## **ABC: Attention with Bounded-Memory Control**

Hao Peng Jungo Kasai♠ Nikolaos Pappas★\* Dani Yogatama\* Zhaofeng Wu<sup>♦</sup>\* **Roy Schwartz**<sup>\(\nabla\)</sup> Noah A. Smith♠♦ Lingpeng Kong<sup>♦</sup>

Paul G. Allen School of Computer Science & Engineering, University of Washington <sup>©</sup>School of Computer Science & Engineering, Hebrew University of Jerusalem ◆Department of Computer Science, The University of Hong Kong

{hapeng, jkasai, npappas, zfw7, nasmith}@cs.washington.edu dyogatama@deepmind.com, lpk@cs.hku.hk

Transformer architectures have achieved stateof-the-art results on a variety of natural language processing (NLP) tasks. However, their attention mechanism comes with a quadratic complexity in sequence lengths, making the computational overhead prohibitive, especially for long sequences. Attention context can be seen as a random-access memory with each token taking a slot. Under this perspective, the memory size grows linearly with the sequence length, and so does the overhead of reading from it. One way to improve the efficiency is to bound the memory size. We show that disparate approaches can be subsumed into one abstraction, attention with bounded-memory control (ABC), and they vary in their organization of the memory. ABC reveals new, unexplored possibilities. First, it connects several efficient attention variants that would otherwise seem distinct. Second, this abstraction gives new insights—an established approach (Wang et al., 2020b) previously thought to not be applicable in causal attention, actually is. Last, we present a new instance of ABC, which draws inspiration from existing ABC approaches, but replaces their heuristic memory-organizing functions with a *learned*, contextualized one. Our experiments on language modeling, machine translation, and masked language model finetuning show that our approach outperforms previous efficient attention models; compared to strong transformer baselines, it significantly improves the inference time and space efficiency with no or negligible accuracy loss.

#### 1 Introduction

Transformer architectures are now central in natural language processing (Vaswani et al., 2017). They rely on the attention mechanism (Bahdanau et al., 2015) to contextualize the input. The context can be seen as a random access memory whose size linearly grows with the sequence length; each query

roy.schwartz1@mail.huji.ac.il
Abstract reads from it using a softmax-normalized linear combination, with overhead linear in the memory size. This amounts to a quadratic complexity overall, making transformers' computational overhead prohibitive, especially for long sequences.

> One way to improve attention's efficiency is to bound its memory size. Imposing a constantsized constraint over the memory ensures that reading from it has constant time and space overhead, yielding a linear overall complexity in sequence lengths. This is in fact a common strategy adopted by several recent works. In this work, we show that some of these works are closely connected in ways that, to date, have gone unremarked. We propose attention with bounded-memory control (ABC), a unified abstraction over them. In ABC, constant-sized memories are organized with various control strategies, e.g., induced from heuristic patterns (Beltagy et al., 2020; Zaheer et al., 2020; Ainslie et al., 2020; Rae et al., 2020, inter alia), locality assumptions (Parmar et al., 2018; Liu et al., 2018), or positions (Wang et al., 2020b).

> These strategies, by and large, are "contextagnostic." In response to this, we propose ABC<sub>MLP</sub>, a particular instance of ABC that learns a contextualized control strategy from data. Specifically, ABC<sub>MLP</sub> uses a neural network to determine how to store each token into the memory (if at all). Compared to previous bounded-memory models, it strikes a better trade-off between accuracy and efficiency: controlling for the accuracy, ABC<sub>MLP</sub> can get away with much smaller memory sizes.

> ABC models (including ABC<sub>MLP</sub>) come with a linear complexity in sequence lengths, and admit recurrent computation graphs in causal attention (self-attention over the prefix). Therefore they are appealing choices in a variety of applications, including text encoding, language modeling and text generation. This leads to a surprising finding. Linformer (Wang et al., 2020b), an established efficient attention method, was previously thought not

<sup>\*</sup>This work was done while Zhaofeng Wu and Nikolaos Pappas were at the University of Washington.

to be applicable in causal attention or autoregressive decoding (Tay et al., 2020). Through the ABC view, we show that it actually *is*, and achieves competitive performance in our machine translation experiments.

ABC connects existing models that would otherwise seem distinct, reveals new insights into established methods, and inspires new efficient attention architectures. We explore its applications in transformers, as a drop-in substitute for the canonical softmax attention. ABC offers a novel lens that can help future research in the analysis of transformers, where the theoretical insights are still catching up with empirical success. Experiments on language modeling, machine translation, and masked language model finetuning show that our ABC<sub>MLP</sub> model outperforms previous ABC approaches in accuracy with a much smaller memory size. Compared to the strong transformer baseline, ABC<sub>MLP</sub> achieves a significant speedup and memory savings at inference time, with no or negligible accuracy loss. The efficiency improvements are more prominent for long sequences, suggesting that the asymptotic savings are even more appealing in applications involving long sequences. We release our code at https://github.com/Noahs-ARK/ABC.

#### 2 An Outer-Product View of Attention

This section presents our outer-product memory perspective of attention, which allows for a smooth transition to later discussion.

In attention, a sequence of **queries**  $\{\mathbf{q}_i\}_{i=1}^N$  attend to a **memory** with N slots, each storing a **key** and **value** pair:  $\mathbf{K} = [\mathbf{k}_1, \dots, \mathbf{k}_N]^\top, \mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_N]^\top \in \mathbb{R}^{N \times d}$ . Query  $\mathbf{q}$  reads from the memory using a softmax-normalized linear combination, producing a d-dimensional vector:

$$\operatorname{attn}(\mathbf{q}, \{\mathbf{k}_i\}, \{\mathbf{v}_i\}) = \mathbf{V}^{\top} \operatorname{softmax}(\mathbf{K}\mathbf{q}).$$
 (1)

This takes  $\mathcal{O}(N)$  time and space. When the attention with N queries can be parallelized (e.g., in text encoding), it takes linear time and quadratic space; when it *cannot* be (e.g., in decoding), it takes quadratic time and linear space.

The memory can be equivalently represented as sums of vector outer products:  $\mathbf{K} = \mathbf{I}\mathbf{K} = \sum_{i=1}^{N} \mathbf{e}_i \otimes \mathbf{k}_i$ ,  $\mathbf{V} = \sum_{i=1}^{N} \mathbf{e}_i \otimes \mathbf{v}_i$ . I is the identity matrix, and  $\otimes$  denotes the outer product:  $[\mathbf{x} \otimes \mathbf{v}_i]$ 

 $\mathbf{y}]_{i,j} = x_i y_j$ . *N*-dimensional vectors  $\{\mathbf{e}_i\}$  form the standard basis:  $\mathbf{e}_i$  has the *i*th element being one and others zeros. We can view  $\mathbf{e}_i$  as **control vectors** that determine where to store  $\mathbf{k}_i$  and  $\mathbf{v}_i$ :

$$\mathbf{e}_{i} \otimes \mathbf{k}_{i} = \left[\underbrace{0, \dots, 0}_{i-1}, 1, \underbrace{0, \dots, 0}_{N-i}\right]^{\top} \otimes \mathbf{k}_{i}$$
$$= \left[\underbrace{\mathbf{0}}_{d \times (i-1)}; \mathbf{k}_{i}; \underbrace{\mathbf{0}}_{d \times (N-i)}\right]^{\top}.$$
 (2)

The N-by-d matrix on the last line has its *i*th row being  $\mathbf{k}_i^{\top}$  and all others zeros; in this sense,  $\mathbf{k}_i$  is stored in the *i*th slot by  $\mathbf{e}_i$ , not affecting others.

## 3 Attention with Bounded Memory

A straightforward way to improve attention's efficiency is to bound its memory size. Our outer-product view of attention provides a straightforward way to devise this, by replacing  $\{\mathbf{e}_i\}$  with control vectors that select  $n \ll N$  vectors to attend to. We dub this approach attention with bounded-memory control (ABC). Concretely, let  $\widetilde{\mathbf{K}}, \widetilde{\mathbf{V}} \in \mathbb{R}^{n \times d}$  denote a constant-size memory with n slots, with n set a priori.

$$\widetilde{\mathbf{K}} = \sum_{i=1}^{N} \phi_i \otimes \mathbf{k}_i, \quad \widetilde{\mathbf{V}} = \sum_{i=1}^{N} \phi_i \otimes \mathbf{v}_i.$$
 (3)

 $\{\phi_i \in \mathbb{R}^n\}_{i=1}^N$  denotes a sequence of **control** vectors. The output is calculated by attending to  $\widetilde{\mathbf{K}}$  and  $\widetilde{\mathbf{V}}$ : ABC  $(\mathbf{q}, \{\mathbf{k}_i\}, \{\mathbf{v}_i\}, \{\phi_i\}) =$ 

$$\widetilde{\mathbf{V}}^{\top} \operatorname{softmax} \left( \widetilde{\mathbf{K}} \mathbf{q} \right).$$
 (4)

We will discuss various ways to construct  $\{\phi_i\}$  in the subsequent sections. Reading from the memory takes a constant  $\mathcal{O}(n)$  time and space; therefore ABC's overall complexity is  $\mathcal{O}(Nn)$ , linear in the sequence length.<sup>2</sup>

Eq. 3 offers an equivalent recurrent computation, which is particularly useful in causal attention where only the prefix is looked at,

$$\widetilde{\mathbf{K}}_{t+1} = \widetilde{\mathbf{K}}_t + \phi_{t+1} \otimes \mathbf{k}_{t+1}, \tag{5}$$

likewise for  $V_t$ .  $K_t$  and  $V_t$  can be seen as the **recurrent hidden state** that encodes the prefix.

In what follows, we study several existing efficient attention approaches and show that they are in fact instances of the ABC abstraction.

<sup>&</sup>lt;sup>1</sup>The number of queries and key-value pairs may differ, e.g., in the cross attention of a sequence-to-sequence model.

<sup>&</sup>lt;sup>2</sup>Using bounded memory distinguishes ABC from softmax attention. If growing-size memory *were* allowed (n=N), an ABC with  $\phi_i=\mathbf{e}_i$  would fall back to softmax attention.

#### 3.1 Linformer

Linformer (Wang et al., 2020b) is an established efficient transformer variant that has proven successful in masked language modeling and text encoding. It assumes fixed-length inputs and learns a low-rank approximation of the attention weights. A learned n-by-N matrix  $\mathbf{W}^{\mathrm{LF}}$  down projects the N-by-d dimensional keys and values along the timestep dimension, to an n-by-d memory:  $\widetilde{\mathbf{K}}^{\mathrm{LF}} = \mathbf{W}^{\mathrm{LF}}\mathbf{K}$ ,  $\widetilde{\mathbf{V}}^{\mathrm{LF}} = \mathbf{W}^{\mathrm{LF}}\mathbf{V}$ ; they are then used for attention computation with Eq. 4. This yields a linear complexity in the input length. Linformer is an ABC instance with  $\phi_i^{\mathrm{LF}} = \mathbf{W}_{:,i}^{\mathrm{LF}}$  (ith column), and in this sense, it learns a control vector for each position.

Previous works have noted that Linformer *cannot* be efficiently applied in causal attention (Table 1 of Tay et al., 2020). Indeed, it is less straightforward to avoid mixing future with the past when projecting along the timestep dimension. ABC reveals that, in fact, Linformer *is* applicable in causal attention. Like all ABC models, it admits a linear-complexity recurrent computation (Eq. 5):  $\widetilde{\mathbf{K}}_{t+1}^{\mathrm{LF}} = \widetilde{\mathbf{K}}_t + \phi_{t+1}^{\mathrm{LF}} \otimes \mathbf{k}_{t+1}$ . This confirms ABC's benefits: it reveals new insights about existing models and reassesses their applications and impact. Our experiments show that Linformer achieves competitive performance in machine translation.

## 3.2 Clustering-Based Attention

Improving attention's efficiency with clustering has received an increasing amount of interest (Kitaev et al., 2020; Roy et al., 2020; Wang et al., 2020a, inter alia). ABC bears interesting connections to clustering-based methods. Here we discuss an approach that closely follows Vyas et al. (2020), except that it clusters keys and values instead of queries, and only attends to the centroids to reduce the effective context size. Formally, keys and values are grouped into n < N clusters  $\{\tilde{\mathbf{k}}_j^{\text{CL}}\}_{j=1}^n$ ,  $\{\tilde{\mathbf{v}}_j^{\text{CL}}\}_{j=1}^n$ . Let an N-by-n binary matrix  $\mathbf{M}$  denote the cluster membership shared between keys and values.  $M_{i,j} = 1$  iff.  $\mathbf{k}_i$  is assigned to cluster  $\tilde{\mathbf{k}}_j^{\text{CL}}$  and  $\mathbf{v}_i$  to  $\tilde{\mathbf{v}}_j^{\text{CL}}$ . The jth centroid for the keys is

$$\widetilde{\mathbf{k}}_{j}^{\text{CL}} = \sum_{i=1}^{N} \frac{M_{i,j}}{\sum_{\ell=1}^{N} M_{\ell,j}} \mathbf{k}_{i};$$
 (6)

likewise for the values. It then attends over the centroids using Eq. 4, with  $\widetilde{\mathbf{K}}^{\text{CL}} = [\widetilde{\mathbf{k}}_1^{\text{CL}}, \dots, \widetilde{\mathbf{k}}_n^{\text{CL}}]^{\top} =$ 

$$\sum_{j=1}^{n} \mathbf{e}_{j} \otimes \widetilde{\mathbf{k}}_{j}^{\text{CL}} = \sum_{j=1}^{n} \mathbf{e}_{j} \otimes \sum_{i=1}^{N} \frac{M_{i,j}}{\sum_{\ell=1}^{N} M_{\ell,j}} \mathbf{k}_{i}$$
$$= \sum_{i=1}^{N} \left( \sum_{j=1}^{n} \mathbf{e}_{j} \frac{M_{i,j}}{\sum_{\ell=1}^{N} M_{\ell,j}} \right) \otimes \mathbf{k}_{i}.$$

The last line indicates that this model is an instance of ABC:  $\phi_i = \sum_{j=1}^n (M_{i,j}/\sum_{\ell=1}^N M_{\ell,j}) \mathbf{e}_j$ . The stack of centroids can be seen as the constant-size memory. Putting aside the clustering overhead (i.e., constructing  $\mathbf{M}$  and computing centroids), it has a linear complexity in the sequence length.

#### 3.3 Sliding-Window Attention

In some applications, being able to remove entries from the memory can be beneficial: clearing up older context frees slots for more recent ones, promoting a locality inductive bias. ABC offers the capability to do so, if augmented with an additional matrix multiplication. We use the sliding-window attention as an example.

Attending to the most recent n input tokens (Beltagy et al., 2020; Zaheer et al., 2020; Sukhbaatar et al., 2021, inter alia) can be seen as a first-in-first-out queue that "pops" out the oldest token while "pushing" in the most recent one:  $\widetilde{\mathbf{K}}_t^{\mathrm{WD}} = [\mathbf{k}_{t-n+1},...,\mathbf{k}_t]^{\top}$ . The pop operation can be achieved by multiplying an n-by-n upper shift matrix:  $U_{i,j} = \delta_{i+1,j}$ , with  $\delta$  being the Kronecker delta (i.e.,  $\mathbf{U}$  has ones only on the superdiagonal and zeros elsewhere). Left-multiplying  $\mathbf{U}$  against  $\widetilde{\mathbf{K}}_t^{\mathrm{WD}}$  shifts its rows one position up, with zeros appearing in the last:

$$\mathbf{U}\widetilde{\mathbf{K}}_{t}^{\mathrm{WD}} = \mathbf{U} \left[ \underbrace{\mathbf{k}_{t-n+1}, \dots, \mathbf{k}_{t}}_{n} \right]^{\top}$$
$$= \left[ \underbrace{\mathbf{k}_{t-n+2}, \dots, \mathbf{k}_{t-1}, \mathbf{k}_{t}}_{n-1}, \mathbf{0} \right]^{\top} \in \mathbb{R}^{n \times d}.$$

Then the most recent token can be put into the slot freed up:  $\widetilde{\mathbf{K}}_{t+1}^{\text{WD}} = \mathbf{U}\widetilde{\mathbf{K}}_{t}^{\text{WD}} + \mathbf{e}_{n} \otimes \mathbf{k}_{t+1}$ . U and  $\phi_{t} = \mathbf{e}_{n}$  ensure a first-in-first-out queue. Dilated and stride convolution patterns (Beltagy et al., 2020) can be similarly recovered (§A.4).

Recurrently multiplying U simulates the discrete pop operation (Grefenstette et al., 2015; Joulin and Mikolov, 2015; Yogatama et al., 2018) in a differentiable way. This is reminiscent of recurrent neural networks, while in this case U is *never* updated as

<sup>&</sup>lt;sup>3</sup>We use  $\widetilde{\mathbf{k}}_{j}^{\text{CL}}$  to denote both the *j*th cluster and its centroid.

parameters. It is exciting to explore learning U, but is beyond the scope of this work.

**Discussion.** Besides the models discussed above, certain variants of Rae et al. (2020) and sparse attention patterns (local-to-global attention; Beltagy et al., 2020; Zaheer et al., 2020; Ainslie et al., 2020) can also be seen as instances of ABC (§A). ABC provides a unified perspective of them, and at the same time points out their limitations: their control strategies are context-agnostic. In response to this, in §4 we propose to learn a contextualized strategy from data. Table 1 analyzes various ABC models, and Table 2 details their complexity.

## 4 Learned Memory Control

The ABC abstraction connects several existing approaches that would otherwise seem distinct. This inspires the design of new architectures. We hypothesize that learning a contextualized strategy can achieve better performance. This section introduces ABC<sub>MLP</sub>. It parameterizes  $\phi$  with a single-layer multi-layer perceptron (MLP) that takes as input the token's representation  $\mathbf{x}_i$ , and determines which slots to write it into and how much.

$$\alpha_i = \exp(\mathbf{W}_{\phi}\mathbf{x}_i), \quad \phi_i = \alpha_i / \sum_{j=1}^N \alpha_j.$$
 (7)

Matrix  $\mathbf{W}_{\phi}$  is learned.  $\exp$  is an elementwise activation function. The motivation is to allow for storing a "fractional" (but *never* negative) amount of input into the memory.<sup>4</sup> Using a non-negative activation, however, has a drawback: the scales of  $\sum_i \phi_i \otimes \mathbf{k}_i$  and  $\sum_i \phi_i \otimes \mathbf{v}_i$  would grow with the sequence lengths, making training less stable. To overcome this, we divide  $\alpha_i$  vectors by their sum. This functions as normalization and aims to offset the impact of varying sequence lengths.<sup>5</sup> It admits the recurrent computation graph as in Eq. 5, and has a linear complexity in the sequence length.

A key design choice of ABC<sub>MLP</sub> is that its  $\phi_i$  depends *only* on current input  $x_i$ . This helps (1) keep the recurrent computation efficient in practice (Lei et al., 2018), and (2) make it applicable

in not only encoder self-attention and cross attention, but also causal attention. Concurrently to this work, Goyal et al. (2021) and Ma et al. (2021) also proposed methods to learn contextualized control. They compute  $\phi_i$  from *previous* layer's memory, revealing the full sequence to the control vectors. As a result, these two approaches are *unsuitable* for causal attention.<sup>6</sup>

ABC<sub>MLP</sub>, as other ABC models, can be used as a drop-in replacement for the canonical softmax attention, and we apply its multihead variant in transformers. With proper parameter sharing, the number of additional parameters ABC<sub>MLP</sub> incurs is small: inspired by Wang et al. (2020b), we tie  $\phi$ -MLP's parameters across different layers, which adds less than 1% parameters to the models.

**ABC<sub>MLP</sub>: context-agnostic then context-dependent attention.** We now dissect ABC<sub>MLP</sub> and show that it can be seen as a cascade of two attention mechanisms: one with a learned context-agnostic "pseudo query" followed by one with a context-dependent query. Our analysis starts with a one-dimensional example; the conclusion generalizes to higher-dimensional cases.

**Example 1.** Consider ABC<sub>MLP</sub> with a *single* memory slot (n = 1). It is parameterized with a learned vector  $\mathbf{w}_{\phi}$ , and  $\phi_i = \exp(\mathbf{w}_{\phi} \cdot \mathbf{x}_i) / \sum_{j=1}^{N} \exp(\mathbf{w}_{\phi} \cdot \mathbf{x}_j)$ . Since  $\phi_i$  is a scalar here,  $\phi_i \otimes \mathbf{k}_i = \phi_i \mathbf{k}_i^{\top}$ .

$$\widetilde{\mathbf{K}}^{\top} = \sum_{i=1}^{N} (\phi_i \otimes \mathbf{k}_i)^{\top}$$

$$= \sum_{i=1}^{N} \frac{\exp(\mathbf{w}_{\phi} \cdot \mathbf{x}_i)}{\sum_{j=1}^{N} \exp(\mathbf{w}_{\phi} \cdot \mathbf{x}_j)} \mathbf{k}_i$$

$$= \operatorname{attn} \left( \mathbf{w}_{\phi}, \{\mathbf{x}_i\}_{i=1}^{N}, \{\mathbf{k}_i\}_{i=1}^{N} \right).$$

In other words,  $\widetilde{\mathbf{K}}$  uses  $\mathbf{w}_{\phi}$  as a "pseudo-query" to attend to  $\{\mathbf{x}_i\}$  and  $\{\mathbf{k}_i\}$ . Likewise,  $\widetilde{\mathbf{V}}^{\top} = \operatorname{attn}(\mathbf{w}_{\phi}, \{\mathbf{x}_i\}_{i=1}^N, \{\mathbf{v}_i\}_{i=1}^N)$ . Despite its similarity to the standard softmax attention, Example 1 has a more efficient linear complexity in sequence lengths.  $\mathbf{w}_{\phi}$ 's being context-independent is the key to the savings. Table 2 details its complexity.

Example 1's conclusion generalizes to higherdimensional cases: the jth dimension of  $\{\phi_i\}$  attends to  $\{\mathbf{x}_i\}$  and  $\{\mathbf{k}_i\}$  using the jth row of  $\mathbf{W}_{\phi}$ as the context-independent pseudo-query; n such attention mechanisms run in parallel, stacking the

<sup>&</sup>lt;sup>4</sup>We experiment with other activations in §C.2.

<sup>&</sup>lt;sup>5</sup>Here encoder self-attention or cross attention is assumed, and the normalization sums over the entire sequence. Causal attention is slightly different, normalizing by the sum over the prefix instead:  $\phi_i = \alpha_i / \sum_{j=1}^i \alpha_j$ . This does *not* require access to future tokens. §B.1 details a linear complexity computation graph of causal  $\phi_i$ .

<sup>&</sup>lt;sup>6</sup>Both are instances of ABC (§A.5). Ma et al. (2021) resorts to a variant of Katharopoulos et al. (2020) for causal attention.

Model	Section	$oldsymbol{\phi}_t$	Mem. Control
Sliding-window	§3.3	$\mathbf{e}_n$	$\widetilde{\mathbf{K}}_{t+1} = \mathbf{U}\widetilde{\mathbf{K}}_t + \phi_{t+1} \otimes \mathbf{k}_{t+1}$
Linformer	§3.1	$\mathbf{W}^{ ext{LF}}_{::t}$	
L2G Pattern	§A.1	$\mathbf{e}_i$ if $\mathbf{x}_t$ is the <i>i</i> th global token	
$ABC_{RD}$	§A.2	$\mathbf{e}_{i_t}$ , where $i_t \sim \mathrm{unif}\{1, n\}$	$\widetilde{\mathbf{K}}_{t+1} = \widetilde{\mathbf{K}}_t + \phi_{t+1} \otimes \mathbf{k}_{t+1}$
Comp. Trans.	§A.3	$\mathbf{e}_{\mid nt/N\mid}$	$\mathbf{K}_{t+1} \equiv \mathbf{K}_t + \mathbf{\varphi}_{t+1} \otimes \mathbf{K}_{t+1}$
Clustering	§3.2	$\sum_{j=1}^{n} \left( M_{t,j} / \sum_{\ell=1}^{N} M_{\ell,j} \right) \mathbf{e}_j$	
$ABC_{MLP}$	§4	$\exp{(\hat{\mathbf{W}}_{oldsymbol{\phi}}\mathbf{x}_t)}/\sum_{i=1}^t \exp{(\mathbf{W}_{oldsymbol{\phi}}\mathbf{x}_t)}$	

Table 1: A comparison of different ABC models. N denotes the sequence length, and n the memory size.  $\phi_t$  denotes the memory control vector for  $\mathbf{k}_t$  and  $\mathbf{v}_t$ , and unif is the discrete uniform distribution.

	Time Complexity			Space Complexity		
Model	Mem.	Per Query	Overall	Mem.	Per Query	Overall
Softmax Attention	-	$\mathcal{O}(N)$	$\mathcal{O}(N^2)$	-	$\mathcal{O}(N)$	$\mathcal{O}(N^2)$
ABC	$\mathcal{O}(N)$	$\mathcal{O}(n)$	$\mathcal{O}(nN)$	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(nN)$

Table 2: ABC's time and space complexity in sequence length against the softmax attention's. "Mem." indicates the time and space needed for calculating and storing memory  $\widetilde{\mathbf{K}}, \widetilde{\mathbf{V}}$ . N denotes the sequence length, and n the memory size. The time complexity analysis assumes that the softmax attention *cannot* be parallelized across the queries. In practice, this is common in autoregressive decoding or for long sequences where the accelerators (e.g., GPUs) do not have enough threads to fully parallelize softmax attention's computation across different queries.

results into n-by-d memory  $\widetilde{\mathbf{K}}$  and  $\widetilde{\mathbf{V}}$ . Intuitively, it is the "real queries"  $\{\mathbf{q}_i\}$  that encode "what information is useful for the prediction task." Without access to them,  $\mathsf{ABC}_{\mathsf{MLP}}$  summarizes the input for n times using different pseudo-queries, aiming to preserve enough information in the memory for onward computation. The attention output is calculated with the context-dependent real queries using Eq. 4. §B.2 presents a detailed derivation.

Connections to other prior works. Although starting from distinct motivations,  $ABC_{MLP}$  closely relates to hierarchical attention (HA; Yang et al., 2016). HA summarizes the context into higher-level representations with a cascade of attention mechanisms, e.g., words to sentences, and then to documents.  $ABC_{MLP}$  applies two types of attention. The first learns context-agnostic pseudo-queries and attends to the same sequence for n times in parallel, while the second retrieves from the memory with real queries. HA, in contrast, summarizes non-overlapping segments at each level.

The learned pseudo-queries closely relate to the inducing point method in set attention (ISA; Lee et al., 2019). ISA applies a non-linear feedforward network between a cascade of two attention mod-

ules. This precludes the outer-product memory computation and efficient recurrences in ABC.

Another line of work "linearizes" attention through kernel tricks and also applies bounded memory: their feature map dimensions are analogous to memory sizes. They substitute the softmax with approximations (Peng et al., 2021; Choromanski et al., 2021), heuristically designed (Katharopoulos et al., 2020; Schlag et al., 2021), or learned (Kasai et al., 2021b) functions. ABC<sub>MLP</sub> keeps the softmax, but over a smaller constant-sized context. This can be useful in practice: (1) ABC provides a unified perspective of several efficient attention methods, allowing for borrowing from existing wisdom to design new architectures; (2) it draws a close analogy to the canonical softmax attention, and is better-suited as its drop-in substitute in various application settings, as we will show in the experiments; (3) empirically, we find that ABCMLP can get away with a much smaller memory size to retain the accuracy. Peng et al. (2021) and Schlag et al. (2021) use gating to promote recency bias. The same technique is equally applicable in ABC models.

The learned contextualized memory control is reminiscent of the content-based addressing in neu-

ral Turing machines (NTM; Graves et al., 2014). ABC<sub>MLP</sub> computes the control vectors  $\{\phi_i\}$  as a function of the input, but *not* of the memory as in NTM. This ensures that the control vectors at different timesteps can be computed in parallel, improving the time efficiency in practice (Lei et al., 2018; Peng et al., 2018). Analogies between memory and neural architectures are also made by other previous works (Hochreiter and Schmidhuber, 1997; Weston et al., 2015; Le et al., 2020, *inter alia*).

## 5 Experiments

We evaluate ABC models on language modeling (§5.1), sentence-level and document-level machine translation (§5.2), and masked language model fine-tuning (§5.3). Dataset statistics and implementation details are summarized in §C.

### 5.1 Language Modeling

Setting. We experiment with WikiText-103, sampled text from English Wikipedia (Merity et al., 2017). The BASE model with standard softmax attention is the strong transformer-based language model by Baevski and Auli (2019). We compare the following ABC variants, which build on BASE, but replace the softmax attention with linear-complexity bounded-memory attention alternatives while keeping other components the same.

- ABC<sub>MLP</sub>, as described in §4, learns a contextualized exp-MLP as the φ function.
- Linformer (§3.1; Wang et al., 2020b).
- ABC<sub>RD</sub> stores each token in a randomly-selected memory slot with  $\phi_t = \mathbf{e}_{i_t}$ .  $i_t$  is uniformly drawn from  $\{1,\ldots,n\}$  at each time step. This helps us quantify the differences between random and learned bounded-memory controls.

We consider two model size settings:

- 16 layers (Baevski and Auli, 2019). All models have around ~242M parameters. They train with 512-token segments, and evaluate with 0 or 480 context sizes: a 0- or 480- length prefix precedes each evaluation segment.
- 32 layers (Kasai et al., 2021b). All models have ~484M parameters. This setting applies layer dropout (Fan et al., 2020), and evaluates with a 256 context size. It aims to compare ABC<sub>MLP</sub> to several kernel-based efficient attention variants: ELU (Katharopoulos et al., 2020), RFA (Peng et al., 2021), and T2R (Kasai et al., 2021b).

**Results.** Table 3a compares ABC variants using Baevski and Auli (2019)'s 16-layer setting. Among

		D	ev.	Test		
Model	n	0	480	0	480	
BASE	-	19.8	18.4	20.5	19.0	
Linformer ABC <sub>RD</sub>	64 64	26.5 23.2	27.1 22.3	27.2 24.0	30.7 23.1	
ABC <sub>MLP</sub>	32 64	21.2 <b>20.4</b>	19.7 <b>18.9</b>	21.9 <b>21.1</b>	20.5 <b>19.5</b>	

(a) 16-layer setting. 0/480 indicate evaluation context sizes.

Model	n	Dev.	Test
†BASE	-	17.9	18.5
†ELU	128	22.0	22.8
†RFA	32	20.4	21.3
†T2R	32	20.1	20.8
ABC <sub>MLP</sub>	32	19.2	19.9

(b) 32-layer setting. A 256-length context is used at evaluation time. † numbers are due to Kasai et al. (2021b).

Table 3: WikiText-103 language modeling perplexity (**lower is better**). n denotes the memory size. Bold numbers perform the best among linear-complexity models.

ABC models, ABC<sub>MLP</sub> achieves the best performance for both context sizes. With a memory size n = 64, ABC<sub>MLP</sub> outperforms both Linformer and ABC<sub>RD</sub> by more than 2.9 test perplexity; and the gap is larger with the longer 480-length context: more than 3.6 test perplexity. ABC<sub>MLP</sub>-32 outperforms its larger-memory ABC counterparts by more than 2.1 test perplexity. These results confirm ABC<sub>MLP</sub>'s advantages of using a contextualized strategy. Surprisingly, Linformer underperforms ABCRD, and its performance drops with the larger 480-length context window. This suggests that, while successful in text encoding, Linformer's position-based strategy is a suboptimal design choice for causal attention, at least for long context. All ABC models underperform the BASE, with ABC<sub>MLP</sub>-64 having the smallest gap of 0.5 perplexity. ABC<sub>MLP</sub>-32 outperforms kernel-based methods by more than 0.9 test perplexity, using Kasai et al. (2021b)'s 32-layer setting (Table 3b).

#### **5.2** Machine Translation

**Datasets.** To assess their performance over various output lengths, we compare ABC models on sentence- and document-level machine translation.

• Sentence-level translation with WMT14 EN-DE

Model	Cross n	Causal n	BLEU
BASE	-	-	27.2
ABC <sub>RD</sub>	32	32	25.7
$ABC_{RD}$	64	64	26.2
Linformer	32	32	26.6
Linformer	64	64	26.7
$ABC_{MLP}$	32	8	27.1
$ABC_{MLP}$	32	32	27.3

(a) Bolded number outperforms BASE.

Model	Cross n	Causal n	BLEU
BASE	-	-	39.9
Linformer	128	64	-
$\overline{ABC_{RD}}$	128	64	38.6
ABC <sub>MLP</sub>	128	64	39.7

(b) Linformer fails to converge even with multiple random seeds. Bold number performs the best among ABC models.

Table 4: Machine translation test SacreBLEU. Left: sentence-level translation with WMT14 EN-DE; right: document-level translation with IWSLT14 ES-EN.

(Bojar et al., 2014). The preprocessing and data splits follow Vaswani et al. (2017).

• Document-level translation with IWSLT14 ES-EN (Cettolo et al., 2014). We use Miculicich et al. (2018)'s data splits and preprocessing. Following standard practice (Voita et al., 2019), a 4-sentence sliding window is used to create the dataset, i.e., each instance has 4 sentences.

**Setting.** We compare ABC variants as in §5.1. §C.2 further compares to the clustering-based (§3.2) and sliding-window (§3.3) ABC variants.

The BASE model they build on is our implementation of transformer-base (Vaswani et al., 2017). ABC variants replace decoder cross attention and causal attention with bounded-memory attention, while keeping softmax attention for the encoder, since its overhead is much less significant (Kasai et al., 2021a); other components are kept the same. §C.2 studies a model that replaces *all* softmax attention with ABC<sub>MLP</sub>. It performs on par with BASE, confirming ABC<sub>MLP</sub>'s broad applicability in various application scenarios. We evaluate with SacreBLEU (Post, 2018).

**Results.** Table 4a summarizes sentence-level machine translation results on the WMT14 EN-DE test set. Overall ABC<sub>MLP</sub> performs on par with BASE, with either 32-32 cross-causal memory sizes or 32-8. Even with smaller memory sizes, it outperforms other ABC variants by more than 1.1 BLEU. Differently from the trend in the language modeling experiment (§5.1), Linformer outperforms ABC<sub>RD</sub> by more than 0.5 BLEU. We attribute this to the smaller sequence lengths of this dataset. ABC<sub>MLP</sub> outperforms other ABC models by more than 0.4 BLEU, even with smaller memory sizes.

The trend is similar on document-level translation with IWSLT14 ES-EN (Table 4b), except that ABC<sub>MLP</sub> slightly *underperforms* BASE by 0.2 BLEU. This suggests that even with longer sequences, ABC<sub>MLP</sub> is effective despite its bounded memory size. Linformer fails to converge even with multiple random seeds, suggesting the limitations of its purely position-based strategy in tasks involving decoding varying-length text.

#### 5.3 Masked Language Model Finetuning

**Setting.** We compare the ABC variants as in §5.1. It is interesting to pretrain ABC from scratch, but we lack the resources to do so. Instead, we warm-start from a pretrained RoBERTa-base (Liu et al., 2019) trained with the softmax transformer, swap its attention with ABC variants, and continue pretraining with the masked language modeling (MLM) objective on a concatenation of BookCorpus (Zhu et al., 2015), English Wikipedia, Open-WebText (Gokaslan and Cohen, 2019), and Real-News (Zellers et al., 2019). Then the models are finetuned and evaluated on downstream classification datasets from the the GLUE benchbark (Wang et al., 2019). This is an appealing setting, since it avoids reinvesting the huge amounts of resources already put into pretraining.8

**Results.** Table 5 compares downstream text classification performance. BASE indicates a baseline that continues pretraining RoBERTa-base on our data. Following standard practice, we report development accuracy. Linformer achieves competitive

<sup>&</sup>lt;sup>7</sup>Our data differs from RoBERTa's, which we do *not* have access to. We replace CC-News (Nagel, 2016) with RealNews, and drop Stories (Trinh and Le, 2018), whose public access is broken at the time of this work.

<sup>&</sup>lt;sup>8</sup>In preliminary experiments, we explored swapping in ABC, and then directly finetuning on downstream tasks *without* continued MLM pretraining; all models fail.

<sup>&</sup>lt;sup>9</sup>BASE slightly *underperforms* RoBERTa-base. This could be due to overfitting, or the pretraining data discrepancy.

Model	n	MNLI	QNLI	QQP	SST	Avg.
BASE	-	87.2	92.4	91.7	94.3	91.4
Linformer	64	85.3	91.8	90.8	92.4	90.1
Linformer	128	86.1	91.9	91.4	93.7	90.8
$ABC_{MLP}$	64	85.6	91.8	<u>91.7</u>	93.8	90.7
$ABC_{MLP}$	128	<b>87.1</b>	<u>92.6</u>	<u>91.8</u>	<u>94.4</u>	<u>91.5</u>

Table 5: Text classification development set accuracy. All models continue pretraining RoBERTa-base on our data with the MLM objective. Bold numbers perform the best among ABC models, and underlined ones perform on par with or better than BASE.

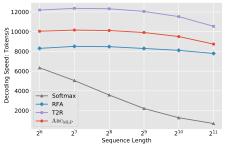
performance, aligned with Wang et al. (2020b)'s results. ABC<sub>MLP</sub> outperforms Linformer, and performs on par with or better than BASE, affirming the benefits of using contextualized memory organization in MLM. ABC<sub>RD</sub> fails to converge in continued pretraining even with multiple seeds.

Based on the above results, we think  $ABC_{MLP}$  can achieve competitive performance when pretrained from scratch, just as Linformer does (Wang et al., 2020b). Further empirical exploration is beyond our budget and left for future work.

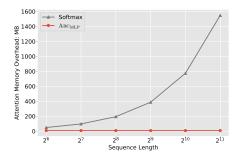
## 6 Analysis

**Decoding efficiency over varying sequence lengths.** ABC's efficiency gains can be more prominent for long sequences. We study ABC<sub>MLP</sub>'s decoding overhead with varying sequence lengths. Following Kasai et al. (2021b), we consider a sequence-to-sequence generation experiment. Three linear-complexity models are compared: RFA (with 256/128 cross/causal memory sizes; Peng et al., 2021), T2R (32/4; Kasai et al., 2021b), and ABC<sub>MLP</sub> (32/8). The sizes are chosen to maximize efficiency *without* accuracy drop. T2R needs to be finetuned from a pretrained transformer to match its performance, while others don't.

All linear-time models achieve consistent decoding speed for different lengths (Figure 1a), substantially outpacing the softmax attention baseline, especially for long sequences. In particular, ABC<sub>MLP</sub> decodes ~1.25 times faster than RFA, another competitive model that can match transformer's accuracy *without* a warm start from a pretrained model. This can be attributed to the fact that ABC<sub>MLP</sub> achieves similar accuracy with a much smaller memory. T2R's memory sizes are similar to ABC<sub>MLP</sub>'s, but it decodes about 20% faster.



(a) Decoding Speed.



(b) Decoding memory overhead.

Figure 1: Sequence-to-sequence decoding speed (top) and memory consumption (bottom) varying sequence lengths. Greedy decoding is used, with batch size 16.

This is because it does *not* compute the softmax when calculating attention output, while ABC<sub>MLP</sub> does (Eq. 4). These results show that ABC<sub>MLP</sub> is an appealing modeling choice for decoding tasks, especially when training from scratch is desired.

ABC<sub>MLP</sub> also achieves significant savings in terms of memory overhead (Figure 1b). ABC<sub>MLP</sub>, RFA, and T2R's curves are similar.

**Text encoding efficiency.** We compare the efficiency of ABC<sub>MLP</sub> against softmax attention and Linformer when used as text encoders. The models' sizes mirror those in the MLM experiment (§5.3). Table 6 summarizes inference time and memory overhead with 512-length inputs, batch size 16. Both ABC<sub>MLP</sub> and Linformer achieve inference speed gains and memory savings over BASE. Linformer is faster, since its linear projection is cheaper to compute than ABC<sub>MLP</sub>'s MLP. Inference speed is measured on the same V100 GPU. The trend in memory overhead is similar.

Although ABC<sub>MLP</sub> slightly underperforms Linformer in terms of inference speed, it can be a more appealing architectural choice in practice: in all of our 5 experiments, ABC<sub>MLP</sub> outperforms other ABC models in accuracy. Linformer, in contrast, fails to converge or yields sub-optimal performance on some tasks. This confirms its flexibility and ap-

	BASE	Linfo	Linformer		MLP
n	-	64	128	64	128
Speed	1.0×	1.7×	1.5×	1.5×	1.3×
Memory	1.0×	0.5×	0.6×	0.5×	0.6×

Table 6: Text encoding inference speed (higher is better) and memory (lower is better). Inputs are text segments with 512 tokens and batch size 16.

		Cross $n$				
		8	16	32	64	
n	8	24.7	25.2	25.6	25.5	
ausal	16	-	25.4	25.7	25.6	
au	32	-	-	25.7	25.8	
C	64	-	-	-	25.8	

Table 7: ABC<sub>MLP</sub>'s SacreBLEU on WMT14 EN-DE development data varying memory sizes.

plicability in various settings.

Memory size's impact on accuracy. Practically, one may want to minimize the memory size to improve efficiency. We use the WMT14 EN-DE experiment to investigate how memory size affects accuracy. Using the §5.2's setup, we vary ABC<sub>MLP</sub>'s cross and causal attention memory sizes and compare their translation quality on the development data. They are selected from  $\{8, 16, 32, 64\}$ , with cross attention's equal to or larger than causal's: cross attention is more important than causal attention in machine translation (Michel et al., 2019). Our results (Table 7) align with this observation: when cross attention memory is large enough, reducing causal attention memory size from 64 to 8 has a minor 0.3 BLEU drop. Surprisingly, ABC<sub>MLP</sub> with 8-8 sized cross-causal memory is only 1.1 BLEU behind the best-performing configuration.

## 7 Conclusion

We presented attention with bounded-memory control (ABC). It provides a unified perspective of several recently-proposed models, and shows that they vary in the organization of the bounded memory. ABC reveals new insights into established methods and inspires new architectures. We proposed ABC<sub>MLP</sub>, a particular instance of ABC that learns a contextualized memory control. On language modeling, machine translation, and masked language model finetuning, ABC<sub>MLP</sub> outperforms previous ABC models. Compared to the strong transformer

baseline, ABC<sub>MLP</sub> achieves substantial efficiency improvements with no or negligible accuracy loss.

### **Acknowledgments**

We would like to thank the ARK group at the University of Washington for their helpful feedback, and the anonymous reviewers for their thoughtful comments. This work was supported in part by NSF grant 2113530 and a Google Fellowship. Nikolaos Pappas was supported by the Swiss National Science Foundation grant P400P2\_183911.

#### References

Joshua Ainslie, Santiago Ontanon, Chris Alberti, Vaclav Cvicek, Zachary Fisher, Philip Pham, Anirudh Ravula, Sumit Sanghai, Qifan Wang, and Li Yang. 2020. ETC: Encoding long and structured inputs in transformers. In *Proc. of EMNLP*.

Alexei Baevski and Michael Auli. 2019. Adaptive input representations for neural language modeling. In *Proc. of ICLR*.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *Proc. of ICLR*.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer.

Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, and Aleš Tamchyna. 2014. Findings of the 2014 workshop on statistical machine translation. In *Proc. of WMT*.

Mauro Cettolo, Jan Niehues, Sebastian Stüker, Luisa Bentivogli, and Marcello Federico. 2014. Report on the 11th IWSLT evaluation campaign. In *Proc. of IWSLT*.

Krzysztof Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Łukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. 2021. Rethinking attention with performers. In *Proc. of ICLR*.

Kornél Csernai. 2017, accessed September 1, 2020. First Quora Dataset Release: Question Pairs.

Angela Fan, Edouard Grave, and Armand Joulin. 2020. Reducing transformer depth on demand with structured dropout. In *Proc. of ICLR*.

Aaron Gokaslan and Vanya Cohen. 2019. Openwebtext corpus. http://Skylion007.github.io/OpenWebTextCorpus.

- Anirudh Goyal, Aniket Didolkar, Alex Lamb, Kartikeya Badola, Nan Rosemary Ke, Nasim Rahaman, Jonathan Binas, Charles Blundell, Michael Mozer, and Yoshua Bengio. 2021. Coordination among neural modules through a shared global workspace.
- Alex Graves, Greg Wayne, and Ivo Danihelka. 2014. Neural turing machines.
- Edward Grefenstette, Karl Moritz Hermann, Mustafa Suleyman, and Phil Blunsom. 2015. Learning to transduce with unbounded memory. In *Proc. of NeurIPS*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*.
- Armand Joulin and Tomás Mikolov. 2015. Inferring algorithmic patterns with stack-augmented recurrent nets. In *Proc. of NeurIPS*.
- Jungo Kasai, Nikolaos Pappas, Hao Peng, James Cross, and Noah A. Smith. 2021a. Deep encoder, shallow decoder: Reevaluating non-autoregressive machine translation. In *Proc. of ICLR*.
- Jungo Kasai, Hao Peng, Yizhe Zhang, Dani Yogatama, Gabriel Ilharco, Nikolaos Pappas, Yi Mao, Weizhu Chen, and Noah A. Smith. 2021b. Finetuning pretrained transformers into RNNs. In *Proc. of EMNLP*.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and Francois Fleuret. 2020. Transformers are RNNs: Fast autoregressive transformers with linear attention. In *Proc. of ICML*.
- Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. 2020. Reformer: The efficient transformer. In *Proc. of ICLR*.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proc. of ACL*.
- Hung Le, Truyen Tran, and Svetha Venkatesh. 2020. Self-attentive associative memory. In *Proc. of ICML*.
- Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. 2019. Set transformer: A framework for attention-based permutation-invariant neural networks. In *Proc. of ICML*.
- Tao Lei, Yu Zhang, Sida I. Wang, Hui Dai, and Yoav Artzi. 2018. Simple recurrent units for highly parallelizable recurrence. In *Proc. of EMNLP*.
- Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Łukasz Kaiser, and Noam Shazeer. 2018. Generating Wikipedia by summarizing long sequences. In *Proc. of ICLR*.

- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized BERT pretraining approach.
- Xuezhe Ma, Xiang Kong, Sinong Wang, Chunting Zhou, Jonathan May, Hao Ma, and Luke Zettlemoyer. 2021. Luna: Linear unified nested attention. In *Proc. of NeurIPS*.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. Pointer sentinel mixture models. In *Proc. of ICLR*.
- Paul Michel, Omer Levy, and Graham Neubig. 2019. Are sixteen heads really better than one? In *Proc. of NeurIPS*.
- Lesly Miculicich, Dhananjay Ram, Nikolaos Pappas, and James Henderson. 2018. Document-level neural machine translation with hierarchical attention networks. In *Proc. of EMNLP*.
- Sebastian Nagel. 2016. News dataset available. https://commoncrawl.org/2016/10/news-dataset-available/.
- Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Łukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. 2018. Image transformer. In *Proc. of ICML*.
- Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah Smith, and Lingpeng Kong. 2021.Random feature attention. In *Proc. of ICLR*.
- Hao Peng, Roy Schwartz, Sam Thomson, and Noah A. Smith. 2018. Rational recurrences. In *Proc. of EMNLP*.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proc. of WMT*.
- Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap. 2020. Compressive transformers for long-range sequence modelling. In *Proc. of ICLR*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proc. of EMNLP*.
- Aurko Roy, Mohammad Taghi Saffar, David Grangier, and Ashish Vaswani. 2020. Efficient content-based sparse attention with routing transformers. *TACL*.
- Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. 2021. Linear transformers are secretly fast weight programmers. In *Proc. of ICML*.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proc. of ACL*.

- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proc. of EMNLP*.
- Sainbayar Sukhbaatar, Da Ju, Spencer Poff, Stephen Roller, Arthur Szlam, Jason Weston, and Angela Fan. 2021. Not all memories are created equal: Learning to forget by expiring. In *Proc. of ICML*.
- Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. 2020. Efficient transformers: A survey.
- Trieu H. Trinh and Quoc V. Le. 2018. A simple method for commonsense reasoning.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proc. of NeurIPS*.
- Elena Voita, Rico Sennrich, and Ivan Titov. 2019. When a good translation is wrong in context: Context-aware machine translation improves on deixis, ellipsis, and lexical cohesion. In *Proc. of ACL*.
- Apoorv Vyas, Angelos Katharopoulos, and François Fleuret. 2020. Fast transformers with clustered attention. In *Proc. of NeurIPS*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proc. of ICLR*.
- Shuohang Wang, Luowei Zhou, Zhe Gan, Yen-Chun Chen, Yuwei Fang, Siqi Sun, Yu Cheng, and Jingjing Liu. 2020a. Cluster-Former: Clustering-based sparse transformer for long-range dependency encoding. *Findings of ACL*.
- Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. 2020b. Linformer: Self-attention with linear complexity.
- Jason Weston, Sumit Chopra, and Antoine Bordes. 2015. Memory networks. In *Proc. of ICLR*.
- Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proc. of NAACL*.
- Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. Hierarchical attention networks for document classification. In *Proc. of NAACL*.
- Dani Yogatama, Yishu Miao, Gabor Melis, Wang Ling,Adhiguna Kuncoro, Chris Dyer, and Phil Blunsom.2018. Memory architectures in recurrent neural network language models. In *Proc. of ICLR*.

- Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. 2020. Big bird: Transformers for longer sequences. In *Proc. of NeurIPS*.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. In *Proc. of NeurIPS*.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proc. of ICCV*.

## **Appendices**

#### A Other ABC Models

## A.1 Sparse Local-to-global Attention

It sparsifies attention pattern to reduce the number of tokens that are attended to (Beltagy et al., 2020; Zaheer et al., 2020, inter alia). All queries attend to a subset of n < N "global tokens," while ignoring others. Therefore the effective context size is reduced to n. The global tokens are usually preselected by positions according to some heuristics. Local-to-global attention is an instance of ABC: it can be recovered by letting  $\phi_t = \mathbf{e}_i$  if  $x_t$  is the ith global token  $(i = 1, \ldots, n)$ , and the zero vectors for others.

### A.2 Random Memory Control

As a baseline,  $ABC_{RD}$  stores each token in a randomly-selected memory slot. This is achieved by letting  $\phi_t = \mathbf{e}_{i_t}$ , where  $i_t$  is uniformly drawn from  $\{1,\ldots,n\}$  for each t. It is designed as a baseline to  $ABC_{MLP}$  and Linformer to quantify the differences between random and learned bounded-memory control.

Random sparse attention patterns are explored by Zaheer et al. (2020), where a subset of n < N tokens are randomly selected to be attended to by all tokens. ABC<sub>RD</sub> is different, and it attends to *all* tokens, but randomly "squash" them into an n-slot memory.

# A.3 Compressive Transformer with Mean Pooling

The compressive transformer (Rae et al., 2020) explores various ways to "squash" long context into smaller and more compact representations. It achieves state-of-the-art performance on several language modeling benchmarks. We show that at least the mean-pooling variant of the compressive transformer can be seen as an ABC instance.

The mean-pooling variant of the compressive transformer compresses the context by

$$\mathbf{K} = \begin{bmatrix} \mathbf{k}_{1}, \dots, \mathbf{k}_{N} \end{bmatrix}^{\top} \in \mathbb{R}^{N \times d}$$

$$\rightarrow \widetilde{\mathbf{K}} = \underbrace{\begin{bmatrix} (\mathbf{k}_{1} + \dots + \mathbf{k}_{c}) / c, \\ c \end{bmatrix}}_{c} / c,$$

$$\underbrace{(\mathbf{k}_{c+1} + \dots + \mathbf{k}_{2c}) / c \dots,}_{c}$$

$$\underbrace{(\mathbf{k}_{N-c+1} + \dots + \mathbf{k}_{N})}_{c} / c \end{bmatrix}^{\top} \in \mathbb{R}^{n \times d}.$$

where c=N/n is the compression ratio. Here  $N \mod n = 0$  is assumed, since otherwise the sequence can be padded to.

The above model is an ABC instance by letting

$$\phi_i = \mathbf{e}_{|(i-1)/c|+1}/c.$$
 (8)

#### A.4 Dilated Convolution Attention Patterns

The dilated attention pattern is similar to the sliding window attention and only considers the context within a predefined window. It differs in that it attends to every other token:

$$\widetilde{\mathbf{K}}_t = [\mathbf{k}_{t-2n+2}, \mathbf{k}_{t-2n+4}, ..., \mathbf{k}_{t-2}, \mathbf{k}_t]^{\top}.$$
 (9)

It can be simulated with two separate queues  $\widetilde{\mathbf{K}}^{\text{odd}}$  and  $\widetilde{\mathbf{K}}^{\text{even}}$ :

$$\begin{split} \widetilde{\mathbf{K}}_t^{\text{odd}} &= \left\{ \begin{array}{ll} \mathbf{U}\widetilde{\mathbf{K}}_{t-1}^{\text{odd}} + \mathbf{e}_n \otimes \mathbf{k}_t, & \text{if $t$ is odd} \\ \widetilde{\mathbf{K}}_{t-1}^{\text{odd}}, & \text{otherwise} \end{array} \right. \\ \widetilde{\mathbf{K}}_t^{\text{even}} &= \left\{ \begin{array}{ll} \mathbf{U}\widetilde{\mathbf{K}}_{t-1}^{\text{even}} + \mathbf{e}_n \otimes \mathbf{k}_t, & \text{if $t$ is even} \\ \widetilde{\mathbf{K}}_{t-1}^{\text{even}}, & \text{otherwise} \end{array} \right. \end{split}$$

Likewise for the values. Depending on t, the query attends to one of the two queues: output =

$$\begin{cases} \left(\widetilde{\mathbf{V}}^{\text{odd}}\right)^{\top} \operatorname{softmax}(\widetilde{\mathbf{K}}^{\text{odd}}\mathbf{q}_{t}), & \text{if } t \text{ is odd} \\ \left(\widetilde{\mathbf{V}}^{\text{even}}\right)^{\top} \operatorname{softmax}(\widetilde{\mathbf{K}}^{\text{even}}\mathbf{q}_{t}), & \text{otherwise.} \end{cases}$$

The above implementation could incur considerable amount of overhead and may be actually more expensive than the the original dilated window formulation. Therefore it has more conceptual value than practical value.

## A.5 Shared Workspace and Linear Unified Nested Attention

Concurrently to this work, shared workspace (SW; Goyal et al., 2021) and linear unified nested attention (LUNA; Ma et al., 2021) also propposed methods to learn contextualized memory control strategies. Both can be seen as instances of ABC. At layer  $\ell$ , their  $\phi_i^{\ell}$  is a function of previous layer's memory  $\widetilde{\mathbf{X}}^{\ell-1} \in \mathbb{R}^{n \times d}$  and current layer's input  $\mathbf{X}^{\ell} \in \mathbb{R}^{N \times d}$ :

$$\phi_i = \left[ \operatorname{softmax} \left( \widetilde{\mathbf{X}}^{\ell-1} {\mathbf{X}}^{\ell^{\top}} \right) \right]_{::i}, \quad (10)$$

where  $[\cdot]_{:,i}$  denotes the *i*th column of a matrix. Query, key, and value projections are suppressed for notation clarity.

SW and LUNA reveal the entire sequence to the control vectors, by constructing  $\phi$  as a function of previous layer's memory. Although both admit the recurrent computation as all ABC models do, they are ill-suited for causal attention and autoregressive decoding, since future information is "leaked" to  $\phi_i$  from the previous layer. LUNA resorts to a variant of Katharopoulos et al. (2020) in causal attention (Ma et al., 2021). In contrast, ABC<sub>MLP</sub> never conditions  $\phi_i$  on previous layer's memory, but only on the current layer's input.

## **B** More Details about ABC-MLP

#### **B.1** Normalization in Causal Attention

An equivalent implementation to Eq. 7 is to normalize  $\widetilde{\mathbf{K}}$  and  $\widetilde{\mathbf{V}}$  instead of  $\phi_i$  vectors:

$$\begin{split} \boldsymbol{\alpha}_i &= \mathbf{exp} \, \left( \mathbf{W}_{\boldsymbol{\phi}} \mathbf{x}_i \right), \quad \boldsymbol{\phi}_i = \boldsymbol{\alpha}_i, \\ \bar{\mathbf{K}} &= \widetilde{\mathbf{K}} \left/ \sum_{j=1}^N \boldsymbol{\alpha}_j. \quad \bar{\mathbf{V}} = \widetilde{\mathbf{V}} \left/ \sum_{j=1}^N \boldsymbol{\alpha}_j. \right. \\ \text{output} &= \bar{\mathbf{V}}^\top \operatorname{softmax}(\bar{\mathbf{K}} \mathbf{q}). \end{split}$$

 $\mathbf{M}/\mathbf{z}$  divides the  $\ell$ th row of matrix  $\mathbf{M}$  by vector  $\mathbf{z}$ 's  $\ell$ th dimension. This admits a linear complexity computation graph for the causal variant of  $ABC_{\mathbf{MLP}}$ .

## **B.2** Higher-Dimensional Case of Example 1

This section generalizes Example 1 to higher dimensional cases. Assume that the constant-sized memory has n slots.  $\phi_i$  is cauculated as in Eq. 7. Then  $\widetilde{\mathbf{K}} = \sum_{i=1}^N \phi_i \otimes \mathbf{k}_i \in \mathbb{R}^{n \times d}$ . Each row of  $\widetilde{\mathbf{K}}$  can be seen as a separate attention mechanism with a pseudo query. Let  $[\cdot]_\ell$  denote the  $\ell$ th row/dimension of a matrix/vector. Then for any  $\ell = 1, \ldots, n$ ,

$$\begin{split} \left[\widetilde{\mathbf{K}}\right]_{\ell} &= \sum_{i=1}^{N} [\phi_{i}]_{\ell} \otimes \mathbf{k}_{i} \\ &= \sum_{i=1}^{N} \frac{\exp([\mathbf{W}_{\phi}]_{\ell} \cdot \mathbf{x}_{i})}{\sum_{j=1}^{N} \exp([\mathbf{W}_{\phi}]_{\ell} \cdot \mathbf{x}_{j})} \mathbf{k}_{i}^{\top} \\ &= \operatorname{attn} \left( [\mathbf{W}_{\phi}]_{\ell}, \{\mathbf{x}_{i}\}_{i=1}^{N}, \{\mathbf{k}_{i}\}_{i=1}^{N} \right)^{\top} \in \mathbb{R}^{1 \times d}. \end{split}$$

In other words, there are n attention mechanisms in total, each with a separately-parameterized pseudoquery  $[\mathbf{W}_{\phi}]_{\ell}$ . They summarize the context for n times in parallel, each producing a d-dimensional vectors. These output vectors are then stacked into n-by-d memory  $\widetilde{\mathbf{K}}$ .  $\widetilde{\mathbf{V}}$  is similar.

### C Experimental Details

### C.1 Language Modeling

We closely build on Baevski and Auli (2019) and Kasai et al. (2021b). The hyperparameters are summarized in Table 10. All models are trained on 4 A100 GPUs.

#### **C.2** Machine Translation

We experiment with a sentence-level (WMT14 EN-DE, Bojar et al., 2014) and a document-level benchmark (IWSLT14 ES-EN, Cettolo et al., 2014) to assess model performance over various sequence lengths. The preprocessing and data splits of WMT14 EN-DE follow Vaswani et al. (2017). A 32,768 byte pair encoding (BPE; Sennrich et al., 2016) vocabulary is shared between source and target languages. For IWSLT14, we follow Miculicich et al. (2018) and use the dev2010 subset for development and tst2010-2012 for testing. The tokenization is also the same as Miculicich et al. (2018): we tokenize and truecase Spanish and English with Moses (Koehn et al., 2007) and run byte-pair encoding with 30k splits, shared between the two languages. The final dataset contains 1421, 8, and 42 documents for training, development, and testing. On average, each document contains 126.7 sentences, and each sentence contains 21.7(ES)/22.5(EN) BPE subwords. We use a sliding window with length-4 and stride-one to generate our dataset. During inference, we use predicted context on the target side.

We average the checkpoints from the last five epochs to obtain the final model (Vaswani et al., 2017). In inference, we apply beam search with size 5 and length penalty 0.6. Other hyperparameters are summarized in Table 11. All models are trained on 4 RTX 2080 Ti GPUs.

Additional machine translation results. In addition to the results presented in §5.2, Table 8 further compares, on the WMT14 EN-DE dataset, the clustering-based (§3.2) and sliding-window (§3.3) models of ABC, as well as ReLU and sigmoid variants of ABC<sub>MLP</sub>. Clustering and sliding-window ABC variants *underperform* ABC<sub>MLP</sub> with the same memory sizes by more than 0.5 BLEU. Both ReLU and sigmoid *underperform* their exp counterpart.

MLP-exp-all replaces the encoder's softmax attention modules with ABC, in addition to the decoder's. It underperforms ABC<sub>MLP</sub> by only 0.3 BLEU.

Model	$\phi$	Cross n	Causal n	<b>Encoder</b> n	BLEU
BASE	-	-	-	-	27.2
	Window	32	32	-	26.3
	Cluster	32	32	-	26.8
	MLP-ReLU	32	8	-	_
ABC	MLP-ReLU	32	32	-	26.4
ABC	MLP-sigmoid	32	8	-	26.8
	MLP-sigmoid	32	32	-	27.0
	MLP-exp	32	8	-	27.1
	MLP-exp	32	32	-	27.3
	MLP-exp-all	32	32	32	27.0

Table 8: ABC variants' performance (SacreBLEU) on the WMT14 EN-DE test set for sentence-level machine translation. MLP-ReLU with 32/8 memory sizes fails to converge. MLP-exp-all applies ABC in both the encoder and the decoder, while others only in the decoders.

Figure 1b compares ABC<sub>MLP</sub>'s (32-8 memory sizes) attention memory overhead with softmax attention's. Following Kasai et al. (2021b), we consider a synthetic sequence-to-sequence generation task with varying sequence lengths. A batch size of 16 and greedy decoding is used. The models are of the same size as those in §5.2.

**C.3** Masked Language Model Finetuning

Our data for continued pretraining is a concatenation of BookCorpus (Zhu et al., 2015), English Wikipedia, OpenWebText (Gokaslan and Cohen, 2019), and RealNews (Zellers et al., 2019). Our data differs from RoBERTa's pretraining data, which we do not have access to. We replace their CC-News (Nagel, 2016) with RealNews, and drop Stories (Trinh and Le, 2018). At the time of this project, the public access to the Stories dataset is broken. 10 Our machine does *not* have a large enough memory to load all the data. We therefore split the training data into 20 shards, after shuffling. Other preprocessing is the same as Liu et al. (2019).<sup>11</sup> The hyperparameters for continued pretraining follow base-sized RoBERTa, part of which are summarized in Table 12. All models are trained on a single TPU v3 accelerator.

For downstream task finetuning, we use the same

hyperparameters as Liu et al. (2019).<sup>12</sup> Table 13 briefly describes the tasks. The readers are referred to Wang et al. (2019) for futher details.

<sup>10</sup>https://console.cloud.google.com/
storage/browser/commonsense-reasoning/
reproduce/stories\_corpus?pli=1

inhttps://github.com/pytorch/fairseq/
blob/master/examples/roberta/README.
pretraining.md

<sup>12</sup>https://github.com/pytorch/fairseq/
blob/master/examples/roberta/README.glue.
md

Data	Train	Dev.	Test	Vocab.	Sent./doc
WikiText-103	103M	218K	246K	268K	_
WMT14 EN-DE	4.5M	3K	3K	32K	_
IWSLT14 ES-EN	1713	8	56	30K	121.5

Table 9: Statistics for the datasets. WikiText-103 split sizes are in number of tokens, WMT14 in number of sentences, and IWSLT14 in number of documents.

Hyperprams.	B&A	Kasai
# Layers	16	32
# Heads	8	8
Embedding Size	1024	1024
Head Size	128	128
FFN Size	4096	4096
Batch Size	64	64
Learning Rate	1.0	1.0
Dropout	0.3	0.3
Layer Dropout	-	0.2
Memory size	[32, 64]	64

Table 10: Hyperparameters used in the language modeling experiments. B&A: Baevski and Auli (2019); Kasai: Kasai et al. (2021b).

Hyperprams.	Values
# Layers	12
# Heads	12
Embedding Size	768
Head Size	64
FFN Size	3072
Dropout	0.1
Memory Size	[64, 128]

Table 12: Hyperparameters for continued pretraining in the masked language model finetuning experiments.

Data	Task	Train	Dev.
MNLI	Entailment	392K	9.8K
<b>QNLI</b>	Entailment	105K	5.5K
QQP	Paraphrase	363K	40K
SST-2	Sentiment	67K	873

Table 13: GLUE datasets and statistics. MNLI: Williams et al. (2018); QNLI is compiled by GLUE's authors using Rajpurkar et al. (2016); QQP: Csernai (2017, accessed September 1, 2020); SST-2: Socher et al. (2013).

Hyperprams.	WMT14	IWSLT14
# Layers	6	6
# Heads	8	8
Embedding Size	512	512
Head Size	64	64
FFN Size	2048	1024
Warmup Steps	6000	4000
Dropout	0.1	0.3
Cross Attn. $n$	32	128
Causal Attn. $n$	8	64

Table 11: Hyperparameters used in the machine translation experiments.