PRE-TRAINED SUMMARIZATION DISTILLATION

Sam Shleifer *
Hugging Face
sam@huggingface.co

Alexander M. Rush Hugging Face and Cornell University sasha@huggingface.co

ABSTRACT

Recent state-of-the-art approaches to summarization utilize large pre-trained Transformer models. Distilling these models to smaller student models has become critically important for practical use; however there are many different distillation methods proposed by the NLP literature. Recent work on distilling BERT for classification and regression tasks shows strong performance using direct knowledge distillation. Alternatively, machine translation practitioners distill using pseudo-labeling, where a small model is trained on the translations of a larger model. A third, simpler approach is to "shrink and fine-tune" (SFT), which avoids any explicit distillation by copying parameters to a smaller student model and then fine-tuning. We compare these three approaches for distillation of Pegasus and BART, the current and former state of the art, pre-trained summarization models, and find that SFT outperforms knowledge distillation and pseudo-labeling on the CNN/DailyMail dataset, but under-performs pseudo-labeling on the more abstractive XSUM dataset. PyTorch Code and checkpoints of different sizes are available through Hugging Face transformers.

1 Introduction

Pre-trained transformer models continue to grow in size (Brown et al., 2020), motivating researchers to try to compress large pre-trained checkpoints into smaller, faster versions that retain strong performance.

Recently, researchers have developed promising methods for utilizing pre-trained models for sequence-to-sequence ("Seq2Seq") language generation tasks, showing particularly large improvements in performance on summarization. BART (Lewis et al., 2019), a Seq2Seq transformer (?) recently achieved state of the art performance on the Extreme Summarization ("XSUM") and CNN/Dailymail("CNN") summarization datasets (Narayan et al., 2018; See et al., 2017), with particularly large improvements on XSUM. A few months later, Pegasus (Zhang et al., 2019) achieved further performance improvements by replacing BART's more general pre-training objective with a pre-training objective specifically tailored to abstractive text summarization and a 25% larger model.

In parallel to the progress on summarization, DistilBERT, BERT of Theseus, TinyBERT, MobileBERT, and MiniLM showed that BERT, a large pre-trained transformer, can be shrunk substantially without much performance degradation on the GLUE suite of non-generative tasks using direct Knowledge Distillation ("KD") (Sanh et al., 2019; Jiao et al., 2019; Sun et al., 2020; Wang et al., 2020; Devlin et al., 2018; Wang et al., 2018; Hinton et al., 2015). On the other hand, past work in machine translation suggests that Seq2Seq models should be compressed with pseudo-labeling ("PL") (Kim and Rush, 2016). PL approaches run beam search with the teacher model on the whole training dataset, then retrain a smaller student model from scratch on those translations. Since BART and Pegasus are both pre-trained, like BERT, and Seq2Seq, like the translation models, it is not clear whether Knowledge Distillation or pseudo-labeling is the best approach. Other approaches are possible as well. Several works suggest that subsets of trained teacher models can be extracted directly (Sanh et al., 2019; Xu et al., 2020; Fan et al., 2019). We therefore propose a "shrink and fine-tune" ("SFT") approach that extracts a student model from the maximally spaced layers of a fine-tuned teacher. Since transformer layers are stacked using residual connections, we hypothesize that removing full layers has a minimal impact on summarization performance. This shrunken student model is then used to re-run the original fine-tuning procedure without modification.

We test all three methods on the CNN and XSUM datasets. On CNN, SFT outperforms the more expensive methods. For both BART and Pegasus, SFT produces distilled models that are 75% faster than their teacher with minimal loss in performance. On the more abstractive XSUM task, KD and PL can generate significant improvements over SFT. For BART, we use KD ¹ to match teacher performance. For Pegasus, no technique matches teacher performance, but PL comes closest. As shown in Figure 1, we manage to find an approach that generates the best available model at its computational budget for each task and teacher model. In the BART case, we generate many such models of various sizes.

The paper is organized as follows: Section 2 discusses related work in further detail. Section 3 describes the specifics of our implementation of the three families of techniques. Section 5 describes summarization speed and quality for various teachers, datasets, student sizes, and distillation methods. Sections 6.2 and 6.3 describe extensions of pseudo-labeling and knowledge distillation which can further improve performance on the XSUM task.

2 RELATED WORK

Knowledge distillation is a compression technique where a smaller student model is trained to reproduce the logits of a larger teacher, rather than simply minimize the cross-entropy between the model's predicted distribution

^{*}Ask questions here, or start your own thread and tag @sshleifer. We are grateful to Stas Bekman, Zoe Shleifer, Patil Suraj and Victor Sanh for comments.

¹More specifically, we use extensions of KD proposed in Jiao et al. (2019), and explained in Section 3.



Figure 1: The best distilled checkpoint from Pegasus (P) and Bart (B) for XSUM and CNN at different sizes. In three out of four settings we are able to distill a student model to the same Rouge-2 score as the teacher with at least a 90% speedup.

Method		Setup			Knowledge Transfer			
wiethou	Task	Pre-Train Teacher	Pre-Train Student	Init.	Logits	Hidden	Gens.	
DistilBERT♡	GLUE	√	✓	√	√			
TinyBERT♣	GLUE	\checkmark	✓		\checkmark	\checkmark		
BERT-of-Theseus♦	GLUE	\checkmark		✓				
Seq-Level KD♠	MT						\checkmark	
KD†	Summ.	✓		√	√	√		
Pseudo-Labels †	Summ.	\checkmark		✓			\checkmark	
SFT †	Summ.	\checkmark		✓				

Table 1: A comparison of the setting studied and knowledge transfer techniques employed by different transformer distillation methods. † indicates our implementation. INIT: Are weights copied from teacher to student? PRETRAIN STUDENT: must the student be pre-trained? LOGITS: does the student learn from the teacher's logits? HIDDEN: does the student learn from the teacher's hidden states? GENS: Does the student learn from the teacher's generations? ♡: Sanh et al. (2019), ♣: Jiao et al. (2019), ♦: Xu et al. (2020), ♠: Kim and Rush (2016).

and the training labels (Bucila et al., 2006; Hinton et al., 2015). In a language modeling context, this allows the student model to learn a full distribution of possible next words in a given context, rather than just the next word in the training data.

Recent research on KD for pre-trained models has overwhelmingly focused on distilling BERT to perform well on GLUE tasks, rather than tasks that require text generation. Sanh et al. (2019) use a weighted average of KD loss and the traditional cross entropy data loss to train DistilBERT, a 6 layer distilled version of BERT, that is 60% faster on CPU and 50% faster on GPU. DistilBERT intializes student models by copying alternating layers. an idea we extend – in all of our experiments, we initialize students by copying maximally spaced layers. In TinyBERT, Jiao et al. (2019) add terms to the KD loss function which enforce student/teacher alignment at intermediate levels and improve performance. ² Bert-of-Theseus (Xu et al., 2020) randomly replaces multiple teacher layers with a single student layer during fine-tuning with probability r, such that each student layer learns to replicate 2 teacher layers. LayerDrop, a related technique, drops random parts of the teacher model during one long training run, allowing a smaller student model to be extracted at inference time (Fan et al., 2019). Distillation for Seq2Seq models has primarily used pseudo-labeling and produces strong results on machine translation, as shown in Kasai et al. (2020), Junczys-Dowmunt (2019), and Sun et al. (2019). Their approach consists of re-generating a new distilled dataset containing original source documents with pseudo-labels. The pseudo-labels are summaries generated by the teacher using beam search. After the long dataset generation process, they train a

²We refer to both of these formulations as KD. DistillBERT can be described as the tinyBERT variant with a zero coefficient on all terms besides the logits loss.

smaller student model on the "distilled" dataset. Kim and Rush (2016) call this type "Sequence-level Knowledge Distillation" in contrast to "Word-Level Knowledge Distillation", where knowledge is transferred through logits.

Recent work from Liu et al. (2020a) presents a new method to further improve fine-tuned summarization models by fine-tuning them on their own logits with added noise. Like quantization (Jacob et al., 2017), this method could be used before or after the other methods in this work.

Table 1 compares the attributes of these methods to our three approaches. Like Theseus, our experiments do not re-run pre-training. SFT is most similar to BERT-of-Theseus, and can even be described as running the Theseus procedure with r fixed at 100%, thereby saving computation. Our KD implementation is most similar to TinyBERT, and Pseudo-labels is most similar to Sequence Level Knowledge Distillation.

3 BACKGROUND AND METHODS

Assume we have a source document $x_1 \dots x_M$ and target document $y_1 \dots y_N$ in the standard sequence-to-sequence setting. A Seq2Seq transformer is composed of a transformer-based encoder (Enc) and decoder (Dec). Enc is trained to map x to contextual embeddings and Dec to map those contextual embeddings and the previously decoded words to a probability distribution for the next word, $p(y_{t+1}|y_{1:t},x)$.

Pre-trained Seq2Seq models such as BART and Pegasus learn parameters which are subsequently fine-tuned on Seq2Seq tasks like summarization. BART is pre-trained to reconstruct corrupted documents. Source documents x are corrupted versions of original target documents y, e.g. spans of words are masked and sentences are shuffled. Pegasus is pre-trained to generate the most important sentences extracted from an unlabeled document (y is the important sentences and x is the original documented with those sentences removed).

To fine-tune the models, we assume a dataset where each example (x, y) is a (document, summary) pair. In the Seq2Seq fine-tuning setting, we train the student model using the standard cross entropy loss:

$$\mathcal{L}_{\text{Data}} = -\sum_{t=1}^{T} \log p(y_{t+1}|y_{1:t}, x)$$
 (1)

where T in the target sequence length p is the model's predicted probability for the correct word. Our distillation experiments start with a teacher model fine-tuned in this manner.

3.1 DISTILLATION

We consider different approaches for compressing these models through distillation. All settings assume that we are learning a student model from a larger teacher. Define the notation Dec^L to represent a decoder with L Transformer layers (and similarly for Enc). Assuming we have a large pre-trained teacher model with decoder Dec^L , we are interested in compressing it to a smaller student model $\operatorname{Dec}^{L'}$. In most experiments, we do not compress the teacher's encoder.

Shrink and Fine-Tune Our most basic method SFT simply *shrinks* the teacher model to student size and re-fine-tunes this student model. Here each $l \in L'$ is copied fully from L; students are initialized by copying full maximally spaced decoder layers from teacher to student. For example, when creating a BART student with 3 decoder layers from the 12 encoder layer 12 decoder layer teacher, we copy the teacher's full Enc^L and decoder layers 0, 6, and 11 to the student. When deciding which layers to copy, we break ties arbitrarily; copying layers 0, 5, and 11 might work just as well. When copy only 1 decoder layer, we copy layer 0. We found this to work better than copying layer 11. The impact of initialization on performance is measured experimentally in Section 6.1. After initialization, the student model continues to fine-tune on the summarization dataset, with the objective of minimizing $\mathcal{L}_{\mathrm{Data}}$. As the initialization approach is simple and effective, it is used to initialize student models for both other methods.

Pseudo-labels In the pseudo-label setting, we replace the ground truth target documents Y with \hat{Y} , the teacher's generations for the source documents X, computed with beam search.

$$\mathcal{L}_{Pseudo} = -\sum_{t=1}^{T} \log p(\hat{y}_{t+1}|\hat{y}_{1:t}, x)$$
 (2)

After this procedure the student model is fine-tuned only on this new pseudo-labeled data.

Direct Knowledge Distillation (KD) In the KD setting, even more information is transferred from teacher to student, by encouraging the student to match the teacher's full probability distribution over possible next words at each position, by minimizing KL-Divergence(Kullback and Leibler, 1951; Sanh et al., 2019):

$$\mathcal{L}_{\text{Logits}} = \sum_{t=1}^{T} KL(Q_{t+1}, P_{t+1}), \tag{3}$$

where Q_{t+1} and P_{t+1} are teacher and student probability distributions over each next possible word at position t+1, and KL is the KL-Divergence.³ Since we use layer based compression, student and teacher layers output the same shape, and we can add another term to the loss function that encourages students to match teacher hidden states.

$$\mathcal{L}_{\text{Hid}} = \sum_{t=1}^{T} \sum_{l=1}^{L'} \text{MSE}(\boldsymbol{H}_{l}^{S}, \boldsymbol{H}_{\phi(l)}^{T})$$

$$\tag{4}$$

³KL-Divergence is implemented in PyTorch (Paszke et al., 2019) and explained well on Wikipedia.

Technique	Extra Supervision	Cost	Loss
SFT PseudoLabeling KD	T's Generations T's Hidden States, Logits	2.5 19 14	$egin{array}{c} \mathcal{L}_{ ext{Data}} \ \mathcal{L}_{ ext{Pseudo}} \ \mathcal{L}_{ ext{KD}} \end{array}$

Table 2: Training time of different distillation approaches. Cost is an estimate of how many hours were required to run the technique for the CNN dataset with BART as a teacher and the 12 encoder layer, 6 decoder layer student on a Titan RTX 2080 GPU.

Data	# Train	Avg. Source Words	Avg. Target Words	Size (MB)
CNN XSUM	262,567 204,017	756 358	56 21	1,331 501
EN-RO	610,319	23	23	178

Table 3: Dataset Statistics.

Here, MSE stands for mean squared error, \boldsymbol{H}_l^S retrieves the hidden state returned by student layer l, and $\phi(l)$ maps student layer l to the teacher layer whose output we would like them to emulate. $\boldsymbol{H}_{\phi(l)}^T$, therefore, is the output of a teacher layer. ⁴ For example, when creating a BART student with 3 decoder layers, we copy the full teacher encoder and decoder layers 0, 6, and 11 to the student. We then choose pairings in ϕ such that each student decoder layer is taught to behave like 4 decoder layers. Student layer 0's hidden state is paired o teacher layer 3, 1 to 7, and 11 to 11 ($\phi = [3,7,11]$). The student layers are therefore trained to perform the work of teacher layers 0-3, 4-7 and 8-11 respectively. ⁵

Our final KD formulation is a weighted average:

$$\mathcal{L}_{KD} = \alpha_{logits} \mathcal{L}_{Logits} + \alpha_{data} \mathcal{L}_{Data} + \alpha_{Hidn} \mathcal{L}_{Hidn}$$
(5)

We set $\alpha_{logits} = 0.8$ and $\alpha_{data} = 1$ following Sanh et al. (2019), and found $\alpha_{Hidn} = 3$ to perform best out of [1, 3, 10, 100] for BART on the XSUM development set.

Training Time Comparison Table 2 compares the training time of these three approaches. Whereas SFT simply requires fine-tuning a small model, computing \mathcal{L}_{KD} requires teacher logits; for each training example, we must run the large teacher model forwards as well as the student model forwards and backwards. Similarly, \mathcal{L}_{Pseudo} requires \hat{Y} , which is computed by running beam search with the teacher on the full training dataset. This large preprocessing cost can dwarf the cost of fine-tuning the student model on the pseudo-labels, as shown in Table 2, where 16.5 of the 19 GPU hour cost of producing a student model is spent generating the pseudo-labels. After initialization, SFT does not use the teacher model, and is therefore much cheaper.

4 EXPERIMENTAL SETUP

We experiment with both the CNN and XSUM abstractive summarization datasets, both of which are based on english language news articles. The CNN summaries are roughly 3 sentences long, and tend to be similar to text from the beginning of the document. The XSUM summaries are the first sentence of a BBC news article, which is then removed from the article, so are both shorter and more abstractive than CNN summaries. The original BART model's improvement over its predecessors was much more significant (roughly 6 ROUGE-2 points) on the more abstractive XSUM dataset than on the CNN dataset (1.5 points). Table 3 shows dataset statistics.

Generation and Evaluation We run beam search on the distilled models and measure summary quality using the Rouge implementation from the rouge_scorer python package. ROUGE scores for teachers and students can be found in Tables 6 and 7. Unlike the BART paper, we do not tokenize or otherwise preprocess summaries before scoring, leading to slightly lower scores. For inference speed comparison, we measure Summaries per Second using a batch size of 32 and mixed precision for BART on 1 GPU. Mixed precision overflows for Pegasus, so we use full precision.

Translation Experiments Translation experiments are included for comparison. These use the English-Romanian WMT 2016 English-Romanian Dataset ("EN-RO") (Bojar et al., 2016), and two teachers. "mBART" is pre-trained on many languages and then fine-tuned on bilingual data. (Liu et al., 2020b), "Marian" is trained from scratch on bilingual data (Tiedemann and Thottingal, 2020). These experiments are evaluated using BLEU with no post-processing (Papineni et al., 2002).

Training We stop training at whichever point comes first: the end of epoch 5 or the validation score not increasing for four consecutive evaluations (a full epoch). We measure "training cost" as the amount of time training takes on one Nvidia-RTX-2080 GPU + the cost of generating pseudo-labels, if applicable.

In experiments with a full sized (completely copied) encoder, we freeze its parameters during training. Initial experiments suggested that this did not impact performance but made fine-tuning faster by a factor of 5.6 We also freeze the positional and token embeddings.

 $^{^4}$ For a more detailed description of \mathcal{L}_{Hidn} , read the Methods Section of the TinyBert paper. Our approach is inspired by theirs, but we do not use per-layer weights, per-layer learning rates, embedding loss, or attention loss.

⁵A complete list of the ϕ mappings we used can be found here

⁶For KD, if the encoder is the same for teacher and student, it only needs to be run once. Back propagation is also much cheaper, as it can stop at the end of the encoder.

Teacher	Dataset	# GPU Hours	% GPU Hours	# Experiments	% Experiments
BART	XSUM	787	30%	102	36%
BART	CNN	365	14%	59	21%
mBART	EN-RO	332	13%	48	17%
Pegasus	XSUM	766	29%	42	15%
Marian	EN-RO	185	7%	26	9%
Pegasus	CNN	196	7%	10	3%
-	TOTALS	2,631		287	

Table 4: Effort calculations. Each row represents the resources spent attempting to distill a teacher to a smaller student model on a given dataset. Experiments were only counted if they lasted 15 minutes or more. % columns divide # columns by their sum.

Teacher	Size	Data	Teacher	SFT	SFT		KD		Pseudo	
			Score	Score	Cost	Score	Cost	Score	Cost	
BART †	12-3	XSUM	22.29	21.08	2.5	21.63	6	21.38	15	
Pegasus	16-4	XSUM	24.56	22.64	13	21.92	22	23.18	34	
BART	12-6	CNN	21.06	21.21	2	20.95	14	19.93	19.5	
Pegasus	16-4	CNN	21.37	21.29	31	-	-	20.1	48	
Marian	6-3	EN-RO	27.69	25.91	4	24.96	4	26.85	28	
mBART	12-3	EN-RO	26.457	25.6083	16	25.87	24	26.09	50	

Table 5: Main results. Score is Rouge-2 for the 2 summarization datasets (first 4 rows), and BLEU for the bottom two rows. Cost measures the GPU hours required to run the approach end to end, which, in the case of Pseudo-labeling, requires running beam search on the full training set. The highest scoring distillation technique is in bold.

Effort We did not spend equal resources on all datasets and models, as shown in Table 4. In particular, we ran fewer CNN experiments because SFT worked well in that case, and fewer Pegasus experiments because Pegasus takes longer to train. Many of the BART experiments on XSUM tested variants and hyperparameters for KD, which has yet to work well for Pegasus. If we had run 60 more Pegasus experiments on XSUM data, we might have found something that works better.

Model Notation We use shorthand notation to describe student models generated with our initialization procedure. For example, dBART-12-3 is a student model extracted from BART with (all) 12 encoder layers and 3 decoder layers. Similarly, all "Size" columns in tables use the Encoder Layers-Decoder Layers convention.

5 RESULTS

Table 5 shows the performance of 3 different approaches for different tasks, teachers and students. No approach dominates the others across all datasets. On CNN, SFT works best for both teachers. On XSUM, BART performs best with KD, while Pegasus performs best with PL. ROUGE-1 and ROUGE-L scores follow a similar pattern to ROUGE-2 in Table 5 on both summarization datasets. We additionally include translation experiments for comparison. On the English-Romanian translation dataset, PL works best for both teacher models.

Tables 6 and 7 show scores and inference times for many different student models on XSUM and CNN, respectively. These tables show the best student of a given size, regardless of distillation method. In 3 out of 4 contexts, distillation leads to relatively minor performance losses and significant speedups. On XSUM, both the 12-3 and 12-6 sized BART students outperform the teacher model at 93% and 43% speedups, whereas the Pegasus student falls more than a full ROUGE-2 point below the teacher model. On CNN, the 12-6 sized BART student outperforms the teacher, and the Pegasus teacher is close.

Note that Table 6 shows a higher score for the BART/XSUM 12-3 student than Table 5 shows. The stronger student was trained on pseudo-labels generated by Pegasus. The result is not included in Table 5's PL column, which shows results for student models trained on pseudo-labels generated by their teacher. We discuss this further in Section 6.2.

6 ANALYSIS

6.1 How does initialization impact distillation?

In Table 8, we show the validation cross entropy loss of dBART-12-3 students trained with the same, frozen encoder, but different decoder layers copied from different sources. The default SFT initialization for 3 layer students, copying layers 0, 6, 11, (the low, blue line in Figure 2) converges more quickly and to a better loss than other initialization strategies. We show that this result holds on the CNN and EN-RO datasets in Table 9.

6.2 When does Pseudo-Labeling help performance?

Table 10 shows results from fine-tuning teacher models on combinations of real labels and pseudo-labels. The Orig and Orig+PL columns show that, for summarization on XSUM, PL can improve over the SFT baseline when the pseudo-labels are added to the original fine-tuning dataset. For translation, (EN-RO), pseudo-labels can

Teacher	Student	MM Params	Time (MS)	Speedup	Rouge-2	Rouge-L
BART	12-1	222	743	2.35	17.98	33.31
	12-3	255	905	1.93	22.40	37.30
	6-6	230	1179	1.48	21.17	36.21
	9-6	268	1184	1.47	22.08	37.24
	12-6	306	1221	1.43	22.32	37.39
	Baseline (12-12)	406	1743	1.00	22.29	37.20
Pegasus	16-4	369	2038	2.40	23.18	38.13
_	16-8	435	2515	1.94	23.25	38.03
	Baseline (16-16)	570	4890		24.46	39.15
BertABS	Baseline (6-6)	110	1120		16.50	31.27

Table 6: Best XSUM results across all methods. Each sub-table is sorted fastest to slowest by inference time. dBART-12-3 and dPegasus-16-4 are trained on Pegasus pseudo-labels. dBART-12-6, dBART-6-6, and dBART-9-6 are trained with KD. dPegasus-16-8 and dBART-12-1 are trained with SFT. For the BART experiments where the encoder is smaller than 12 layers, we do not freeze it during training.

Teacher	Student	MM Params	Inference Time (MS)	Speedup	Rouge-2	Rouge-L
BART	12-3	255	1483	1.66	20.57	40.31
	6-6	230	1684	1.46	20.17	39.55
	12-6	306	1709	1.44	21.19	41.01
	Baseline (12-12)	406	2461	1.00	21.08	40.89
Pegasus	16-4	369	3728	2.67	21.29	40.34
	Baseline (16-16)	570	9965		21.37	41.04
BertABS	Baseline (6-6)	110	1582		19.6	39.18

Table 7: Best CNN/Daily Mail Results across all methods, which is always SFT.

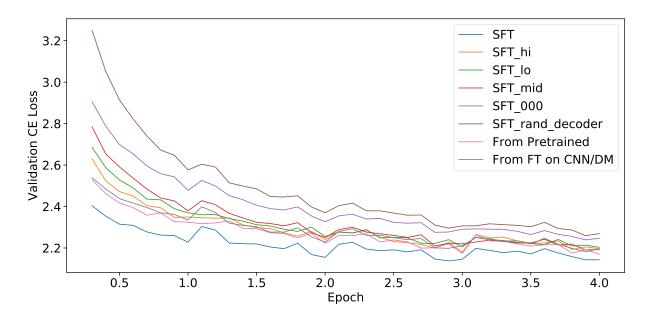


Figure 2: Training curves for different initialization strategies. Each line represents one fine-tuning run for a BART student on XSUM using a different initialization strategy. Initialization strategies are described in Table 8.

simply replace the original training data. For CNN (not shown), PL performs worse than SFT.⁷ For BART on XSUM, fine-tuning on the original dataset (SFT) generates a student that is 1.2 ROUGE-2 points worse than the teacher, fine-tuning on the original dataset and Pseudo-labels generates a better student, that is only 0.8 points behind the teacher. Adding Pseudo-labels generated by Pegasus, (the Orig+PL+PL column), generates a substantial improvement: the finetuned student is 0.1 points better than the teacher.

For Pegasus on XSUM, however, there is no benefit to adding pseudo-labels generated by BART. Comparing Orig+PL to Orig+PL+PL in Table 10 Row 2 shows that a student trained on the original data and Pegasus pseudo-labels is 1.2 ROUGE-2 below the teacher, whereas a student trained on the original data, Pegasus pseudo-labels, and BART pseudo-labels is 1.6 ROUGE-2 below the teacher.

The quality of the pseudo-labels may be driving this pattern. If we take the ROUGE-2 of pseudo-labels (against the training set labels) as proxy for their quality, the quality of the Pegasus pseudo-labels is 4 points higher than BART. Additionally, we did not find that pseudo-labels helped on CNN, where ROUGE scores are lower for both teachers, supporting the quality hypothesis.

⁷All pseudo-labels are made available for download here.

Name	Layers Copied	From	Min Loss
SFT	0,6,11	XSUM	2.14
SFT hi	9,10,11	XSUM	2.18
SFT lo	0,1,2	XSUM	2.20
SFT mid	5,6,7	XSUM	2.19
SFT 000	0,0,0	XSUM	2.24
SFT rand decoder	-	-	2.26
From Pre-trained	0,6,11	PT	2.17
From FT on CNN	0,6,11	CNN	2.18

Table 8: Loss for different initialization strategies on XSUM. Each row represents one fine-tuning run for a BART student on XSUM using a different initialization strategy (and one line in Figure 2.) LAYERS COPIED indicates which decoder layers were copied from the teacher. FROM indicates the BART model the layers were copied from, where XSUM is the BART teacher fine-tuned on the correct dataset, CNN is the teacher fine-tuned on the wrong dataset, and PT is the pre-trained (but not fine-tuned) BART checkpoint. MIN LOSS is cross entropy on the XSUM dev set. Validation loss was checked 10 times every epoch. This table corresponds to the figure above it.

Name	Layers Copied	From	Min Loss
SFT	0,6,11	CNN	2.02
From Pre-trained	0,6,11	PT	2.17
SFT hi	9,10,11	CNN	2.31
SFT rand	-	-	3.91
SFT	0, 2, 5	Marian	1.66
SFT Back	3,4,5	Marian	1.70
SFT Front	0,1,2	Marian	1.88
SFT rand	-	-	2.2

Table 9: Loss for different initialization strategies. See Table 8 for column descriptions. The top half of the table uses BART as a teacher and CNN as a dataset, the bottom half uses the fine-tuned Marian MT model as a teacher and EN-RO as a dataset. In the FROM column, CNN is BART fine-tuned on CNN, and PT is the pre-trained (but not fine-tuned) BART checkpoint, and Marian is the fine-tuned Marian MT checkpoint, which uses 6 encoder layers and 6 decoder layers.

6.3 Do Changes to \mathcal{L}_{kd} improve performance?

Except for BART on XSUM, KD did not generate improvements over SFT, and, as previously discussed, is always more expensive. This was not for lack of effort. Here are few modifications that did not improve performance:

- 1. Removing \mathcal{L}_{Hidn} , which encourages student layer l to produce the same hidden state as teacher layer ϕL , hurt performance for BART on XSUM. In the other settings, removing \mathcal{L}_{Hidn} had a negligible affect on performance.
- 2. Adding TinyBERT's \mathcal{L}_{Attn} , which encourages student layer l to produce the same attention weights as teacher layer $M_{-}\phi L$, further slowed training without improving performance. (Jiao et al., 2019)
- 3. Adding the cosine loss used in DistillBERT to \mathcal{L}_{kd} did not impact performance. (Sanh et al., 2019)

This suggests that more work is needed for adapting KD approaches that work on BERT to Seq2Seq tasks, and that practitioners should try SFT first, followed by pseudo-labeling.

6.4 INFERENCE TIME ANALYSIS

To further understand why the 6-6 models ran slower than 12-3 models in Tables 6 and 7, we ran a single forward pass on 12,000 different randomly initialized BART configurations in a GPU half-precision environment, and estimated the effects of changing the number of encoder layers, feed forward dimensions, number of decoder layers, and embedding size (width) on inference time with a linear regression. The results suggest that adding a decoder layer would slow down inference by 8%, while adding an encoder layer would slow down inference by

Teacher	Size	Dataset	Teacher Score	Orig	PL	Orig+PL	Orig+PL+PL*
BART Pegasus BART Pegasus	12-3 16-4 12-3 16-4	XSUM XSUM CNN CNN	22.3 24.5 21.1 21.37	-1.2 -1.9 -1.4 -0.1	-0.9 -2.2 -2.0 -1.4	-0.8 - 1.2 -2.0	+0.1 -1.6 -
Marian mBART	6-3 12-3	EN-RO EN-RO	27.7 26.5	-1.8 -0.8	-0.8 -0.4	-1.8 -0.6	- -

Table 10: Pseudo-labeling Strategies. Columns (Orig, PL, Orig+PL, and Orig+PL+PL*) report student scores relative to their teacher using (the original training data, pseudo-labels generated by the Teacher, both, and all pseudo-labels available for a given dataset + the original data). The score units are ROUGE-2 for the top four rows, BLEU for the two bottom rows, with the score for each student subtracted from the teacher score. All students are initialized by copying maximally spaced layers from the teacher and trained for 2 epochs.

only 4%. We also observed that changing width or feed forward dimensions had negliglible impact on run time. This difference is exacerbated during beam search, where the decoder is run beam_size times per example.

7 Conclusion

In this paper, we show that for summarization tasks, removing carefully chosen decoder layers from a Seq2Seq transformer and then continuing fine-tuning generates high quality student models quickly, and that in some situations more expensive training techniques with the same initialization strategy can generate additional quality improvements.

Future experiments could (1) evaluate these techniques on other summarization datasets, other tasks, and other teachers, like $T5^9$. (2) Explore distilling the knowledge in pre-trained, but not fine-tuned, Seq2Seq models. (3) Explore more of the large KD hyper-parameter space. (4) Explore strategies to improve pseudo-label quality. (5) Our experiments target speedups on GPU, but SqueezeBERT 10 suggests that reducing the width of each student layer is key to unlocking more efficient CPU inference. 11

⁸Sanh et al. (2019) found similar results with respect to the BERT architecture; (Kasai et al., 2020) found similar results for MT.

⁹(Raffel et al., 2020)

¹⁰(Iandola et al., 2020)

¹¹Discussion here

REFERENCES

- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, 2019.
- Canwen Xu, Wangchunshu Zhou, Tao Ge, Furu Wei, and Ming Zhou. Bert-of-theseus: Compressing bert by progressive module replacing, 2020.
- Forrest N. Iandola, Albert E. Shaw, Ravi Krishna, and Kurt W. Keutzer. Squeezebert: What can computer vision teach nlp about efficient neural networks?, 2020.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding, 2019.
- Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. Mobilebert: a compact task-agnostic bert for resource-limited devices. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020. doi: 10.18653/v1/2020.acl-main.195. URL http://dx.doi.org/10.18653/v1/2020.acl-main.195.
- 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, 2020.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2018.
- 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. *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, 2018. doi: 10.18653/v1/w18-5446. URL http://dx.doi.org/10.18653/v1/W18-5446.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv* preprint arXiv:1910.13461, 2019.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *ArXiv*, abs/1808.08745, 2018.
- Abigail See, Peter J. Liu, and Christopher D. Manning. Get to the point: Summarization with pointer-generator networks. *CoRR*, abs/1704.04368, 2017. URL http://arxiv.org/abs/1704.04368.
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization, 2019.
- Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *ArXiv*, abs/1503.02531, 2015.
- Yoon Kim and Alexander M. Rush. Sequence-level knowledge distillation. *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016. doi: 10.18653/v1/d16-1139. URL http://dx.doi.org/10.18653/v1/D16-1139.
- Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In KDD, 2006.
- Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with structured dropout, 2019.
- Jungo Kasai, Nikolaos Pappas, Hao Peng, James Cross, and Noah A. Smith. Deep encoder, shallow decoder: Reevaluating the speed-quality tradeoff in machine translation, 2020.
- Marcin Junczys-Dowmunt. Microsoft translator at WMT 2019: Towards large-scale document-level neural machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 225–233, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-5321. URL https://www.aclweb.org/anthology/W19-5321.
- Meng Sun, Bojian Jiang, Hao Xiong, Zhongjun He, Hua Wu, and Haifeng Wang. Baidu neural machine translation systems for WMT19. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 374–381, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-5341. URL https://www.aclweb.org/anthology/W19-5341.
- Yang Liu, Sheng Shen, and Mirella Lapata. Noisy self-knowledge distillation for text summarization, 2020a.
- Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference, 2017.

- Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/W16-2301. URL https://www.aclweb.org/anthology/W16-2301.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation, 2020b.
- Jörg Tiedemann and Santhosh Thottingal. OPUS-MT Building open translation services for the World. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation (EAMT)*, Lisbon, Portugal, 2020.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2020.