

---

# Alternating Updates for Efficient Transformers

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 It has been well established that increasing scale in deep transformer networks leads  
2 to improved quality and performance. However, this increase in scale often comes  
3 with prohibitive increases in compute cost and inference latency. We introduce  
4 Alternating Updates (AltUp), a simple-to-implement method to increase a model’s  
5 capacity without the computational burden. AltUp enables the widening of the  
6 learned representation, i.e., the token embedding, while only incurring a negligible  
7 increase in latency. AltUp achieves this by working on a subblock of the widened  
8 representation at each layer and using a predict-and-correct mechanism to update  
9 the inactivated blocks. We present extensions of AltUp, such as its applicability  
10 to the sequence dimension, and demonstrate how AltUp can be synergistically  
11 combined with existing approaches, such as Sparse Mixture-of-Experts models, to  
12 obtain efficient models with even higher capacity. Our experiments on benchmark  
13 transformer models and language tasks demonstrate the consistent effectiveness  
14 of AltUp on a diverse set of scenarios. Notably, on SuperGLUE and SQuAD  
15 benchmarks, AltUp enables up to 87% speedup relative to the dense baselines at  
16 the same accuracy.

## 17 1 Introduction

18 Contemporary machine learning models have been remarkably successful in many domains, ranging  
19 from natural language [6, 20] to computer vision [53, 38]. Many of these successes have come in  
20 part through sheer scale. A vast amount of empirical studies justify the conventional wisdom that  
21 bigger (models and data sets) is better [19, 24]. Accordingly, state-of-the-art Transformer [46] models  
22 often contain billions of parameters and are trained for weeks on enormously large data sets using  
23 thousands of AI accelerators. Their immense size leads to prohibitive compute and energy costs [34]  
24 and prevents their deployment to resource-constrained applications [30].

25 To alleviate these costs and enable scalability of modern Transformers, a recent line of works  
26 have proposed techniques to increase the capacity of models without drastically increasing the  
27 computational costs via conditional computation. A notable paradigm is sparsely-activated networks,  
28 such as Mixture-of-Experts (MoE) models [10, 56, 33, 1, 27, 41, 43]. The main idea of MoE is to  
29 effectively *widen* each network layer by accessing dynamically invoked parameters, i.e., experts,  
30 where each expert corresponds to a small subset of disjoint parameters that can be acted on by the  
31 input. During training and inference, a given input to the network is routed to a small subset of  
32 experts (parameters) to compute the output. As a result, the computation cost remains small relative  
33 to the total number of parameters. This scheme enables models with higher capacity with only a  
34 relatively small increase in computation.

35 While prior approaches in conditional computation have predominantly focused on the *processing*  
36 *power* of transformers, there is a research gap in efficiently incorporating *widened learned represen-*  
37 *tations*. Recent works have empirically and theoretically established that a wider token representation  
38 (i.e., a larger model dimension) helps in learning more complicated functions by enabling more

39 information to be packed in the representation vectors [19, 24, 52]. This phenomenon is also evident  
 40 in modern architectures of increasing capability. For instance, the representation dimension grows  
 41 from 512 (small) to 768 (base) and 1024 (large, 3B, and 11B) in T5 models [35], and from 4096  
 42 (8B) to 8192 (64B) and 18432 (540B) in PaLM models [6]. A widened representation dimension  
 43 also significantly improves performance for dual encoder retrieval models [31, 32]. However, naively  
 44 widening the learned representation requires accordingly increasing the model dimension (see Fig. 1),  
 45 which quadratically increases the amount of computation in the feedforward computation. In light  
 46 of the above, a natural question arises: can we leverage the benefit of wider representations without  
 47 incurring the additional cost of wider transformer layers?

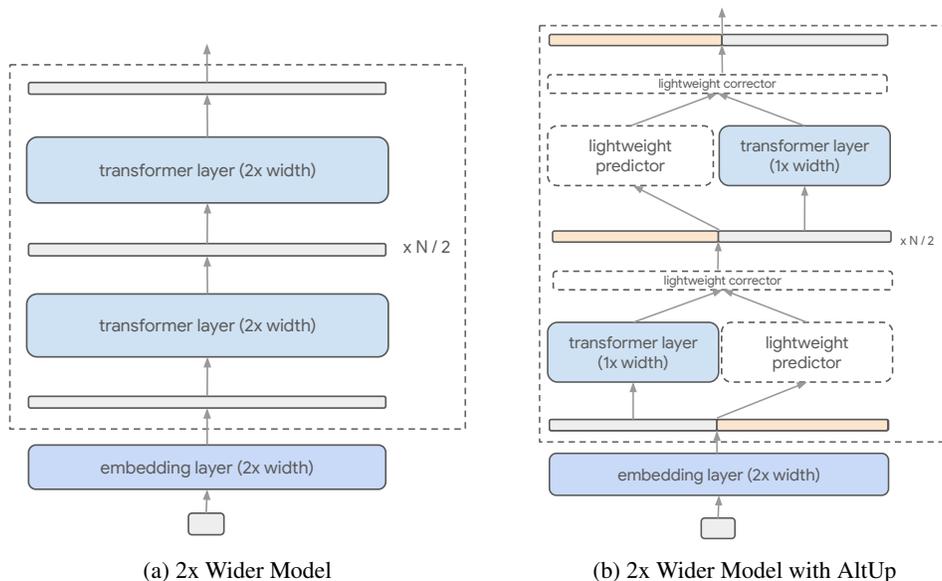


Figure 1: An illustration of widening the token representation without (left) and with Alternating Updates (right). This widening causes a near-quadratic increase in computation in a vanilla transformer due to the increased layer width. In contrast, Alternating Updates keeps the layer width constant and efficiently computes the output by operating on a sub-block of the representation at each layer.

48 In this paper, we address this research gap by introducing *Alternating Updates* (AltUp), a technique  
 49 to incorporate wider representations in a simple and efficient way. AltUp operates by partitioning the  
 50 widened representation vector into blocks, processing only a single block at each layer, and using  
 51 an efficient prediction mechanism to infer the outputs of the other blocks (see Fig. 1). Processing a  
 52 single block in each transformer layer enables AltUp to simultaneously keep the model dimension,  
 53 hence the computation cost, constant and take advantage of using an increased token dimension.  
 54 Unlike prior approaches, e.g., Sparse Mixture of Experts, AltUp is easy to implement, requires  
 55 minimal hyperparameter tuning, and does not necessitate sharding. Moreover, since AltUp focuses on  
 56 increasing the representation dimension, it can be applied synergistically with orthogonal techniques  
 57 like MoE [57] to obtain complementary performance gains.

58 In particular, our contributions are:

- 59 1. We introduce *Alternating Updates* (AltUp) to bridge the research gap in efficiency techniques  
 60 by enabling wider representations with little additional computation cost. AltUp is simple-  
 61 to-implement, requires minimal hyperparameter tuning, and does not necessitate sharding.
- 62 2. We develop two notable extensions of AltUp: (i) *Recycled-AltUp*, a faster variant of AltUp  
 63 that requires virtually no additional learnable parameters and (ii) *Sequence-AltUp*, an  
 64 extension of the AltUp idea to the sequence dimension.
- 65 3. We present an extensive evaluation of AltUp on T5 models on various benchmark language  
 66 tasks. Our experimental results show that AltUp and its variants uniformly lead to models  
 67 with improved speed-accuracy trade-offs. Notably, on SuperGLUE and SQuAD benchmarks,  
 68 AltUp enables up to 87% speedup relative to the dense baselines at the same accuracy.

## 69 2 Related Work

70 Prior work is rich with a diverse set of techniques to increase the efficiency of contemporary  
71 transformer models. Here, we cover the most relevant subset of state-of-the-art techniques and refer  
72 the interested reader to [45] for a more comprehensive survey.

73 Recent works have introduced extremely large, yet scalable models with the use of conditional routing  
74 of inputs to a learnable subset of parameters. These sparsely-activated models have achieved state-of-  
75 the-art performance on various benchmarks [10] and exhibit favorable theoretical properties [4, 1].  
76 Notably, the Sparse Mixture of Experts (SMoE) [43, 54, 22] family of models use a learned softmax  
77 probability distribution to conditionally direct the computation to *experts*, i.e., subsets of network  
78 parameters. By routing the computation to a small subset of parameters on an input-dependent  
79 basis, SMoE leads to higher capacity models with a relatively small and controllable increase in  
80 computation. Switch Transformers [11] show that routing to a single expert on an input-dependent  
81 basis reduces computation and outperforms prior SMoE approaches on language tasks. Follow-up  
82 work on SMoE include those that improve the load balancing of experts [57, 28], use reinforcement  
83 learning to learn the routing function [7], and leverage smooth top- $k$  expert selection [18] (see [10]  
84 for a survey). Other choices for the routing function include non-learnable ones such as Locality  
85 Sensitivity Hashing (LSH) [33] which generally maps similar inputs to the same expert, Hash Layers  
86 that use token-based hashing [41], and language-specific deterministic routing [9]. Residual Mixture  
87 of Experts [49] separates the expert weights into input-independent and input-dependent components.

88 Conditionally accessing *external memory* is another related approach to vastly increase model  
89 capacity at the cost of a relatively small increase in computation [14, 13]. For examples, Memorizing  
90 Transformers [51], Memformer [50], and Product key memory [26] leverage dynamic memory to  
91 encode and retrieve relevant information. Additional works include those that use an immensely  
92 large untrainable corpus, such as Wikipedia, REALM [17], or a 2 trillion token database, RETRO [2].  
93 These prior works that focus on routing(expert)-based mechanisms often necessitate complicated,  
94 sharded implementations due to the sheer number of additional parameters that they introduce — often  
95 on the order of billions. Our work, on the other hand, is simple-to-implement and requires virtually  
96 no hyperparameter tuning. Moreover, AltUp can be synergistically combined with sparsely-activated  
97 models like MoE to obtain complementary improvements in efficiency.

98 Additional relevant works in the realm of efficient transformers include Funnel transformers [8],  
99 Reformers [25], Performers [5], Big-Bird [55], and LongT5 [16], among others. These works notably  
100 present methods to reduce the quadratic cost of the attention mechanism of transformers. Another  
101 flavor of methods complementary to our work is that of adaptive computation, e.g., CALM [42],  
102 DynaBERT [21] CascadeBERT [29] and DeeCap [12], where different amounts of computational  
103 power is allotted on an example-specific basis via some type of early-exit strategy. AltUp achieves its  
104 efficiency via the orthogonal direction of conditionally leveraging wider token representations, and  
105 hence can be easily combined with these techniques.

## 106 3 Alternating Updates

107 In this section, we introduce the method of *Alternating Updates* (AltUp), an approach to enable  
108 increased token dimension with little additional computation cost.

### 109 3.1 Background

110 At a high level, a standard transformer with  $L$  layers generates a  $d$ -dimensional representation by  
111 applying a sequence of layers transformer layers  $\mathcal{L}_1, \dots, \mathcal{L}_L$  as follows. For a particular input token  
112 within a sequence of length  $N$ , the initial token representation  $x_1 \in \mathbb{R}^d$  is computed by an embedding  
113 table lookup. Subsequently, this  $d$ -dimensional representation is refined across the transformer layers  
114 by iteratively computing  $x_{i+1} = \mathcal{L}_i(x_i)$  for each layer  $i \in [L]$ ; here and throughout  $[N]$  denotes the  
115 set  $\{1, \dots, N\}$  for  $N \in \mathbb{N}$ . Each transformer layer  $\mathcal{L}_i : \mathbb{R}^{d_{\text{model}}} \rightarrow \mathbb{R}^{d_{\text{model}}}$  has width  $d_{\text{model}}$  (with  
116  $d_{\text{model}} = d$  in the standard setting) and contains an attention block and a FeedForward (FFN) block.  
117 The width of the layer  $d_{\text{model}}$  controls the dimensions of the matrices involved in the attention and  
118 FFN blocks. Consequently, the computation cost of attention and FFN scales with  $\mathcal{O}(N^2 d_{\text{model}})$   
119 and  $\mathcal{O}(N d_{\text{model}}^2)$ , respectively. The output,  $x_{L+1}$  is the output token representation generated by  
120 the transformer. This computation is usually followed by a linear layer operation that maps from

121 the  $d$ -dimensional representation  $x_{L+1}$  to  $|\mathcal{V}|$ -dimensional logits (in  $\mathcal{O}(|\mathcal{V}|d)$  time), followed by a  
 122 softmax non-linearity to generate the probabilities over the vocabulary  $\mathcal{V}$ .

123 Increasing the representation dimension  $d$  is a way to enhance the capacity of the transformer model,  
 124 as a wider representation enables the transformer to store richer information about the input and helps  
 125 in learning more complicated functions [19, 24, 52]. Naively widening the token representation  $d$   
 126 requires widening each layer as well, since  $d_{\text{model}}$  must match  $d$  in a standard transformer model.  
 127 However, the computation time of each transformer layer grows roughly quadratically with  $d_{\text{model}}$ ,  
 128 notably for relatively short input sequences. This means that, growing the token dimension from  $d$  to  
 129  $2d$ , for example, leads to a model that is at least 2 times (and closer to 4 times for small  $N$ ) slower  
 130 than the original model with a  $d$ -dimensional representation.

### 131 3.2 Alternating Updates

132 The core idea of Alternating Updates is to *widen the representation vector, but perform computation*  
 133 *with a  $d$ -dimensional sub-block*, and estimate the updated representation using a Predict-Compute-  
 134 Correct algorithm, as illustrated in Figure 1, right. More specifically, AltUp expands the representation  
 135 width from  $d$  to  $Kd$ , for integers  $K > 1$ ,  $d > 0$  (for example,  $K = 2$  in Fig. 1), but uses layers of  
 136 width  $d_{\text{model}} = d$  (not  $d_{\text{model}} = Kd$ ) to transform the representation vector. By keeping the width of  
 137 each transformer layer constant, AltUp avoids incurring the quadratic increase in computation cost  
 138 that would otherwise be present with a naive expansion of the representation.

139 Alg. 1 depicts the details of the per-layer computation involved in a transformer with AltUp with a  $Kd$ -  
 140 dimensional representation vector. The input to the AltUp layer is assumed to be the concatenation  
 141 of  $d$ -dimensional contiguous subblocks  $x_{old} = \text{concat}(x_{old}^1, x_{old}^2, \dots, x_{old}^K) \in \mathbb{R}^{dK}$ . Inspired by  
 142 predictor-corrector methods used to solve ordinary differential equations [3], AltUp first generates a  
 143 prediction  $\hat{x}^i$  for each of the subblocks  $i \in [K]$  (Line 1). This prediction takes the form of a mixture  
 144 of subblocks  $\hat{x}^i = \sum_{j=1}^K p_{i,j} x_{old}^j$ , where  $p_{i,j} \in \mathbb{R}$  for  $i, j \in [K]$  are learnable scalars. Subsequently,  
 145 one of the  $K$  sub-blocks is chosen and the computation with the unexpanded transformer layer of  
 146 width  $d_{\text{model}} = d$  is performed on this sub-block (Line 2). Finally, the result of this computation is  
 147 used in the correction step to generate the updated representation for each sub-block (Line 3).

---

#### Algorithm 1 Alternating Updates (AltUp) Layer

---

**Input:**  $x_{old} = \text{concat}(x_{old}^1, x_{old}^2, \dots, x_{old}^K) \in \mathbb{R}^{dK}$ :  $dK$ -dimensional input representation vector  
 to the layer, where  $x_{old}^j \in \mathbb{R}^d$ ,  $j = 1, 2, \dots, K$  are contiguous sub-blocks of  $x_{old}$ .

**Output:**  $x_{new} \in \mathbb{R}^{dK}$ : The layer’s  $dK$ -dimensional output representation.

1: **Predict:** for each  $i \in [K]$ , predict the updated representation with a trainable linear map:

$$\hat{x}^i = \sum_{j=1}^K p_{i,j} x_{old}^j,$$

where  $p_{i,j} \in \mathbb{R}$ ,  $i, j \in [K]$  are trainable scalars.

2: **Compute:** select a sub-block  $j^* \in [K]$  and update this block with  $\mathcal{L}$ :

$$\tilde{x}^{j^*} = \mathcal{L}(x_{old}^{j^*}).$$

3: **Correct:** for each  $i \in [K]$ , correct the prediction with the computation result:

$$x_{new}^i = \hat{x}^i + g_i(\tilde{x}^{j^*} - \hat{x}^{j^*}),$$

where  $g_i \in \mathbb{R}$ ,  $i \in [K]$  are trainable scalars.

---

148 **Computation time** We see from Alg. 1 that AltUp introduces negligible amount of additional  
 149 computation per layer, as the prediction and correction steps involve only vector addition and scalar-  
 150 vector multiplications ( $\mathcal{O}(d)$  operations). Thus, relative to the computation cost of a transformer layer  
 151 with width  $d$  (which we incur on Line 2 in AltUp), the cost of AltUp is only an additional  $\mathcal{O}(dK^2)$  per  
 152 token, where  $d$  is the original model dimension and  $K$  is the factor of increase in the representation  
 153 dimension (typically  $K = 2$  or  $K = 4$ , see Sec. 5). This additional  $\mathcal{O}(dK^2)$  cost per token is a factor

154 of  $d$  smaller than the  $\mathcal{O}(d^2 K^2)$  per token cost of the FFN block alone in a  $K$ -times wider transformer  
 155 layer. In fact, an AltUp layer does not lead to an increased computation time relative to a  $d$ -width  
 156 transformer layer asymptotically, since the  $\mathcal{O}(dK^2)$  additional cost per token per layer is dominated  
 157 by the cost of the FFN block as  $K \ll d$  in practice. At a higher level, AltUp requires using an  
 158 embedding table with width  $Kd$  and invoking the final linear operation with  $Kd$ -dimensional vectors.  
 159 The initial embedding lookup using a wider table and the linear + softmax operation with  $Kd$  (instead  
 160 of  $d$ ) dimensional vectors may lead to a perceptible increase in computation time. However, since we  
 161 only incur this additional cost in the beginning and the end, these factors are often inconsequential,  
 162 and increasingly so for deeper transformers. Nevertheless, we present an extension to AltUp in Sec. 4  
 163 that avoids this slowdown altogether for specialized applications.

164 **Parameter count** AltUp introduces  $K^2 + K$  additional learnable parameters per layer, where  $K^2$  is  
 165 due to  $p_{i,j}, i, j \in [K]$  and  $K$  is a result of  $g_i, i \in [K]$ . Since  $K \ll d$ , this is an imperceptible amount  
 166 of additional parameters per layer in practice. Zooming out, AltUp with an expansion factor of  $K$   
 167 requires a  $Kd$ -width embedding table, and consequently requires  $(K - 1)|\mathcal{V}|d$  additional parameters,  
 168 where  $|\mathcal{V}|$  is the vocabulary size. In Sec. 4, we present a variant that requires no additional embedding  
 169 parameters to be added to the model.

170 **Selection of sub-blocks** The selection of the sub-block  $j^*$  for the computation step in Algorithm 1  
 171 can be any user-specified technique. We consider two simple, deterministic selection methods in this  
 172 paper and leave more sophisticated methods for future work: (i) **same**: choose the same sub-block  
 173 for all the layers in a neural network and (ii) **alternating** (default method): for a sequence of layers,  
 174 alternating through the sub-blocks, that is, if the sub-blocks are indexed with zero-based index, then  
 175 sub-block  $\ell \bmod K$  is selected for the computation step for layer  $\ell \in [L]$ . This alternating selection  
 176 is the default for Algorithm 1 (hence the name Alternating Updates). We compare the two selection  
 177 methods empirically in the supplementary material and find that using alternating blocks performs  
 178 better empirically.

## 179 4 AltUp Extensions

180 In this section, we present extensions of the core AltUp idea introduced in the previous section.

### 181 4.1 Recycled-AltUp: Faster AltUp via embedding recycling

182 The AltUp formulation presented in Sec. 3 adds  
 183 an insignificant amount of per-layer computa-  
 184 tion, however, it does require using a  $K$ -times  
 185 wider embedding table. In certain scenarios  
 186 where the vocabulary  $\mathcal{V}$  is very large, this may  
 187 lead to a non-trivial amount of added computa-  
 188 tion for the initial embedding lookup and the fi-  
 189 nal linear + softmax operation. A colossal voca-  
 190 bulary may also lead to an undesirable amount of  
 191 added embedding parameters. *Recycled-AltUp*  
 192 is an extension of AltUp that avoids these compu-  
 193 tational and parameter costs by keeping the em-  
 194 bedding table’s width  $d$ -dimensional.

195 Figure 2 depicts an example application of Recy-  
 196 cled AltUp with  $K = 2$ . The general idea is to  
 197 *recycle* the initial  $d$ -dimensional lookup by  
 198 replicating the  $d$ -dimensional lookup  $K$  times.  
 199 Hence, Recycled-AltUp virtually adds no addi-

200 tional parameters relative to the baseline width  $d$  model. Subsequently, the regular AltUp layers  
 201 (Alg. 1) are applied until the last linear + softmax operation. To avoid the computational cost of this fi-  
 202 nal operation, Recycled AltUp downprojects the  $Kd$ -dimensional representation vector  $x_{L+1} \in \mathbb{R}^{dK}$   
 203 to a  $d$ -dimensional representation by simply elementwise-adding the  $d$ -dimensional contiguous sub-  
 204 blocks in  $\mathcal{O}(Kd)$  time. Applying the linear + softmax operation on this down-projected vector implies  
 205 that the computation cost of this operation is now  $\mathcal{O}(|\mathcal{V}|d)$  rather than  $\mathcal{O}(K|\mathcal{V}|d)$ , effectively reducing

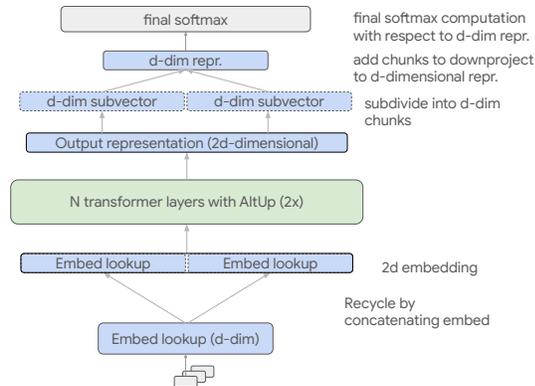


Figure 2: Recycled-AltUp with  $K = 2$ .

206 the amount of computation by  $\mathcal{O}((K-1)|\mathcal{V}|d)$ . Our results in Sec. 5 show that Recycle-AltUp’s  
 207 improved speed and reduced parameter count may make it more appealing for certain applications.

208 **4.2 Sequence-AltUp: Extension of AltUp to the sequence dimension**

209 Here, we introduce *Sequence-AltUp*, a natural extension of Alternating Updates to reduce the apparent  
 210 sequence length. This extension is motivated by the computation cost associated with the cost of the  
 211 attention mechanism for long input sequence lengths. Our approach is similar in its goal to that of prior  
 212 techniques focused on designing efficient attention mechanisms to reduce the quadratic dependency  
 213 of attention cost on the sequence length: Funnel transformers [8], Reformers [25], Performers [5],  
 214 Big-Bird [55], and LongT5 [16]. Similar to the Funnel transformer [8] and Conformer [15], Sequence-  
 215 AltUp uses a simple striding operation to reduce the sequence length. Only sampled tokens are  
 216 processed by the transformer layer while the rest of the tokens require little computation, leading to a  
 computation cost reduction by a factor of  $k$ , where  $k$  is the stride parameter.

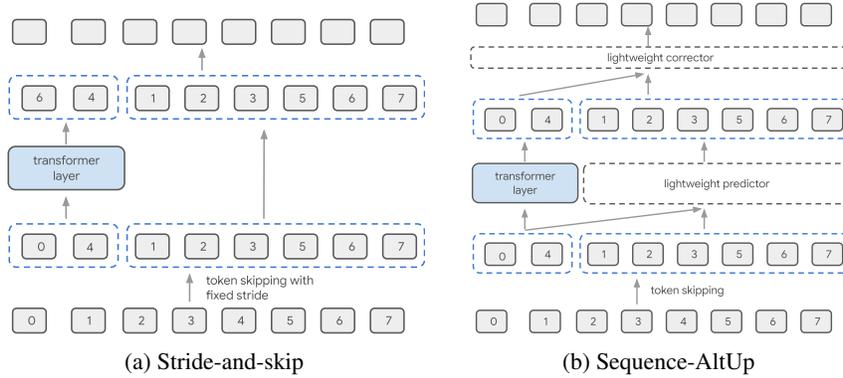


Figure 3: An illustration of Sequence-AltUp (right) and the baseline Stride-and-skip method (left). Sequence-AltUp has virtually the same computation cost as Stride-and-skip, but enables contextual information passing to the skipped tokens.

217

218 Figure 3 depicts the baseline stride-and-skip technique (left) and the proposed Sequence-AltUp  
 219 method (right). Given an input sequence of vectors  $(x_0, x_1, \dots, x_{T-1})$  to a transformer layer  $\mathcal{L}$ ,  
 220 we propose the following extension of Algorithm 1 to reduce the effective sequence length by a  
 221 factor of  $k$ . First, we apply a lightweight predictor on the full sequence to obtain a predicted output  
 222  $\hat{y} = (\hat{y}_0, \dots, \hat{y}_{T-1})$ . Next, we subsample the input with a fixed stride  $k$  and apply  $\mathcal{L}$  on the subsampled  
 223 input and get computed output  $\tilde{y} = (\tilde{y}_0, \tilde{y}_k, \tilde{y}_{2k}, \dots, \tilde{y}_{\lfloor (T-1)/k \rfloor * k}) = \mathcal{L}(x_0, x_k, \dots, x_{\lfloor (T-1)/k \rfloor * k})$ .  
 224 Finally, we use a lightweight corrector to combine  $\tilde{y}$  and  $\hat{y}$  to form the final output sequence. This  
 225 design allows unsampled token vectors to obtain contextual information, even though they are not  
 226 processed by the transformer layer directly—analogueous to the inactivated sub-blocks in original AltUp.  
 227 In contrast, a simple stride-and-skip approach (Figure 3, left) lacks the ability to bring contextual  
 228 information to the skipped tokens. We present the full algorithm pseudocode and implementation  
 229 details of Sequence-AltUp in the supplementary material.

230 **5 Results**

231 In this section, we apply AltUp and its variants to benchmark language models and tasks. We proceed  
 232 by outlining the experimental setting below. In Secs. 5.1 and 5.2 we present the performance of  
 233 AltUp on standard benchmarks with varying configurations and model sizes; in Secs. 5.3 and 5.4  
 234 we evaluate the performance of AltUp extensions and demonstrate their effectiveness. We present  
 235 the full details of our evaluations and additional experimental results in the supplementary; namely,  
 236 the supplementary contains additional evaluations that demonstrate the synergistic combination of  
 237 AltUp with other conditional compute techniques, additional finetune results, and complementary  
 238 ablation studies. Overall, our results consistently show that AltUp and its variants enable sizeable  
 239 performance gains, e.g., up to 87% faster models, across all evaluations on standard benchmarks.

240 **Setting** We performed all of our experiments using T5-model architectures [35] of varying sizes  
 241 (small, base, large, and 3B) which we pretrained on the C4 dataset for 500,000 steps with a batch size  
 242 of 256. The pretrained models were then finetuned on either the GLUE [48], SuperGLUE (SG) [47],  
 243 SQuAD [37] or Trivia-QA (closed-book) [23, 40] benchmark tasks for a further 50,000 steps with a  
 244 batch-size of 256. The pretraining task is to predict corrupted text spans, and the finetuning tasks are  
 245 re-cast into text generation tasks. We report both pretraining and finetuning metrics: for pretraining,  
 246 we report span prediction accuracy on a hold-out validation set, and for finetuning, we follow the  
 247 same recipe as the T5 models, see [35] for more details. The supplementary contains the full details  
 248 of our evaluations and hyperparameters.

## 249 5.1 Alternating updates on benchmarks

250 First, we investigate whether incorporating AltUp on a baseline model leads to an unambiguously  
 251 better model when we consider the predictive performance and *actual observed latency* (not theoretical  
 252 FLOPS). To this end, we compare the dense T5-Base/Large/XL models to models augmented with  
 253 AltUp with  $K = 2$  on GLUE, SuperGLUE, SQuAD, and TriviaQA (closed-book) finetuning tasks.  
 254 Figure 4 plots the performance and normalized speed of the evaluated models.

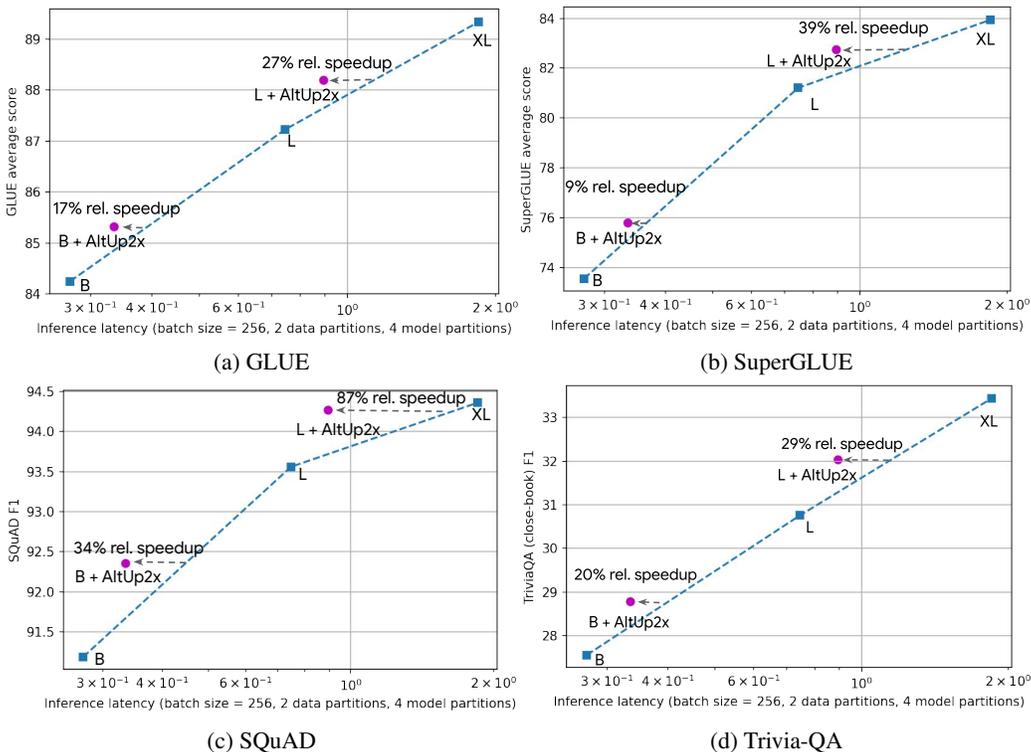


Figure 4: Evaluations of AltUp on T5 models of various sizes and popular benchmarks. AltUp consistently leads to sizeable speedups relative to baselines at the same accuracy. Latency is measured on TPUv3 with 8 cores. Relative speedup is defined as latency delta divided by AltUp latency.

255 As the figure depicts, models augmented with AltUp are uniformly faster than the extrapolated  
 256 dense models at the same accuracy. For example, we observe that a T5 large model augmented with  
 257 AltUp leads to a 27%, 39%, 87%, and 29% speedup on GLUE, SuperGLUE, SQuAD, and Trivia-QA  
 258 benchmarks, respectively. Moreover, we see that AltUp’s relative performance improves as we apply  
 259 it to larger models (compare relative speedup of T5 Base + AltUp to that of T5 Large + AltUp). This  
 260 demonstrates the scalability of AltUp to and its improved performance on even larger models. Overall,  
 261 **AltUp consistently leads to models with better predictive performance than the corresponding**  
 262 **baseline models with the same speed on all model sizes and benchmarks.**

263 **5.2 AltUp with varying representation size**

264 In the previous subsection, we had used a value of  $K = 2$  for the runs with AltUp. As discussed in  
 265 Sec. 3,  $K$  controls the width of the representation vector, and is the only hyperparameter required by  
 266 AltUp. Can we obtain even more performant models by using a larger expansion factor  $K$ ? Here, we  
 267 compare the performance of AltUp with  $K = 2$  to AltUp with a larger expansion factor  $K = 4$ .

Table 1: Performance of AltUp with varying representation dimension scaling parameter  $K$  on T5.

Model	Pretrain Accuracy	GLUE	SG	SQuAD (EM/F1)	TriviaQA (EM/F1)
S	61.21	75.83	59.52	76.44/84.97	19.03/22.83
S + AltUp (K=2)	61.86	<b>76.82</b>	<b>59.60</b>	<b>77.51/85.79</b>	<b>19.27/22.95</b>
S + AltUp (K=4)	<b>62.00</b>	76.40	59.54	76.38/84.86	19.07/22.84
B	66.42	84.25	73.56	83.78/91.19	23.1/27.56
B + AltUp (K=2)	66.96	<b>85.32</b>	75.80	<b>85.24/92.36</b>	24.35/28.78
B + AltUp (K=4)	<b>67.18</b>	84.95	<b>78.91</b>	84.82/92.07	<b>24.41/28.90</b>
L	69.13	87.23	81.21	86.77/93.56	26.15/30.76
L + AltUp (K=2)	69.32	88.20	82.75	<b>87.81/94.29</b>	27.10/32.04
L + AltUp (K=4)	<b>69.55</b>	<b>88.42</b>	<b>82.94</b>	87.59/94.02	<b>27.36/32.42</b>

268 Table 1 summarizes the results with AltUp instantiated on T5 small, base, and large sized models  
 269 with hyperparameter  $K = 2$  and  $K = 4$ . We observe that a larger value of  $K = 4$  leads to strict  
 270 improvements in pretrain accuracy over AltUp with  $K = 2$  for all models (Table 1, column 2). This  
 271 is perhaps intuitive, as a wider representation vector enables more information to be learned during  
 272 the pretraining stage. Interestingly, however, a larger  $K$  does not always lead to better finetune  
 273 performance, especially for smaller models. For example, despite having a worse pretrain accuracy,  
 274 AltUp with  $K = 2$  is better than AltUp with  $K = 4$  on all finetune tasks GLUE, SuperGLUE, and  
 275 SQuAD. We see a similar phenomenon occur for the Base model, but here  $K = 4$  is better on GLUE;  
 276 and on Large, the trend reverses:  $K = 4$  is better on every metric except for SQuAD. Our results  
 277 indicate that a larger value of  $K$  has potential to increase the performance of models on pretrain and  
 278 fine-tune metrics when AltUp is applied to larger models. *We note that there is an inherent trade-off*  
 279 *between a larger factor  $K$  and trainability, however, as a larger value of  $K$  leads to less frequent*  
 280 *activation of each sub-block which may impair performance.* We envision that practitioners can pick  
 281 a value of  $K$  other than the default  $K = 2$  to optimize performance on an application-specific basis.

282 **5.3 Recycled-AltUp**

283 Next, we consider the performance of the lightweight extension of AltUp, Recycled-AltUp, introduced  
 284 in Sec. 4. We apply Recycled-AltUp with  $K = 2$  to T5 base, large, and XL models and compare its  
 285 pretrain accuracy and speed to those of baselines. We record both the training speed and inference  
 286 speed of the resulting models. Since Recycled-AltUp does not require an expansion in the embedding  
 287 table dimension (see Sec. 4), we remark that the models augmented with it have virtually the same  
 288 number of trainable parameters as the baseline models.

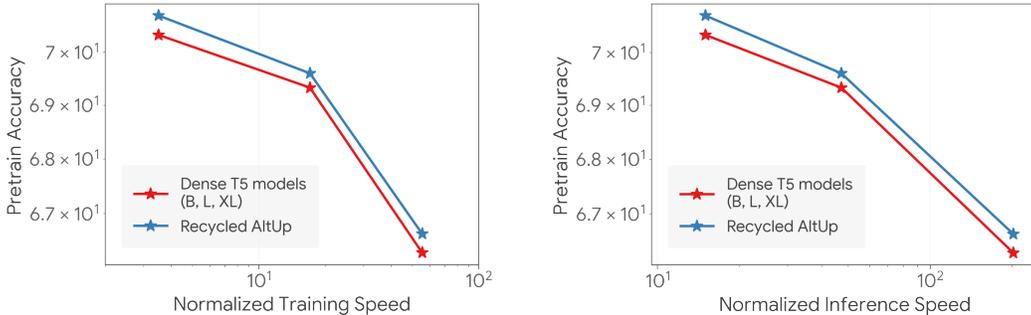


Figure 5: Recycled-AltUp on T5-B/L/XL compared to baselines. Recycled-AltUp leads to strict improvements in pretrain performance without incurring any perceptible slowdown.

289 The results of our experiments are shown in Fig. 5. The figures for both the training and inference  
 290 speed show that models with Recycled-AltUp clearly improve over baselines in pretrain accuracy,

291 without any perceptible slowdown. While Recycled-AltUp’s predictive strength generally falls below  
 292 that of standard AltUp (cf., pretrain values for AltUp in Table 1), its improved speed and reduced  
 293 parameter count may make it more suitable for certain applications. We present additional fine-tuning  
 294 results with Recycled-AltUp in the supplementary material; overall, our results demonstrate that  
 295 Recycled-AltUp is similarly effective on fine-tuning tasks.

## 296 5.4 Sequence-AltUp

297 Here, we evaluate Sequence-AltUp (from Sec. 4) to reduce the apparent sequence length for the  
 298 T5 base model. In particular, we apply average pooling, stride-and-skip, and Sequence-AltUp to  
 299 the encoder layers to reduce the apparent input sequence length. We apply stride-and-skip and  
 300 Sequence-AltUp to layers  $2, \dots, L - 1$  of the encoder, rather than all the layers, with stride length 4  
 301 as we found that this configuration results in a better accuracy/speed trade-off for both techniques. For  
 302 average pooling, the sequence length is immutably reduced from the start according to the method.

Table 2: Performance and pretrain speed of different methods for sequence length reduction on T5.

Model	Pretrain Accuracy	Finetune GLUE	Finetune SG	Speed
S	61.21	59.52	76.44/84.97	166.1
B (Baseline)	66.42	73.56	83.78/91.19	52.4
Average pooling	63.89	57.85	71.37/81.87	91.9
Stride-and-Skip	65.02	65.98	79.72/87.64	79.4
Sequence-AltUp	<b>65.39</b>	<b>66.94</b>	<b>81.67/89.37</b>	74.9

303 Table 2 presents the comparisons on pretrain and finetune metrics (GLUE and SuperGLUE) and  
 304 pretrain speed (measured by the number of sequences per second per core). The table additionally  
 305 lists the relevant metrics for T5 Base (which is the baseline model) and T5 Small as reference  
 306 points in the table. We observe that average pooling gives a large speed-up, but suffers from severe  
 307 quality degradation, especially on the finetune metrics where it performs even worse than T5 small.  
 308 Stride-and-skip and Sequence-AltUp, on the other hand, offer an improved quality and speed trade-off  
 309 relative to T5 Base. In particular, Sequence-AltUp is only slightly slower than stride-and-skip (yet,  
 310 still  $\approx 40\%$  faster than the baseline), but is much closer to the baseline model’s quality.

## 311 6 Conclusion

312 We propose the method of *Alternating Updates* (AltUp) to increase the capacity of modern transformer  
 313 models without incurring a significant increase in latency. Our approach bridges the research gap  
 314 in efficient transformers by enabling the use of wider token representations without widening the  
 315 transformer layers. AltUp utilizes lightweight prediction and correction steps to update a wider  
 316 representation vector without increasing the transformer layer’s computation cost. As a result,  
 317 we achieve strong performance improvements on language modeling and language understanding  
 318 benchmarks. We present extensions of AltUp that enable additional gains in efficiency. Given its  
 319 orthogonal scope, AltUp can be synergistically applied with existing techniques like MoE. On popular  
 320 language understanding and QA benchmarks, AltUp enables up to 87% speedup relative to the dense  
 321 baselines at the same accuracy.

322 **Limitations and future work** A current limitation of the technique we propose is the lack of a  
 323 deep theoretical understanding of its properties due to the complicated nature of rigorously analyzing  
 324 transformer models. An interesting open question is whether it would be possible to analyze AltUp  
 325 by relating its performance to a block compressed layer, and transitively relating that to a wide layer  
 326 without block compression. A deeper understanding of AltUp may also shed light on the optimal  
 327 hyperparameter  $K$  on an application-specific basis. In future work, we plan to conduct a theoretical  
 328 investigation of alternating updates to develop a deeper understanding of its effectiveness across  
 329 differing applications. We also plan to experiment with the use of a very large expansion factor  $K$ .

330 **Broader Impact** Training and deploying modern neural network models consumes colossal  
 331 amounts of resources. This leads to detrimental effects on the environment and hampers the  
 332 widespread applicability and democratization of AI. We envision that AltUp can serve as a valuable  
 333 component of efficient architectures of the future and help alleviate these negative impacts.

334 **References**

- 335 [1] Cenk Baykal, Nishanth Dikkala, Rina Panigrahy, Cyrus Rashtchian, and Xin Wang. A theoretical  
336 view on sparsely activated networks. *arXiv preprint arXiv:2208.04461*, 2022. 1, 3
- 337 [2] Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie  
338 Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark,  
339 et al. Improving language models by retrieving from trillions of tokens. In *International  
340 Conference on Machine Learning*, pages 2206–2240. PMLR, 2022. 3
- 341 [3] John Charles Butcher. *Numerical methods for ordinary differential equations*. John Wiley &  
342 Sons, 2016. 4
- 343 [4] Zixiang Chen, Yihe Deng, Yue Wu, Quanquan Gu, and Yuanzhi Li. Towards understanding  
344 mixture of experts in deep learning. *arXiv preprint arXiv:2208.02813*, 2022. 3
- 345 [5] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane,  
346 Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking  
347 attention with performers. *arXiv preprint arXiv:2009.14794*, 2020. 3, 6
- 348 [6] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
349 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:  
350 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 1, 2
- 351 [7] Aidan Clark, Diego de Las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan  
352 Hoffmann, Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, et al. Unified  
353 scaling laws for routed language models. In *International Conference on Machine Learning*,  
354 pages 4057–4086. PMLR, 2022. 3
- 355 [8] Zihang Dai, Guokun Lai, Yiming Yang, and Quoc Le. Funnel-transformer: Filtering out  
356 sequential redundancy for efficient language processing. *Advances in neural information  
357 processing systems*, 33:4271–4282, 2020. 3, 6
- 358 [9] Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal,  
359 Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. Beyond english-  
360 centric multilingual machine translation. *J. Mach. Learn. Res.*, 22(107):1–48, 2021. 3
- 361 [10] William Fedus, Jeff Dean, and Barret Zoph. A review of sparse expert models in deep learning.  
362 *arXiv preprint arXiv:2209.01667*, 2022. 1, 3, 15
- 363 [11] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion  
364 parameter models with simple and efficient sparsity, 2021. 3, 16
- 365 [12] Zhengcong Fei, Xu Yan, Shuhui Wang, and Qi Tian. Deecap: dynamic early exiting for efficient  
366 image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
367 Recognition*, pages 12216–12226, 2022. 3
- 368 [13] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint  
369 arXiv:1410.5401*, 2014. 3
- 370 [14] Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-  
371 Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou,  
372 et al. Hybrid computing using a neural network with dynamic external memory. *Nature*,  
373 538(7626):471–476, 2016. 3
- 374 [15] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han,  
375 Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented  
376 transformer for speech recognition. *arXiv preprint arXiv:2005.08100*, 2020. 6
- 377 [16] Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung,  
378 and Yinfei Yang. Longt5: Efficient text-to-text transformer for long sequences. *arXiv preprint  
379 arXiv:2112.07916*, 2021. 3, 6

- 380 [17] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval  
381 augmented language model pre-training. In *International Conference on Machine Learning*,  
382 pages 3929–3938. PMLR, 2020. 3
- 383 [18] Hussein Hazimeh, Zhe Zhao, Aakanksha Chowdhery, Maheswaran Sathiamoorthy, Yihua Chen,  
384 Rahul Mazumder, Lichan Hong, and Ed Chi. Dselect-k: Differentiable selection in the mixture  
385 of experts with applications to multi-task learning. *Advances in Neural Information Processing*  
386 *Systems*, 34:29335–29347, 2021. 3
- 387 [19] Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. Scaling laws for  
388 transfer. *arXiv preprint arXiv:2102.01293*, 2021. 1, 2, 4
- 389 [20] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
390 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.  
391 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022. 1
- 392 [21] Lu Hou, Zhiqi Huang, Lifeng Shang, Xin Jiang, Xiao Chen, and Qun Liu. Dynabert: Dynamic  
393 bert with adaptive width and depth. *Advances in Neural Information Processing Systems*,  
394 33:9782–9793, 2020. 3
- 395 [22] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures  
396 of local experts. *Neural computation*, 3(1):79–87, 1991. 3
- 397 [23] Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large  
398 scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint*  
399 *arXiv:1705.03551*, 2017. 7
- 400 [24] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
401 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
402 models. *arXiv preprint arXiv:2001.08361*, 2020. 1, 2, 4
- 403 [25] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer.  
404 *arXiv preprint arXiv:2001.04451*, 2020. 3, 6
- 405 [26] Guillaume Lample, Alexandre Sablayrolles, Marc’Aurelio Ranzato, Ludovic Denoyer, and  
406 Hervé Jégou. Large memory layers with product keys. *Advances in Neural Information*  
407 *Processing Systems*, 32, 2019. 3
- 408 [27] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang,  
409 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with condi-  
410 tional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020. 1
- 411 [28] Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. Base layers:  
412 Simplifying training of large, sparse models. In *International Conference on Machine Learning*,  
413 pages 6265–6274. PMLR, 2021. 3
- 414 [29] Lei Li, Yankai Lin, Deli Chen, Shuhuai Ren, Peng Li, Jie Zhou, and Xu Sun. Cascadebert:  
415 Accelerating inference of pre-trained language models via calibrated complete models cascade.  
416 *arXiv preprint arXiv:2012.14682*, 2020. 3
- 417 [30] Lucas Liebenwein, Cenk Baykal, Brandon Carter, David Gifford, and Daniela Rus. Lost in  
418 pruning: The effects of pruning neural networks beyond test accuracy. *Proceedings of Machine*  
419 *Learning and Systems*, 3:93–138, 2021. 1
- 420 [31] Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. Sparse, dense, and  
421 attentional representations for text retrieval. *Transactions of the Association for Computational*  
422 *Linguistics*, 9:329–345, 2021. 2
- 423 [32] Aditya Menon, Sadeep Jayasumana, Ankit Singh Rawat, Seungyeon Kim, Sashank Reddi, and  
424 Sanjiv Kumar. In defense of dual-encoders for neural ranking. In *International Conference on*  
425 *Machine Learning*, pages 15376–15400. PMLR, 2022. 2
- 426 [33] Rina Panigrahy, Xin Wang, and Manzil Zaheer. Sketch based memory for neural networks. In  
427 *International Conference on Artificial Intelligence and Statistics*, pages 3169–3177. PMLR,  
428 2021. 1, 3, 15

- 429 [34] David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel  
430 Rothchild, David So, Maud Texier, and Jeff Dean. Carbon emissions and large neural network  
431 training. *arXiv preprint arXiv:2104.10350*, 2021. 1
- 432 [35] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,  
433 Yanqi Zhou, Wei Li, Peter J Liu, et al. Exploring the limits of transfer learning with a unified  
434 text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020. 2, 7, 14
- 435 [36] Samyam Rajbhandari, Conglong Li, Zhewei Yao, Minjia Zhang, Reza Yazdani Aminabadi,  
436 Ammar Ahmad Awan, Jeff Rasley, and Yuxiong He. Deepspeed-moe: Advancing mixture-of-  
437 experts inference and training to power next-generation ai scale. In *International Conference*  
438 *on Machine Learning*, pages 18332–18346. PMLR, 2022. 15
- 439 [37] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions  
440 for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016. 7
- 441 [38] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André  
442 Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts.  
443 *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021. 1
- 444 [39] Adam Roberts, Hyung Won Chung, Anselm Levskaya, Gaurav Mishra, James Bradbury, Daniel  
445 Andor, Sharan Narang, Brian Lester, Colin Gaffney, Afroz Mohiuddin, Curtis Hawthorne, Aitor  
446 Lewkowycz, Alex Salcianu, Marc van Zee, Jacob Austin, Sebastian Goodman, Livio Baldini  
447 Soares, Haitang Hu, Sasha Tsvyashchenko, Aakanksha Chowdhery, Jasmijn Bastings, Jannis  
448 Bulian, Xavier Garcia, Jianmo Ni, Andrew Chen, Kathleen Kenealy, Jonathan H. Clark, Stephan  
449 Lee, Dan Garrette, James Lee-Thorp, Colin Raffel, Noam Shazeer, Marvin Ritter, Maarten  
450 Bosma, Alexandre Passos, Jeremy Maitin-Shepard, Noah Fiedel, Mark Omernick, Brennan  
451 Saeta, Ryan Sepassi, Alexander Spiridonov, Joshua Newlan, and Andrea Gesmundo. Scaling  
452 up models and data with  $\tau 5x$  and `seq.io`. *arXiv preprint arXiv:2203.17189*, 2022. 14
- 453 [40] Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the  
454 parameters of a language model? *arXiv preprint arXiv:2002.08910*, 2020. 7
- 455 [41] Stephen Roller, Sainbayar Sukhbaatar, Jason Weston, et al. Hash layers for large sparse models.  
456 *Advances in Neural Information Processing Systems*, 34:17555–17566, 2021. 1, 3
- 457 [42] Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Q Tran, Yi Tay, and  
458 Donald Metzler. Confident adaptive language modeling. *arXiv preprint arXiv:2207.07061*,  
459 2022. 3
- 460 [43] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,  
461 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts  
462 layer. *arXiv preprint arXiv:1701.06538*, 2017. 1, 3, 15
- 463 [44] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory  
464 cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR, 2018. 14
- 465 [45] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey.  
466 *ACM Computing Surveys (CSUR)*, 2020. 3
- 467 [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
468 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural*  
469 *Information Processing Systems*, volume 30, 2017. 1
- 470 [47] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix  
471 Hill, Omer Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose  
472 language understanding systems. *Advances in neural information processing systems*, 32, 2019.  
473 7
- 474 [48] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.  
475 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv*  
476 *preprint arXiv:1804.07461*, 2018. 7

- 477 [49] Lemeng Wu, Mengchen Liu, Yinpeng Chen, Dongdong Chen, Xiyang Dai, and Lu Yuan.  
478 Residual mixture of experts. *arXiv preprint arXiv:2204.09636*, 2022. [3](#)
- 479 [50] Qingyang Wu, Zhenzhong Lan, Jing Gu, and Zhou Yu. Memformer: The memory-augmented  
480 transformer. *arXiv preprint arXiv:2010.06891*, 2020. [3](#)
- 481 [51] Yuhuai Wu, Markus N Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing trans-  
482 formers. *arXiv preprint arXiv:2203.08913*, 2022. [3](#)
- 483 [52] Greg Yang and Edward J Hu. Feature learning in infinite-width neural networks. *arXiv preprint*  
484 *arXiv:2011.14522*, 2020. [2](#), [4](#)
- 485 [53] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui  
486 Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint*  
487 *arXiv:2205.01917*, 2022. [1](#)
- 488 [54] Seniha Esen Yuksel, Joseph N Wilson, and Paul D Gader. Twenty years of mixture of experts.  
489 *IEEE transactions on neural networks and learning systems*, 23(8):1177–1193, 2012. [3](#)
- 490 [55] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santi-  
491 ago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers  
492 for longer sequences. *Advances in neural information processing systems*, 33:17283–17297,  
493 2020. [3](#), [6](#)
- 494 [56] Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew Dai,  
495 Zhifeng Chen, Quoc Le, and James Laudon. Mixture-of-experts with expert choice routing.  
496 *arXiv preprint arXiv:2202.09368*, 2022. [1](#)
- 497 [57] Barret Zoph, Irwan Bello, S ameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer,  
498 and William Fedus. Designing effective sparse expert models. *arXiv preprint arXiv:2202.08906*,  
499 2022. [2](#), [3](#), [16](#)

## 500 **Supplementary Material for *Alternating Updates for Efficient Transformers***

501 In this supplementary, we present the full details and hyperparameters of our evaluations, provide  
502 details of our algorithmic contributions, and present complementary empirical results that support the  
503 effectiveness the presented work.

### 504 **A Experimental Setup**

505 We experiment with various T5 baseline model sizes: small (S), base (B), large (L), and XL. The base  
506 (B), large (L), and XL models follow the same model configurations as in the T5 paper, while the small  
507 model is shallower than the T5 paper [35] to cover a larger range of model sizes (4 encoder/decoder  
508 layers instead of 8 encoder/decoder layers). In particular, we use the T5 version 1.1 models with gated  
509 GELU feedforward network and pre layer norm. The models are implemented on top of the T5X [39]  
510 codebase. During pretraining, we use 256 batch size, Adafactor optimizer [44] with base learning rate  
511 1.0 and reciprocal square-root decay with 10000 warmup steps, and zero dropout. During finetuning,  
512 we use 256 batch size, Adafactor optimizer with constant learning rate of 0.001 and 0.1 dropout.  
513 Unless explicitly stated, we pretrain for 500k steps and finetune for 50k steps. Our experiments were  
514 implemented in Python and run on TPUv3 with 8 cores.

### 515 **B Parameter Counts and Speed**

516 Here, we present the number of additional parameters needed by adding AltUp, its speed, and its  
517 pretrain accuracy on T5 models of varying sizes. In the following tables, the embedding param-  
518 eters include input embedding table parameters (shared between encoder and decoder) and output  
519 embedding table. Non-embedding parameters include all the transformer blocks. Train speed is  
520 measured by number of examples per second per core. Table 3 documents the parameter count and  
521 training speed comparison. Note that Alternating Updates increases the number of embedding param-  
522 eters while leaving the non-embedding parameters roughly the same. Since the narrow transformer  
523 layer computation is not changed by alternating updates and since the predict and correct steps are  
524 lightweight (see Sec. 3), we incur a relatively small increase in the computation cost compared to a  
525 dense 2x width model.

Model	# emb params	# non-emb params	train speed
S	3.29E+07	3.78E+07	166.1
S + AltUp	6.58E+07	3.99E+07	119.4
B	4.93E+07	1.98E+08	52.4
B + AltUp	9.87E+07	2.12E+08	42.3
L	6.58E+07	7.17E+08	17.1
L + AltUp	1.32E+08	7.68E+08	14.4

Table 3: Model size and train speed comparisons on T5X models with AltUp instantiated with  $K = 2$ .

526 Table 4 documents the parameter count, training speed and pretrain accuracy comparison when the  
527 representation dimension is scaled up with AltUp or dense scaling. Note that Alternating Updates  
528 increases the number of embedding parameters while leaving the non-embedding parameters roughly  
529 the same, providing an efficient way to scale up the representation dimension relative to a  $K$ -times  
530 wider model.

531 Table 5 contains pretrain performances for T5 XL sized models. We note the AltUp technique  
532 continue to offer quality boost at the billion parameters scale (note that T5XL has roughly 3B  
533 parameters), suggesting that AltUp is a robust technique for increasing model capacity for modern  
534 large language models.

### 535 **C Combination with MoE**

536 Here, we investigate whether AltUp can be combined with orthogonal techniques, namely MoE,  
537 to obtain additive performance gains in pretrain accuracy for T5 small, base, and large models. In

Model	# emb params	# non-emb params	Train speed	Pretrain accuracy
T5 Base	4.93E+07	1.98E+08	52.4	65.29
T5 Base + AltUp2x	9.87E+07	2.12E+08	42.3	65.78
T5 Base + Dense2X	9.87E+07	3.97E+08	32.9	66.45
T5 Base + AltUp4x	1.97E+08	2.41E+08	28.1	66.00
T5 Base + Dense4X	1.97E+08	7.93E+08	12.6	67.01

Table 4: AltUp compared with dense scaling evaluated at 250k pretrain steps.

Model	# emb params	# non-emb params	Train speed	Pretrain accuracy
T5 XL	1.32E+08	2.72E+09	3.6	70.01
T5 XL + AltUp2x	2.63E+08	2.92E+09	3.0	70.61

Table 5: Pretrain performances for T5 XL sized models. Pretrain accuracy is measured at 400k steps. AltUp continues to offer a performance boost even on the scale of models with billions of parameters.

538 particular, we consider the *partial experts* setting similar to [36, 33], where at each layer, in addition to  
 539 the layer’s module, we route the input to a smaller expert module and combine the outputs of the  
 540 main and auxiliary modules as the input to the subsequent layer (see Fig. 6).

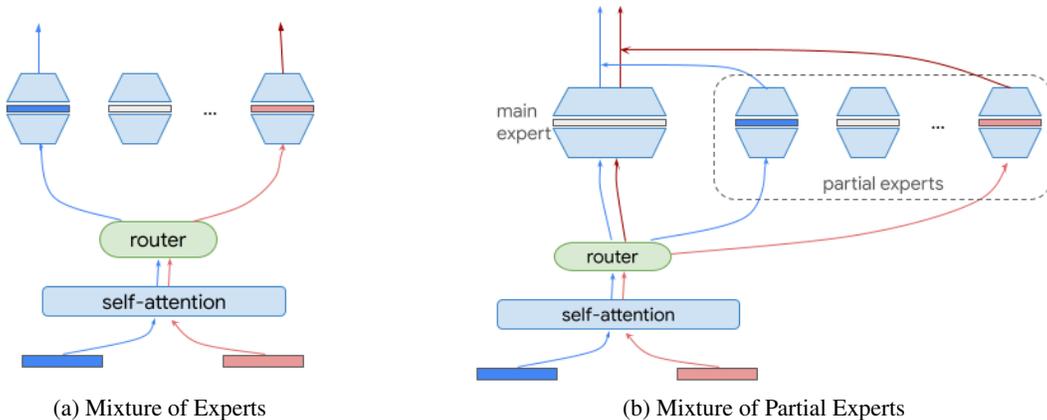


Figure 6: The partial experts setting in the context of the evaluations in Sec. C. The standard MoE model (left) routes the inputs to one or more of  $n$  experts based on a routing function. Mixture of Partial Experts (right) always routes the input to the main expert and additionally routes the input to one or more partial experts; the output is a function of the main expert’s and partial experts’ outputs.

The MoE layer routes an input token  $x$  to  $k$  of  $n$  experts where each expert is itself a parametrized subnetwork (e.g., a fully-connected layer). Following [10], we let  $\{E_i(\cdot)\}_{i \in [n]}$  and  $E_i(x)$  denote the set of experts and the output of lookup the input token  $x$  to expert  $i$ , respectively. For an input token  $x$ , a learnable weight matrix  $W$  is applied to obtain the logits  $h(x) = Wx$ . The lookup probabilities are computed by taking the softmax of  $h(x)$

$$p_i(x) = \frac{\exp(h_i(x))}{\sum_{j \in [n]} \exp(h_j(x))} \quad \forall i \in [n].$$

The token  $x$  is routed to the expert(s)  $\mathcal{T} \subset [n]$  with the top- $k$  probabilities  $p(x)$ . Since this operation is not differentiable, the output  $y$  is computed as a probability weighted combination of the experts’ outputs to enable gradients to propagate back to the router parameters [43], i.e.,

$$y = \sum_{i \in \mathcal{T}} p_i(x) E_i(x).$$

541 For MoE, we used the simplified implementation of the top-1 softmax routing of [11]. We use 128  
 542 experts each per encoder and decoder layer, with each expert representing a 2 layer fully-connected  
 543 neural network with hidden dimension 16. For sake of the synergistic demonstration even with the  
 544 core MoE implementation, we did not incorporate a sophisticated mechanism for load balancing such  
 545 as load balancing loss [11] or router z loss [57]. We use multiplicative jitter noise sampled from a  
 546 uniform distribution over  $[1 - \varepsilon, 1 + \varepsilon]^{d_{in}}$  with  $\varepsilon = 0.01$ . The router matrix  $W$  was initialized by  
 547 drawing from a zero mean Normal distribution with standard deviation  $2 \times 10^{-2}$ .

Method	T5 Small	T5 Base	T5 Large
Baseline	59.10	63.35	65.58
MoE [57]	59.42	63.62	65.71
AltUp (K=2)	59.67	63.97	65.73
AltUp (K=2) + MoE	<b>59.91</b>	<b>64.13</b>	<b>65.95</b>

Table 6: Pretrain accuracy at 100k steps of T5 models augmented with Alternating Updates (see Sec. 3) and MoE. MoE synergistically combines with Alternating Updates and enables further increases in model capacity.

548 Table 6 synthesizes the pretraining performance on the C4 dataset at 100k training steps of AltUp  
 549 and compared techniques on T5 Small, Base, and Large models. Perhaps most notably, we show that  
 550 combining AltUp and MoE leads to even further sizable improvements in the pretrain performance  
 551 (last row of Table 6). For example, the combination of MoE and AltUp improves over the baseline by  
 552 0.81, over AltUp alone by 0.24, and over MoE alone by 0.49. For all model sizes, the combination  
 553 of AltUp and MoE is synergistic and leads to significant improvements compared to not only the  
 554 baseline, but also to each approach in isolation.

## 555 D Alternating updates with varying block selection

556 In this section, we present empirical results on the alternating updates technique and comparison  
 557 with other techniques that widen token representation vectors. We increase the token representation  
 558 dimension by a factor of 2 (corresponding to  $K = 2$  in Algorithm 1) unless otherwise specified. The  
 559 model capacity increase comes from a wider embedding table at the bottom layer of the model while  
 560 the transformer layers remain the same, which results in minimal additional computation cost.

Model	Pretrain Accuracy	Finetune GLUE	Finetune SG	Finetune SQuAD (EM/F1)
S (baseline)	61.21	75.83	59.52	76.44/84.97
S + Sum	61.67	77.54	59.63	75.06/83.82
S + SameUp	<b>61.91</b>	<b>77.75</b>	<b>60.81</b>	76.85/85.51
S + AltUp	61.86	76.82	59.60	<b>77.51/85.79</b>
B (baseline)	66.42	84.25	73.56	83.78/91.19
B + Sum	66.82	84.85	75.2	84.36/91.36
B + SameUp	66.82	84.06	74.15	84.41/91.76
B + AltUp	<b>66.96</b>	<b>85.32</b>	<b>75.80</b>	<b>85.24/92.36</b>
L (baseline)	69.13	87.23	81.21	86.77/93.56
L + Sum	69.09	86.18	78.93	86.19/93.08
L + SameUp	<b>69.45</b>	87.95	82.72	<b>87.65/94.13</b>
L + AltUp	69.32	<b>88.20</b>	<b>82.75</b>	87.58/ <b>94.27</b>

Table 7: Comparison of Algorithm 1 with various sub-block selection methods on T5-S/B/L.

561 In Table 7, we compare the summation method (Sum) in which additional embedding vectors are  
 562 added to the token representation vector, Algorithm 1 with same block selection (SameUp), and  
 563 Algorithm 1 with alternating block selection (AltUp), all on top of the T5 version 1.1 small (S), base  
 564 (B), and large (L) models. We observe the Prediction-Compute-Correct scheme as described in Sec. 3  
 565 with *same* and *alternating* block selection methods outperforms the summation method. For the  
 566 small models, *same* block selection method performs better in most tasks, while for the base and

567 large models, *alternating* block selection method performs better in most tasks. We note all three  
 568 methods bring improvements in both pretraining and finetuning, and AltUp is generally the most  
 569 effective one. While pretraining accuracies for all three methods are mostly similar, differences in  
 570 finetuning metrics are large, and AltUp generally achieves roughly twice the gains of the other two  
 571 variants. Moreover, we observe that the gains of AltUp in pretraining accuracies show diminishing  
 572 returns when model sizes grows, but the gains in finetuning metrics do not.

## 573 E Sequence-AltUp Details

574 Here, we provide the full pseudocode of Sequence-AltUp from Sec. 4 (see Alg. 2).

---

### Algorithm 2 AltUp extension to sequence dimension

---

**Input:** A sequence of vectors  $x = (x_0, x_2, \dots, x_{T-1})$ , where  $x_i \in \mathbb{R}^d$ . Transformer layer  $\mathcal{L}$  and stride parameter  $k$ .

**Output:** A sequence of vectors  $y = (y_0, y_2, \dots, y_{T-1})$ .

1: **Prediction:** predict the output sequence with a trainable linear map:

$$\hat{y}_i = a_1 x_i + a_2 x_{\lfloor i/k \rfloor * k}$$

for  $i = 0, 1, \dots, T - 1$ , where  $a_1, a_2 \in \mathbb{R}$  are trainable scalars;

2: **Computation:** subsample the input sequence with stride  $k$  and apply the transformer layer on the subsampled sequence:

$$(\tilde{y}_0, \tilde{y}_k, \dots, \tilde{y}_{\lfloor (T-1)/k \rfloor * k}) = \mathcal{L}(x_0, x_k, \dots, x_{\lfloor (T-1)/k \rfloor * k});$$

3: **Correction:** correct the prediction with the computation result:

$$y_i = \hat{y}_i + b(\tilde{y}_{\lfloor i/k \rfloor * k} - \hat{y}_{\lfloor i/k \rfloor * k})$$

for  $i = 0, 1, \dots, T - 1$ , where  $b \in \mathbb{R}$  is a trainable scalar.

---

## 575 F Recycled-AltUp Fine-tune Evaluations

576 We conclude the supplementary material by presenting evaluations with Recycled-AltUp as described  
 577 in Sec. 4. Table 8 presents the results of our pretrain and fine-tune evaluations on T5 Small, Base,  
 578 and Large. The pretrain accuracy is the one reported at 500k steps, and we fine-tune for an additional  
 579 50k steps for the fine-tune evaluations.

Model	Pretrain Acc.	GLUE	SG	SQuAD (EM/F1)	TriviaQA (EM/F1)
S	61.21	75.83	59.52	76.44/84.97	19.03/22.83
S + Recycled-AltUp	61.33	<b>77.24</b>	59.12	<b>77.76</b> /85.64	19.06/22.77
S + AltUp	<b>61.86</b>	76.82	<b>59.60</b>	77.51/ <b>85.79</b>	<b>19.27/22.95</b>
B	66.42	84.25	73.56	83.78/91.19	23.1/27.56
B + Recycled-AltUp	66.63	<b>85.60</b>	74.83	84.81/91.93	22.72/27.15
B + AltUp	<b>66.96</b>	85.32	<b>75.80</b>	<b>85.24/92.36</b>	<b>24.35/28.78</b>
L	69.13	87.23	81.21	86.77/93.56	26.15/30.76
L + Recycled-AltUp	69.30	87.91	82.53	87.37/93.88	<b>27.51/32.38</b>
L + AltUp	<b>69.32</b>	<b>88.20</b>	<b>82.75</b>	<b>87.81 / 94.29</b>	27.10/32.04

Table 8: The performance of baseline, Recycled-AltUp, and AltUp on pretrain and fine-tune evaluation metrics. Recycled-AltUp and AltUp were instantiated with  $K = 2$  for all evaluations.

580 Consistent with our results presented in the main body of our paper, we observe that AltUp and  
 581 Recycled-AltUp both provide clear and consistent gains over the baseline on virtually all pretrain  
 582 and fine-tune metrics. As conjectured in Sec. 4 Recycled-AltUp generally does not provide the full  
 583 benefits of AltUp in terms of pretrain and fine-tune accuracies, however, this gap seems to shrink  
 584 for larger models. Moreover, Recycled-AltUp has the appeal that it practically adds no additional

585 parameters to the model and, as a result, has roughly the same speed as the baseline model (see  
586 Fig. 5). We envision that Recycled-AltUp's improved speed and reduced parameter count may make  
587 it more appealing for certain applications.