

Impact of Transformer Depth on Small Datasets for Chatbot Models

Anonymous EMNLP submission

Abstract

This study investigates the impact of transformer layer variations (2, 4, 8, and 16 layers) on the performance and efficiency of chatbot models trained on a small curated dataset. From the original DailyDialog dataset containing 76,052 samples, 10,000 instances were selected after preprocessing, including duplicate removal and filtering by token count (3–60). Evaluation metrics include BLEU, ROUGE, METEOR, and qualitative assessments. Results show that the 2-layer model balances efficiency and semantic performance during early training. The 16-layer model performs best in BLEU but offers diminishing returns in other metrics. These findings emphasize the importance of optimizing transformer depth for small datasets in low-resource scenarios.

1 Introduction

Modern NLP models, particularly transformers, have demonstrated exceptional performance across tasks with large-scale datasets. However, challenges arise when applying such architectures to small datasets. This study evaluates the impact of varying transformer depth (2, 4, 8, and 16 layers) on model performance and computational efficiency in a small-data setting.

2 Background and Related Works

TBU

3 Method

3.1 Dataset Preprocessing

The experiments were conducted using the DailyDialog dataset, a widely-used datasets for dialog generation tasks. The original dataset consists of 76,052 multi-turn conversational samples. To dapt the data for this study, the following preprocessing steps were applied:

- 1. Duplicate Removal:** All duplicate samples were removed to ensure diversity in the training data and minimize inference of perplexity in model training.
- 2. Token-length Filtering:** Dialog samples were filtered to retain only those with token counts between 3 and 60, ensuring meaningful yet manageable input lengths. 95% of total training datasets are included in this range.
- 3. Sample Selection:** A subset of 10,000 instances was randomly selected from the pre-processed dataset to simulate a small-scale training scenario.
- 4. Tokenization:** Tokenization was performed using the NLTK punkt tokenizer and a vocabulary size of 15,000 tokens was constructed.

3.2 Transformer Model Configuration

The configuration of the model follows the standard transformer design with a slight modification in embedding dimension size and variation in encoder and decoder layer numbers, ranging from 2, 4, 8 to 16.

- **Embedding Dimension (d_{model}):** 256
The weights of the embedding layers in the encoder and decoder are shared, since the source and target languages are the same (English).
- **Number of Attention Heads (n_{heads}):** 8.
- **Feedforward Layer Dimension (d_{ff}):** 1024.
- **Dropout Rate:** 0.1.
- **Positional Encoding:** Maximum sequence length of 128 tokens.

3.3 Training Setup

The models were trained using the Adam optimizer with the following hyperparameters:

071	• Epochs per model: 10	2. Convergent Speed:	113
072	• Learning Rate Scheduler: Warm-up steps of	• The 2-layer model achieved competitive	114
073	4,000, followed by inverse square root decay.	performance in early epochs (e.g., within	115
074	• Optimizer Parameters:	10 epochs), particularly in METEOR	116
075	$\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1e - 9$	and ROUGE, which suggests that shal-	117
076	• Batch Size: 64	lower models may balance efficiency and	118
077	• Loss Function: Sparse Categorical Crossen-	performance during early-stage training	119
078	tropy, with a padding mask applied to ig-	with small dataset.	120
079	gnore <pad> tokens in the input and target se-	• Deeper models, especially the 16-layer	121
080	quences.	model, require more epochs to converge	122
081	3.4 Evaluation Metrics	effectively, as indicated by their higher	123
082	To assess model performance and training effi-	loss values.	124
083	ciency, following evaluations were performed:	• Yet, the fluctuation of model perfor-	125
084	1. Quantitative Metrics for Text Generation	mance observed within the 10 epochs	126
085	The 3 following evaluation metrics have	indicates that the training epochs set for	127
086	scores between 0 to 1, and score above 0.2	the experiment may not be sufficient for	128
087	is suggested for generative models.	all models.	129
088	• BLEU: Measures n-gram overlap be-	4.2 Model Performance	130
089	tween generated responses and ground	1. BLEU:	131
090	truth.	• Across all epochs, the 2-layer model con-	132
091	• ROUGE: Evaluates overlap of longer se-	sistently recorded a BLEU score of 0,	133
092	quences, and ROUGE-L was used for the	indicating its inability to generate text	134
093	experiments.	with meaningful n-gram overlap.	135
094	• METEOR: Consider semantic similar-	• Among the other models, the ranking of	136
095	ity, penalizing over-reliance on exact	BLEU performance was $16 > 8 > 4 >$	137
096	matches	2, with the 16-layer model achieving the	138
097	2. Computational Efficiency Metrics:	highest BLEU score.	139
098	• Training Time Per Epoch: Time taken to	• Yet, BLEU scores of all models exceed-	140
099	complete one training epoch	ing low (near-zero values), suggesting	141
100	• Convergence Speed: Number of epochs	that BLEU might not be a reliable metric	142
101	required to achieve optimal performance	for assessing the quality of responses in	143
102	4 Result	this experiment.	144
103	4.1 Training Efficiency	2. ROUGE:	145
104	1. Time per Epoch: "Image"	• The ranking for ROUGE performance	146
105	• The 2-layer model exhibited the fastest	was $8 > 2 > 16 > 4$, with the 8-layer	147
106	training time (91 seconds per epoch),	model achieving the highest ROUGE	148
107	making it highly efficient for small	score. However, the 2-layer model per-	149
108	datasets.	formed better than deeper models like 16	150
109	• Deeper models (e.g. 16-layer) had sig-	and 4 layers in ROUGE.	151
110	nificant higher training times (658 sec-	• Significant fluctuations in ROUGE	152
111	onds per epoch), increasing computa-	scores were observed within the first 10	153
112	tional costs substantially.	epochs, suggesting instability in captur-	154
		ing sequence-level overlap during train-	155
		ing.	156
		• Similar to BLEU, ROUGE values were	157
		all below 0.08, indicating that none of	158
		the models could consistently generate	159
		meaningful textual content.	160

Table 1: Performance Metrics by Model Layer Depth

Layer	BLEU Avg	ROUGE Avg	METEOR Avg	Observations
2	0	0.0737	0.0413	Strong semantic similarity and sequence overlap in early training.
4	$4.58e - 157$	0.0330	0.0433	Moderate performance; outperformed by 2-layer in METEOR and ROUGE.
8	$2.55e - 157$	0.0276	0.0427	Best ROUGE performance but struggled in METEOR.
16	$8.90e - 157$	0.0248	0.0434	Best BLEU performance but less semantic relevance.

3. METEOR:

- The 2-layer model outperformed all other models in METEOR at maximum, followed by $4 > 16 > 8$, showcasing its ability to capture similarity better than deeper models in the early stages of training.
- Like BLEU and ROUGE, all METEOR scores were below 0.08, further supporting the conclusion that the generated text lacked meaningful content.

- Across all layers, BLEU, ROUGE, and METEOR scores remained below 0.08, indicating the difficulty of generating meaningful text with the given dataset size and experimental setup.

In conclusion, while deeper models such as the 16-layer transformer excel in BLEU, the results highlight the suitability of shallow models (e.g., 2 layers) for efficient and effective training in low-resource scenarios. For more balanced performance, the 8-layer configuration offers the best trade-off between quality and computational cost. Future work could explore techniques like data augmentation or transfer learning to improve the overall performance of transformers on small datasets.

5 Conclusion

This study analyzed the impact of transformer layer depth on model performance and training efficiency for a chatbot model trained on a small dataset. Experiments with models having 2, 4, 8, and 16 layers revealed several key insights:

- The **2-layer model** demonstrated strong semantic similarity and sequence overlap in the early stages of training, as reflected by the highest ROUGE and METEOR scores among shallow models. Its computational efficiency makes it well-suited for low-resource or early-training scenarios.
- The **8-layer model** achieved the best overall balance of performance and efficiency, with strong ROUGE scores and moderate METEOR values. This configuration is recommended for generating more stable results over prolonged training on small datasets.
- The **16-layer model** performed best in terms of BLEU score but showed limited semantic relevance and higher computational costs, making it less practical for small datasets or limited epochs.

6 Acknowledgements

TBU

7 References

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. *Attention is all you need*. Advances in neural information processing systems, 30, 2017.