Task-agnostic Distillation of Encoderdecoder Language Models

Chen Zhang, Yang Yang, Quichi Li, Jingang Wang, Dawei Song



Language Model Distillation **Teacher-student Paradigm**

Teacher





Language model (LM) distillation aims at reducing inference compute by distilling the large LM into a small LM under a teacher-student paradigm.

Language Model Distillation Existing Methods

- Task-specific distillation with finetuning data (e.g., MRPC).
 - KD (Hinton, et al.)
 - MiniDisc (Zhang, et al.)
- Task-agnostic distillation with pretraining data (e.g., Wikipedia).
 - MiniLM (Wang, et al.): attention distribution-based
 - DistilGPT2 (Sanh et al.): logit distribution-based
- Task-agnostic distillation is commonly viewed as a better choice.

Failure of Distillation Encoder-decoder Language Models

encoder-decoder language models (e.g., T5).



In distilling a base-scale teacher to a 6-layer student, we have found that previous studies that are applicable either to encoder-only language models (e.g., BERT) or to decoder-only language models (e.g., GPT2) fail to handle

Encoder-decoder Interplay Gradient Perspective

 We suspect previous studies lack a component accounting for encoderdecoder interplay with only (MiniLM-like) logit distribution-based or previous implicit objectives, we propose explicit objectives that involve encoder-decoder interplay.

$$\begin{split} \mathcal{L}^{\mathsf{SelfAttn}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}) &= \mathbb{E}_{\mathbf{Z} \sim \mathcal{D}_{\mathbf{Z}}} \sum_{k=1}^{R} \\ \mathsf{KL}(\mathsf{Reln}(\mathbf{Z}; {}^{\mathcal{T}}\mathbf{W}_{k}^{\mathsf{Q}}), \mathsf{Reln}(\mathbf{Z}; {}^{\mathcal{S}}\mathbf{W}_{k}^{\mathsf{Q}})) \\ &+ \mathsf{KL}(\mathsf{Reln}(\mathbf{Z}; {}^{\mathcal{T}}\mathbf{W}_{k}^{\mathsf{K}}), \mathsf{Reln}(\mathbf{Z}; {}^{\mathcal{S}}\mathbf{W}_{k}^{\mathsf{K}})) \\ &+ \mathsf{KL}(\mathsf{Reln}(\mathbf{Z}; {}^{\mathcal{T}}\mathbf{W}_{k}^{\mathsf{V}}), \mathsf{Reln}(\mathbf{Z}; {}^{\mathcal{S}}\mathbf{W}_{k}^{\mathsf{V}})) \end{split}$$

 $\mathcal{L}^{\mathsf{CrossAttn}}(\mathcal{S};\mathcal{T},\mathcal{D}$

 $KL(Reln(\mathbf{Z}; ^{\mathcal{T}}\mathbf{W}_{k}^{O}))$ $+ \mathsf{KL}(\mathsf{Reln}(\mathbf{E}; \mathcal{T}))$ $+ \mathsf{KL}(\mathsf{Reln}(\mathbf{E}; \mathcal{T}))$

 $\operatorname{\mathsf{Reln}}(\mathbf{Z}; \mathcal{T}\mathbf{W}_{k}^{\mathsf{Q}})$

 $= \operatorname{Softmax}(\mathbf{Z}^{\mathcal{T}}\mathbf{W}_{k}^{\mathsf{Q}\mathcal{T}}\mathbf{W}_{k}^{\mathsf{Q}\mathcal{T}}\mathbf{Z}^{\mathcal{T}}/d^{\mathsf{R}}),$

(DistilGPT2-like) attention distribution-based distillation. Distinguishing from

$$egin{aligned} & \mathcal{D}_{\mathbf{Z}}, \mathcal{D}_{\mathbf{E}} \end{pmatrix} = \mathbb{E}_{\mathbf{Z} \sim \mathcal{D}_{\mathbf{Z}}, \mathbf{E} \sim \mathcal{D}_{\mathbf{E}}} \sum_{k=1}^{R} & \sum_{k=1}^{Q'} & \sum_{$$

$$\mathcal{L}^{\mathsf{Imp}} = \mathcal{L}^{\mathsf{Logit}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}) + \mathcal{L}^{\mathsf{SelfAttn}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}),$$

$$\begin{split} \mathcal{L}^{\mathsf{Exp}} = & \mathcal{L}^{\mathsf{Logit}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}) + \mathcal{L}^{\mathsf{SelfAttn}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}) \\ & + \mathcal{L}^{\mathsf{CrossAttn}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}, \mathcal{D}_{\mathbf{E}}) \end{split}$$

Encoder-decoder Interplay Gradient Perspective

stable.



• We have uncovered than implicit objectives would impose unstable gradient norms, thereby unstable training. In contrast, explicit objectives are much more





Encoder-Decoder Interplay MiniEnD

other is based on both encoder and decoder self-attention.

$$\begin{split} \mathcal{L}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}, \mathcal{D}_{\mathbf{E}}) &= \mathcal{L}^{\mathsf{Logit}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}) + \\ \mathcal{L}^{\mathsf{SelfAttn}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}) + \mathcal{L}^{\mathsf{CrossAttn}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}, \mathcal{D}_{\mathbf{E}}). \end{split}$$

 $\mathcal{L}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}, \mathcal{D}_{\mathbf{X}}) = \mathcal{L}^{\mathsf{Logit}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}) +$ $\mathcal{L}^{\mathsf{SelfAttn}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{Z}}) + \mathcal{L}^{\mathsf{EncSelfAttn}}(\mathcal{S}; \mathcal{T}, \mathcal{D}_{\mathbf{X}}),$



 Based on the observation, we put forward two implementations of the explicit objectives. One is based directly on encoder-decoder cross-attention, and the

Experiments Setup

- Distillation on C4 for T5, on OpenWebText for BART.
- Finetuning on GLUE (sequence and sequence-pair classification), CNN/ DailyMail and XSum (summarization).
- T5-based and BART-base as teachers.

Experiments GLUE

MiniEnD.

Method	GFLOP	s SST-2 Acc	MRPC F1	STS-B SpCorr	QQP F1	MNLI-m/mm Acc	QNLI Acc	RTE Acc	GLUE Score
T5 _{base}	25.4	<u>×</u> 94.6	93.0	90.0	88.9	86.7/86.8	92.9	74.7	88.5
T5 _{6L;384H}	3.18	92.2	90.2	86.0	87.3	81.2/81.7	88.2	70.0	84.6
MiniDisc _{5%} ^①	7.80	9 3.8	89.8	85.3	86.7	82.9 /82.7	89.2	64.6	84.4
MImKD _{6L;384H}	3.18	92.3	88.7	86.2	87.5	81.6/82.1	88.2	67.9	84.3
MiniLM _{6L;384H}	3.18	92.1	89.6	85.2	87.0	81.2/81.5	88.0	68.6	84.1
MImKD+MiniLM _{6L;384H}	3.18	92.4	89.2	86.0	87.3	81.7/82.1	89.1	67.9	84.5
$\begin{array}{l} MINIEND-D_{6L;384H} \\ w/o \ \mathcal{L}^{Logit} \\ MINIEND-E_{6L;384H} \\ w/o \ \mathcal{L}^{Logit} \end{array}$	3.18	92.1	90.6	85.8	87.7	81.8/82.3	89.0	68.6	84.7
	3.18	92.2	90.1	86.6	87.6	82.2/82.8	89.1	68.6	84.9
	3.18	92.7	90.0	86.1	87.4	81.8/82.1	88.8	69.3	84.8
	3.18	92.3	89.9	86.6	87.7	82.5/ 83.1	89.2	69.0	85.0

^① MiniDisc is distilled from T5_{xlarge}, and owns larger GFLOPs.

The distillation cannot surpass pretraining-from-scratch baseline on GLUE until

Experiments **CNN/DailyMail & XSum**

Similar results are observed on CNN/DailyMail and XSum.

Method	GFLOPs		CN	N/Dailyl	Mail	XSum			
			Rg-1	Rg-2	Rg-L	Rg-1	Rg-2	Rg-L	
T5 _{base}	25.4	1	40.1	19.4	31.5	34.7	12.4	29.7	
T5 _{6L:384H}	3.18		35.7	16.8	28.4	28.6	8.9	24.8	
MImKD _{6L:384H}	3.18	×	36.0	17.0	28.7	28.9	9.2	25.0	
MiniLM _{6L;384H}	3.18	ŵ	35.0	16.5	28.0	25.9	7.5	22.5	
MImKD+MiniLM _{6L;384H}	3.18		35.8	17.0	28.7	29.0	9.1	25.1	
МіліЕлD-D _{6L;384Н}	3.18		36.2	17.2	28.9	29.5	9.2	25.4	
w/o \mathcal{L}^{Logit}	3.18	×	35.7	17.0	28.6	27.3	8.2	23.7	
MiniEnD-E _{6L;384H}	3.18	8	36.1	17.3	28.9	28.9	9.1	24.9	
w/o \mathcal{L}^{Logit}	3.18		35.8	17.1	28.7	27.2	8.0	23.6	
BART _{base}	12.7	1 ×	39.4	18.5	30.6	36.9	14.7	31.9	
LogitKD _{3/11.768H} ^①	4.23	\times	38.0	16.0	25.2	32.9	12.4	26.9	
DQ-BART _{8bit} ®	12.7	1~3	42.4	19.3	28.8	38.2	15.7	30.7	
MiniEnD-D _{6L;384H}	3.18	4 ×	38.5	18.5	29.7	33.6	12.9	29.2	

^① LogitKD is distilled with an asymmetric layer setting, i.e., more encoder layers than decoder layers, for saved performance decline.

⁽²⁾ DQ-BART only quantizes parameter precision to lower one, i.e., 8 bit, but does not reduce parameter amount. Quantization would not give any speedup in GFLOPs though nice reduction in model size.

Experiments Data Scaling

• MiniEnD is also applicable when half data is used.





Experiments **Model Scaling**

scaled up.

distillation and $\dots \Rightarrow \{\dots \Rightarrow \dots\}$ indicates progressive distillation.

	GLUE	CNN/DailyMail			XSum			
Methoa	Score	Rg-1	Rg-2	Rg-L	Rg-1	Rg-2	Rg-L	
T5 _{6L;384H}	84.6	35.7	16.8	28.4	28.6	8.9	24.8	
T5 _{12L;384H}	85.0	37.2	17.9	29.6	31.2	10.5	27.0	
T5 _{base}	88.5	40.1	19.4	31.5	34.7	12.4	29.7	
T5 _{large}	90.7	40.6	19.4	31.7	38.2	15.1	32.9	
T5 _{xlarge}	92.0	40.8	19.7	32.1	41.1	17.6	35.5	
T5 _{base} ⇒T5 _{6L;384H}	84.7	36.2	17.2	28.9	29.5	9.2	25.4	
T5 _{large} ⇒T5 _{6L;384H}	84.5	36.4	17.4	29.0	29.4	9.3	25.3	
$T5_{xlarge} \Rightarrow T5_{6L;384H}$	84.2	36.1	17.2	28.8	29.1	9.1	25.1	
$T5_{xlarge} \Rightarrow T5_{12L;384H} \Rightarrow T5_{6L;384H}$	84.6	36.6	17.5	29.2	29.2	9.1	25.1	
T5 _{large} ⇒T5 _{12L;384H}	85.5	38.3	18.4	30.4	32.4	11.2	27.9	
T5 _{xlarge} ⇒T5 _{12L;384H}	85.2	38.0	18.4	30.3	32.2	11.1	27.7	
$T5_{xlarge} \Rightarrow \{T5_{large} \Rightarrow T5_{12L;384H}\}$	85.8	38.4	18.5	30.6	32.9	11.5	28.3	

Special methods should be attached to MiniEnd (e.g., TA) when teacher is

Table 5: The results of model scaling using MINIEND-D. \Rightarrow denotes a distillation step, which should be operated sequentially otherwise $\{\}$ is prioritized. $\dots \Rightarrow \dots \Rightarrow \dots$ indicates teacher assistant-based

Conclusion

- and two abstractive summarization datasets.
- We further scale the distillation of encoder-decoder LMs to a 3B teacher that requires efforts to exploring how large language models can be distilled.
- arXiv: <u>https://arxiv.org/abs/2305.12330</u>
- GitHub: <u>https://github.com/GeneZC/MiniEnD</u>
- Slides: <u>https://genezc.github.io/assets/files/COLING2024_MiniEnD.pdf</u>
- Thank you all!

• We find through a pilot study that the encoder-decoder interplay is a key component that should be aligned in the distillation so that the distilled encoder-decoder LMs are promising. Based on the idea, we propose two directions that the encoder-decoder interplay alignment can be incorporated and verify their effectiveness on a language understanding benchmark

additional distillation steps. In this sense, we recommend future research to devote more