

Exploring Versatile Generative Language Model Via Parameter-Efficient Transfer Learning

Zhaojiang Lin*, Andrea Madotto*, Pascale Fung

Center for Artificial Intelligence Research (CAiRE)

Department of Electronic and Computer Engineering

The Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong

{zlinao, amadotto}@connect.ust.hk,

pascale@ece.ust.hk

Abstract

Fine-tuning pre-trained generative language models to down-stream language generation tasks has shown promising results. However, this comes with the cost of having a single, large model for each task, which is not ideal in low-memory/power scenarios (e.g., mobile). In this paper, we propose an effective way to fine-tune multiple down-stream generation tasks simultaneously using a single, large pre-trained model. The experiments on five diverse language generation tasks show that by just using an additional 2-3% parameters for each task, our model can maintain or even improve the performance of fine-tuning the whole model¹.

1 Introduction

Large-scale language models (Radford et al., 2019; Dai et al., 2019) have shown to be effective in learning highly transferable embedding, which can be used in several down-stream tasks. For instance, bidirectional models (Peters et al., 2018; Devlin et al., 2019) are fine-tuned to improve classification tasks (Wang et al., 2019), while, unidirectional language models (Radford et al., 2019) are more effective in language generation tasks. In this work, we focus on the latter, and show that it is possible to dynamically steer the output of a language model (e.g., GPT-2) towards a specific task (e.g., summarization) without modifying the original model parameters.

Feature-based transfer (Howard and Ruder, 2018; Fan et al., 2020a,b) and fine-tuning (Devlin et al., 2019) are the most commonly used methods for transfer learning of a language. The former freezes the pre-trained model and uses it as a feature extractor for training a new classifier, and the

latter uses the pre-trained weight as a weight initialization for the model to be trained for downstream tasks. The feature-based transfer strategy has not shown promising results (Devlin et al., 2019), while fine-tuning, on the other hand, can achieve state of the art performance in multiple tasks (Dong et al., 2019). However, the downside of the latter is the need for a separate model for each of the fine-tuned tasks. This is especially relevant for on-device applications, where a limited amount of computation/memory is available.

Therefore, we study how to effectively use a single pre-trained model as the backbone for multiple language generation tasks, such as conversational question answering, summarization, machine translation, multi-turn chat dialogue, and task-oriented natural language generation. This is a particular parameter-sharing schema, where we constrain the shared parameters to be the ones in the pre-trained model, and we learn task-specific parameters for each of the considered datasets.

In this paper, we propose to use residual adapter layers (Houlsby et al., 2019) and task embeddings for modelling the aforementioned task-specific parameters, and we explore different training strategies such as distillation (Hinton et al., 2015; Kim and Rush, 2016). We also analyse the trade-off between freezing or not freezing the language model parameters by leveraging two learning settings, multi-task (MT) (Caruana, 1997) and continual learning (CL) (Thrun and Pratt, 2012). With our experiments, we empirically demonstrate that by adding less than 3% task-specific parameters, our model can maintain or even achieve better performance than fine-tuning the whole model.

2 Related work

Pre-trained generative language models (Radford et al., 2019, 2018; Dai et al., 2019; Yang et al.,

* Equal contributions.

¹Code available in <https://github.com/zlinao/VGLM>

2019; Peters et al., 2018) have shown to be very effective in language generation, whereas, bidirectional pre-trained models (Devlin et al., 2019; Liu et al., 2019; Sanh et al., 2019) significantly improve the performance of several down-stream classification tasks. Fine-tuning large pre-trained models has shown positive results in dialogue tasks (Wolf et al., 2019b; Budzianowski and Vulić, 2019) and other language generation tasks (Dong et al., 2019). However, all of the previous works only consider fine-tuning on each generation task individually, which requires a separate model for each task. In this work, we use only a **single** model, for multiple generation tasks.

Residual adapters, derived from residual networks (He et al., 2016), were first introduced by Rebuffi et al. (2017) for multiple visual domain learning. Houlsby et al. (2019) proposed low-rank residual adapters to improve the scalability of the adapter module, and effectively transfer BERT (Devlin et al., 2019) to multiple text classification tasks simultaneously, while Bapna and Firat (2019) applied an adapter layer to language/domain adaptation for neural machine translation. On the other hand, Dathathri et al. (2019) proposed a plug and play method to control the language model generation without finetuning the model. Differently, in this paper, we extend the idea of adapters to a large variety of language generation tasks, which has not been considered before, and we compare the idea of a fixed pre-trained back-bone for continual learning with multi-task training (Stickland and Murray, 2019).

3 Methodology

The Versatile Language Model (VLM) is composed of three components: a pre-trained language model back-bone (e.g., GPT-2), and two kinds of specialized parameters for each generation tasks such as low-rank residual adapters and task embedding. Figure 1 shows the VLM architecture with the specialized parameters in different colours.

Residual Adapters These are trainable modules which steer the pre-trained model to different down-stream tasks. We adapt the design of the feed-forward Transformer sub-layer following Bapna and Firat (2019). To elaborate, the adapter block consists of 1) a layer normalization (Ba et al., 2016) for an efficient adaptation and 2) a following autoencoder (Hinton and Zemel, 1994), with a residual connection. Formally, given the hidden rep-

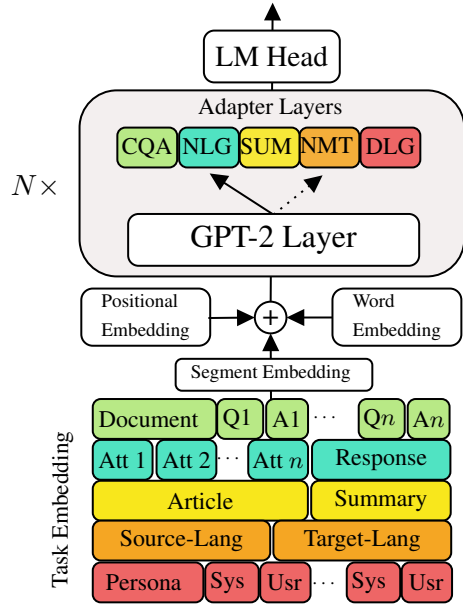


Figure 1: Simplified illustration of the Versatile Language Model. A detailed illustration is reported in Appendix A1.

resentation $H_i \in \mathbb{R}^{t \times d}$ from the language model layer i , where d is the hidden dimension and t is the current generation step, the residual adapter computes the following:

$$\mathbf{Adapter}(H_i) = (\mathbf{ReLU}(\mathbf{LN}(H_i)W_i^E))W_i^D + H_i$$

where W_i^E and W_i^D are parameters of dimension $d \times m$ and $m \times d$ respectively, and $\mathbf{LN}(\cdot)$ denotes layer normalization. The bottleneck dimension m is tunable and it allows to adjust the capacity of the adapter according to complexity of the target task.

Task Embedding. To adapt unconditional generative language models to different conditional language generation tasks (e.g., CoQA, Summarization), we construct a set of task-specific segment embeddings. For example, in multi-turn dialogue, we alternate between *System* and *User* embeddings to help the model to capture the hierarchical structure of dialogues. Figure 1 shows the task embedding for each task, and more details are available in Appendix A2.

Knowledge Distillation In tasks with a large distributional shift from the original pre-trained language model (e.g., Machine Translation), we expect a larger performance gap between VLM and full fine-tuning. To cope with this issue, we propose to use sentence-level knowledge distillation (Kim and Rush, 2016), to help the task-specific parameters to better adapt to the task. Specifically,

	Param.	Persona (DLG)		NMT	SUM	CoQA	NLG	
		ppl. ↓	BLEU ↑	BLEU ↑	ROUGE 2 ↑	F1 ↑	BLEU ↑	AVG ↑
GPT-2 Finetune	5×	13.13	2.17	25.45	18.1	67.7	66.4	57.77
w/o Pre-Train	5×	37.77	0.99	16.52	17.0	15.1	60.5	53.51
w/o Task Emb.	5×	13.24	0.00	0.61	15.0	35.2	53.1	47.25
LM Head	2.55×	17.58	1.34	12.05	15.8	47.0	65.2	55.25
VLM MT	1.13×	13.15	0.84	22.49	17.7	69.3	65.6	57.08
VLM	1.13×	14.06	1.99	24.19*	18.0*	66.2	67.1	57.97
w/o Task Emb.	1.13×	14.31	0.00	0.95	15.0	32.2	58.3	50.99
Reference	-	38.08 [¶]	-	29.2 [§]	17.20 ^{¶¶}	45.4 ^{††}	65.9 ^{‡‡}	57.54
SOTA	-	17.51 [†]	-	35.2 [‡]	21.53 ^{§§}	82.5	66.2 ^{‡‡}	57.44

Table 1: Results of VLM versus other fine-tuning techniques on the five evaluated datasets. Param. refers to the number of parameters that need to be stored after training. We use the adapter with distillation* for translation and summarization. The Reference and SOTA results are: Profile Memory[¶](Zhang et al., 2018), TransferTransfo[†] (Wolf et al., 2019b), DynamicConv[‡](Wu et al., 2019), Transformer[§](Vaswani et al., 2017), PG^{¶¶} (See et al., 2017), T5-11B^{§§}(Raffel et al., 2019), UniLM^{||}(Dong et al., 2019), PG^{††} (Reddy et al., 2019) and SOTA system^{‡‡} in Dušek et al. (2019)

we first fully fine-tune a GPT-2 model on the training set of a task (e.g., Machine Translation). Then we replace the gold target (e.g., gold translation) in the training set with the greedy decoded output from the full fine-tuned model. Finally, the new constructed training set is used to fine-tune the student VLM.

4 Experiments

4.1 Datasets & Evaluation Metrics

We conduct our experiment on five diverse datasets covering multiple generation tasks: Persona-Chat (Zhang et al., 2018; Dinan et al., 2019) for chit-chat based dialogue (DLG), IWSLT (Cettolo et al., 2016) German-English neural machine translation (NMT), CNN/Daily-Mail (Hermann et al., 2015; Nallapati et al., 2016) for text-summarization (SUM), CoQA (Reddy et al., 2019) for generative conversational question answering (CQA), and E2E NLG-challenge (Dušek et al., 2019) for task-oriented natural language generation (NLG).

We use a large variety of evaluation metrics, such as perplexity, F1 score, BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), NIST (Lin and Och, 2004), METEOR (Denkowski and Lavie, 2014) and CiDER (Vedantam et al., 2015). Each task uses the appropriate measure, as reported in Table 1, where in NLG we report the normalized average score of multiple metrics, as in Dušek et al. (2019). More information about task description and the metrics used in each task are reported in Appendix A2.

4.2 Implementation and model comparison

We implement VLM based on GPT-2-small (124M) (Wolf et al., 2019a), and experiment with varying adapter bottleneck dimensions in {10, 50, 100, 300} and pick the best one in each task to trade-off the performance with the parameter efficiency. Specifically, we choose bottleneck sizes 100, 300, 100, 300 and 10 for DLG, NMT, SUM, QA, and NLG, respectively, which results in 13% additional parameters in total. We ablate the adapter training with and without knowledge distillation and task embeddings. We also test the performance of a frozen back-bone (VLM) to show the ability to continuously learn tasks, and multi-task fine-tune (VLM MT) with a trainable backbone to show possible positive transferring among tasks as in Stickland and Murray (2019). More training details and the dataset pre-processing are reported in Appendix A2.

To show the effectiveness of the proposed methodology of learning a versatile generative model, we compare (i) fine-tuning the whole GPT-2 model for each task separately (GPT-2 Finetune), (ii) fine-tuning the language model head of GPT-2 for each task (LM-Head), (iii) existing baseline models reported (Reference), and (iv) the state-of-the-art models for all the tasks (SOTA).

4.3 Results and Analysis

Table 1 shows the experimental results of the aforementioned models. Appendix A3 and A4 report detailed results and generated samples for all the

datasets. Our findings can be summarized as follows:

Fine-tuning GPT-2 vs Baseline & SOTA. Fine-tuning the whole GPT-2-small in each task can generally improve on the performance of competitive baselines such as Pointer-Generator (See et al., 2017) in summarization (SUM) and CoQA. In both the Persona-Chat and the NLG tasks GPT-2 fine-tuning slightly outperforms the current SOTA, whereas, we observe a performance gap between GPT-2 and SOTA in NMT and SUM. Notably, the advantage of GPT-2 pre-training is limited in NMT: 1) no or little German text is present in the pretraining corpus; 2) the GPT-2 BPE (Sennrich et al., 2016) tokenizer is optimized for English text, and not for multiple languages. Finally, in SUM and CoQA, the SOTA models use $100\times$ bigger models (Raffel et al., 2019) and bidirectional attention (Dong et al., 2019), where instead, GPT-2 uses unidirectional attention.

Adapter vs Fine-tuning GPT-2 & LM Head. Fine-tuning only the adapter layers introduces 13% additional parameters to the model with a minimal loss in performance (0.4%) compared to fine-tuning a separate GPT-2 model. Moreover, the adapter layers are more effective, both in terms of performance and number of additional parameters, compared to fine-tuning LM-Head.

Knowledge Distillation (KD) Using KD in the training procedure is especially useful in tasks such as NMT and SUM, where the gap between fine-tuning the whole model and adapter is large. This is because KD reduces the complexity of training targets (by replacing the gold training target with a teacher model generated target), which helps with low-capacity adapter (with 4% parameter) by providing an easier translation/summarization task (Zhou et al., 2019). Figure 2 shows the effect of using distillation training when the gap with the full fine-tuning is more substantial. On the other hand, when the adapter performance is very close to that of the fine-tuning baseline, or better (i.e. NLG), distillation has no impact on the final performance.

Task Embedding The specialized segment embedding (a.k.a. task embedding) is very important for achieving competitive performance, independently of the adapter. In Table 1, we can observe a substantial drop in performance when the task

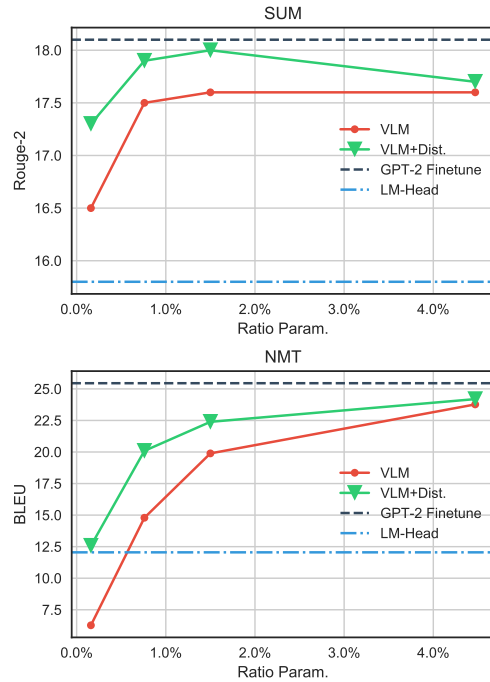


Figure 2: Performance comparison among different ratios of additional parameters for the SUM and NMT tasks.

embedding is not deployed. Indeed, without a task embedding the model struggles to learn when the input sequence ends, and how to distinguish the different parts of the input sequence (e.g., attributes in NLG, document and question in CoQA etc.).

Frozen Backbone vs Trainable Backbone As previously mentioned, VLM can be trained either by freezing the weight of the GPT-2 model, i.e., independently and continually learning one task at a time, or by multitasking all the tasks and thus fine-tuning both the GPT-2 model and the adapters. The latter model has the advantage of being able to transfer knowledge among tasks, as we can observe in Table 1 for the CoQA task, where VLM Multi-Task improve the F1 score by 3%. On the other hand, the frozen back-bone model has the big advantage of learning tasks sequentially, since the original GPT-2 weights remain untouched.

5 Conclusion

In this paper, we have presented a Versatile Language Model which learns five diverse natural language generation tasks in a single model. We found that a residual adapter is more effective than fine-tuning other parts of the model (e.g., LM-Head), and that distillation helps in reducing the gap in

performance in hard to fine-tune tasks, such as summarization and translation. Finally, we show the trade-off between a frozen and trainable back-bone, showing that the former has a competitive performance, with the advantage of being extendable to future tasks without full re-training.

References

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450*.
- Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2016. An actor-critic algorithm for sequence prediction. *arXiv preprint arXiv:1607.07086*.
- Ankur Bapna and Orhan Firat. 2019. Simple, scalable adaptation for neural machine translation. 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 1538–1548.
- Paweł Budzianowski and Ivan Vulić. 2019. Hello, it’s gpt-2—how can i help you? towards the use of pre-trained language models for task-oriented dialogue systems. *arXiv preprint arXiv:1907.05774*.
- Rich Caruana. 1997. Multitask learning. *Machine learning*, 28(1):41–75.
- Mauro Cettolo, Niehues Jan, Stüker Sebastian, Luisa Bentivogli, Roldano Cattoni, and Marcello Federico. 2016. The iwslt 2016 evaluation campaign. In *International Workshop on Spoken Language Translation*.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. [Transformer-XL: Attentive language models beyond a fixed-length context](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988, Florence, Italy. Association for Computational Linguistics.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2019. Plug and play language models: a simple approach to controlled text generation. *arXiv preprint arXiv:1912.02164*.
- Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the EACL 2014 Workshop on Statistical Machine Translation*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, et al. 2019. The second conversational intelligence challenge (convai2). *arXiv preprint arXiv:1902.00098*.
- Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. Unified language model pre-training for natural language understanding and generation. *arXiv preprint arXiv:1905.03197*.
- Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2019. [Evaluating the state-of-the-art of end-to-end natural language generation: The E2E NLG Challenge](#). *arXiv preprint arXiv:1901.11528*.
- Yingruo Fan, Jacqueline CK Lam, and Victor On Kwok Li. 2020a. Facial action unit intensity estimation via semantic correspondence learning with dynamic graph convolution. In *AAAI*, pages 12701–12708.
- Yingruo Fan, Victor Li, and Jacqueline CK Lam. 2020b. Facial expression recognition with deeply-supervised attention network. *IEEE Transactions on Affective Computing*.
- Sebastian Gehrmann, Yuntian Deng, and Alexander Rush. 2018. Bottom-up abstractive summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4098–4109.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Tianyu He, Xu Tan, Yingce Xia, Di He, Tao Qin, Zhibo Chen, and Tie-Yan Liu. 2018. Layer-wise coordination between encoder and decoder for neural machine translation. In *Advances in Neural Information Processing Systems*, pages 7944–7954.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in neural information processing systems*, pages 1693–1701.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Geoffrey E Hinton and Richard S Zemel. 1994. Autoencoders, minimum description length and helmholtz free energy. In *Advances in neural information processing systems*, pages 3–10.

- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799.
- Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 328–339.
- Po-Sen Huang, Chong Wang, Sitao Huang, Dengyong Zhou, and Li Deng. 2017. Towards neural phrase-based machine translation. *arXiv preprint arXiv:1706.05565*.
- Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Chin-Yew Lin and Franz Josef Och. 2004. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 605. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Andrea Madotto, Zhaojiang Lin, Chien-Sheng Wu, and Pascale Fung. 2019. Personalizing dialogue agents via meta-learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5459.
- Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Caglar Gulcehre, and Bing Xiang. 2016. Abstractive text summarization using sequence-to-sequence rnns and beyond. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 280–290.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of NAACL-HLT*, pages 2227–2237.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. URL https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language_understanding_paper.pdf.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. In *ArXiv*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. *arXiv preprint arXiv:1511.06732*.
- Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. 2017. Learning multiple visual domains with residual adapters. In *Advances in Neural Information Processing Systems*, pages 506–516.
- Siva Reddy, Danqi Chen, and Christopher D Manning. 2019. Coqa: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725.
- Asa Cooper Stickland and Iain Murray. 2019. Bert and pals: Projected attention layers for efficient adaptation in multi-task learning. In *International Conference on Machine Learning*, pages 5986–5995.
- Sebastian Thrun and Lorien Pratt. 2012. *Learning to learn*. Springer Science & Business Media.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.

- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *International Conference on Learning Representations*.
- Yijun Wang, Yingce Xia, Li Zhao, Jiang Bian, Tao Qin, Guiquan Liu, and Tie-Yan Liu. 2018. Dual transfer learning for neural machine translation with marginal distribution regularization. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R’emi Louf, Morgan Funtowicz, and Jamie Brew. 2019a. Huggingface’s transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.
- Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. 2019b. [Transfertransfo: A transfer learning approach for neural network based conversational agents](#). *CoRR*, abs/1901.08149.
- Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. 2019c. Transfertransfo: A transfer learning approach for neural network based conversational agents. *arXiv preprint arXiv:1901.08149*.
- Felix Wu, Angela Fan, Alexei Baevski, Yann N Dauphin, and Michael Auli. 2019. Pay less attention with lightweight and dynamic convolutions. *arXiv preprint arXiv:1901.10430*.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*.
- Semih Yavuz, Abhinav Rastogi, Guan-Lin Chao, and Dilek Hakkani-Tur. 2019. Deepcopy: Grounded response generation with hierarchical pointer networks. *arXiv preprint arXiv:1908.10731*.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213.
- Chunting Zhou, Graham Neubig, and Jiatao Gu. 2019. Understanding knowledge distillation in non-autoregressive machine translation. *arXiv preprint arXiv:1911.02727*.

A Appendices

A.1 Model details

Figure 3 illustrates a detailed version of VLM. VLM shares a GPT-2 back-bone and for each task, the model looks up a set of task embeddings for modeling different input structures and chooses the corresponding adapter.

A.2 Experiment details

In this section, we will describe the dataset, evaluation metrics, dataset preprocessing and training details for each task.

Conversational Question Answering (CQA)

CoQA (Reddy et al., 2019) is a free-form conversational question answering dataset. The task is to answer the questions in a conversation. Each turn in the conversation contains a question, and we need to answer the questions based on conversation histories and documents. We use *document*, *question*, and *answer* segment embedding to help the model to distinguish the document and alternating questions and answers in the input sequence. We fine-tune the full GPT2-small or VLM (trainable adapter with a fixed GPT2-small) for five epochs with the Adam optimizer. For distillation we only fine-tune VLM for three epochs. We set the batch size to 16 and limit the maximum length of the document to 400 tokens and only retain the last two turns of questions and answers in the dialogue history. Following Reddy et al. (2019) we use the F1 score as evaluation metrics.

Summarization (SUM) CNN/Daily-Mail is a benchmark (Hermann et al., 2015; Nallapati et al., 2016) for text summarization. We use *article*, *summary* segment embedding to divide the article and the summary. We fine-tune the full GPT2-small and VLM for 10 epochs with the Adam optimizer. For distillation, we only fine-tune VLM for five epochs. We set the batch size to 32 and limit the maximum length of the article to 400 tokens and that of the summary to 130 tokens. We use the ROUGE-1, ROUGE-2, and ROUGE-L scores (Lin, 2004) as evaluation metrics.

Neural Machine Translation (NMT) We use the spoken German-English translation dataset IWSLT (Cettolo et al., 2016) as our NMT benchmark. We use *source*, *target* segment embedding to divide the source language and the target language.

We fine-tune the full GPT2-small, VLM and distilled VLM for 8 epochs with the Adam optimizer. We set the batch size to 32 and limit the maximum length of the source and target sequence to 100 tokens. We use BLEU (Papineni et al., 2002) as the evaluation metric.

Persona Dialogue (DLG) The Persona-Chat dataset (Zhang et al., 2018) is a persona-grounded multi-turn conversation dataset. We use *persona*, *system*, *user* segment embedding to help the model to distinguish the persona, alternating system utterance and user utterance in an input sequence. We fine-tune the full GPT2-small or VLM for three epochs with the Adam optimizer. We set the batch size to 16 and only retain the last five utterances in the dialogue history. We use perplexity, BLEU, and Consistency score (Madotto et al., 2019) as evaluation metrics.

Natural Language Generation (NLG) The natural language generation challenge (Dušek et al., 2019) is a dataset for building a response generation module for task-oriented dialogue systems. Given a set of response attributes, the model needs to generate responses. For example, when the input attribute is *name[The Wrestlers]*, *priceRange[cheap]*, *customerRating[low]*, the output should be *The wrestlers offers competitive prices, but is not highly rated by customers*. We use a set of attribute segment embedding to segment the input attributes. We fine-tune the full GPT2-small and VLM for 10 epochs with the Adam optimizer. We set the batch size to 32 and use BLUE (Papineni et al., 2002), ROUGE (Lin, 2004), NIST (Lin and Och, 2004), METEOR (Denkowski and Lavie, 2014) and CiDER (Vedantam et al., 2015) as evaluation metrics.

Computational Cost Fine-tuning VLM requires around 80%-90% GPU memory compared to full-finetune the whole GPT-2 model, as it only updates the small ratio of parameters. And both models have similar training cost, we report the training speed with single GTX 1080 Ti:

Task	Training Speed	Training set size
SUM	7.5h/epoch	300, 000
NMT	1.6h/epoch	200, 000
DLG	1.5h/epoch	130, 000
QA	5.0h/epoch	100, 000
NLG	0.2h/epoch	42, 000

A.3 Detailed Results

In this section, we report the detailed results for each task in Tables 2-6. We use a greedy decoding strategy for all the tasks.

A.4 Example

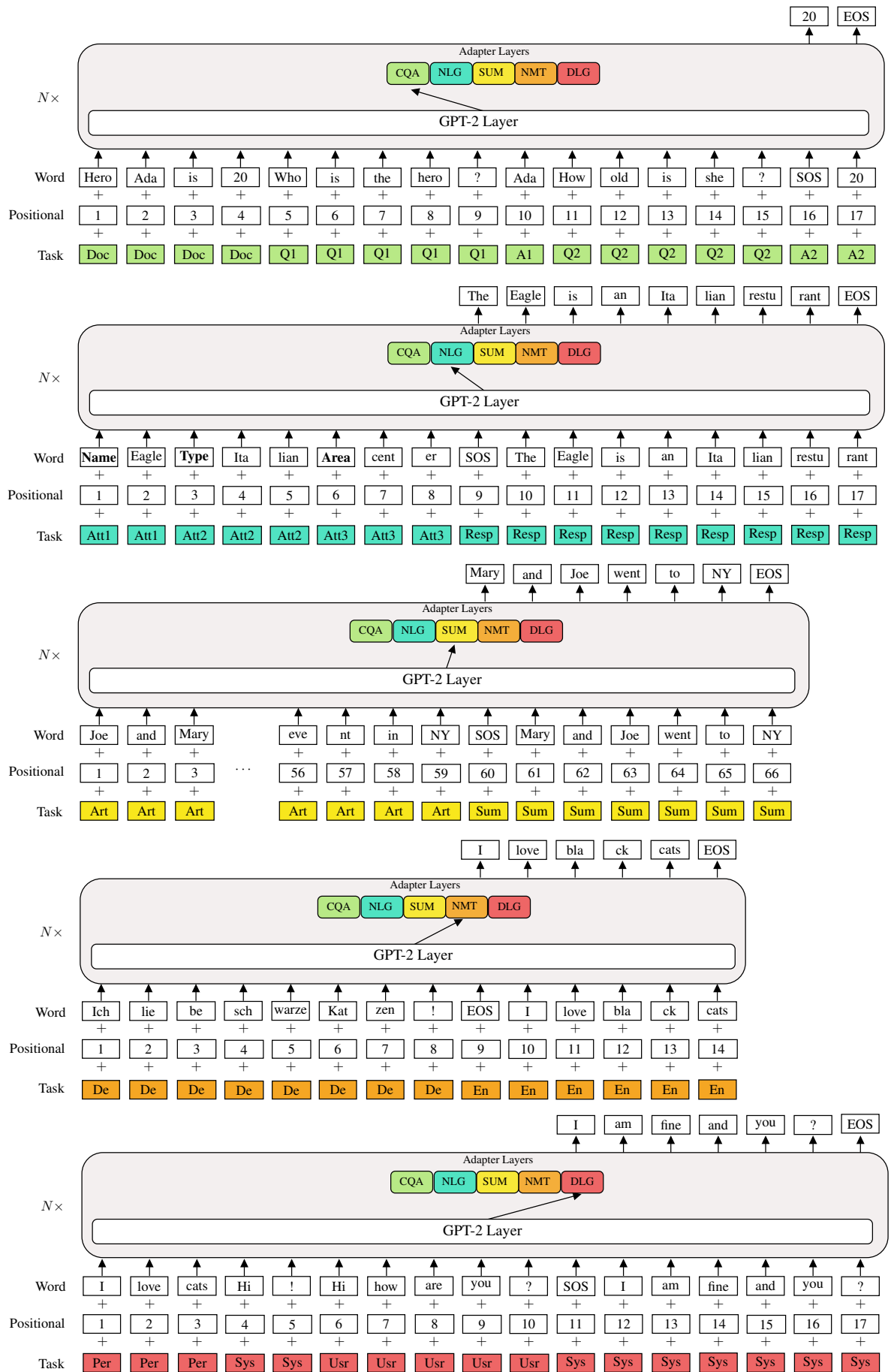


Figure 3: A detailed version of VLM. VLM shares a GPT-2 back-bone and for each task, the model looks up a set of task embeddings and chooses the corresponding adapter.

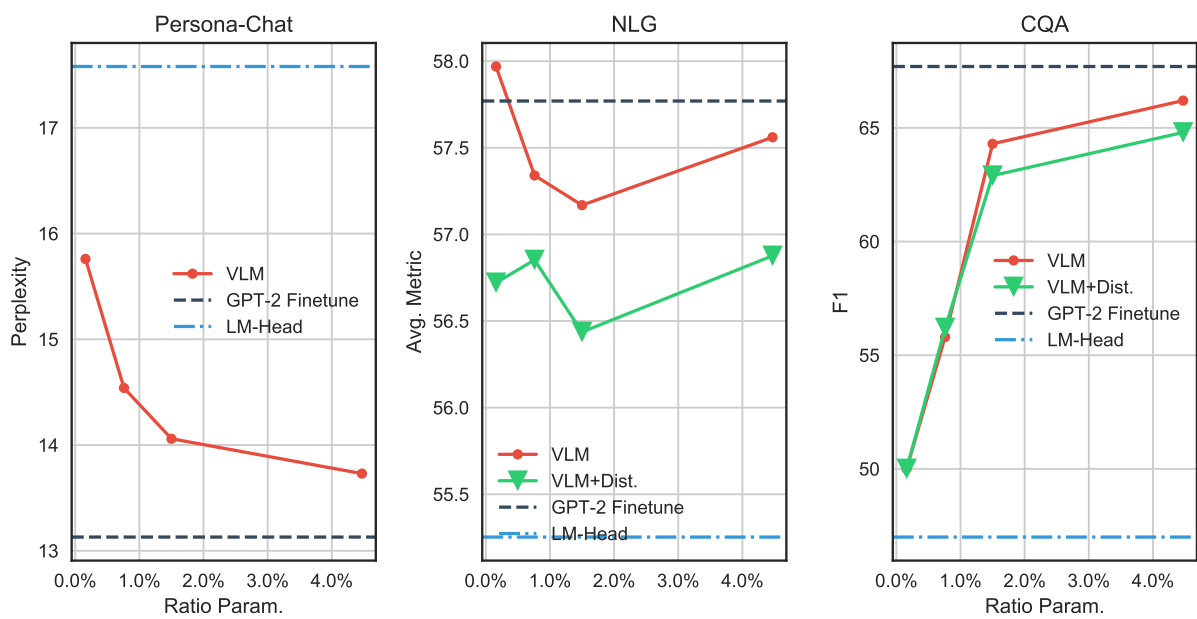


Figure 4: Performance comparison among different ratios of additional parameters. Here we can see that knowledge distillation does not improve the performance of the NLG task because of the small gap between VLM and the full fine-tuned GPT-2. Instead for the dialogue and QA tasks, the gold target is always better than the distilled target.

CNN / Daily Mail			
Models	<i>ROUGE 1</i>	<i>ROUGE 2</i>	<i>ROUGE L</i>
GPT Finetune	37.4	18.1	27.7
w/o Pre-Train	35.5	17	26.2
VLM mutli-task	36.6	17.7	27
VLM-10 (+ Dlst.)	35.0 (36.2)	16.5 (17.3)	25.0 (25.7)
VLM-50 (+ Dlst.)	36.4 (36.8)	17.5 (17.9)	26.6 (26.8)
VLM-100 (+ Dlst.)	36.5 (37.0)	17.6 (18.0)	27.0 (27.0)
VLM-300 (+ Dlst.)	36.6 (36.7)	17.6 (17.7)	26.6 (26.7)
PGNet (See et al., 2017)	39.53	17.28	36.38
Bottom-Up (Gehrmann et al., 2018)	41.22	18.68	38.34
UniLM (Dong et al., 2019)	43.33	20.21	40.51
T5-11B (Raffel et al., 2019)	43.52	21.55	40.69

Table 2: Summarization results.

Persona			
Models	<i>Perplexity</i>	<i>BLEU</i>	<i>Consistency (C)</i>
GPT Finetune	13.13	2.17	0.71
w/o Pre-Train	37.77	0.99	0.12
VLM mutli-task	13.15	0.84	0.27
VLM-10	15.76	1.63	0.86
VLM-50	14.54	1.84	0.72
VLM-100 (+ Dlst.)	14.06 (89.34)	1.99 (2.15)	0.76 (0.72)
VLM-300	13.73	1.98	0.74
Deep Copy (Yavuz et al., 2019)	54.58	4.09	-
PAML-TRS (Madotto et al., 2019)	30.42	1.0	0.07
ADAPT Centre (ConvAI2) (Dinan et al., 2019)	29.85	-	-
Persona-Chat (Zhang et al., 2018)	35.07	-	-
TransferTransfero (Wolf et al., 2019c)	17.51	-	-

Table 3: Persona Chat results.

CoQA	
Models	<i>F1</i>
GPT Finetune	67.7
w/o Pre-Train	15.1
VLM mutli-task	69.3
VLM-50 (+ Dlst.)	55.8 (56.2)
VLM-100 (+ Dlst.)	64.3 (62.9)
VLM-300 (+ Dlst.)	66.2 (64.8)
Seq2Seq (Reddy et al., 2019)	27.5
PGNet (Reddy et al., 2019)	45.4
DrQA (Reddy et al., 2019)	54.7
UNILM (Dong et al., 2019)	82.5
Human (Reddy et al., 2019)	89.8

Table 4: CoQA results.

NMT	
Models	BLUE
GPT Finetune	25.45
w/o Pre-Train	16.52
VLM mutli-task	22.49
VLM-10 (+ DIst.)	6.27(12.57)
VLM-50 (+ DIst.)	14.79(20.09)
VLM-100 (+ DIst.)	19.89(22.39)
VLM-300 (+ DIst.)	23.77(24.19)
Transformer (Vaswani et al., 2017)	29.2
DynamicConv (Wu et al., 2019)	35
MIXER (Ranzato et al., 2015)	21.83
AC+LL (Bahdanau et al., 2016)	28.53
NPMT (Huang et al., 2017)	28.96
Dual Transfer Learning (Wang et al., 2018)	32.35
LYC Transforemer (He et al., 2018)	35.07

Table 5: NMT results.

NLG						
Models	BLEU	NIST/10	METEOR	ROUGE L	CIDEr/10	norm. avg.
GPT Finetune	66.44	0.85279	0.4548	0.6911	0.22546	57.771
w/o Pre-Train	60.54	0.81697	0.4152	0.6471	0.19086	53.5106
VLM mutli-task	65.63	0.8342	0.4525	0.6889	0.22213	57.0806
VLM-10	67.1	0.85046	0.4545	0.6935	0.229	57.9692
VLM-50	66.01	0.84124	0.4568	0.6876	0.22128	57.3404
VLM-100	65.38	0.83922	0.4564	0.6893	0.21972	57.1688
VLM-300	66.18	0.84876	0.4539	0.6897	0.22387	57.5606
VLM-10 + DIst.	65.03	0.83199	0.456	0.6849	0.21286	56.721
VLM-50 + DIst.	65.23	0.83326	0.4576	0.6866	0.21287	56.8526
VLM-100 + DIst.	64.35	0.82485	0.4584	0.6852	0.20989	56.4368
VLM-300 + DIst.	65.19	0.83481	0.4575	0.6878	0.21182	56.8766
TGEN baseline (Dušek et al., 2019)	65.93	0.86094	0.4483	0.685	0.22338	57.5384
SLUG (Dušek et al., 2019)	66.19	0.8613	0.4454	0.6772	0.22615	57.44

Table 6: NLG results.

NMT IWSLT 2014	
<i>Source</i>	Wenn ihr mit jemanden in den 20ern arbeitet, einen liebt, wegen einem in den 20ern Schlaf verliert, ich möchte euch seh en... O.k. Großartig. Leute in den 20ern sind wirklich wichtig.
<i>GPT-2 Finetune</i>	If you work with somebody in the '20s, you love them because you lost a loved one in the '20s, I want to see you – – great. People in the '20s are really important.
<i>VLM</i>	If you work with somebody in the '20s, because of a love lost in the '20s, I want to see you – OK. Great. People in the '20s are really important.
<i>LM-Head</i>	If you work with someone in the 20ern, you love,
<i>Target</i>	If you work with twentysomethings, you love a twentysomething, you're losing sleep over twentysomethings, I want to see — Okay. Awesome, twentysomethings really matter.
<i>Source</i>	Ja, die Leute lassen sich später häuslich nieder als früher, aber das machte Alex' 20er nicht zum Entwicklungsausfall.
<i>GPT-2 Finetune</i>	Yes, people will be more domestic in the future than they used to be, but that didn't make Alex' 20s for failure.
<i>VLM</i>	Yes, people would come up later than they used to, but that didn't make Alex' 20s a disaster.
<i>LM-Head</i>	Yes, people are later going to come back as former former, but that doesn't make Alex' 20s anymore.
<i>Target</i>	Yes, people settle down later than they used to, but that didn't make Alex's 20s a developmental downtime.
<i>Source</i>	Leute in den 20ern wie Alex und ich hatten nichts als Zeit.
<i>GPT-2 Finetune</i>	People in the '20s like Alex and I didn't have time for time.
<i>VLM</i>	People in the '20s like Alex and I had nothing but time.
<i>LM-Head</i>	People like Alex and I had nothing as a time.
<i>Target</i>	Twentysomethings like Alex and I had nothing but time.

Table 7: NMT examples.

Summurization CNN Daily Mail	
<i>Source</i>	los angeles -lrb- cnn -rrb- it's more than just one state's internal problem. the historic california drought hurts the,rest of the union, too. that's because california is a breadbasket to the nation, growing more than a third of its vegetables and nearly two-thirds of its fruits and nuts. here's why we should heed the ongoing drought in the most populous state, a slowly expanding natural disaster now in its fourth year that this week prompted gov. jerry brown to announce a,mandatory 25 %
<i>GPT-2 Finetune</i>	the drought hurts the rest of the union, too, because california is a breadbasket to the nation, the drought is, causing the nation's price to rise, thanks to a strong dollar.
<i>VLM</i>	california drought hurts rest of the union, too, says jerry brown. brown: prices rose last year for these items on your kitchen table. he says prices expected to rise this year because of inflation.
<i>LM-Head</i>	new : the drought has hurt the rest of the union, the drought has hurt the rest of the union, too, the drought has hurt the rest of the union, too, the drought has hurt the rest of the union, too.
<i>Target</i>	americans paid more for some fruits and vegetables last year because of the drought, tourists will now have to ask for a glass of water at a california restaurant, perhaps the only good thing is another " great " wine grape harvest last year.
<i>Source</i>	-lrb- cnn -rrb- governments around the world are using the threat of terrorism – real or perceived – to advance executions, amnesty international alleges in its annual report on the death penalty. " the dark trend of governments using the death penalty in a futile attempt to tackle real or imaginary threats to state security and public safety was stark last year, " said salil shetty, amnesty's secretary general in a release. " it is shameful that so many states around the world are essentially playing with people's lives – putting people to death for ' terrorism' or to quell internal instability on the ill-conceived premise of deterrence.
<i>GPT-2 Finetune</i>	report: governments are using the threat of terrorism to advance executions, china is mentioned as having used, the death penalty as a tool in its " strike hard " campaign., the number of executions worldwide has gone down by almost 22 % on the previous year.
<i>VLM</i>	amnesty international releases its annual report on the death penalty, the report cites the example of pakistan, lifting a six-year moratorium on the execution of civilians.
<i>LM-Head</i>	amnesty international : death penalty is a tool in state-sanctioned killing campaigns, china is also mentioned, as having used the death penalty as a tool in its " strike hard " campaign against terrorism.,china is also mentioned, as having used the death penalty as a tool in its " strike hard " campaign against terrorism
<i>Target</i>	amnesty's annual death penalty report catalogs encouraging signs, but setbacks in numbers of those sentenced to death.,organization claims that governments around the world are using the threat of terrorism to advance executions., the number of executions worldwide has gone down by almost 22 % compared with 2013, but death sentences up by 28 %.

Table 8: SUM examples.

Conversational QA (CoQA)	
<i>Source</i>	<p>(CNN) – Dennis Farina, the dapper, mustachioed cop turned-actor best known for his tough-as-nails work in such TV series as "Law Order," "Crime Story," and "Miami Vice," has died. He was 69.</p> <p>"We are deeply saddened by the loss,of a great actor and a wonderful man," said his publicist, Lori De Waal, in a statement Monday. "</p> <p>Dennis Farina was always warmhearted and professional, with a great sense of humor and passion for his profession. He will be greatly missed by his family, friends and colleagues." Farina, who had a long career as a police officer in Chicago, got into acting through director Michael Mann, who used him as a consultant and cast him in his 1981 movie,"Thief." That role led to others in such Mann-created shows as "Miami Vice" (in which Farina played a mobster) and "Crime Story" (in which he starred as Lt. Mike Torello). Farina also had roles, generally as either cops or gangsters, in a number of movies, including "Midnight Run" (1988), "Get Shorty" (1995), "The Mod Squad" (1999) and "Snatch" (2000). In 2004, he joined the cast of the long-running "Law Order" after Jerry Orbach's departure, playing Detective Joe Fontana, a role he reprised on the spinoff "Trial by Jury."</p> <p>Fontana was known for flashy clothes and an expensive car, a distinct counterpoint to Orbach's rumpled Lennie Briscoe. Farina was on "Law Order" for two years, partnered with Jesse L. Martin's Ed Green. Martin's character became a senior detective after Farina left the show.</p> <p>Q1: Is someone in showbiz? A1: yes Q2: Whom?</p>
<i>GPT-2 Finetune</i>	Dennis Farina
<i>VLM</i>	Dennis Farina
<i>LM-Head</i>	Michael Mann
<i>Target</i>	Dennis Farina

Table 9: CQA examples.

Conversational QA (CoQA)	
<i>Source</i>	<p>Document:</p> <p>Once upon a time, in a barn near a farm house, there lived a little white kitten named Cotton. Cotton lived high up in a nice warm place above the barn where all of the farmer's horses slept. But Cotton wasn't alone in her little,home above the barn, oh no. She shared her hay bed with her mommy and 5 other sisters. All of her sisters were cute and,fluffy, like Cotton. But she was the only white one in the bunch. The rest of her sisters were all orange with beautiful white tiger stripes like Cotton's mommy. Being different made Cotton quite sad. She often wished she looked like the rest of her family. So one day, when Cotton found a can of the old farmer's orange paint, she used it to paint herself like them. When her mommy and sisters found her they started laughing. "What are you doing, Cotton?!" "I only wanted to be,more like you". Cotton's mommy rubbed her face on Cotton's and said "Oh Cotton, but your fur is so pretty and special, like you. We would never want you to be any other way". And with that, Cotton's mommy picked her up and dropped her into,a big bucket of water. When Cotton came out she was herself again. Her sisters licked her face until Cotton's fur was all all dry. "Don't ever do that again, Cotton!" they all cried. "Next time you might mess up that pretty white fur of yours and we wouldn't want that!" Then Cotton thought, "I change my mind. I like being special".</p> <p>Q1: Whose paint was it? A1: the farmer Q2: What did Cotton's mother and siblings do when they saw her painted orange?</p>
<i>GPT-2 Finetune</i>	laughed
<i>VLM</i>	laughed
<i>LM-Head</i>	They put her in a big bucket of water
<i>Target</i>	laughed

Table 10: CQA examples.

Dialogue Persona-Chat	
<i>Source</i>	<p>Persona: i read twenty books a year i'm a stunt double as my second job. i only eat kosher. i was raised in, a single parent household.</p> <p>Context: USR: hello what are doing today? SYS: i am good, i just got off work and tired, i have two jobs. USR: i just got done watching a horror movie</p>
<i>GPT-2 Finetune</i>	i love horror movies, i'm a stunt double.
<i>VLM</i>	i love horror movies, i read twenty books a year.
<i>LM-Head</i>	i love horror movies, i love them.
<i>Target</i>	i rather read, i've read about 20 books this year.
<i>Source</i>	<p>Persona: i'm a little girl. i've superpowers. i like to make time stop. i live in the cloud.</p> <p>Context: USR: hi i dye my hair 4 times a year</p>
<i>GPT-2 Finetune</i>	that is cool. i like to dye my hair.
<i>VLM</i>	wow that's a lot of dyeing. i love to dye.
<i>LM-Head</i>	i love to make time stop
<i>Target</i>	do you live in the sky? i do. in the clouds.

Table 11: DLG examples.

Natural Language Generation NLG	
<i>Source</i>	name = Blue Spice eatType = coffee shop near = Crowne Plaza Hotel customer rating = 5 out of 5
<i>GPT-2 Finetune</i>	Blue Spice is a coffee shop near Crowne Plaza Hotel with a 5 out of 5 customer rating.
<i>VLM</i>	Blue Spice is a coffee shop near Crowne Plaza Hotel. It has a customer rating of 5 out of 5.
<i>LM-Head</i>	Blue Spice is a coffee shop near Crowne Plaza Hotel with a customer rating of 5 out of 5.
<i>Target</i>	<ul style="list-style-type: none"> - The coffee shop Blue Spice is based near Crowne Plaza Hotel and has a high customer rating of 5 out of 5. - The Blue Spice coffee shop, near Crowne Plaza Hotel, has a customer rating of 5 out of 5. - If you want a coffee shop rated 5 out of 5 pick Blue Spice. It is located near Crowne Plaza Hotel.
<i>Source</i>	name = Giraffe eatType = pub familyFriendly = yes food = English near = Rainbow Vegetarian Café area = riverside
<i>GPT-2 Finetune</i>	Giraffe is a family friendly pub located near Rainbow Vegetarian Café. It serves English food and is located in the riverside area
<i>VLM</i>	Giraffe is a family friendly pub located near Rainbow Vegetarian Café in the riverside area.
<i>LM-Head</i>	Giraffe is a pub near Rainbow Vegetarian Café that serves English food. It is family friendly and located in the riverside area.
<i>Target</i>	<ul style="list-style-type: none"> - The Giraffe, a child-friendly, English pub, is located in riverside area, near Rainbow Vegetarian Café. - Giraffe is a pub that is children friendly and provides English food. It is located on the riverside near Rainbow Vegetarian Café. - There is an English based pub called Giraffe. It is located in the riverside area near the Rainbow Vegetarian Café and, yes, it is kid friendly.

Table 12: NLG examples.