus on an efficient way to learn what to forget, thereby reducing the computational burden of the model and easing the challenges of the search problem.

We propose EXPIRE-SPAN, a straightforward extension to attention mechanisms that learns when to *expire* unneeded memories. By expiring memories that are no longer useful, EXPIRE-SPAN enables scaling to tens of thousands of timesteps into the past. This learnable mechanism allows the model to adjust the span size as needed, selecting which information is critical to retain and forgetting the rest. More concretely, we augment the self-attention with a simple predictor that outputs an expiration value for each hidden state that determines how long a memory should be retained and accessible to the model. After the EXPIRE-SPAN runs out, the memory will be forgotten, but in a gradually differentiable way to retain end-to-end training with backpropagation. This process is done independently for each layer, allowing different layers to specialize at different time-scales. As EXPIRE-SPAN can flexibly adjust its span based on context, it is more efficient in terms of memory and training time compared to existing long memory approaches.

¹Facebook AI Research ²LORIA. Correspondence to: Sainbayar Sukhbaatar <sainbar@fb.com>.

We demonstrate that EXPIRE-SPAN can distinguish between critical and irrelevant information on several illustrative tasks in natural language processing and reinforcement learning that are specifically designed to test this ability. We then show we can achieve state-of-the-art results on long-context language modeling benchmarks, and EXPIRE-SPAN can scale to memories in the tens of thousands on a frame-by-frame colliding objects task — by expiring irrelevant information, capacity is freed to have even larger memory. Then, we compare the efficiency of our method to competitive baselines and show EXPIRE-SPAN is faster and has a smaller memory footprint. Finally, we analyze the information retained and expired by EXPIRE-SPAN models, to understand the importance of long context memory.

2. Related Work

Memory is crucial for many tasks and has been studied in recurrent networks (Elman, 1990; Hochreiter & Schmidhuber, 1997; Mikolov et al., 2010) for a long time. The development of memory augmented networks (Graves et al., 2014; Sukhbaatar et al., 2015b) made it possible to store large quantities of information and selectively access them using attention (Bahdanau et al., 2015). The Transformer (Vaswani et al., 2017) took full advantage of this approach. Processing long sequences with Transformers is an active area with applications in language understanding (Brown et al., 2020), reinforcement learning (Parisotto et al., 2020), video processing (Wu et al., 2019), and protein folding (Rives et al., 2019; Choromanski et al., 2020). However, extending the memory span is computationally expensive due to the quadratic time and space complexity of self-attention. Other work focuses on benchmarking long memories (Tay et al., 2021), but focuses on encoder-only tasks, whereas we focus on decoder-only Transformers.

Various work has focused on reducing this complexity and increasing memory capacity (Schlag et al., 2021). Dynamic attention spans, such as Adaptive-Span (Sukhbaatar et al., 2019a) and Adaptively Sparse Transformer (Correia et al., 2019), focus on learning which attention heads can have shorter spans, but can only extend to spans of a few thousand. Other work sparsifies attention by computing fewer tokens (Fan et al., 2019b), often by using fixed attention masks (Child et al., 2019) or sliding windows and dilation (Beltagy et al., 2020). The BP Transformer (Ye et al., 2019) structures tokens as a tree, so some tokens have coarse attention. These works focus on learning what to attend to, but searching larger and larger memories is very difficult. In contrast, we focus on learning to expire what is irrelevant. Compressive Transformer (Rae et al., 2020) reduces the number of memories by replacing every few memories with a single compressed one. A disadvantage of this is that all memories have the same compression ratio, so relevant

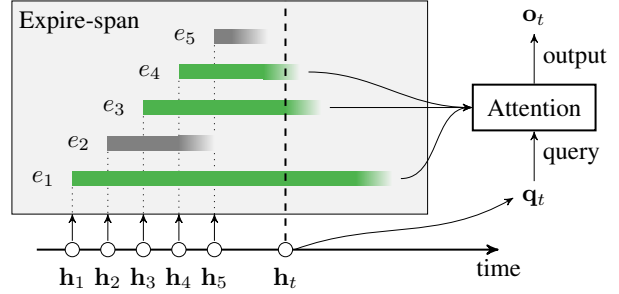


Figure 1. **Expire-Span.** For every memory h_i , we compute an EXPIRE-SPAN e_i that determines how long it should stay in memory. Here, memories h_2, h_5 are already expired at time t , so the query q_t can only access $\{h_1, h_3, h_4\}$ in self-attention.

memories are equally compressed.

Another line of work investigates linear-time attention mechanisms. Wu et al. (2018) replace self-attention with convolutions that run in linear time, but the scalability to long context tasks remains limited. Wang et al. (2020) propose linear time attention by decomposing attention into multiple smaller attentions, that recombine to form a low-rank factorization of the original attention. Katharopoulos et al. (2020) propose linear attention by expressing self-attention as instead a linear dot-product of feature maps. Peng et al. (2021) propose Random Feature Attention, used to approximate the softmax. Those methods, however, focus on making attention more efficient without reducing the number of memories. Further, as our goal is to reduce the number of memories that feed to self-attention by learning to expire, EXPIRE-SPAN can be easily combined with these efficiency improvements. For a review of further recent Transformer variants, see Tay et al. (2020).

3. Background

Transformer architectures have been widely used as decoder-only auto-regressive models for sequential tasks. A Transformer decoder is made of a stack of identical layers, composed of a multi-head self-attention sublayer followed by a feedforward sublayer. The output of each timestep is the hidden state $h_t^l \in \mathbb{R}^d$ at layer l , which is then projected to key k , value v , and query q vectors:

$$q_t^l = W_q^l h_t^l, \quad k_t^l = W_k^l h_t^l, \quad v_t^l = W_v^l h_t^l. \quad (1)$$

Going forward, we focus on a single layer and omit the layer index l for brevity. Information from previous timesteps is accessed through attention a_{ti} to create output o_t :

$$a_{ti} = \text{Softmax} \left(q_t^\top k_i \right), \quad o_t = W_o \sum_{i \in C_t} a_{t,i} v_i. \quad (2)$$

The set C_t indicates which memories can be accessed at time t , which is the focus of this work. The space and time complexity of self-attention is linearly correlated to the size of this set $|C_t|$, making it an important metric of efficiency. For the rest of the paper, we will refer to $|C_t|$ as the *memory size*.

Including all previous timesteps in self-attention by setting $C_t = \{1, \dots, t-1\}$ results in a quadratic complexity $\mathcal{O}(T^2)$ to compute the full attention over a sequence of length T . *Fixed-spans* (Dai et al., 2019) take a more scalable approach such that $C_t = \{t-L, \dots, t-1\}$ so the attention is restricted to previous L steps. The total complexity in this case is $\mathcal{O}(TL)$, where L is the attention span.

Adaptive-Span (Sukhbaatar et al., 2019a) further improves upon upon this by learning an optimal span L per attention head from data, which results in small L values for many heads. *Compression* approaches (Rae et al., 2020) reduce memory size by compressing multiple timesteps into a single memory, with complexity $\mathcal{O}(TL/c)$, where c is the compression rate. However, in all these approaches, all memories are treated equally without regards to their importance to the task. In this work, we focus on distinguishing between relevant and irrelevant memories by learning to expire unneeded information — by expiring, the remaining attention on relevant information can scale beyond existing long context memory approaches.

4. Expire-Span

We describe EXPIRE-SPAN and how to integrate it into Transformers to focus on relevant information and expire the rest, meaning memories can be permanently deleted. We describe how to scale EXPIRE-SPAN and practically train with drastically longer memory spans.¹

4.1. Method

EXPIRE-SPAN, depicted in Figure 1, allows models to selectively forget memories that are no longer relevant. We describe it in the context of a single Transformer layer and omit the layer index l for brevity. Our goal is to reduce the size of C_t defined in Section 3 for more efficiency without performance degradation. For each memory $\mathbf{h}_i \in \mathbb{R}^d$, we will compute a scalar EXPIRE-SPAN $e_i \in [0, L]$:

$$e_i = L\sigma(\mathbf{w}^\top \mathbf{h}_i + b). \quad (3)$$

¹The full implementation can be found at <https://github.com/facebookresearch/transformer-sequential>.

Here $\mathbf{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ represent trainable parameters, σ is the sigmoid function, and L is the maximum span. This expire-span e_i determines how long \mathbf{h}_i should be kept and included in C_t . At time t , the remaining span of \mathbf{h}_i is $r_{ti} = e_i - (t - i)$. When r_{ti} becomes negative, it indicates the memory \mathbf{h}_i is expired and can be removed from C_t . This can be implemented by updating attention weights a_{ti} with a binary masking function $m_{ti} = \mathbf{1}_{r_{ti} > 0}$:

$$a'_{ti} = \frac{m_{ti}a_{ti}}{\sum_j m_{tj}a_{tj}}, \quad \mathbf{o}_t = \sum_i a'_{ti}\mathbf{v}_i. \quad (4)$$

However, with such discrete masking, the Expire-Span e_i will not receive any gradient for training. Instead, we use a soft masking function from Sukhbaatar et al. (2019a) that smoothly transitions from 0 to 1 (see Figure 2):

$$m_{ti} = \max(0, \min(1, 1 + r_{ti}/R)), \quad (5)$$

where R is a hyperparameter that determines the length of a ramp that is bounded between 0 to 1. This function has non-zero gradient for values in $[-R, 0]$ to train e_i , but also can take a value of 0, which is necessary for expiring memories. Thus $C_t = \{i \mid m_{ti} > 0\}$. Since m_{ti} is a monotonically decreasing function of t , once a memory is expired, it can be permanently deleted.

Our goal is to reduce the average memory size, which is directly related with the average EXPIRE-SPAN:

$$\begin{aligned} \frac{1}{T} \sum_t |C_t| &= \frac{1}{T} \sum_t \sum_{i < t} \mathbf{1}_{m_{ti} > 0} \\ &= \frac{1}{T} \sum_i \left(R + \sum_{t > i} \mathbf{1}_{r_{ti} > 0} \right) \\ &= \frac{1}{T} \sum_i \left(R + \sum_{t > i} \mathbf{1}_{e_i > t-i} \right) \\ &= R - 1 + \frac{1}{T} \sum_i [e_i] \end{aligned} \quad (6)$$

Therefore, we add an auxiliary term to the loss function to penalize the L1-norm of EXPIRE-SPAN:

$$L_{\text{total}} = L_{\text{task}} + \alpha \sum_i e_i/T, \quad (7)$$

where $\alpha > 0$ is a hyperparameter. This term decreases the span of memories that contribute less to the main task, resulting in a small memory that focuses only on relevant information. Note the new parameters, \mathbf{w} and b , and the computations of EXPIRE-SPANS are negligible in size compared to the total number of parameters and computations.

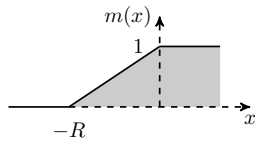


Figure 2. Soft Mask

4.2. Adding Expire-Span to Transformers

We describe how EXPIRE-SPAN can be utilized within Transformer self-attention layers to decrease the memory size and focus on salient information. This section describes each modification clearly, to facilitate easier reproduction. We discuss practical training concerns, such as efficiency and regularization. Additional details can be found in the appendix.

Modifications to Multi-Head Attention Self-attention consists of multiple heads that have different keys, values, and queries. However, they all share one underlying memory, so a memory cannot be removed if it is used by any of the heads. Thus, we compute an EXPIRE-SPAN at each layer that is shared amongst the heads.

Block Parallel This modification allows memories to be permanently deleted in EXPIRE-SPAN. We use the caching mechanism (Dai et al., 2019), where a block of timesteps $B = [t, \dots, t + K - 1]$ is processed in parallel for efficiency — once a block is computed, its hidden states $[\mathbf{h}_t, \dots, \mathbf{h}_{t+K-1}]$ are cached so that future blocks can attend to them. This means a memory can be deleted only if it is not used by any of the queries in B . Concretely, \mathbf{h}_i will be deleted when $m_{ti} = 0$ where t is the first token of B . However, this is not a concern for very long-term memories where $L \gg K$.

Loss Computation The L1-norm loss for EXPIRE-SPAN must be computed for every memory \mathbf{h}_i . A straightforward way is to compute it for the current block B . This empirically results in poor performance — a possible explanation is that the time between positive and negative gradients on e_i may become too distant. Negative gradients that increase e_i only come from the main loss L_{task} through the masking function m_{ti} , which has non-zero gradients only when memory \mathbf{h}_i is about to expire with $0 < m_{ti} < 1$ for $t \in B$. For a large $L \gg K$, \mathbf{h}_i may have been computed many blocks before and since then the model weights would have changed. In contrast, the positive gradients that decrease e_i are computed on the current block $i \in B$. To remove this discrepancy, we compute the auxiliary loss on e_i at the same time as negative gradients when $0 < m_{ti} < 1$ for $t \in B$.

Regularization A potential challenge in exceptionally long memory is greater capacity to overfit. As EXPIRE-SPAN can scale to memories in the tens of thousands, it can overfit to learning specific span sizes on the training set that do not generalize. As a form of regularization, we propose to randomly shorten the memory during training. For each batch, we sample $l \sim \mathcal{U}(0, L)$ and set $a_{ti} = 0$ for all $t - i > l$ only during training. This way, the model cannot assume the memory will always contain specific

information, as the memory is randomly shortened.

Stable Training with Extremely Large Spans Multiplier L in Eq. 3 is the maximum span, so it can take very large values, exceeding tens of thousands. This is a potential problem because small changes in \mathbf{h}_i or \mathbf{w} will be amplified in EXPIRE-SPAN e_i , and subsequently have dramatic effects on the model behaviour. As a straightforward remedy, for very large L values, we replace Eq. 3 with

$$e_i = L\sigma((\mathbf{w}^\top \mathbf{h}_i + b)/R). \quad (8)$$

5. Experiments and Results

We show that EXPIRE-SPAN focuses on salient information on various constructed and real-world tasks that necessitate expiration. First, we describe baselines and efficiency metrics for comparing various models. Second, we illustrate the importance of expiration on various constructed tasks. Then, we highlight the scalability of EXPIRE-SPAN when operating on extremely large memories. Additional experiments and details are in the appendix.

5.1. Baselines

We compare our method against several baselines from Section 3 that take different approaches to access information in the past. We compare the performance of these methods, along with two efficiency metrics: GPU memory and training speed for a fixed model size and batch size. First, we compare to *Transformer-XL* (Dai et al., 2019), which corresponds to the fixed-span approach where simply the last L memories are kept. Our Transformer-XL implementation also serves as a base model for all the other baselines to guarantee that the only difference among them is how memories are handled. The other baselines are *Adaptive-Span* (Sukhbaatar et al., 2019a) and *Compressive Transformer* (Rae et al., 2020), two popular approaches for long memory tasks. For Compressive Transformer, we implemented the mean-pooling version, which was shown to have strong performance despite its simplicity.

5.2. Importance of Expiration: Illustrative Tasks

Remembering One Key Piece of Information To illustrate a case where proper expiration of unnecessary memories is critical, we begin with an RL gridworld task: walking down a corridor. In this *Corridor* task, depicted in Figure 3 (left), the agent is placed at one end of a very long corridor, next to an object that is either red or blue. The agent must walk down the corridor and go to the door that corresponds to the color of the object that it saw at the beginning to receive +1 reward. The requirement on the memory is very low: the agent only needs to remember the object color so it can walk through the correct door.

Model	Maximum span	Accuracy (%)
Transformer-XL	2k	26.7
EXPIRE-SPAN	16k	29.4
EXPIRE-SPAN	128k	52.1

Table 1. **Copy Task.** We report accuracy on the test set.

multiple instructions can be in queue for execution.

We experiment with a dataset where the average distance between receiving and executing instructions is around 950 distractor words. Models are trained as language models, but evaluated only on their success in executing the instruction. Task details and model architecture are provided in the appendix. We illustrate in Figure 4 (right) that EXPIRE-SPAN is much more successful at this task than Transformer-XL and Adaptive-Span (see the appendix), as it can focus on the specific instruction lines.

5.3. Scalability of Expire-Span

We analyze the scalability of EXPIRE-SPAN. On a copy task, we train models with spans up to 128k timesteps. Then, we show the utility of EXPIRE-SPAN on character-level language modeling — Enwik8 and PG-19 — and a moving objects task that is processed frame by frame. For these tasks, we also analyze the efficiency of EXPIRE-SPAN compared to existing methods, and demonstrate that our method has a smaller memory footprint and faster processing speed. We quantify efficiency with two metrics: (1) peak GPU memory usage and (2) training time per batch (comparing fixed batch size for similar size models).

Extremely Long Copy To illustrate the scalability of EXPIRE-SPAN, we construct a copy task where the model sees a sequence of A very far in the past. The rest of the characters are B . The model must copy the correct quantity of A . We design the task such that a long span (up to 128k) can be required, as the A tokens are very far into the past. In Table 1, we show that only by scaling the maximum span to 128k it is possible to achieve improved performance. We compare to a Transformer-XL baseline with 2k attention span and a EXPIRE-SPAN model with smaller span.

Character Level Language Modeling: Enwik8 We subsequently experiment on Enwik8 for character level language modeling (Mahoney, 2011). We compare the performance of EXPIRE-SPAN with Adaptive-Span and Transformer-XL, varying the average span size (see Figure 5). Models with EXPIRE-SPAN achieve stronger results — when comparing at any given memory size, EXPIRE-SPAN outperforms both baselines. Further, the performance of EXPIRE-SPAN does not vary much even if the memory size is drastically reduced, indicating the model retains a

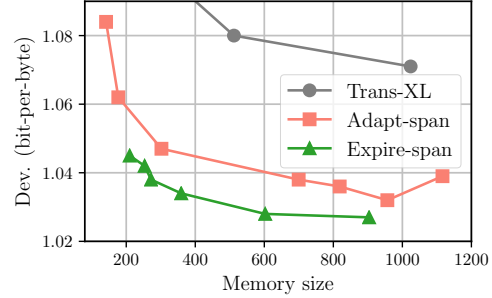


Figure 5. **Performance as a Function of Memory Size on Enwik8.** Lower bpb and smaller memory size is better.

Model	Params	Test
<i>Small models</i>		
Trans-XL 12L (Dai et al., 2019)	41M	1.06
Adapt-Span 12L (Sukhbaatar et al., 2019a)	39M	1.02
Our Trans-XL 12L baseline	38M	1.06
EXPIRE-SPAN 12L	38M	0.99
Trans-XL 24L (Dai et al., 2019)	277M	0.99
Sparse Trans. (Child et al., 2019)	95M	0.99
Adapt-Span 24L (Sukhbaatar et al., 2019a)	209M	0.98
All-Attention (Sukhbaatar et al., 2019b)	114M	0.98
Compressive Trans. (Rae et al., 2020)	277M	0.97
Routing Trans. (Roy et al., 2020)	-	0.99
Feedback Trans. (Fan et al., 2020b)	77M	0.96
EXPIRE-SPAN 24L	208M	0.95

Table 2. **Enwik8 Results.** We report bit-per-byte (bpb) on test and the number of parameters.

small quantity of salient information for good performance.

Next, we compare EXPIRE-SPAN to existing work in Table 2. A small EXPIRE-SPAN model with the maximum span $L = 16k$ outperforms similarly sized baselines by a large margin. We also trained a larger EXPIRE-SPAN model with $L = 32k$ and LayerDrop (Fan et al., 2020a), which outperforms the Compressive Transformer and sets a new state of the art on this task. This indicates that models can learn to expire relevant information and encode long context effectively, even on very competitive language modeling benchmarks.

Finally, we compare the efficiency of EXPIRE-SPAN with the Transformer-XL, Adaptive-Span and Compressive Transformer baselines. We find that EXPIRE-SPAN models achieve much better performance, as shown in Table 4 with substantially less GPU memory and faster training time per batch.

Character Level Language Modeling: PG-19 We use the PG-19 (Rae et al., 2020) benchmark and convert it to character-level language modeling with a vocabulary size of

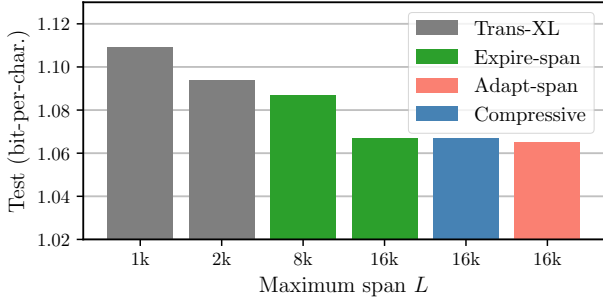


Figure 6. **Performance on Character-level PG-19.** We report bit-per-character on test.

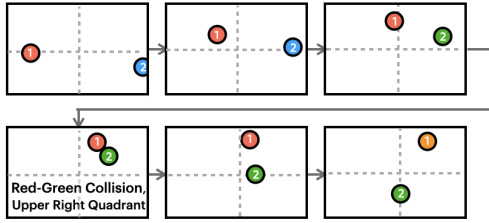


Figure 7. **Object Collision** task tests if models can remember the location of specified colored collisions.

3506. We train several baselines: Transformer-XL with maximum spans of 1k and 2k, and Adaptive-Span and Compressive Transformers with 16k span. We train EXPIRE-SPAN with maximum spans of 8k and 16k. We present results in Figure 6, where we show that EXPIRE-SPAN is substantially better than Transformer-XL, and matches the performance of Adaptive-Span and Compressive Transformer.

However, EXPIRE-SPAN uses its available memory very effectively. The 16k maximum span EXPIRE-SPAN model has an average memory size of 860. In comparison, the Adaptive-Span model has an average memory size of 2440, almost 3x that of the 16k EXPIRE-SPAN model. This indicates that EXPIRE-SPAN enables models to identify the critical bits of information and expire the rest, reaching the same performance with a much smaller memory.

Finally, comparing efficiency (Table 4), EXPIRE-SPAN trains at double the speed of Compressive Transformer. EXPIRE-SPAN is faster than Adaptive-Span, though uses slightly more memory. The memory usage of EXPIRE-SPAN is usually lower, around 12GB, but spikes for some sentences. Lastly, while the average span size of EXPIRE-SPAN is lower than Adaptive-Span, the computation requires additional tensors allocated in memory, which can potentially be addressed by an optimized implementation.

Model	Maximum Span	Test Error (%)
Transformer-XL	1k	73.3
Compressive	8k	63.8
Adaptive-Span	16k	59.8
EXPIRE-SPAN	16k	52.2
	32k	36.7
	64k	26.7

Table 3. **Results on Object Collision.** We report the error on the test set comparing to various baselines.

Frame-by-Frame Processing: Object Collision An important setting where learning which long context may be important is in video understanding, a field with increasing focus as model architectures provide the capability to process long sequences. Despite video data being memory intensive, salient events might be localized in space and time. We test our model on a task where two objects move around and collide, and the goal is to reason about the location of specified-color collisions. Objects have a color that can randomly change. We divide the grid into four quadrants and the model is asked to recall the quadrants of the last collision of a specific color pair. Because the collisions are rare, and collisions of specific colors are even rarer, the model must process a large quantity of frames.

We illustrate the task in Figure 7 and results in Table 3. The task requires many frames, so long context is very beneficial — as the EXPIRE-SPAN maximal span increases, performance steadily rises. Our largest span, 64k, matches the size of the largest attention limit reported to date (Kitaev et al., 2019) and has the strongest performance. This model is trained with the random drop regularization method described in Section 4.2. Compared to Compressive Transformer and Adaptive-Span baselines, our EXPIRE-SPAN model has the strongest performance.

Comparing efficiency, EXPIRE-SPAN trains almost 3x faster than both baselines (see Table 4) while having much stronger performance. Further, expiration is critical to this performance — a Adaptive-Span model with $L = 32k$ runs out of memory in the same setting where we trained our EXPIRE-SPAN model with $L = 64k$. Through expiration, our model can keep the GPU memory usage reasonable and train with the longer spans necessary for strong performance.

6. Analysis and Discussion

EXPIRE-SPAN creates the phenomena of *selective forgetting*: it allows memories to be permanently deleted if the model learns they are not useful for the final task. In this section, we analyze the information retained and expired by EXPIRE-SPAN models to better understand how models use

	Model	Performance	GPU Memory (GB)	Time/Batch (ms)
Enwik8	Transformer-XL	1.06 bpb	27	649
	Compressive Transformer	1.05 bpb	21	838
	Adaptive-Span	1.04 bpb	20	483
	EXPIRE-SPAN	1.03 bpb	15	408
Char-level PG-19	Compressive Transformer	1.07 bpc	17	753
	Adaptive-Span	1.07 bpc	13	427
	EXPIRE-SPAN	1.07 bpc	15	388
Object Collision	Compressive Transformer	63.8% Error	12	327
	Adaptive-Span	59.8% Error	17	365
	EXPIRE-SPAN	52.2% Error	12	130

Table 4. **Efficiency of EXPIRE-SPAN.** We report peak GPU memory usage and per-batch training time, fixing the batch size.

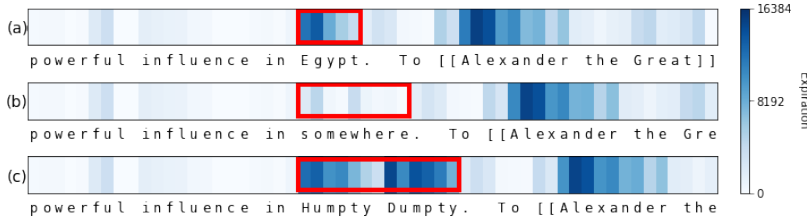


Figure 8. **Expiration in EXPIRE-SPAN on Enwik8.** In (a), the model strongly memorizes two areas, “Egypt” and “Alexander”. In (b), if we replace “Egypt” with “somewhere”, then it’s forgotten fast. In (c), we insert “Humpty Dumpty” and the model retains these rare words in memory.

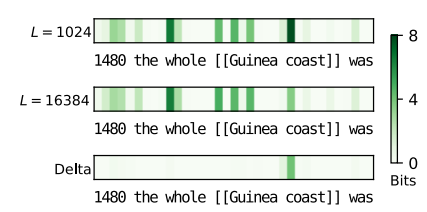


Figure 9. **Accuracy Needs Memory.** As the maximum span is artificially decreased at inference time from 16k to only 1k, the prediction is less accurate.

the ability to forget. Additional analyses are in the appendix.

Retaining Salient Information We analyze what is retained by an EXPIRE-SPAN model on Enwik8 to understand how models utilize the ability to forget. In Figure 8 (a), we show that the model retains information about named entities such as *Egypt* and *Alexander the Great* by giving them longer spans (darker color). Next, we analyze how expire-spans changes when we artificially edit the past text. In Figure 8 (b), we replace the entity *Egypt* with the generic text *somewhere*, and this generic word is quickly expired. In Figure 8 (c), we edit *Egypt* to *Humpty Dumpty*, which is a very rare entity, and the model retains it in memory without expiring. In addition to entities, EXPIRE-SPAN memorizes spaces, newlines, and section titles, all of which retain information about words, sentences, or sections. The model’s expiration choices vary by layer, indicating that EXPIRE-SPAN models use the memory at each layer to remember different information.

Importance of Long Term Memory Next, we analyze which predictions benefit the most from memory capacity. We take an EXPIRE-SPAN model trained on Enwik8 and decrease the maximum span size to 1024 at inference time, even though the model was trained with a maximum span of 16k. We then compare which predictions decreased in accuracy. In Figure 9, we see that models have a much higher

loss when predicting the named entity *Guinea coast* compared to having the full 16k maximal span. *Guinea coast* was mentioned 3584 tokens earlier, which indicates that long attention is often necessary to predict words mentioned in far away context. In general, we found that rare tokens and structural information about documents, such as section headings or document titles, required longer attention span to accurately predict.

Efficiency Advantages of Expire-Span Finally, we end with a brief discussion about why EXPIRE-SPAN is more efficient compared to existing architectures that focus on long context. First, Transformer-XL cannot adapt to the data at all, so it becomes slow and inefficient quite quickly as the span size increases. Adaptive-Span can adapt to the data and adjust its memory, but this memory size is fixed after training and does not have the dynamic adjustment of Expire-Span (where memory depends on local context even at inference time). Finally, the Compressive Transformer compresses past memories, but it compresses always at a fixed rate. The compression rate is an adjustable parameter, but aggressive compression potentially hurts performance. In contrast, EXPIRE-SPAN can expire irrelevant content, which both improves performance by focusing on salient information, and reduces the load on GPU memory and allows for faster processing per batch.

7. Conclusion

We present EXPIRE-SPAN, an operation that can be added to any attention mechanism to enable models to learn what to forget. By expiring irrelevant information, models can scale attention to tens of thousands of past memories. We highlight the strong performance of EXPIRE-SPAN in language modeling, reinforcement learning, object collision, and algorithmic tasks, and use it to attend over tens of thousands of past memories. The scalability and much greater efficiency of our proposed EXPIRE-SPAN method has strong potential for allowing models to be applied to more challenging, human-like tasks that would require expiration.

References

- Baevski, A. and Auli, M. Adaptive input representations for neural language modeling. In *ICLR*, 2019.
- Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate. In *ICLR*, 2015.
- Bahrack, H. P., Hall, L. K., and Da Costa, L. A. Fifty years of memory of college grades: Accuracy and distortions. *Emotion*, 8(1):13, 2008.
- Bai, S., Kolter, J. Z., and Koltun, V. Deep equilibrium models. In *Advances in Neural Information Processing Systems*, pp. 688–699, 2019.
- Beltagy, I., Peters, M. E., and Cohan, A. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In *NeurIPS*, 2020.
- Child, R., Gray, S., Radford, A., and Sutskever, I. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Choromanski, K., Likhoshesterov, V., Dohan, D., Song, X., Davis, J., Sarlos, T., Belanger, D., Colwell, L., and Weller, A. Masked language modeling for proteins via linearly scalable long-context transformers. *arXiv preprint arXiv:2006.03555*, 2020.
- Correia, G. M., Niculae, V., and Martins, A. F. Adaptively sparse transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2174–2184, 2019.
- Dai, Z., Yang, Z., Yang, Y., Carbonell, J. G., Le, Q. V., and Salakhutdinov, R. Transformer-xl: Attentive language models beyond a fixed-length context. In *ACL (1)*, pp. 2978–2988. Association for Computational Linguistics, 2019.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, 2019.
- Elman, J. Finding structure in time. *Cogn. Sci.*, 14:179–211, 1990.
- Fan, A., Gardent, C., Braud, C., and Bordes, A. Using local knowledge graph construction to scale seq2seq models to multi-document inputs. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4177–4187, 2019a.
- Fan, A., Lewis, M., and Dauphin, Y. Strategies for structuring story generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2650–2660, 2019b.
- Fan, A., Grave, E., and Joulin, A. Reducing transformer depth on demand with structured dropout. In *International Conference on Learning Representations*, 2020a.
- Fan, A., Lavril, T., Grave, E., Joulin, A., and Sukhbaatar, S. Addressing some limitations of transformers with feedback memory. *arXiv preprint arXiv:2002.09402*, 2020b.
- Glorot, X. and Bengio, Y. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256, 2010.
- Grave, E., Joulin, A., Cissé, M., and Jégou, H. Efficient softmax approximation for gpus. In *ICML*, 2017.
- Graves, A., Wayne, G., and Danihelka, I. Neural turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Izacard, G. and Grave, E. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*, 2020.
- Katharopoulos, A., Vyas, A., Pappas, N., and Fleuret, F. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, pp. 5156–5165. PMLR, 2020.

- Kitaev, N., Kaiser, L., and Levskaya, A. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2019.
- Lample, G., Sablayrolles, A., Ranzato, M., Denoyer, L., and Jégou, H. Large memory layers with product keys. In *Advances in Neural Information Processing Systems*, pp. 8548–8559, 2019.
- Loshchilov, I. and Hutter, F. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- Mahoney, M. Large text compression benchmark. URL: <http://www.mattmahoney.net/text/text.html>, 2011.
- Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016.
- Mikolov, T., Karafiát, M., Burget, L., Černocký, J., and Khudanpur, S. Recurrent neural network based language model. In *Eleventh annual conference of the international speech communication association*, 2010.
- Murre, J. and Dros, J. Replication and analysis of ebbinghaus’ forgetting curve. *PLoS ONE*, 10, 2015.
- Parisotto, E., Song, H. F., Rae, J. W., Pascanu, R., Gülçehre, Ç., Jayakumar, S. M., Jaderberg, M., Kaufman, R. L., Clark, A., Noury, S., Botvinick, M., Heess, N., and Hassel, R. Stabilizing transformers for reinforcement learning. In *ICML*, 2020.
- Peng, H., Pappas, N., Yogatama, D., Schwartz, R., Smith, N. A., and Kong, L. Random feature attention. *arXiv preprint arXiv:2103.02143*, 2021.
- Rae, J. W., Potapenko, A., Jayakumar, S. M., Hillier, C., and Lillicrap, T. P. Compressive transformers for long-range sequence modelling. In *International Conference on Learning Representations*, 2020.
- Rives, A., Goyal, S., Meier, J., Guo, D., Ott, M., Zitnick, C. L., Ma, J., and Fergus, R. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *bioRxiv*, 2019.
- Roller, S., Dinan, E., Goyal, N., Ju, D., Williamson, M., Liu, Y., Xu, J., Ott, M., Shuster, K., Smith, E. M., et al. Recipes for building an open-domain chatbot. *arXiv preprint arXiv:2004.13637*, 2020.
- Roy, A., Saffar, M., Vaswani, A., and Grangier, D. Efficient content-based sparse attention with routing transformers. *arXiv preprint arXiv:2003.05997*, 2020.
- Schlag, I., Irie, K., and Schmidhuber, J. Linear transformers are secretly fast weight memory systems. *arXiv preprint arXiv:2102.11174*, 2021.
- Shaw, P., Uszkoreit, J., and Vaswani, A. Self-attention with relative position representations. In *NAACL-HLT (2)*, 2018.
- Sukhbaatar, S., Szlam, A., Synnaeve, G., Chintala, S., and Fergus, R. Mazebase: A sandbox for learning from games. *ArXiv*, abs/1511.07401, 2015a.
- Sukhbaatar, S., Szlam, A., Weston, J., and Fergus, R. End-to-end memory networks. In *NIPS*, 2015b.
- Sukhbaatar, S., Grave, É., Bojanowski, P., and Joulin, A. Adaptive attention span in transformers. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 331–335, 2019a.
- Sukhbaatar, S., Grave, E., Lample, G., Jégou, H., and Joulin, A. Augmenting self-attention with persistent memory. *arXiv preprint arXiv:1907.01470*, 2019b.
- Tay, Y., Dehghani, M., Bahri, D., and Metzler, D. Efficient transformers: A survey. *arXiv preprint arXiv:2009.06732*, 2020.
- Tay, Y., Dehghani, M., Abnar, S., Shen, Y., Bahri, D., Pham, P., Rao, J., Yang, L., Ruder, S., and Metzler, D. Long range arena: A benchmark for efficient transformers. In *International Conference on Learning Representations*, 2021.
- Urbanek, J., Fan, A., Karamcheti, S., Jain, S., Humeau, S., Dinan, E., Rocktäschel, T., Kiela, D., Szlam, A., and Weston, J. Learning to speak and act in a fantasy text adventure game. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 673–683, 2019.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Wang, S., Li, B., Khabza, M., Fang, H., and Ma, H. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- Wixted, J. T. The psychology and neuroscience of forgetting. *Annu. Rev. Psychol.*, 55:235–269, 2004.
- Wu, C.-Y., Feichtenhofer, C., Fan, H., He, K., Krahenbuhl, P., and Girshick, R. Long-term feature banks for detailed video understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 284–293, 2019.

Wu, F., Fan, A., Baevski, A., Dauphin, Y., and Auli, M. Pay less attention with lightweight and dynamic convolutions. In *International Conference on Learning Representations*, 2018.

Ye, Z., Guo, Q., Gan, Q., Qiu, X., and Zhang, Z. Bp-transformer: Modelling long-range context via binary partitioning. *arXiv preprint arXiv:1911.04070*, 2019.

Zhang, J., Luan, H., Sun, M., Zhai, F., Xu, J., Zhang, M., and Liu, Y. Improving the transformer translation model with document-level context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 533–542, 2018.

A. Appendix

A.1. Additional Method Details

Position Embedding Relative position embeddings (Shaw et al., 2018) make it possible to condition on the order of inputs by modifying the attention to $a_{ti} = \text{Softmax}(\mathbf{q}_t^\top \mathbf{k}_i + \mathbf{q}_t^\top \mathbf{p}_{t-i})$. However, because this second term is computed for the whole block in parallel for efficiency, it can become expensive for a large L even when the average memory size $|C_t|$ is small. Our solution is to remove position embeddings from older memories $i < t - K$ (where K is the block size), which empirically does not affect performance. The computational complexity of the position embeddings is then $\mathcal{O}(K)$, thus allowing us to increase the maximum span L . This modification makes training EXPIRE-SPAN more efficient, but does not improve accuracy.

Training with Small Initial Spans EXPIRE-SPAN scales to long attention spans as it quickly learns to expire irrelevant content. However, at the beginning of training, the long span can use large quantities of GPU memory. To circumvent this, we initialize the bias term b with a negative value. This prevents large memory usage at the beginning of training, after which the model quickly learns to expire and the memory usage is no longer problematic.

A.2. Additional Experimental Results

Efficiency for Instruction Task We include a comparison of EXPIRE-SPAN to Adaptive-Span and Compressive Transformer in Table 5 and show that EXPIRE-SPAN has stronger performance, is faster, and saves GPU memory.

Wikitext-103 Language Modeling The Wikitext-103 word-level language modeling benchmark (Merity et al., 2016) consists of a collection of Wikipedia articles and a fixed vocabulary size of 270K. We set the max attention span for EXPIRE-SPAN to 8K. We compare EXPIRE-SPAN to existing work in Table 6 and show that even fairly small models trained with EXPIRE-SPAN achieve competitive results. Next, we analyze the performance of EXPIRE-SPAN on Wikitext-103 as the memory size increases. We compare to a Transformer-XL model in Figure 10 — even with far smaller memory, EXPIRE-SPAN performs much better.

Expire-span Performance and Analysis on Enwik8 In Figure 11, we analyze multiple layers of a trained model and show that different layers memorize different types of information. Several layers retain summarizing information about sentences or sections by increasing the expire-spans of spaces, new lines, and section titles.

Additionally, we did an ablation by running our large Expire-Span model without LayerDrop. Its validation performance

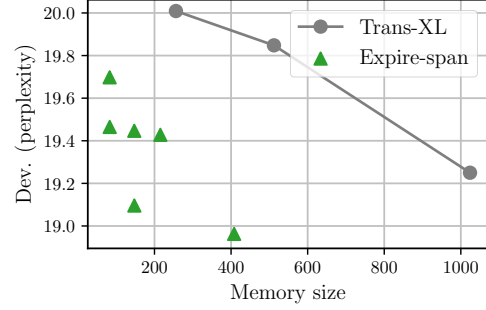


Figure 10. Performance as a function of Memory Size on Wikitext-103

dropped from 0.98bpb to 1.00bpb.

Importance of Structured Dropout for Regularization

We analyze the importance of structured dropout to regularize the large memory capacity provided by EXPIRE-SPAN. In an experiment on enwiki8, Figure 12 shows that loss on a portion of validation data was incredibly large. This part corresponds to a 66K token long table. We hypothesize that the model likely never encountered such a table during training. During validation, this caused all non-table tokens to expire. Without regularizing the model memory size during training, the model can easily overfit.

Colliding Objects, An Easy Version

We experiment with an easier version of the Colliding Objects task where objects do not have colors. The model has to predict either the last collision, or a mapping of the last 3 collisions. In contrast to the harder task, there are no color switches and any collision prediction is valid. As this version is less memory intensive, the EXPIRE-SPAN model almost solves it with a shorter maximum span, as shown in Table 7.

A.3. Additional Implementation Details

A.3.1. REINFORCEMENT LEARNING TASKS

We used MazeBase (Sukhbaatar et al., 2015a) to construct tasks in grid world. Agents can observe its surrounding 3×3 area and move in the four cardinal directions. Every objects and their properties are described by words such as “agent”, “block”, “blue”, etc. Thus, the input to the model is a binary tensor of size $3 \times 3 \times \text{vocabulary-size}$.

We train 2-layer Transformers with 64 hidden units using actor-Critic algorithm. We used a BPTT length of 100, and an entropy cost of 0.0005.

Corridor Task The corridor length is sampled from $\mathcal{U}(3, 200)$. All models are trained for 100M steps. We used RMSProp optimizer with a learning rate of 0.0001 and a batch size of 64. For the expire-span models, we set the

	Model	Performance	GPU Memory (GB)	Time/Batch (ms)
Instruction Task	Compressive Transformer	71% Acc	10	210
	Adaptive-Span	64% Acc	14	240
	EXPIRE-SPAN	74% Acc	8	90

Table 5. **Efficiency of EXPIRE-SPAN.** We report peak GPU memory usage and per-batch training time, fixing the batch size. We evaluate the mean pooling version of the Compressive Transformer.

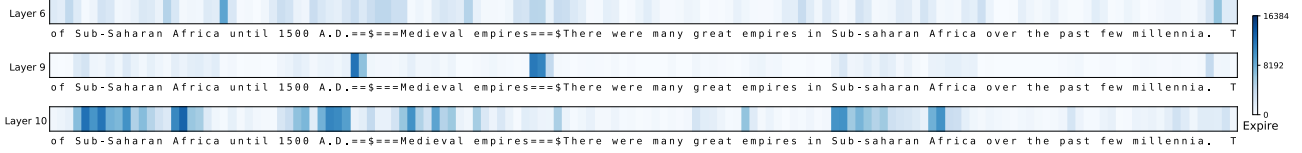


Figure 11. **Per-Layer EXPIRE-SPAN values on Enwik8.** We visualize the expire-spans of different layers: layer 6 gives long span to spaces, layer 9 memorizes special tokens like newlines and section titles, and layer 10 retains named entities in memory.

Model	Params	Test
DEQ-Trans. (Bai et al., 2019)	110M	23.3
Trans-XL (Dai et al., 2019)	257M	18.3
Feedback Trans. (Fan et al., 2020b)	77M	18.3
Trans.+LayerDrop (Fan et al., 2020a)	423M	17.7
Compressive Trans. (Rae et al., 2020)	277M	17.1
Routing Trans. (Roy et al., 2020)	-	15.8
EXPIRE-SPAN	140M	19.6

Table 6. **Wikitext-103 Results.** We report perplexity on test.

Model	Maximum Span	Test Error (%)
Transformer-XL	1k	39.1
	1k	19.5
EXPIRE-SPAN	2k	9.1
	4k	3.2

Table 7. **Colliding Objects Results.** We report test error.

maximum span L to 200, the loss coefficient α to $5e-6$, and the ramp length R to 16.

Multi-Room Portal In this task, there are 50 rooms sequentially connected together. Each room is 5×5 in size, and have two doors with different colors. If agent go to the correct door, it will be teleported to the next room, but if it is the wrong door, the agent will be teleported back to the first room and have to start over. Which of the two doors is correct in each room is randomly decided and fixed throughout the episode. This information is not visible to the agent, thus can only be discovered by trial and error within each episode. The current room number is visible to the agent.

When the agent successfully transitions from the k -th room to the next, it receives a reward of $0.1k$. The episode ends if the agent makes two mistakes in the same room, reaches the last room, or when the number of steps reach 1000. A reward discount of 0.98 is used. All models are trained with Adam optimizer with a learning rate of $5e-4$, and a batch size of 1024, with gradients are clipped at 0.1. We set $L = 100$, $R = 16$ and $\alpha = 1e-6$ for the expire-span models.

A.3.2. INSTRUCTION TASK IN LIGHT

We train 6-layer models with a hidden size of 512 and 8 attention heads. To train, we use the Adam optimizer with a learning rate of $7e-4$ and 8000 warmup updates. We set the expire-span ramp R to 64 and the expire-span loss α to $2e-6$.

A.3.3. COLLISION TASK

At the start of the simulation, each particle samples a Gaussian Normal velocity and position uniform inside a 16×16 grid. At each time step the particles' position is updated by adding its velocity (unless it would go off the grid, in which case its velocity is re-sampled). There are 5 different colors, and a particle can change its color randomly at each step with 0.05 probability. A collision happens when the two particles have the same rasterized locations, but it does not affect the movement.

Given a question specifying two colors, the task is to report in which of the four quadrants of the grid the last collision of the specified-colors occurred. To make the task easier to learn, 40% of the queries will have the matching colors as the last collision.

The model is given an input sequence of tokens that has 8 entries per timestep. The first 4 are the rounded and rasterized (x, y) locations of the two particles, and next 2

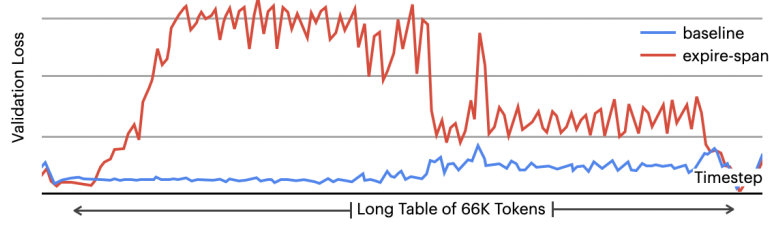


Figure 12. **Extreme Overfitting** on part of validation occurs without proper regularization.

are tokens representing the colors of the particles. The last 2 entries are “question” tokens that specify the colors of the collision. The output sequence has a token for each quadrant. We generate 50M steps for training, which equals to 400M entries.

Easy Version: The particles have no color in this version. There are two types of questions, in which the task is to report either: (i) in which of the four quadrants of the grid the last collision occurred, or (ii) the label mapping of the last 3 collisions.

A.3.4. LANGUAGE MODELING DETAILS

Enwik8 Our small model has 12 layers with a hidden size of 512 and 8 attention heads. To train, we use Adam optimizer with a learning rate of $7e-4$, a batch size of 512, a block size of 512 and 8000 warmup updates. All models are trained for 100k updates. The model in Table 2 is further fine-tuned for another 10k updates with a 10x smaller LR. The baseline models used for comparison are the same size model following the same training protocol.

The large EXPIRE-SPAN model Table 2 is a 24-layer model with a hidden size of 768 and 4096 feedforward units. It is trained with a learning rate of $4e-4$ and 16k warmup steps. In addition to 0.5 dropout, we also add 0.2 layer-drop. The EXPIRE-SPAN parameters are $L = 32k$, $\alpha = 3e-7$, and $R = 128$. We used the version of Eq. 6 due to the very long maximum span.

Character-level PG-19 Besides the maximum span, all model parameters and training parameters were held constant. Each model had 12 layers, a hidden size of 512, a feedforward size of 2048, 8 attention heads, and processed a block of 512 characters at a time. We initialized the weights using a uniform distribution as described by (Glorot & Bengio, 2010), used dropout of 0.2, clipped the gradients at 0.3, warmed up the learning rate linearly for 8000 steps, and used cosine annealing to decay the learning rate after warmup (Loshchilov & Hutter, 2016). For the EXPIRE-SPAN models, we used a ramp of $R = 128$ and an expiration loss coefficient of $\alpha = 1e-6$ ($3e-7$) for $L = 8k$ ($16k$).

Wikitext-103 All models have 8 layers and 1024 hidden units (4096 in feedforward layers). In addition to the dropout of 0.3 applied to attention and ReLU activation, outputs from the embedding layer and the last layer had a dropout of 0.2. We used the adaptive input (Baevski & Auli, 2019) and the adaptive softmax (Grave et al., 2017) for reducing the number of parameters within word embeddings. The models are trained for 300k updates with a block size of 256, and gradients are clipped at 0.1. The other hyperparameters are the same as the small Enwik8 experiments.