Few-shot Task-agnostic Neural Architecture Search for Distilling Large Language Models

Anonymous Author(s) Affiliation Address email

Abstract

Traditional knowledge distillation (KD) methods manually design student archi-1 tectures to compress large models given pre-specified computational cost. This 2 3 requires several trials to find viable students, and repeating the process with change in computational budget. We use Neural Architecture Search (NAS) to automat-4 ically distill several compressed students with variable cost from a large model. 5 Existing NAS methods train a single SuperLM consisting of millions of sub-6 networks with weight-sharing, resulting in interference between subnetworks of 7 different sizes. Additionally, many of these works are task-specific requiring task 8 labels for SuperLM training. Our framework AutoDistil addresses above chal-9 lenges with the following steps: (a) Incorporates inductive bias and heuristics to 10 partition Transformer search space into K compact sub-spaces (e.g., K=3 can 11 generate typical student sizes of base, small and tiny); (b) Trains one SuperLM for 12 each sub-space using task-agnostic objective (e.g., self-attention distillation) with 13 weight-sharing of students; (c) Lightweight search for the optimal student without 14 re-training. Task-agnostic training and search allow students to be reused for fine-15 tuning on any downstream task. Experiments on GLUE benchmark demonstrate 16 AutoDistil to outperform state-of-the-art KD and NAS methods with upto 3x 17 additional reduction in computational cost and negligible loss in task performance. 18

19 **1** Introduction

While large pre-trained language models (e.g., BERT [1], GPT-3 [2]) are effective, their huge size 20 poses significant challenges for downstream applications in terms of energy consumption and cost 21 of inference [3] limiting their usage in on the edge scenarios and under constrained computational 22 inference budgets. Knowledge distillation [4, 5, 6, 7] has shown strong results in compressing pre-23 trained language models into small student models. However, these works require pre-specification of 24 the student architecture and computational cost (e.g., number of parameters, FLOPs) for distillation. 25 This poses two significant challenges: (i) it requires several trials to come up with viable architectures 26 as they are hand-engineered and to define several hyper-parameters (e.g., number of layers and 27 attention heads, hidden dimension, etc.); (ii) one has to re-run distillation with any change in 28 specification for the student architecture or computational cost for using it in a target environment. 29 Neural Architecture Search (NAS) [8, 9, 10, 11] provides a natural solution to automatically search 30

through a large space of candidate models. The dominant NAS paradigm consists of two main steps:
 (a) Training a Super model combining all possible architectures into a single graph and jointly training
 them via weight-sharing; (b) Searching for optimal architecture from Super model with best accuracy
 on a downstream task, satisfying user-specified latency constraint for target device. Parallel to above

³⁵ computer vision (CV) works, NAS has shown strong results in recent works like DynaBERT [12],

36 AutoTinyBERT [13] and NAS-BERT [14] for natural language understanding (NLU).

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

37 Drawbacks of existing NAS methods.

38 [D1: Co-adaptation in weight-sharing] Above

³⁹ works train one single large Super Language

40 Model (SuperLM) consisting of millions of di-

verse student architectures. This results in someundesirable effects of co-adaptation [15] like

conflicts in weight-sharing where bigger student

44 models converge faster in contrast to smaller

⁴⁵ ones converging slower [16, 11].

46 [D2: Multi-stage training] A single SuperLM

47 may not have sufficient capacity to encode a
48 large search space. As a result, these works
49 use multi-stage training process, where they
50 first conduct NAS to identify candidate students
51 and then perform further pre-training [13] and

⁵² knowledge distillation [14] of the candidates.

[D3: Task-specific training] NAS works in 53 the CV domain (e.g., AutoFormer [17], Once-54 for-all [10], One-Shot NAS [11, 18]) leverage 55 hard class labels from a given task (e.g., im-56 age classification) or soft labels from ImageNet 57 pre-trained models (e.g., MobileNet [7], Reg-58 Net [19]) for task-specific optimization with ac-59 curacy as an evaluation metric. Different from 60

61 CV domain, NLU tasks have different objec-



Figure 1: AutoDistil uses few-shot taskagnostic NAS to distill several compressed students with variable #FLOPs (x-axis) from K=3SuperLMs (corresponding to each point cloud) trained on K sub-spaces of Transformer search space. Each student (blue dot) extracted from the SuperLM is fine-tuned on MNLI with accuracy on y-axis. The best student from each SuperLM is marked in red. Given any state-of-the-art distilled model, AutoDistil generates a better candidate with less #FLOPs and improved task performance from corresponding search space.

tives and evaluation metrics for classification (e.g., MNLI), regression (e.g., STS-B) and correlation 62 (e.g., CoLA). Correspondingly, pre-trained language models like BERT [1] are also trained in self-63 supervised fashion without using task labels. This makes it challenging to adapt existing NAS works 64 to the NLU domain in a task-agnostic setting. Recent NAS works in the NLU domain are not fully 65 task-agnostic. For instance, DynaBERT [12] accesses both task labels for knowledge distillation 66 and task development set for network rewiring. NAS-BERT [14] performs two-stage knowledge 67 distillation with pre-training and fine-tuning of the candidates for best performance. While AutoTiny-68 BERT [13] also explores task-agnostic training, we demonstrate better performance from few-shot 69 NAS and much cheaper cost from single stage training without additional pre-training and distillation. 70

Contributions. We address above challenges with fully task-agnostic few-shot NAS consisting of 71 three steps. (S1) Search space design. We partition the Transformer search space into K sub-spaces 72 considering important architectural hyper-parameters like the network depth, width and attention 73 heads. We further leverage inductive bias and heuristics to limit the number of student architectures 74 in each sub-space. (S2) Fully task-agnostic SuperLM training. We train K SuperLM overall, one 75 for every sub-space. This allows each SuperLM more capacity to encode a sub-space as opposed to a 76 single large one. We train each SuperLM with a fully task-agnostic objective (without accessing any 77 task labels) like deep self-attention distillation, where we transfer knowledge from the self-attention 78 79 module (including keys, queries and values) of a pre-trained teacher (e.g., BERT) to the student and 80 use weight-sharing to train the SuperLM. (S3) Lightweight optimal student search. We obtain 81 optimal student(s) directly from well-trained SuperLM(s) without any re-training that can be *simply* 82 fine-tuned on downstream tasks. Our contributions over existing NAS works can be summarized as: • In contrast to prior works (e.g., DynaBERT, AutoTinyBERT, NAS-BERT), we do a single-stage 83 training combining NAS and distillation with no further pre-training or augmentation and demonstrate 84 superior performance of the NAS process itself with significantly reduced training cost. Obtained 85 subnetworks are simply fine-tuned on downstream tasks. 86 • Fully task-agnostic training with subnetwork attention state alignment for self-attention relation 87

distillation and search in contrast to prior works in NLU (e.g., DynaBERT, NAS-BERT) and CV (e.g.,
 AutoFormer, BigNAS, Once-For-All).

• Few-shot NAS to mitigate gradient conflicts in SuperNet training compared to prior One-shot NAS

⁹¹ works in NLU (e.g., DynaBERT, AutoTinyBERT, NAS-BERT). AutoFormer in the CV domain is an

92 exception to this point which also uses few-shot NAS but accesses task labels during training.

• Strong results over all the above NAS and distillation works in NLU with 3x additional compression

over best performing distillation technique with negligible drop in task performance.



Figure 2: Overview of AutoDistil. It considers K partitions of the Transformer architecture subspace to train one SuperLM for each partition with weight-sharing of the constituent subnetworks trained via task-agnostic deep self-attention distillation. Optimal compressed subnetworks can be easily extracted from the SuperLMs without additional training or distillation.

2 Background 95

- We present an overview of Transformers [20], especially its two main sub-layers, multi-head self-96
- attention (MHA) and feed-forward network (FFN). Transformer layers are stacked to encode contex-97
- tual information for input tokens as: $\mathbf{X}^{l} = \text{Transformer}_{l}(\mathbf{X}^{l-1}), l \in [1, L]$ where L is the number 98
- of Transformer layers, $\mathbf{X}^l \in \mathbb{R}^{s*d_{hid}}$, s is the sentence length, and d_{hid} is the hidden dimension. In 99 the following, we omit the layer indices for simplicity. 100
- Multi-Head Self-Attention (MHA). Given previous Transformer layer's output X, MHA computes: 101

Attention
$$(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h) = \operatorname{softmax}(\frac{\mathbf{Q}_h \mathbf{K}_h^{\top}}{\sqrt{d_{head}}}) \mathbf{V}_h; \ \mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h = \mathbf{X} \mathbf{W}_h^Q, \mathbf{X} \mathbf{W}_h^K, \mathbf{X} \mathbf{W}_h^V, \ (1)$$

$$MHA(\mathbf{X}) = Concat(head_1, \cdots, head_H) \mathbf{W}^O,$$
(2)

where W_h^Q , W_h^K , $W_h^V \in \mathbb{R}^{d_{hid}*d_{head}}$, $W^O \in \mathbb{R}^{d_{hid}*d_{hid}}$ are linear transformations. \mathbf{Q}_h , \mathbf{K}_h , $\mathbf{V}_h \in \mathbb{R}^{s*d_{head}}$ are called queries, keys, and values, respectively. H is the number of heads. head_h = 102 103 Attention $(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h)$ denotes the *h*-th attention head. Concat is the concatenating operation. 104 $d_{head} = d_{hid}/H$ is the dimension of each head. 105 Feed-Forward Network (FFN). Each Transformer layer contains an FNN sub-layer, which is stacked 106

on the MHA. FFN consists of two linear transformations with a ReLU activation as: 107

$$FFN(x) = \max(0, xW^{1} + b_{1})W^{2} + b_{2},$$
(3)

where $W^1 \in \mathbb{R}^{d_{hid}*d_f}$, $W^2 \in \mathbb{R}^{d_f*d_{hid}}$, $b_1 \in \mathbb{R}^{d_f}$, and $b_2 \in \mathbb{R}^{d_{hid}}$. In addition, there are residual 108 connection and layer normalization on top of MHA and FFN (denoted by \oplus in Figure 2), which are 109 formulated as LayerNorm(x + MHA(x)) and LayerNorm(x + FFN(x)), respectively. 110

Few-shot Task-agnostic NAS 3 111

Given a large pre-trained language model (e.g., BERT) as teacher, AutoDistil distills several 112 compressed models with variable computational cost with the following major components. 113

3.1 Search Space Design 114

Searchable transformer components. From Transformers overview (Section 2) and our framework 115 (Figure 2), we observe four important hyper-parameters for the Transformer blocks to include: 116

- (1) Feed-forward network (FFN) dimension we encode this by the MLP (multi-layer perceptron) 117
- ratio defined as $r = \frac{d_f}{d_{hid}}$, with d_f and d_{hid} representing the intermediate dimension of the FFN and hidden dimension respectively; (2) Number of layers (L) to capture the network depth; (3) Hidden 118
- 119
- dimension (d_{hid}) to encode input; (4) Attention heads (H) for multi-head self-attention. 120

All of the above factors are important for model capacity and have a significant impact on the model 121 size and computational cost. For instance, different layers have different feature representation 122 capabilities. Recent works show that Transformer models are overparameterized [21, 22], such as 123 the feed-forward layer (FFN), which is one of the most computation intensive components [23]. 124 Therefore, we search for optimal MLP ratio and hidden dimension that reduce computational cost 125 resulting from FFN layers. Furthermore, studies [24, 25] show that attention heads can be redundant 126 when they learn to encode similar relationships for each word. Thus, we make the number of attention 127 heads searchable as well. 128

Inductive bias. Prior work [26] demonstrate that thinner and deeper neural networks with improved representation capacity perform better than wider and shallower ones. We incorporate this as an inductive bias to decide the number of layers to consider for the students in each of our K sub-spaces (base, small, tiny), where we prefer deeper students in terms of the number of layers. Furthermore, we constrain all the Transformer layers in a given student model to share identical and homogeneous structures, i.e., the same number of attention heads, hidden dimension, etc. This not only reduces the size of the search space, it is also more friendly to hardware and software frameworks [13].

Search space partition. Existing works [13, 14] train a single large SuperLM containing millions of 136 student architectures by weight-sharing. This leads to performance degradation due to optimization 137 interference and convergence of subnetworks with very different sizes [11]. To mitigate such 138 interference, we employ a few-shot learning strategy [17, 16] as follows. We partition the whole 139 Transformer search space into K sub-spaces such that each sub-space covers different sizes of student 140 models given by the number of parameters. For instance, K = 3 can cover typical student sizes, 141 namely base, small and tiny versions. Table 1 shows the parameter ranges for the K sub-spaces, 142 along with the student configurations contained in each. 143

We now encode each sub-144 space into a SuperLM, 145 where each student model 146 in the space is a subnetwork 147 of the SuperLM. Further-148 more, all the student subnet-149 works share the weights of 150 their common dimensions, 151 with the SuperLM being the 152 largest one in the search 153 space. Considering K in-154 dependent SuperLMs, each 155 one now has more capacity 156 to encode a sub-space, in 157 contrast to a limited capac-158

	$SuperLM_{\rm Tiny}$	$SuperLM_{\rm Small}$	$SuperLM_{\rm Base}$	BERT
#Subnets	256	256	256	N/A
#Layers	(4, 7, 1)	(9, 12, 1)	(9, 12, 1)	12
#Hid_dim	(128, 224, 32)	(256, 352, 32)	(544, 640, 32)	768
MLP Ratio	(2.0, 3.5, 0.5)	(2.5, 4.0, 0.5)	(2.5, 4.0, 0.5)	4.0
#Heads	(7, 10, 1)	(7, 10, 1)	(9, 12, 1)	12
#FLOPs	40-367 <i>M</i>	0.5-2.1 <i>G</i>	2.1-7.9 <i>G</i>	11.2 <i>G</i>
#Params	4-10 <i>M</i>	12-28 <i>M</i>	39-79 <i>M</i>	109 <i>M</i>

Table 1: The search space of AutoDistil with K=3 partitions, each consisting of 256 subnets with variable computational cost. We train one SuperLM with weight-sharing for *each partition* with child models sharing transformer blocks. Each tuple represents the lowest value, highest value, and steps for each factor.

159 ity single SuperLM as in prior works. Furthermore, our choices for the heuristic partition and 160 inductive bias result in less number of student models of comparable size in each sub-space which

161 alleviates conflicts in weight-sharing.

We extract student subnetworks from the SuperLM by a simple truncation strategy like bottom-left extraction. We defer more sophisticated extraction strategies to future work. In the above strategy, given a specific architecture $\alpha = \{l, d_{hid}, r, h\}$, (i) we first extract alternate *l* Transformer layers from the SuperLM; (ii) then extract bottom-left sub-matrices in terms of d_{hid} and *r* from the original matrices that represent the hidden dimension and the MLP ratio respectively; (iii) finally, for the attention heads, we extract the leftmost *h* heads and retain the dimension of each head as the SuperLM.

168 **3.2 Task-agnostic SuperLM Training**

We illustrate SuperLM training process in Algorithm 1. Given a large pre-trained language model (e.g., 169 BERT) as teacher, we initialize the SuperLM with the weights of teacher. In each step of SuperLM 170 training, we randomly sample several student subnetworks from the search space; apply knowledge 171 distillation between sampled subnetworks and the teacher to accumulate gradients; and then update 172 the SuperLM. During sampling, we employ Sandwich rule [27], also used in BigNAS [11], that 173 samples the smallest subnetwork, the largest subnetwork and M random ones for updating SuperLM. 174 The motivation is to improve the performance of all subnetworks by increasing the performance 175 lower bound (smallest subnetwork) and upper bound (largest one) across all subnetworks. 176

We leverage deep self-attention distillation [4] for task-agnostic training. To this end, we employ 177 multi-head self-attention relation distillation to align the attention distributions as well as scaled 178 dot-product of keys, queries and values of the teacher and sampled student subnetworks. Consider 179 A_1 , A_2 , A_3 to denote the queries, keys and values of multiple relation heads of teacher model, 180 and \mathbf{B}_1 , \mathbf{B}_2 , \mathbf{B}_3 respectively for a sampled subnetwork. Mean squared error (MSE(·)) between 181 multi-head self-attention relation of teacher and sampled subnetwork is used as distillation objective: 182

$$\mathcal{L} = \sum_{i=1}^{3} \beta_i \mathcal{L}_i, \ \mathcal{L}_i = \frac{1}{H} \sum_{k=1}^{H} \text{MSE}(\mathbf{R}_{ik}^T, \mathbf{R}_{ik}^S),$$
(4)

where $\mathbf{R}_{i}^{T} = \operatorname{softmax}(\mathbf{A}_{i}\mathbf{A}_{i}^{\top}/\sqrt{d_{k}}), \mathbf{R}_{i}^{S} = \operatorname{softmax}(\mathbf{B}_{i}\mathbf{B}_{i}^{\top}/\sqrt{d_{k}}), H$ is the number of attention heads; \mathbf{R}_{i}^{T} represents the teacher's Q - Q, K - K, or V - V relation; \mathbf{R}_{i}^{S} represents the same for student. \mathbf{R}_{ik}^{T} is the relation information based on one attention head, and d_{k} is the attention head size. 183 184 185

Relation knowledge distillation avoids the introduction of additional parameters to transform the 186 student's representations with different dimensions to align to that of the teacher. For the teacher 187 model and subnetworks with different number of attention heads, we first concatenate the self-188 attention vectors of different attention heads of the subnetwork and then split them according to the 189 number of relation heads of the teacher model. Then, we align their queries with the same number of 190 relation heads for distillation. In addition, we only transfer the self-attention knowledge from the 191 last layer of the teacher model to the last layer of the student model. Automatically selecting which 192 layers to align is an interesting research direction that we defer to future work. 193

The SuperLM for sub-space A_k is trained as: 194 Algorithm 1 Few-shot Task-agnostic Knowledge $\boldsymbol{W}_{\mathcal{A}_{k}}^{*} = argmin_{\boldsymbol{W}} \mathbb{E}_{\alpha \in \mathcal{A}}[\mathcal{L}(\boldsymbol{W}_{\alpha}; \boldsymbol{U}; \mathcal{D}_{train})], \underbrace{\text{Distillation with AutoDistil}}_{\boldsymbol{U}_{k}}]$ **Input:** Partitioned K sub-spaces \mathcal{A}_k ; initialized K SuperLMs S_k on \mathcal{A}_k ; pre-trained teacher T; unlabeled where, K is the number of sub-space parti-195 data D; training epochs E; sampling steps Mtions; W are the weights of the SuperLM; 196 **Output:** Trained SuperLMs $\{S_k\}$ W_{α} are the weights in W specified by the ar-197 for k = 1 to K do chitecture α ; U are the weights of the teacher for i = 1 to E do 198 model including the self-attention module Get a batch of data from D199 used for distillation; \mathcal{D}_{train} is the training data for *batch* in D do 200 Clear gradients in SuperLM S_k set, and $\mathcal{L}(\cdot)$ is the self-attention loss function 201 for m = 1 to M do from Eqn. (4). 202 Randomly sample a subnetwork s from S_k Calculate self-attention distil. loss between 3.3 Lightweight Optimal Student Search 203 subnetwork s and teacher T with Eqn. (4) Accumulate gradients We outline two search strategies for selecting 204 end for the optimal student subnetwork. 205 Update S_k with the accumulated gradients end for Task-agnostic search. We adopt this to be our 206

primary strategy to compare against all base-207 lines since it does not access any task label 208 information. We compute the task-agnostic 209

end for end for self-attention distillation loss for all student subnetworks using Eqn. (4) on a heldout validation set

210 from the unlabeled training corpus. The student subnetworks are directly obtained by bottom-left 211 extraction from the well-trained SuperLM (outlined in Section 3.1). This process is lightweight since 212 it does not require any training or adaptation of the student and number of subnetworks is limited. The 213 optimal student is given by the subnetwork with least validation loss subject to following constraint. 214

$$\alpha_{\mathcal{A}}^* = \operatorname{argmin}_{\alpha \in \mathcal{A}_{1,2,\cdots K}} \mathcal{L}(\boldsymbol{W}_{\alpha}^*; \mathcal{D}_{val}), \quad s.t. \quad g(\alpha) < c, \tag{6}$$

where W^*_{α} is the weights of architecture α obtained from $W^*_{\mathcal{A}_k}$, \mathcal{D}_{val} is the validation data set, \mathcal{L} 215

is the self-attention distillation loss, and $q(\cdot)$ is a function to calculate the computational cost (e.g., 216

#FLOPs, #parameters) of the subnetwork subject to a given user-specified resource constraint c. 217

Task-proxy search. We compare our task-agnostic search against another strategy that considers a 218 proxy task (e.g., MNLI [28]) with label information to fine-tune the 256 candidate subnetworks in 219 each sub-space. The optimal student in each sub-space is given by the one with the best downstream 220 task performance (e.g., accuracy). Note that, for this strategy, the proxy task is used only during 221 search while the NAS training is still fully task-agnostic. 222

223 4 Experiments

Datasets. We conduct experiments on the General Language Understanding Evaluation (GLUE) 224 benchmark [29]. We compare our method with the baseline methods on two single-sentence classifi-225 cation tasks (CoLA [30], SST-2 [31]), two similarity and paraphrase tasks (MRPC [32], QQP [33]), 226 and three inference tasks (MNLI [28], QNLI [34], RTE [35, 36, 37, 38])¹. We report accuracy for 227 MNLI, QNLI, QQP, SST-2, RTE, report f1 for MRPC, and report Matthew's correlation for CoLA. 228 Baselines. We compare against several *task-agnostic methods*² generating compressed models from 229 BERT_{base} teacher, using (i) knowledge distillation like BERT_{SMALL} [39], Truncated BERT [28], 230 DistilBERT [5], TinyBERT [6], MINILM [4]; as well as those based on Neural Architecture Search, 231 like AutoTinyBERT [13], DynaBERT [12], and NAS-BERT [14]. 232 AutoDistil configuration. We use uncased $\text{BERT}_{\text{BASE}}$ as the teacher consisting of 12 Trans-233 former layers, 12 attention heads; with the hidden dimension and MLP ratio being 768 and 4, 234 respectively. It consists of 109M parameters with 11.2G FLOPs. We use English Wikipedia and 235

BookCorpus data for SuperLM training with WordPiece tokenization. We use a batch size of 128 and 4*e*-5 as the peak learning rate for 10 epochs. The maximum sequence length is set to 128. The coefficients in distillation objective (Eqn. (4)), β_1 , β_2 , and β_3 , are all set to 1. We distill the self-attention knowledge of the last layer to train the SuperLM. Both the teacher and SuperLM are initialized with pre-trained BERT_{BASE}. Other hyper-parameter settings are shown in Appendix. We

use 16 V100 GPUs to train the SuperLM with 336 GPU-hours as the training cost.

242 4.1 Finding the Optimal Compressed Models

AutoDistil_{Agnostic} is obtained by fully task-agnostic training and task-agnostic search without using any task label information. We set a constraint in Eqn. (6) such that the #FLOPs of the optimal compressed model is atleast 50% less than the teacher model. We rank all the subnetworks contained in all the partitions of the trained SuperLM by their self-attention distillation loss on the heldout validation set, and select the one that meets the constraint with the minimum loss.

AutoDistil_{Proxy} uses MNLI [28] as a proxy to estimate downstream task performance of different 248 subnetworks. Prior work [40] has demonstrated performance improvements in MNLI to be correlated 249 to other GLUE tasks. To this end, we fine-tune all the 256 subnetworks in each partition of the trained 250 superLMs, and select corresponding subnetworks with the best trade-off between task performance 251 (accuracy) and computational cost (#FLOPs). This results in K=3 optimal students, corresponding to 252 $AutoDistil_{Proxy_B}$, $AutoDistil_{Proxy_S}$ and $AutoDistil_{Proxy_T}$ obtained from the corresponding 253 sub-spaces of SuperLM_{Base}, SuperLM_{Small} and SuperLM_{Tiny}, respectively. Notably all students are 254 obtained from the AutoDistil SuperLM still trained in a fully task-agnostic fashion. 255

4.1.1 Comparison with Traditional Knowledge Distillation Baselines

We compare AutoDistil against state-of-the-art KD models distilled from the same teacher 257 BERT_{BASE} in Table 2 with respect to the following measures: computational cost in the form of (i) 258 FLOPs and (ii) parameters, along with (iii) improvement in the average task performance aggregated 259 over all the GLUE tasks. We observe that the compressed model AutoDistilAgnostic generated 260 261 via our task-agnostic SuperLM training leads to up to 3x reduction in FLOPs over state-of-the-art 262 distilled models (e.g., MINILM [4], TinyBERT [6], DistilBERT [5]) that are hand-engineered while matching the overall task performance. The most aggressive compressed version corresponding to 263 AutoDistil_{Proxy_T} obtains a massive 41x reduction in FLOPs over BERT_{BASE} while incurring 264 5 point accuracy drop in GLUE (excluding CoLA) and 10 point drop (including CoLA). Notably 265 CoLA is a syntactic task in contrast to other semantic tasks in the benchmark like natural language 266 inference, paraphrase detection and sentiment classification. This depicts an interesting impact of 267 massive model compression on varying task types. 268

4.1.2 Comparison with Neural Architecture Search Baselines

We report the performance of several NAS-generated student models of comparable FLOPs and parameters from corresponding papers in Table 2.

¹We ignore STS-B for a fair comparison with our strongest KD baseline MINILM [4] that do not report it. ²For a fair comparison, we do not include MobileBERT [7] that uses BERT_{large} as teacher.

Table 2: Performance comparison between students from traditional task-agnostic distillation; multistage one-shot NAS with additional pre-training, distillation; and single-stage few-shot AutoDistil Our results are averaged over 5 runs with baselines reported from corresponding papers.

Model (Metric)	#FLOPs (G)	#Para (M)	MNLI-m (Acc)	QNLI (Acc)	QQP (Acc)	SST-2 (Acc)	CoLA (Mcc)	MRPC (Acc)	RTE (Acc)	Average
$BERT_{BASE}$ [1] (teacher)	11.2	109	84.5	91.7	91.3	93.2	58.9	87.3	68.6	82.2
Base-sized	d Models f	rom Tas	sk-agnosti	c KD M	ethods	and Au	itoDis	stil		
BERT _{SMALL} [39]	5.66	66.5	81.8	89.8	90.6	91.2	53.5	84.9	67.9	80.0
Truncated BERT [28]	5.66	66.5	81.2	87.9	90.4	90.8	41.4	82.7	65.5	77.1
DistilBERT[5]	5.66	66.5	82.2	89.2	88.5	91.3	51.3	87.5	59.9	78.6
TinyBERT [6]	5.66	66.5	83.5	90.5	90.6	91.6	42.8	88.4	72.2	79.9
MINILM [4]	5.66	66.5	84.0	91.0	91.0	92.0	49.2	88.4	71.5	81.0
$AutoDistil_{Proxy_B}$	4.40	50.1	83.8	90.8	91.1	91.1	55.0	88.8	71.9	81.7
Small-sized Me	odels from	Multi-	stage One-	-shot N	AS Met	hods ar	ud Auto	oDisti	1	
AutoTinyBERT-KD-S1 [13]	1.69	30.0	82.3	89.7	89.9	91.4	47.3	88.5	71.1	80.0
DynaBERT [12]	1.81	37.7	82.3	88.5	90.4	92.0	43.7	81.4	63.2	77.4
NAS-BERT ₁₀ [14]	2.30	10.0	76.4	86.3	88.5	88.6	34.0	79.1	66.6	74.2
AutoDistil _{Proxys}	2.02	26.1	83.2	90.0	90.6	90.1	48.3	88.3	69.4	79.9
AutoDistil _{Agnostic}	2.13	26.8	82.8	89.9	90.8	90.6	47.1	87.3	69.0	79.6
Tiny-sized Models from Multi-stage One-shot NAS Methods and AutoDistil										
AutoTinyBERT-KD-S4 [13]	0.30	10.1	76.0	85.5	86.9	86.8	20.4	81.4	64.9	71.7
NAS-BERT ₅ [14]	0.87	5.00	74.4	84.9	85.8	87.3	19.8	79.6	66.6	71.2
$AutoDistil_{Proxy_T}$	0.27	6.88	79.0	86.4	89.1	85.9	24.8	78.5	64.3	72.6

AutoDistil outperforms all competing methods on aggregate for all sizes; except for small-sized model; where it has marginally lower performance (0.1 points on avg) compared to AutoTinyBERT.

It is worthwhile to note that computational cost of training process is another important dimension for comparing methods. This is especially important when comparing to NAS methods

that use multi-stage training; where additional pre-training and
distillation is applied to NAS-generated candidates.

	-	-	
Cost (GPU hours)	AutoTiny BERT	AutoTiny BERT-Fast	Auto Disti
SuperNet Training	NR	NR	336
Search	150	12	<1
Further Training	870	290	0

To better understand the impact of single-stage vs. multi-stage

methods on the training cost, we compare the overall cost of 280 NAS for AutoDistil and that reported in AutoTinyBERT³ 281 for the small model segment in Table 3. AutoDistil is much 282 cheaper due to its single-stage training protocol; where no ad-283 ditional pre-training or distillation is needed. It is worth noting 284 that the overall SuperNet training cost of AutoDistil (the 285 most expensive component of NAS) is less or comparable to 286 the additional training cost of re-training candidate models 287 for AutoTinyBERT. Note that AutoTinyBERT does not report 288 their SuperNet training cost. Additionally, AutoDistil has 289 a much faster search mechanism due to (1) inductive biases 290

Table 3: Cost (V100 GPU hours) comparison for generating students of similar FLOPs. ^{NR}AutoTiny-BERT does not report the cost of SuperNet training - typically the most expensive step. Further Training refers to additional pre-training applied to NAS-generated candidates

²⁹¹ built into the search space definition to limit the number of student architectures and (2) task-agnostic ²⁹² search that only requires computing self-attention validation loss without the need for any training.

Finally, we show the pareto frontier of student subnetworks generated by several KD and NAS methods in Figure 3 for the MNLI task. The blue points represent all the subnetworks extracted from AutoDistil and red points denote the optimal ones, all fine-tuned on the MNLI task. We observe the optimal AutoDistil models to outperform several competing methods.

297 4.1.3 Task-agnostic Training Strategies

We study different task-agnostic strategies for SuperLM training in AutoDistil . Specifically, we compare three strategies in Table 4. (i) We replacing the KD loss in Eqn. (4) with masked language modeling (MLM) loss [1] to calculate gradients which is the most widely used task-agnostic pre-training and distillation strategy. (ii) KD_{att} +Cont further continues training the searched compressed models on the large language

Table 4: Comparing task-agnostic
SuperLM training strategies.

Strategy	MRPC	RTE
MLM	89.4	68.2
KD _{att} +Cont.	91.0	71.8
KD_{att}	91.2	71.5

corpus (iii) KD_{att} is the strategy adopted in AutoDistil for self-attention distillation. We evaluate subnetworks with the same architecture (6 layers, 768 hidden, 12 heads, MLP ratio 4) from the

³Other NAS methods either use a different hardware for training or do not report the cost.



(d) Acc vs #Para (SuperLM_{Base}). (e) Acc vs #Para (SuperLM_{Small}). (f) Acc vs #Para (SuperLM_{Tiny}).

Figure 3: Computational cost vs. task (MNLI) performance trade-off for all 256 subnetworks contained in each of K SuperLMs (base, small and tiny). 3(a)-3(c) show the trade-off between accuracy (Y-axis) and #FLOPs (X-axis), and 3(d)-3(f) show the trade-off between accuracy (Y-axis) and #Para (X-axis). We show the optimal compressed AutoDistil student for each SuperLM marked in red, along with other state-of-the-art KD and NAS techniques for comparison.

trained SuperLM. We fine-tune the subnetworks on RTE and MRPC tasks, and report accuracy and f1 respectively. First, we observe self-attention distillation to perform better than MLM, for SuperLM training. Second, we observe limited performance gains with continued training of the optimal subnetworks from NAS as done in existing works **demonstrating the effectiveness of our single-stage training protocol**.

312 4.1.4 One-shot vs. Few-shot NAS with Varying K

For few-shot NAS, we choose K=3 (i.e. 3 sub-313 spaces) for following reasons: (1) the 3 sub-spaces 314 correspond to base, small and tiny model sizes; as 315 316 used in prior work in CV; e.g. AutoFormer [17], (2) searching over different values of K is a very 317 resource-extensive process since it requires training 318 K SuperLMs for each choice of K, (3) As K in-319 creases, the search process becomes similar to the 320

Table 5: Search space design strategies.

Search Space Size (#subnetworks) One-shot $(K = 1)$ Few-shot $(K = 3)$					
27	864	11232	256*3		
88.2 67.2	87.5 64.5	85.1 62.8	91.2 71.8		
	Sea One-s 27 38.2 57.2	Search Sp One-shot (<i>K</i> 27 864 38.2 87.5 57.2 64.5	Search Space Siz One-shot $(K = 1)$ 27 864 88.2 87.5 85.1 57.2 64.5 62.8		

undesirable brute-force discrete search that trains all models in search space individually.

To understand the effect of few-shot NAS vs. one-shot NAS, we compare the performance of a 322 single space (K = 1) to multiple sub-spaces (K = 3). We extract subnetworks with the same 323 architecture (6 layers, 768 hidden, 12 heads, MLP ratio 4) from trained SuperLMs for each strategy 324 for evaluation with results presented in Table 5. For one-shot NAS, we consider a single search space 325 containing different numbers of subnetworks (e.g., 27, 864, 11232). The few-shot NAS contains 256 326 subnetworks in each partition. We fine-tune the subnetworks on RTE and MRPC tasks, and report 327 accuracy and f1 respectively. We observe fewer subnetworks contained in a single search space for 328 one-shot NAS result in a better performance. This results from optimization interference and gradient 329 conflicts as the number and size of subnetworks increase in the space. Finally, we observe our design 330 strategy performs the best while containing lesser number of subnetworks **demonstrating the benefit** 331 of few-shot NAS for language model distillation. 332

333 4.1.5 Comparing Search Strategies and Optimal Architectures

From Table 2, we observe that the student models $AutoDistil_{Agnostic}$ and $AutoDistil_{Proxys}$ obtained from SuperLM_{small} by task-agnostic and task-proxy search strategies respectively obtain a similar trade-off between performance and cost. The task-proxy search results in a minor performance gain 0.3 over the fully task-agnostic search mechanism. Table 6 shows the configuration of searched optimal architectures from AutoDistil with corresponding computational cost. For reference, we also show the architecture of the teacher BERT_{BASE} and a state-of-the-art distilled model MINILM [4] that are hand-engineered. We observe that the obtained architectural hyper-parameters are quite

non-standard and difficult to obtain by 341 trial and error considering the large 342 space of Transformer architectures. 343 We also observe that optimal com-344 pressed models have thin-and-deep 345 346 structure consistent with findings that thinner and deeper models perform 347 better [26] than wider and shallower 348 ones. While we use this as an induc-349 tive bias for sub-space partitioning, 350 our search space (Table 1) also con-351 tains diverse subnetworks with differ-352

Table 6: Architecture comparison between the optimal compressed students searched by AutoDistil with state-ofthe-art hand-engineered students distilled from BERT_{BASE}.

Model	#Layers	#Hid	Ratio	#Heads	#FLOPs	#Para
BERTBASE	12	768	4	12	11.2G	109M
MINILM	6	768	4	6	5.66G	66.5M
AutoDis.Agnostic	11	352	4	10	2.13G	26.8M
AutoDis.ProxyB	12	544	3	9	4.40G	50.1M
AutoDis. Proxys	11	352	4	8	2.02G	26.1M
$AutoDis{Proxy_T}$	7	160	3.5	10	0.27G	6.88M

ent depth and width. Non-maximal MLP ratio and attention heads for optimal compression indicate that self-attention and feed-forward layers of Transformers are overparameterized [21, 22].

355 5 Related Work

356 Task-specific knowledge distillation. Knowledge distillation (KD) [41] is a widely used technique 357 for model compression, which transfers knowledge from a large teacher to a smaller student. Taskspecific KD aims to generate smaller students by using downstream task label information. Typical 358 task-specific KD works include BERT-PKD [42], BERT_{SMALL} [39], TinyBERT [6], DynaBERT [12], 359 and SparseBERT [43]. While task-specific KD often achieves good task performance, a typical 360 drawback is that it is resource-consuming to run KD for each and every task, and also not scalable. 361 **Task-agnostic knowledge distillation.** In contrast to task-specific KD, we explore task-agnostic KD 362 that does not use any task label information. The distilled task-agnostic models can be re-used by 363 simply fine-tuning on downstream tasks. Task-agnostic KD leverages knowledge from soft target 364 probabilities, hidden states, layer mappings and self-attention distributions of teacher to train student 365 models. Typical task-agnostic KD works include DistilBERT [5] MobileBERT [7], and MiniLM [4]. 366 MobileBERT assumes that students have the same number of layers as the teacher for layer-by-layer 367 distillation. MiniLM transfers self-attention knowledge from the last layer of the teacher to that of the 368 student. These works rely on hand-designed architecture for the student models for KD that requires 369 several trials, and needs to be repeated for a new student with a different cost. In contrast, we develop 370 techniques to automatically design and distill several student models with variable cost using NAS. 371 Neural Architecture Search. While NAS has been extensively studied in computer vision [8, 372 9, 10, 11], there has been relatively less exploration in natural language processing. Evolved 373 Transformer [44] and HAT [45] search for efficient sub-networks from the Transformer architecture 374 for machine translation tasks. Some recent approaches closest to our method include, DynaBERT [12], 375 AutoTinyBERT [13] and NAS-BERT [14]. DynaBERT performs task-specific distillation. NAS-376 BERT performs two-stage knowledge distillation with pre-training and fine-tuning of candidates. 377 Similar to above approaches, AutoTinyBERT also employs one-shot NAS with a single large search 378 space containing millions of subnetworks that result in co-adaption and weight-sharing challenges 379 for SuperLM training. Further it also uses a multi-stage training protocol for further pre-training 380 and distillation of the NAS-generated candidates. In contrast, AutoDistil employs few-shot NAS 381 with a compact search space design with a single-stage task-agnostic training protocol. This further 382 allows us to do a lightweight search for the optimal student without re-training. 383

384 6 Conclusion

We develop a few-shot task-agnostic NAS framework, namely AutoDistil for distilling large 385 386 language models into compressed students with variable computational cost. To address the coadaption and weight-sharing challenges for SuperLM training, we partition the Transformer search 387 space into K compact sub-spaces covering important architectural components like the network 388 depth, width, and number of attention heads. We leverage deep self-attention distillation for fully 389 task-agnostic SuperLM training and lightweight optimal student search without any re-training. 390 This allows our students to be re-used by simply fine-tuning on downstream tasks. AutoDistil 391 generates students with 3x less computational cost (FLOPs) than state-of-the-art task-agnostic 392 distillation methods while obtaining a similar downstream task performance in the GLUE benchmark. 393

394 **References**

- [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirec tional transformers for language understanding. In *NAACL*, pages 4171–4186, Minneapolis, Minnesota,
 June 2019. Association for Computational Linguistics.
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind 398 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, 399 Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens 400 Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, 401 Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models 402 are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, 403 Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, 404 Inc., 2020. 405
- [3] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep
 learning in NLP. In *ACL*, pages 3645–3650, Florence, Italy, July 2019. Association for Computational
 Linguistics.
- [4] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention
 distillation for task-agnostic compression of pre-trained transformers. In H. Larochelle, M. Ranzato,
 R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*,
 volume 33, pages 5776–5788. Curran Associates, Inc., 2020.
- [5] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert:
 smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [6] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu.
 Tinybert: Distilling bert for natural language understanding. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 4163–4174, 2020.
- [7] Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. Mobilebert: a
 compact task-agnostic bert for resource-limited devices. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2158–2170, 2020.
- [8] Hieu Pham, Melody Guan, Barret Zoph, Quoc Le, and Jeff Dean. Efficient neural architecture search via
 parameters sharing. In *International Conference on Machine Learning*, pages 4095–4104. PMLR, 2018.
- [9] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V
 Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2820–2828, 2019.
- [10] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once for all: Train one network and
 specialize it for efficient deployment. In *International Conference on Learning Representations*, 2020.
- [11] Jiahui Yu, Pengchong Jin, Hanxiao Liu, Gabriel Bender, Pieter-Jan Kindermans, Mingxing Tan, Thomas
 Huang, Xiaodan Song, Ruoming Pang, and Quoc Le. Bignas: Scaling up neural architecture search with
 big single-stage models. In *European Conference on Computer Vision*, pages 702–717. Springer, 2020.
- [12] Lu Hou, Zhiqi Huang, Lifeng Shang, Xin Jiang, Xiao Chen, and Qun Liu. Dynabert: Dynamic bert with
 adaptive width and depth. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors,
 Advances in Neural Information Processing Systems, volume 33, pages 9782–9793. Curran Associates,
 Inc., 2020.
- [13] Yichun Yin, Cheng Chen, Lifeng Shang, Xin Jiang, Xiao Chen, and Qun Liu. AutoTinyBERT: Automatic
 hyper-parameter optimization for efficient pre-trained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5146–5157. Association for Computational
 Linguistics, August 2021.
- [14] Jin Xu, Xu Tan, Renqian Luo, Kaitao Song, Jian Li, Tao Qin, and Tie-Yan Liu. NAS-BERT: task-agnostic
 and adaptive-size BERT compression with neural architecture search. In Feida Zhu, Beng Chin Ooi, and
 Chunyan Miao, editors, *KDD '21: The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, Singapore, August 14-18, 2021*, pages 1933–1943. ACM, 2021.
- 444 [15] Gabriel Bender, Pieter-Jan Kindermans, Barret Zoph, Vijay Vasudevan, and Quoc Le. Understanding
 445 and simplifying one-shot architecture search. In *International Conference on Machine Learning*, pages
 446 550–559. PMLR, 2018.

- [16] Yiyang Zhao, Linnan Wang, Yuandong Tian, Rodrigo Fonseca, and Tian Guo. Few-shot neural architecture 447 search. In International Conference on Machine Learning, pages 12707–12718. PMLR, 2021. 448
- 449 [17] Minghao Chen, Houwen Peng, Jianlong Fu, and Haibin Ling. Autoformer: Searching transformers for 450 visual recognition. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 12270-12280, 2021. 451
- [18] Tianyi Chen, Bo Ji, Tianyu Ding, Biyi Fang, Guanyi Wang, Zhihui Zhu, Luming Liang, Yixin Shi, Sheng 452 Yi, and Xiao Tu. Only train once: A one-shot neural network training and pruning framework. Advances 453 in Neural Information Processing Systems, 34, 2021. 454
- [19] Jing Xu, Yu Pan, Xinglin Pan, Steven Hoi, Zhang Yi, and Zenglin Xu. Regnet: Self-regulated network for 455 image classification. IEEE Transactions on Neural Networks and Learning Systems, 2022. 456
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz 457 458 Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in neural information processing systems, pages 5998-6008, 2017. 459
- [21] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In NeurIPS, 460 pages 14014-14024, 2019. 461
- [22] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-462 attention: Specialized heads do the heavy lifting, the rest can be pruned. arXiv preprint arXiv:1905.09418, 463 2019. 464
- [23] Prakhar Ganesh, Yao Chen, Xin Lou, Mohammad Ali Khan, Yin Yang, Deming Chen, Marianne Winslett, 465 Hassan Sajjad, and Preslav Nakov. Compressing large-scale transformer-based models: A case study on 466 467 bert. arXiv preprint arXiv:2002.11985, 2020.
- [24] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In H. Wallach, 468 H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, Advances in Neural 469 Information Processing Systems, volume 32. Curran Associates, Inc., 2019. 470
- [25] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-471 attention: Specialized heads do the heavy lifting, the rest can be pruned. In Proceedings of the 57th Annual 472 Meeting of the Association for Computational Linguistics, pages 5797–5808, Florence, Italy, July 2019. 473 Association for Computational Linguistics. 474
- [26] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua 475 476 Bengio. Fitnets: Hints for thin deep nets. In Yoshua Bengio and Yann LeCun, editors, 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference 477 Track Proceedings, 2015. 478
- [27] Jiahui Yu and Thomas S Huang. Universally slimmable networks and improved training techniques. In 479 Proceedings of the IEEE/CVF international conference on computer vision, pages 1803–1811, 2019. 480
- [28] Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence 481 understanding through inference. In Proceedings of the 2018 Conference of the North American Chapter of 482 the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), 483
- 484 pages 1112–1122, New Orleans, Louisiana, June 2018. Association for Computational Linguistics.
- 485 [29] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In Proceedings of 486 the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, pages 487 353–355, Brussels, Belgium, November 2018. Association for Computational Linguistics. 488
- [30] Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. Neural network acceptability judgments, 2018.
- [31] Richard Socher et al. Recursive deep models for semantic compositionality over a sentiment treebank. 490 In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, pages 491
- 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. 492

489

- [32] William B Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In 493 494 Proceedings of the Third International Workshop on Paraphrasing (IWP2005), 2005.
- 495 [33] Zihan Chen, Hongbo Zhang, Xiaoji Zhang, and Leqi Zhao. Quora question pairs. URL https://www.kaggle. com/c/quora-question-pairs, 2018. 496

- 497 [34] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for
 498 machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- [35] Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment challenge.
 In *Machine Learning Challenges Workshop*, pages 177–190. Springer, 2005.
- [36] R Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor.
 The second pascal recognising textual entailment challenge. In *Proceedings of the Second PASCAL Challenges Workshop on Recognising Textual Entailment*, 2006.
- [37] Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and William B Dolan. The third pascal recognizing
 textual entailment challenge. In *Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing*, pages 1–9, 2007.
- [38] Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. The fifth pascal recognizing textual
 entailment challenge. In *TAC*, 2009.
- ⁵⁰⁹ [39] Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better: On the ⁵¹⁰ importance of pre-training compact models. *arXiv preprint arXiv:1908.08962*, 2019.
- [40] Tianlong Chen, Jonathan Frankle, Shiyu Chang, Sijia Liu, Yang Zhang, Zhangyang Wang, and Michael
 Carbin. The lottery ticket hypothesis for pre-trained bert networks. In H. Larochelle, M. Ranzato,
 R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*,
 volume 33, pages 15834–15846. Curran Associates, Inc., 2020.
- [41] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015.
- [42] Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. Patient knowledge distillation for bert model compression.
 In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th
 International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 4323–4332,
 2019.
- [43] Dongkuan Xu, Ian EH Yen, Jinxi Zhao, and Zhibin Xiao. Rethinking network pruning–under the pre-train
 and fine-tune paradigm. In *Proceedings of the Human Language Technology Conference of the NAACL*,
 2021.
- [44] David So, Quoc Le, and Chen Liang. The evolved transformer. In Kamalika Chaudhuri and Ruslan
 Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97
 of *Proceedings of Machine Learning Research*, pages 5877–5886. PMLR, 09–15 Jun 2019.
- [45] Hanrui Wang, Zhanghao Wu, Zhijian Liu, Han Cai, Ligeng Zhu, Chuang Gan, and Song Han. Hat:
 Hardware-aware transformers for efficient natural language processing. In *Annual Conference of the Association for Computational Linguistics*, 2020.

530 Checklist

531	1. For all authors
532 533	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
534	(b) Did you describe the limitations of your work? [Yes] Appendix
535	(c) Did you discuss any potential negative societal impacts of your work? [N/A]
536 537	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
538	2. If you are including theoretical results
539	(a) Did you state the full set of assumptions of all theoretical results? [Yes]
540	(b) Did you include complete proofs of all theoretical results? [Yes]
541	3. If you ran experiments
542 543	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [Yes]
544 545	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
546 547	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No]
548 549	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
550	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
551	(a) If your work uses existing assets, did you cite the creators? [Yes]
552	(b) Did you mention the license of the assets? [N/A]
553 554	(c) Did you include any new assets either in the supplemental material or as a URL? $[N/A]$
555 556	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
557 558	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
559	5. If you used crowdsourcing or conducted research with human subjects
560 561	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
562 563	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
564 565	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]