

---

# Efficient Sequence Packing without Cross-contamination: Accelerating Large Language Models without Impacting Performance

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Effective training of today’s large language models (LLMs) depends on large  
2        batches and long sequences for throughput and accuracy. To handle variable-length  
3        sequences on hardware accelerators, it is common practice to introduce padding  
4        tokens, so that all sequences in a batch have the same length. We show in this paper  
5        that the variation in sequence lengths in common NLP datasets is such that up to  
6        50% of all tokens can be padding. In less common, but not extreme, cases (e.g.  
7        GLUE-cola with sequence length 128), the ratio is up to 89%. Existing methods  
8        to address the resulting inefficiency are complicated by the need to avoid ‘cross-  
9        contamination’ in self-attention, by a reduction in accuracy when sequence ordering  
10       information is lost, or by customized kernel implementations only valid for specific  
11       accelerators. This paper introduces a new formalization of sequence packing in  
12       the context of the well-studied bin packing problem, and presents new algorithms  
13       based on this formulation which, for example, confer a 2x speedup for phase 2  
14       pre-training in BERT. We show how existing models can be adapted to ensure  
15       mathematical equivalence between the original and packed models, meaning that  
16       packed models can be trained with existing pre-training and fine-tuning practices.

## 17    1 Introduction

18    Many language datasets, including the de-facto pre-training dataset for BERT—Wikipedia, have  
19    a skewed distribution of sequence lengths (see Figure 1). However, typical machine learning  
20    accelerators, and their corresponding libraries, exhibit poor performance when processing variable-  
21    length workloads. A simple mitigation is to set a maximum sequence length, and to pad shorter  
22    sequences with padding tokens. This naive batching is widely used and provided in the vanilla BERT  
23    implementation as well as the Hugging Face framework [32]. Its effect is enhanced by the offline  
24    dataset generation process which, in BERT, attempts to “pack” together sentences so as to fill the  
25    sequence length as completely as possible [8]. We improve this process at a whole-dataset level.

26    We show that, even after this pre-processing, padding tokens represent 50% of all tokens of the  
27    Wikipedia pre-training dataset at sequence length 512. Thus, by avoiding processing the padding  
28    tokens one can get a 2x speed-up for phase 2. Overall, the lengths range between 5 tokens up to 512.  
29    Samples of length 512 represent only 23.5% of the dataset,

30    Beyond the simple batching, other solutions have been addressed in the literature, and in open-source  
31    software implementations. When processing sequences, most libraries and algorithms mention  
32    packing as reference to concatenating sentences from the same document (BERT) or from different  
33    documents (BERT, T5 [24], GPT-3 [4], and RoBERTa [16]) as they arrive (GREEDY) from the  
34    source dataset to generate the training dataset. None of the respective papers addresses the packing

35 efficiency, i.e., remaining fraction of padding. To “separate” sequences from different documents, a  
 36 separator token is introduced. However, this is not sufficient and can have a significant impact on  
 37 performance. This is discussed only in the RoBERTa paper which shows that downstream F1 scores  
 38 get consistently reduced on average by 0.35%. Alternative common approaches to overcome the large  
 39 amount of padding in many datasets are “un-padding” as in Effective Transformer [5] and sorted  
 40 batching (SORT) as in Faster Transformer [21], lingvo [28] fairseq [22], and RoBERTa. However, for  
 41 running efficiently on arbitrary accelerators, these approaches require substantial hardware-specific  
 42 low-level code optimizations only available on GPUs. Further details are in Sections C [1] and 4.4.

43 Beyond language models, packing has been also present in other areas of machine learning, however  
 44 with little to no exploration in the literature and mostly hidden in some libraries without any further  
 45 discussion. For example, PyG (PyTorch Geometric) combines multiple small graphs in a batch to  
 46 account for the large variation in size and to optimize the hardware usage when training a Graph  
 47 Neural Network (GNN). Another example is the RNN implementation in PyTorch which introduces a  
 48 “PackedSequence” object and states that “All RNN modules accept packed sequences as inputs” but  
 49 does not address how sequences are packed efficiently and how the processing of packed sequences  
 50 is implemented in an efficient manner while avoiding interaction between sequences. Even though  
 51 we focus on BERT [6] and other transformers in this paper, the general principles can be transferred  
 52 to many more machine learning algorithms with differently sized data samples.

53 In this paper, we formally frame the packing problem in transformer based models, and provide some  
 54 solutions, showing that sequences can be packed efficiently, separator tokens are not required, and  
 55 cross-contamination can be avoided with little overhead.

56 In summary, the contributions of the paper are as follows. In Section 2 we produce histograms of a  
 57 variety of datasets showing the high percentage of padding tokens. In Section 3.1, we present two new  
 58 deterministic and efficient packing algorithms based on established solvers which efficiently pack  
 59 datasets with millions of sequences in a matter of seconds (or less). In Section 3.2 and Section 3.3, we  
 60 describe ‘cross-contamination’ —the cause of the accuracy reduction which separator tokens do not  
 61 mitigate— and show how the BERT model can be adjusted to show the same convergence behavior  
 62 on packed and unpacked sequences. We empirically show that the proposed packing algorithms  
 63 produce a nearly-optimal packing scheme for Wikipedia pre-training dataset (Section 4.1) and more  
 64 in the Appendix. In Section 4.2, we demonstrate that the convergence of the BERT large model on  
 65 the packed dataset is equivalent to that on the un-packed dataset with 2x throughput increase on the  
 66 Wikipedia sequence length 512 pre-training dataset. Further experiments underline the necessity and  
 67 efficiency of our changes.

## 68 2 Sequence length distributions

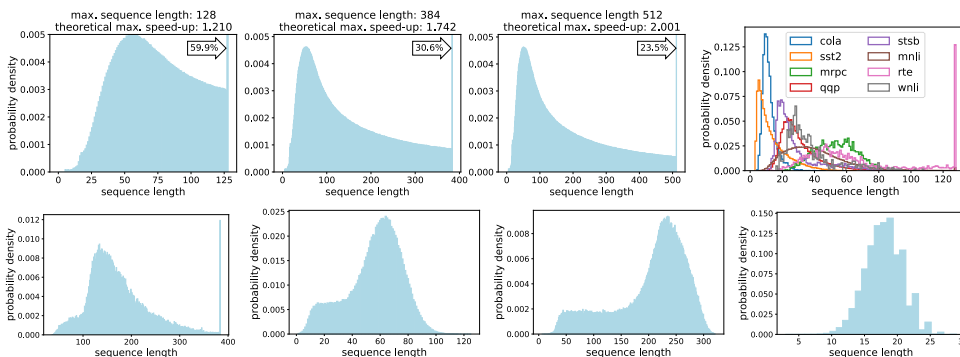


Figure 1: Sequence length distributions for different datasets. The three graphics at the top left show Wikipedia BERT pre-training dataset sequence length histograms (token count excluding padding) for different maximum sequence lengths based on the Wikipedia article dump from October 1st 2020. The theoretical speed-up relates to not using any padding tokens and not having any overhead from processing the different lengths. Top right: GLUE datasets. Bottom from left to right: SQuAD 1.1, LibriSpeech text labels, LibriSpeech audio token sequence, and QM9 molecules of a graph in a sequence.

69 BERT is pre-trained using masked-language modelling and next-sentence prediction on a large  
70 corpus of Wikipedia articles. Each sequence is composed of one  $\langle \text{CLS} \rangle$  token followed by the  
71 first “segment” of sentences, followed by a  $\langle \text{SEP} \rangle$  token, and then finally the second “segment” of  
72 sentences. Because these “segments” are created in sentence-level increments there is no token-level  
73 control of sequence length. Furthermore 10% (default value, [7]) of sequences are intentionally  
74 cut short. This leads to significant levels of padding, especially for longer maximum sequence  
75 lengths (see Figure 1 and Section 3.1). At sequence length 128 (commonly used in phase 1 of  
76 pre-training) the theoretical speed-up is around 1.2, at sequence length 384 this increases to 1.7, and  
77 finally at sequence length 512 (commonly used for phase 2 of pre-training) it is 2.0. Despite the  
78 widespread use of the Wikipedia dataset for pre-training BERT such histograms have, to the best  
79 of our knowledge, not been published previously. This has perhaps led to the underestimation of  
80 the speed-up opportunity available. To put things into perspective, the sequence length 512 dataset  
81 contains 8.33 billion tokens, of which 4.17 billion are padding tokens.

82 Note that the skewed sequence length distributions are neither limited to Wikipedia, as shown with  
83 GLUE [30, 31] from Section 3.1 and SQuAD 1.1 [25] from Section 3.2 (2.2x speed up), to BERT  
84 training, as shown with LibriSpeech text distributions [23] from Section 3.1, nor to text itself,  
85 given the LibriSpeech audio data distributions, and the QM9 molecular data [27, 26] (1.6x speed-up,  
86 Section 3.1). All distributions can be found in Figure 1. Since LibriSpeech audio data is skewed to  
87 longer sequences, only 1.3x speed-up could be achieved despite the theoretical maximum of 1.6x.  
88 For all other cases, the algorithms presented in Section 3.1 lead to close to optimal packing.

## 89 3 Methods

90 Our approach consists of three distinct components. Firstly, we pack the  $n$  data samples efficiently  
91 during pre-processing to make full use of the maximum sequence length,  $s_m$  (Sections 3.1 and F).  
92 Secondly, we introduce a series of model changes in Section 3.2 that preserve the equivalence with  
93 the original BERT implementation. The changes include a self-attention mask to prevent the model  
94 from attending between different sequences in the same pack (Section 3.2.2), and an adjustment  
95 of the positional embeddings (Section 3.2.1) to handle packs of sequences. Other components  
96 of the model, such as the feed-forward layer [29], operate on a per-token basis and do not require  
97 modification for pre-training. In Section 3.2.3, we also demonstrate how to compute a per-sequence  
98 loss and accuracy for NSP and downstream fine-tuning tasks. Thirdly, we provide suggestions for  
99 hyperparameter adjustment (Section 3.3) that lead to analogous convergence behavior between the  
100 packed and un-packed BERT implementations. Additional videos and animations are provided as  
101 supplemental material.

### 102 3.1 Packing algorithms

103 The widely studied and well established bin packing problem deals with the assignment of items into  
104 bins of a fixed capacity such that the number of utilized bins is minimized. It has been known for  
105 decades if not centuries. Since an exact solution is strongly NP-complete [14], numerous approximate  
106 solutions have been proposed [12, 15, 13, 36]. Since most existing approximations have a high  
107 complexity of at least  $O(n \log n)$ , we propose two new heuristic offline algorithms that are tailored  
108 to the NLP setting applied to the whole dataset. For a detailed introduction to packing see Section F.

#### 109 3.1.1 Shortest-pack-first histogram-packing (SPFHP)

110 Shortest-pack-first histogram-packing (SPFHP) works on the bins in the sequence length histogram  
111 (with bin size 1) rather than the individual samples. The histogram is traversed in sorted order from  
112 longest to shortest sequences. Then, to pack the data during the traversal, we apply the worst-fit  
113 algorithm [12, 36] such that the histogram bin being processed goes to the “**pack**”<sup>1</sup> that has the most  
114 space remaining (“shortest-pack-first”). If the histogram bin does not fit completely, a new pack is  
115 created. We also limit the **packing depth**, in other words the maximum number of sequences that  
116 are allowed in a pack. Therefore, an existing pack is only extended if it is not already at maximum  
117 packing depth. The detailed code for the algorithm is provided in Listing 3. The time and space  
118 complexity of the algorithm are  $O(n + s_m^2)$  and  $O(s_m^2)$  (Section G.2.1).

<sup>1</sup>We avoid the ambiguous terms “bin” and “sample/sequence” and use “pack” instead to refer to the multiple sequences concatenated during packing.

### 119 3.1.2 Non-negative least squares histogram-packing (NNLSHP)

120 The proposed NNLSHP algorithm is based on re-stating the packing problem as a (weighted) non-  
121 negative least squares problem (NNLS) [3] of the form  $wAx = wb$  where  $x \geq 0$ . The vector  $b$  is the  
122 histogram containing the counts of all the sequence lengths in the dataset. Next, we define the  $A$   
123 matrix (the “packing matrix”) by first generating a list of all possible sequence length combinations  
124 (“strategies”) that add up exactly to the maximum sequence length. We focus specifically on strategies  
125 that consist of at most 3 sequences per pack (independent of  $b$ ) and encode each strategy as a column  
126 of the sparse matrix  $A$ . For example, a strategy consisting of the sequence length 128, 128, and  
127 256 is represented a column vector that has the value 2 at the 128th row, the value 1 at the 256th  
128 row, and zero at all other rows. The variable  $x$  describes the *non-negative* repetition count for each  
129 strategy. So a 24 in the  $i$ th row of  $x$  means that the strategy represented by the  $i$ th column of  $A$  should  
130 repeat 24 times. Moreover, in the un-weighted setting,  $Ax = b$  states that we would like to “mix” the  
131 pre-defined strategies (columns of  $A$ ) such that the number of samples matches the histogram  $b$ , and  
132 where each strategy is used  $x \geq 0$  times. We use the residual weight  $w$  to control the penalization  
133 of the  $Ax - b$  residual on different sequence lengths (different rows of  $b$ ). Heuristically, we set  
134 the weight of 0.09 for all sequences of length 8 or smaller because they are considered acceptable  
135 padding sequences while all other sequence lengths get weight 1. We discuss this heuristic choice of  
136 parameters in Section F.4.5 and F.5 [1]. The overall efficiency of the packing is not greatly influenced  
137 by the weighing (less than 1% extra speed-up).

138 After solving  $wAx = wb$  for  $x \geq 0$  using an off-the-shelf solver, we obtain a floating point solution,  
139 which means that the repetition counts are not necessarily integers. Since we cannot use a non-natural  
140 number of strategies, we round the solution  $\hat{x}$  to the nearest integer. The error introduced by this  
141 rounding is found to be negligible (a few hundred sequences in the worst case) compared to the size  
142 of the dataset (millions of sequences). The time complexity and space complexity of the algorithm  
143 are  $O(n + s_m^5)$  and  $O(s_m^3)$ . Further details are provided in Section F.4 [1].

## 144 3.2 packedBERT: model changes

145 This section describes how any vanilla BERT implementation should be modified for packed sequence  
146 processing, such that the behavior of the model is the same as when processing unpacked sequences.  
147 Preserving the mathematical equivalence is necessary to ensure existing BERT pre-training and  
148 fine-tuning practices remain valid, as well as being required by benchmarks such as MLPerf™ [17].  
149 The presented approaches and principles apply to a variety of other models.

### 150 3.2.1 Adjust positional embeddings

151 The BERT model uses three types of embeddings: token, segment, and positional embeddings. The  
152 latter is canonically implemented as a bias add operation, rather than a full embedding look-up. This  
153 is possible because the positional indices increase linearly for every sequence. However, when using  
154 the packed data format the position index needs to be reset with each new packed sequence. For  
155 instance, when packing two sequences one of length 2 and one of length 3, the positional embedding  
156 indexes that need to be picked up are  $[0, 1, 0, 1, 2]$ . To achieve this, the bias add needs to be replaced  
157 by an embedding look-up to extract the correct positional embedding for each token in the pack. This  
158 also requires keeping an extra input which specifies the position of each token in its sequence. This  
159 required adjustment has only a minor impact on absolute accuracy/loss (see Section 4.2 and 4.2.1).

### 160 3.2.2 Adjust attention masking

```
# input
mask = np.array([[1, 1, 1, 2, 2]])
# 0, 1 mask
zero_one_mask = tf.equal(mask, mask.T)
# for use with softmax:
softmax_mask = tf.where(
    zero_one_mask, 0, -1000)
```

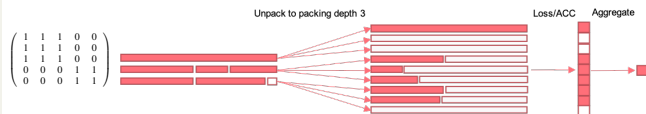


Figure 2: Attention mask code [left], respective zero-one mask [middle], and vectorized unpacking of the sequence loss [right]. White rectangles correspond to padding.

161 To maintain an implementation that is consistent with the un-packed version, tokens from different  
162 sequences within a pack should not be able to attend to each other. This is typically achieved in  
163 other implementations by unpacking the sequences using custom attention kernels and then doing  
164 the attention per-sequence [5]. Instead, we propose directly masking the attention matrix with a  
165 block-diagonal mask before the attention softmax. This is straightforward to implement in modern  
166 frameworks (see Figure 2). Naturally, there is a cost to both the mask construction and applying  
167 it to the attention matrix. However, it is required to keep the accuracy (see Table 1, Section 4.1  
168 Section 4.2). See also the code of the deprecated `tensor2tensor` library and our own provided code.

### 169 3.2.3 Adjust per-sequence loss and accuracy

170 Canonical implementations of BERT compute the cross-entropy loss for the masked language model  
171 on a per-token basis. However other NLP tasks, such as SQuAD, compute the loss and accuracy on  
172 a per-sequence basis. This section discusses how to handle such tasks when training with packed  
173 sequences. Simply feeding packs of sequences to the same implementation of cross-entropy would  
174 result in a per-pack weighted loss. In other words, the overall loss on the micro-batch would sum-up  
175 the losses on the individual packs, rather than individual sequences. As a result, the model would  
176 converge to a different optimum than when running with the un-packed implementation. For instance,  
177 a pack of a single sequence would contribute to the loss with the same weight as a pack of three  
178 sequences.

179 To recover the per-sequence averaging behavior of the canonical un-packed BERT implementation,  
180 we effectively “unpack” the incoming logits and labels. Once the sequences have been unpacked,  
181 we can compute the loss on each sequence separately as usual and then add up the losses. However,  
182 rather than looping through the sequences index, we compute on all indexes in parallel (see Figure 2).  
183 This minimizes the latency overhead of un-packing the loss calculation. As an example, we show how  
184 per-sequence loss can be implemented for the pre-training task. We use the “masked lm weight” [7]  
185 input tensor to represent which sequence a given masked token belongs to (0, 1, 2 and so on). This  
186 is consistent with the canonical BERT implementation where this input takes a value of either 1  
187 (belonging to the sequence) or 0 (belonging to padding). The full methodology is detailed in Listing 5  
188 and can be applied to other classification or pre-training tasks.

### 189 3.3 Adjust hyperparameters

190 In terms of convergence behavior, the primary consequence of packing is an increase in the effective  
191 batch size (with respect to number of sequences and real tokens) with some added variation over  
192 different iterations. If we look on the sentence level, the number of sentences in one batch increases  
193 by the packing factor. Similarly, the number of tokens in one batch increases. Hence, hyperparameters  
194 that are sensitive to these numbers need to be adjusted.

195 A direct solution is to reduce the computational batch size by the packing factor (average number of  
196 sequences per pack) and keep all other hyperparameters the same. For example, if the packing factor  
197 is 2, cutting the gradient accumulation count by half is sufficient. The advantage of this strategy is that  
198 no fine-tuning of hyperparameters is required and performance curves are comparable. However, this  
199 approach might be not desirable as it might imply under-utilizing the memory/compute, especially if  
200 the micro batch size needs to be reduced.

201 Hence to preserve batch size and optimize hardware utilization, we additionally propose an approxi-  
202 mate heuristic for updating the decay parameters of the LAMB optimizer [35]. For a packed dataset  
203 with a packing factor  $p$ , we update the decay parameters as:  $\beta_1 := \beta_1^p$ ,  $\beta_2 := \beta_2^p$ . For  $p = 2$ , this  
204 corresponds to the exact parameters for calculating momentum and velocity, when updating with the  
205 same gradient twice (Section D). A common approach is to scale the learning rate with the batch size.  
206 However, our experiments in Section 4.2 show that this reduces convergence speed.

207 Since these adjustments are only heuristics the convergence of the model will be comparable but not  
208 identical. In particular, it is unlikely that simply adjusting the hyperparameters will fully undo the  
209 impact of the increased batch size. However, with these adjustments, researchers should be able to  
210 continue to use existing configurations.



## 211 4 Experiments

### 212 4.1 Bin packing algorithm comparison

213 We evaluate our algorithms using the following metrics: **number of packs**, **number of all tokens**,  
214 **number of padding tokens**, **solution time of the packing algorithm** (after histogram and strategy  
215 creation), **number of strategies used**, **packing efficiency** (the fraction of non-padding tokens in the  
216 packed dataset), the **speed-up** achieved compared to not packing (depth 1), and the average number  
217 of sequences per sample (**packing factor**). For SPFHP, we analyse different (maximum) packing  
218 depth, since packing is less efficient with smaller depth and we want to get a general understanding  
219 on how the packing depth influences the processing time. For NNLSHP, we focus on packing  
220 depth 3 because it packs the data sufficiently well. For the speed-up analysis, we focus on the  
221 intelligence processing unit (IPU) [11] (IPU-M2000, 16 accelerator chips), BERT phase 2 pretraining  
222 setup as in Section 4.2. A GPU dynamically loads the code into the accelerator; in contrast, the  
223 IPU works with a static pre-compiled engine that gets loaded onto the chip at the start of the run.  
224 While other approaches result in excessive padding or continuous changes of the code, our approach  
225 can work with the same code for the whole dataset. So in this setting the IPU architecture would  
226 especially benefit from our approach since it avoids code changes. Nevertheless, it can be applied  
227 to any implementation on GPU or TPU. For determining the speed-up, we take advantage of the  
228 precompiled kernel. Since time measurements are quite noisy, we can profile the kernel and how  
229 many cycles it takes for processing a batch. That way, we can determine the **overhead** (in cycles)  
230 from processing the additional attention masking and for unpacking the loss. Combining **overhead**  
231 and **packing factor**, we get the **speed-up** estimate. No experiment repetitions are required since the  
232 algorithms and measurements are deterministic.

Table 1: Key performance results of proposed packing algorithms (SPFHP and NNLSHP) on IPU.

pack. depth	packing algorithm	EFF (%)	p	OH (%)	realized speed-up
1	NONE	50.0	1.00	0.000	1.000
1	SORT	99.9	2.00	$\gg 100$	$\ll 1.000$
$\approx 10$	GREEDY	$\approx 78$	$\approx 1.6$	$\approx 4.48$	$\approx 1.5$
2	SPFHP	80.5	1.61	4.283	1.544
3	SPFHP	89.4	1.79	4.287	1.716
3	NNLSHP	99.7	2.00	4.287	<b>1.913</b>
4	SPFHP	93.9	1.88	4.294	1.803
8	SPFHP	98.9	1.98	4.481	1.895
max	SPFHP	99.6	1.99	4.477	1.905

**Packing depth** describes the maximum number of packed sequences. NONE is the baseline BERT implementation, whereas SORT corresponds to sorted batching, and GREEDY concatenates sequences as they arrive until they would exceed 512 tokens. Setting no limit resulted in a maximum packing depth of 16. **EFF**iciency is the percentage of real tokens in the packed dataset. The **packing factor** describes the resulting potential speed-up compared to packing depth 1. With **overhead (OH)**, we denote the percentage decrease in throughput due to changes to the model to enable packing (such as the masking scheme introduced in Section 3.2.2). The **realized speed-up** is the combination of the speed-up due to packing (the **packing factor**) and the decrease in throughput due to the overhead on the IPU. It is used to measure the relative speed-up in throughput and the overhead from masking and loss adjustment. SORT can be only efficient on GPUs (see Section 4.4).

233 The main results for the performance metric evaluation are displayed in Table 1. The processing  
234 time for SPFHP on an Intel(R) Xeon(R) Gold 6138 CPU with 2.00GHz, 80 nodes, and 472G RAM  
235 was around 0.03s and independent from the packing depth. Classical First-Fit-Decreasing requires  
236 87-120s, a lot of memory, and scales almost linear with the number of samples. We see that the  
237 overhead slightly increases with packing depth but that the benefits of packing outweigh the cost. The  
238 best speed-up is obtained with NNLSHP at depth 3 which required 28.4s on the CPU for processing  
239 and ran out of memory for larger depth. With a value of 1.913, it is close to the theoretical upper  
240 bound of 2.001. The results show that efficiency, packing factor, and speed-up can be viewed inter-  
241 changeably. The amount of time needed to process a sample (a pack of sequences) is barely changed  
242 relative to the un-packed implementation. The packing factor, or the improvement in efficiency,

effectively provide an accurate estimate of the speed-up. GREEDY packing as used in T5 shows to be quite inefficient and sorted batching (SORT) is highly efficient in avoiding padding but the resulting different computational graphs cause a major overhead on the IPU that exceeds the benefits of avoiding the padding. Since we made our algorithm and code public available, results have been reproduced with a different framework on the Habana Gaudi accelerator [10] and confirmed that our approach is hardware and software independent giving it a huge advantage over existing approaches.

## 4.2 MLPerf™ phase 2 pretraining setup: learning curves and hyperparameter adjustment

For depth 1 (classic BERT) and NNLSHP with depth 3, we additionally evaluate on the MLPerf™ version 0.7 BERT pre-training benchmark [17]. Briefly, this involves training from a standard checkpoint to a masked-language model accuracy of 71.2% using 3 million sequences with a maximum length of 512 tokens (refer to [19] for details). Following this standardized benchmark supports reproduction of results even on other systems and makes sure that the reproduction effort is moderate and setup rules are clearly documented. We compare the resulting speed-up as well as the respective learning curves by evaluating the data on a held-out validation dataset. The objective of this additional evaluation is to analyse if convergence behavior is changed by the packing strategy and if the theoretical speed-up can be achieved in practice.

With packing, we effectively increase the average batch size by the packing factor ( $\approx 2$ ). However, with a different batch size, different hyperparameters are required (see Section 3.3) and there is no mapping that will generate exact matching of results but only heuristics. In a first comparison, we use the same hyperparameters when comparing packed and unpacked training except for cutting the accumulation count by half. This way, we make sure that the batch size is constant on **average** and we have the same amount of training steps. In the second comparison, we evaluate our heuristics and how they compensate the difference in batch size. This setup is more desirable because it is beneficial to use the hardware to its full potential and cutting the batch size by half usually reduces throughput. In the third comparison, we compare two optimized setups. In these two cases, packing takes half the amount of training steps.

The learning curves are displayed in Figure 3. In the first setup, we see the curves almost matching perfectly when normalizing by the numbers of samples processed. Differences can be explained by the variation of the number of sequences in the packing batch, and general noise in the training process. Especially after the initial phase, the curves show a near-identical match. The second setup shows bigger differences since changing the batch size and hyperparameters changes the training dynamics. We observe slower convergence early on in training due to the increased batch size. This is expected. The adjustment of the learning rate actually decreases performance probably because we correct for the increased number of sequences already in the modified loss. With the adjustment of the decay parameter of LAMB, we see matching performance at the later training stages. However, it is not feasible to completely recover the early convergence behavior of the smaller batch size by adjusting the hyperparameters. For instance doubling the batch size of unpacked BERT to 3000 and adjusting the LAMB decay parameters leads to more of a slow down in convergence than when running packed BERT with a batch size of 1500 and a packing factor of 2. In practice, our implementations exceed the estimated 1.913 maximum speed-up. This estimate is based on the reduction in the computational work needed to process the dataset. However, packing the data also reduces the latency of the transferring the data to the device. Figure 3 shows that the realized total speed-up from packing exceeds  $2x$ .

### 4.2.1 Ablation study

So far, we have shown that with the introduced adjustments, we can match the accuracy of unpacked BERT. In the following, we analyze in how far the masking adjustment is required. In Figure 4, we can see that without our adjustments, training loss and accuracy worsen drastically and a longer training time does not lead to a recovery. When not adjusting the positional embedding, the loss and accuracy almost match. However, the accuracy stalls at 71.8% and does not reach the target accuracy of 72.1%. So overall, both adjustments are crucial to avoid a reduction in performance.

When running packed BERT without the NSP loss but keeping everything else the same in a full training setup, we observed that downstream performance on SQuAD reduced the F1 measure by 1.31% and EM by 1.15%. Hence, we do not consider removing NSP as done in approaches like RoBERTa and T5 as discussed in Section II.

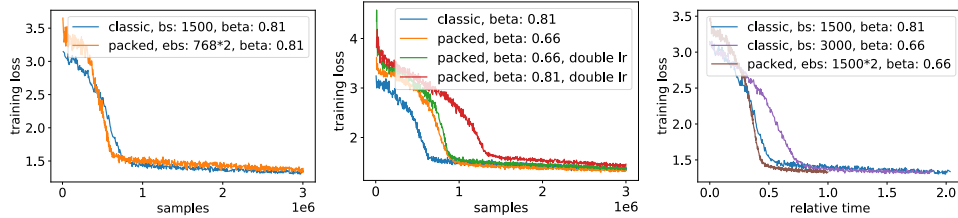


Figure 3: Comparison of learning curves for packed and unpacked processing, where all experiments converged to the target accuracy within the same number of training samples(3 million). [left] same effective batch size (**ebs** is batch size times packing factor), [middle] different heuristic adjustments of the hyperparameters (batch size 1500 for all runs, such that **ebs** for packed runs is  $1500 * 2$ ), and [right] realized speed-up from packing (in excess of desired 2x). Further learning curves are provided in Section [Q](#).

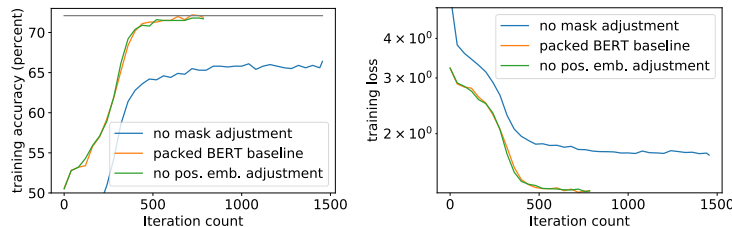


Figure 4: Comparison of learning curves with and without mask or positional embedding adjustment in our packed BERT approach. The grey accuracy baseline to reach is 72.1%.

### 297 4.3 Full pretraining and SQuAD finetuning

298 Packing slightly violates the i.i.d. assumption of data. Thus, we have to check that downstream  
 299 performance is not impacted by packing. This is especially relevant in a full training setup without  
 300 a starting checkpoint. To this aim, we show that the packed and unpacked SQuAD 1.1 scores are  
 301 comparable after a full-pretraining of BERT base and large plus fine-tuning. During pre-training,  
 302 in order to avoid giving an advantage to packing by further hyperparameter tuning, we reduce the  
 303 gradient accumulation count for the packed BERT training for phase 1 and phase 2 to match, on  
 304 average, the total number of sequences that get processed before each weight update. With this  
 305 approach, we can use the same hyperparameters and number of training steps but process each batch  
 306 faster by avoiding the processing of padding. This gives a slight disadvantage to the packed run in  
 307 terms of machine utilization, as explained in Section [3.3](#) and is different to the speedup analysis in  
 308 Section [4.2](#). For Phase 2, we use sequence length 384 since longer range attention is not relevant for  
 309 SQuAD 1.1. The respective speed-ups from packing for BERT base and large are shown in  
 310 Table [2](#); the realized speed-up, measured as the quotient of the throughputs between the packed  
 311 and unpacked runs, is slightly lower to the theoretical throughput (i.e. the packing factor) due to  
 312 the packing overhead. Further learning curves with the loss function and accuracy are provided in  
 313 Section [P](#). For the fine-tuning training on SQuAD 1.1, we do not use packing. The scores, computed  
 314 as the median of 10 different seeds, are displayed in Table [3](#). They are comparable to the reference  
 315 ones in [\[6\]](#): for BERT base (resp. large) the F1 score is reduced by 0.2% (resp. 0.3%) and the EM  
 316 score increases by 0.3% (resp. 0.02%).

Table 2: Measured speed-ups in BERT pretraining with packing.

Model size	Sequence length	Packing factor	Realized speed-up
base	128	1.17	1.15
	384	1.70	1.68
large	128	1.17	1.15
	384	1.70	1.69

Table 3: SQuAD 1.1 scores after BERT pretraining with packing.

Model size	Configuration	F1	Exact match
base	<a href="#">[6]</a>	88.5	80.8
	Packed	88.32	81.03
large	<a href="#">[6]</a>	90.9	84.1
	Packed	90.65	84.12



317 **4.4 Scaling analysis: Impact of accelerators count**

318 A further advantage of packing over competing un-padding approaches is the inherent load balancing  
 319 provided by packing. So called un-padding approaches rely on dynamically launching custom kernels  
 320 that ignore padding. A stated advantage of such implementations is the ability to avoid computing  
 321 the complete (512 x 512) attention matrix. This provides additional computational savings compared  
 322 to packing, where the attention matrix is computed in its entirety and then masked. Because of  
 323 these additional savings, un-padding can exceed the theoretical upper bound for speed-up from  
 324 packing (2.013 on Wikipedia). As a result of the dynamic nature of the approach, the processing  
 325 time with un-padding is different for each sequence in the batch, and the amount of time required to  
 326 process a batch of sequences will be determined by the processing time of the longest sequence in  
 327 the batch (with the sequences being processed in parallel). Furthermore, in the multiple accelerator  
 328 setting the processing time on each device will vary depending on the sequences in the batch that it  
 329 receives. Devices which finish early have to wait for the slowest device to finish before exchanging  
 330 gradients. This load-imbalance between the devices (and inside the batch) leads to a considerable  
 331 decrease in the speed-up from un-padding as the number of accelerators is increased (see Figure 5  
 332 and Section E (Q)). In contrast, packing (our approach) is inherently load-balanced. The processing  
 333 time on each accelerator is independent of the content inside the batch received by the device. Any  
 334 number of accelerators can therefore operate in unison without having to wait for the slowest batch to  
 335 process (all per-device batches are equally fast).

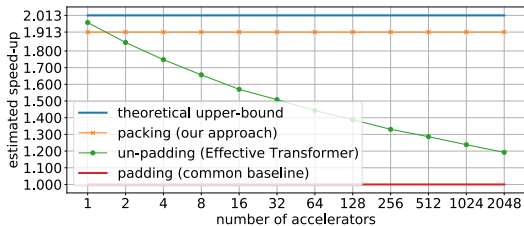


Figure 5: Comparison of the theoretical speed-up as the number of accelerators is increased.

336 **5 Conclusion**

337 Whereas packing is a well known concept, this paper sheds a new light onto it in multiple aspects.  
 338 First, we visualize the sequence length distributions of multiple datasets not just from language  
 339 domains but also audio and molecular domains to emphasize that packing is beneficial for a lot of  
 340 datasets and that in many cases, more than 2x acceleration can be achieved by removing 50% or  
 341 more padding. Second, we provide two new highly efficient packing approaches based on established  
 342 solvers that leave almost no padding and that can tackle arbitrarily large datasets in a matter of  
 343 seconds, in contrast to existing approaches that are slow and suboptimal. Third, we demonstrate that  
 344 without adjusting the sequence processing algorithm (e.g., BERT) to the packed sequences, predictive  
 345 performance is reduced. Thus, we propose several model adjustments that are all necessary to keep  
 346 predictive performance. Last but not least, we prove that, thanks to such adjustments, predictive  
 347 performance is preserved as if no packing was used — but speed significantly increases, especially  
 348 since the adjustments come with an overhead of less than 5%. We prove in our experiments that  
 349 downstream performance is not impacted by packing and that the anticipated 2x acceleration can be  
 350 achieved.

351 In the future, an interesting direction is the packing of images of different sizes to help accelerate  
 352 computer-vision applications. This is especially relevant given the recent advances in the use of  
 353 transformer-based approaches in the computer vision domain, for example the visual transformer [33].  
 354 Note that many images come in different shapes and resolutions and packing them can be a new  
 355 approach to tackle this diversity instead of casting them all to the same resolution and shape. Masking  
 356 out the self-attention within transformers is easier to implement than avoiding cross-contamination of  
 357 convolutions applied to packed images. Future work should explore improving the performance of  
 358 other models (RoBERTa, GPT-3, T5) by avoiding contamination between non-contiguous segments  
 359 from different documents. Even BERT itself might benefit from avoiding contamination between the  
 360 two concatenated segments.

## References

- 361
- 362 [1] ANONYMOUS. Supplemental Material for “Efficient Sequence Packing without Cross-contamination:  
363 Accelerating Large Language Models without Impacting Performance”, 2022.
- 364 [2] BOTTOU, L., CURTIS, F. E., AND NOCEDAL, J. Optimization Methods for Large-Scale Machine Learning.  
365 *SIAM Review* 60, 2 (jan 2018), 223–311.
- 366 [3] BRO, R., AND DE JONG, S. A fast non-negativity-constrained least squares algorithm. *Journal of*  
367 *Chemometrics* 11, 5 (sep 1997), 393–401.
- 368 [4] BROWN, T. B., MANN, B., RYDER, N., SUBBIAH, M., KAPLAN, J., DHARIWAL, P., NEELAKANTAN,  
369 A., SHYAM, P., SASTRY, G., ASKELL, A., AGARWAL, S., HERBERT-VOSS, A., KRUEGER, G.,  
370 HENIGHAN, T., CHILD, R., RAMESH, A., ZIEGLER, D. M., WU, J., WINTER, C., HESSE, C., CHEN,  
371 M., SIGLER, E., LITWIN, M., GRAY, S., CHES, B., CLARK, J., BERNER, C., MCCANDLISH, S.,  
372 RADFORD, A., SUTSKEVER, I., AND AMODEI, D. Language Models are Few-Shot Learners. In *Advances*  
373 *in Neural Information Processing Systems 33 pre-proceedings (NeurIPS 2020)* (may 2020).
- 374 [5] BYTEDANCE INC. Effective Transformer. [https://github.com/bytedance/effective\\_](https://github.com/bytedance/effective_)  
375 [transformer](https://github.com/bytedance/effective_transformer), 2021.
- 376 [6] DEVLIN, J., CHANG, M. W., LEE, K., AND TOUTANOVA, K. BERT: Pre-training of deep bidirectional  
377 transformers for language understanding. *NAACL HLT 2019 - 2019 Conference of the North American*  
378 *Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings*  
379 *of the Conference 1* (oct 2019), 4171–4186.
- 380 [7] DEVLIN, J., CHANG, M. W., LEE, K., AND TOUTANOVA, K. BERT: Pre-training of Deep Bidirectional  
381 Transformers for Language Understanding. <https://github.com/google-research/bert>, 2019.
- 382 [8] DEVLIN, J., CHANG, M. W., LEE, K., AND TOUTANOVA, K. Pre-training data creation script  
383 for BERT. [https://github.com/google-research/bert/blob/master/create\\_pretraining\\_](https://github.com/google-research/bert/blob/master/create_pretraining_data.py#L243)  
384 [data.py#L243](https://github.com/google-research/bert/blob/master/create_pretraining_data.py#L243), 2019.
- 385 [9] FEDUS, W., ZOPH, B., AND SHAZEER, N. Switch Transformers: Scaling to Trillion Parameter Models  
386 with Simple and Efficient Sparsity. *arXiv* (jan 2021).
- 387 [10] INTEL, 2021.
- 388 [11] JIA, Z., TILLMAN, B., MAGGIONI, M., AND SCARPAZZA, D. P. Dissecting the Graphcore IPU  
389 architecture via microbenchmarking. *ArXiv abs/1912.03413* (2019).
- 390 [12] JOHNSON, D. S. *Near-optimal bin packing algorithms*. PhD thesis, Massachusetts Institute of Technology,  
391 1973.
- 392 [13] JOHNSON, D. S., AND GAREY, M. R. A 7160 theorem for bin packing. *Journal of Complexity* 1, 1 (oct  
393 1985), 65–106.
- 394 [14] KORTE, B., AND VYGEN, J. *Combinatorial Optimization*, vol. 21 of *Algorithms and Combinatorics*.  
395 Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.
- 396 [15] LEE, C. C., AND LEE, D. T. A Simple On-Line Bin-Packing Algorithm. *Journal of the ACM (JACM)* 32,  
397 3 (jul 1985), 562–572.
- 398 [16] LIU, Y., OTT, M., GOYAL, N., DU, J., JOSHI, M., CHEN, D., LEVY, O., LEWIS, M., ZETTMLOYER,  
399 L., AND STOYANOV, V. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv* (jul 2019).
- 400 [17] MATTSON, P., REDDI, V. J., CHENG, C., COLEMAN, C., DIAMOS, G., KANTER, D., MICIKEVICIUS,  
401 P., PATTERSON, D., SCHMUELLING, G., TANG, H., WEI, G., AND WU, C. MLPerf: An Industry  
402 Standard Benchmark Suite for Machine Learning Performance. *IEEE Micro* 40, 2 (2020), 8–16.
- 403 [18] MENG, Q., CHEN, W., WANG, Y., MA, Z. M., AND LIU, T. Y. Convergence analysis of distributed  
404 stochastic gradient descent with shuffling. *Neurocomputing* 337 (apr 2019), 46–57.
- 405 [19] MLCOMMONS. v0.7 Results. <https://mlcommons.org/en/training-normal-07/>, 2020. Result  
406 not verified by MLPerf. Throughput/speedup is not the primary metric of MLPerf. MLPerf name and logo  
407 are trademarks. See [www.mlperf.org](http://www.mlperf.org) for more information.

- 408 [20] NVIDIA. Reference numbers for BERT un-padding results. [https://github.com/mlcommons/  
409 training\\_results\\_v0.7/blob/master/NVIDIA/results/dgxa100\\_ngc20.06\\_pytorch/bert/  
410 result\\_0.txt](https://github.com/mlcommons/training_results_v0.7/blob/master/NVIDIA/results/dgxa100_ngc20.06_pytorch/bert/result_0.txt), 2020. Throughput/speedup is not the primary metric of MLPerf. MLPerf name and logo  
411 are trademarks. See [www.mlperf.org](http://www.mlperf.org) for more information.
- 412 [21] NVIDIA. Faster Transformer. [https://github.com/NVIDIA/DeepLearningExamples/tree/  
413 master/FasterTransformer/v1](https://github.com/NVIDIA/DeepLearningExamples/tree/master/FasterTransformer/v1), 2021.
- 414 [22] OTT, M., EDUNOV, S., BAEVSKI, A., FAN, A., GROSS, S., NG, N., GRANGIER, D., AND AULI,  
415 M. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019:  
416 Demonstrations* (2019).
- 417 [23] PANAYOTOV, V., CHEN, G., POVEY, D., AND KHUDANPUR, S. Librispeech: an asr corpus based on  
418 public domain audio books. In *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International  
419 Conference on* (2015), IEEE, pp. 5206–5210.
- 420 [24] RAFFEL, C., SHAZEER, N., ROBERTS, A., LEE, K., NARANG, S., MATENA, M., ZHOU, Y., LI, W.,  
421 AND LIU, P. J. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal  
422 of Machine Learning Research* 21 (oct 2019).
- 423 [25] RAJPURKAR, P., ZHANG, J., LOPYREV, K., AND LIANG, P. SQuAD: 100,000+ questions for machine  
424 comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language  
425 Processing* (Austin, Texas, Nov. 2016), Association for Computational Linguistics, pp. 2383–2392.
- 426 [26] RAMAKRISHNAN, R., DRAL, P. O., RUPP, M., AND VON LILIENFELD, O. A. Quantum chemistry  
427 structures and properties of 134 kilo molecules. *Scientific Data* 1 (2014).
- 428 [27] RUDDIGKEIT, L., VAN DEURSEN, R., BLUM, L. C., AND REYMOND, J.-L. Enumeration of 166 billion  
429 organic small molecules in the chemical universe database gdb-17. *Journal of Chemical Information and  
430 Modeling* 52, 11 (2012), 2864–2875. PMID: 23088335.
- 431 [28] SHEN, J., NGUYEN, P., WU, Y., CHEN, Z., ET AL. Lingvo: a modular and scalable framework for  
432 sequence-to-sequence modeling, 2019.
- 433 [29] VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER, U.,  
434 AND POLOSUKHIN, I. Attention is all you need. In *Proceedings of the 31st International Conference on  
435 Neural Information Processing Systems* (Red Hook, NY, USA, 2017), NIPS’17, Curran Associates Inc.,  
436 p. 6000–6010.
- 437 [30] WANG, A., SINGH, A., MICHAEL, J., HILL, F., LEVY, O., AND BOWMAN, S. GLUE: A multi-task  
438 benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP  
439 Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP* (Brussels, Belgium, Nov.  
440 2018), Association for Computational Linguistics, pp. 353–355.
- 441 [31] WARSTADT, A., SINGH, A., AND BOWMAN, S. R. Neural network acceptability judgments. *arXiv  
442 preprint arXiv:1805.12471* (2018).
- 443 [32] WOLF, T., DEBUT, L., SANH, V., CHAUMOND, J., DELANGUE, C., MOI, A., CISTAC, P., RAULT, T.,  
444 LOUF, R., FUNTOWICZ, M., DAVISON, J., SHLEIFER, S., VON PLATEN, P., MA, C., JERNITE, Y., PLU,  
445 J., XU, C., SCAO, T. L., GUGGER, S., DRAME, M., LHOEST, Q., AND RUSH, A. M. Transformers: State-  
446 of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in  
447 Natural Language Processing: System Demonstrations* (Online, Oct. 2020), Association for Computational  
448 Linguistics, pp. 38–45.
- 449 [33] WU, B., XU, C., DAI, X., WAN, A., ZHANG, P., YAN, Z., TOMIZUKA, M., GONZALEZ, J., KEUTZER,  
450 K., AND VAJDA, P. Visual transformers: Token-based image representation and processing for computer  
451 vision, 2020.
- 452 [34] XLA, T. XLA: Optimizing Compiler for Machine Learning. <https://www.tensorflow.org/xla>,  
453 2021.
- 454 [35] YOU, Y., LI, J., REDDI, S., HSEU, J., KUMAR, S., BHOJANAPALLI, S., SONG, X., DEMMEL, J.,  
455 KEUTZER, K., AND HSIEH, C.-J. Large Batch Optimization for Deep Learning: Training BERT in 76  
456 minutes. *arXiv* (apr 2019).
- 457 [36] YUE, M., AND ZHANG, L. A simple proof of the inequality  $MFFD(L) \leq 71/60OPT(L) + 1$ ,  $L$  for  
458 the MFFD bin-packing algorithm. *Acta Mathematicae Applicatae Sinica* 11, 3 (jul 1995), 318–330.

459 **Checklist**

- 460 1. For all authors...
- 461 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
462 contributions and scope? [Yes] Our paper has four main claims. First, in Figure I  
463 we show the sequence length distribution of Wikipedia and many other datasets and  
464 the excessive padding that they require. Second, in Section 4.1 we show that we  
465 can efficiently pack the data which can be easily reproduced with the shared data and  
466 code II. Third, in Figure 3[right], we clearly show the 2x performance gain from  
467 packing and the related hyperparameter adjustment scheme. Fourth, multiple additional  
468 experiments on downstream tasks, ablation studies, and packing variants further verify  
469 the validity of our proposed approaches.
- 470 (b) Did you describe the limitations of your work? [Yes] We see three potential limitations  
471 that we discuss in the paper. First, as stated in Section A “Broader Impact” in the  
472 appendix II, our approach is clearly dependent on the sequence length distribution  
473 of the dataset. However, we looked into several other datasets beyond Wikipedia and  
474 observed even higher potential for acceleration and document this in multiple sections  
475 throughout the paper as well as in the appendix II. Second, we explain our focus  
476 on the IPU hardware in Section 4.1. Our theoretical analysis in Section 4.4 indicates  
477 that our approach benefits also GPUs. We also cite other work, that shows that our  
478 approach is hardware independent. Third, our changes to the network with a modified  
479 attention mask and loss calculation come with some overhead. This is addressed in  
480 Table I [overhead column] in Section 4.1.
- 481 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We address  
482 this point in Section A “Broader Impact”, third paragraph, in the appendix II.
- 483 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
484 them? [Yes]
- 485 2. If you are including theoretical results...
- 486 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Detailed  
487 algorithm explanations, clarifications of assumptions, and proofs are provided in the  
488 supplemental material II.
- 489 (b) Did you include complete proofs of all theoretical results? [Yes] Sections D, E, and  
490 G in the supplemental material II provide the necessary derivations on theoretical  
491 results.
- 492 3. If you ran experiments...
- 493 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
494 mental results (either in the supplemental material or as a URL)? [Yes] All packing  
495 code is provided in the paper. The packing results on BERT got verified by multiple  
496 independent parties. One party used a draft of this paper to successfully reproduce its  
497 main findings. Links to implementations in three different frameworks will be provided  
498 after acceptance, to avoid violating the blind submission rules.
- 499 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
500 were chosen)? [Yes] In the first part, we follow the MLPerf 0.7 benchmark rules.  
501 We document the parameters that we changed and why we change them. For the  
502 downstream tasks, we follow the reference and report where, how and why we change  
503 hyperparameters.
- 504 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
505 ments multiple times)? [No] The packing algorithms are deterministic and have no  
506 error. Other experiments are only executed once to compare convergence curves. For  
507 downstream tasks, we report repetition details and the median as in the reference.
- 508 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
509 of GPUs, internal cluster, or cloud provider)? [Yes] We used 16 Graphcore Mk2 IPUs  
510 for acceleration on an internal cluster.
- 511 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 512 (a) If your work uses existing assets, did you cite the creators? [Yes] Appropriate references  
513 to the BERT authors, all datasets, and the code snippet from the HuggingFace inc. are  
514 appropriately referenced with citations and links.

- 515 (b) Did you mention the license of the assets? [Yes] For the only taken code snippet, the  
516 license is part of the file [Listing 7 in 11]. Dataset licenses like Wikipedia’s “Creative  
517 Commons Attribution-ShareAlike 3.0 License” are covered by the references. New  
518 materials like packing code and histograms will be provided under an MIT license  
519 which will be added over a link to the resources in the final paper version.
- 520 (c) Did you include any new assets either in the supplemental material or as a URL?  
521 [Yes] New materials like packing code and histograms are included in the supplement  
522 document as well as separate file. To avoid violating the blind submission rules, they  
523 will be linked in the final version like many other assets which are already publicly  
524 available under MIT license.
- 525 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
526 using/curating? [N/A] We did not curate other people’s data. We only provide a very  
527 high level aggregate of the used data.
- 528 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
529 information or offensive content? [N/A] We did not curate other people’s data.
- 530 5. If you used crowdsourcing or conducted research with human subjects...
- 531 (a) Did you include the full text of instructions given to participants and screenshots, if  
532 applicable? [N/A] Our experiments did not include crowdsourcing or human subjects.
- 533 (b) Did you describe any potential participant risks, with links to Institutional Review Board  
534 (IRB) approvals, if applicable? [N/A] Our experiments did not include crowdsourcing  
535 or human subjects.
- 536 (c) Did you include the estimated hourly wage paid to participants and the total amount  
537 spent on participant compensation? [N/A] Our experiments did not include crowd-  
538 sourcing or human subjects.