

A Framework for PromptOps in GenAI Application Development Lifecycle

Anonymous ACL submission

Abstract

The use of "prompts" in the creation process of Generative Artificial Intelligence (GenAI) systems is receiving increasing interest. The significance of these prompts throughout the development cycle, however, is not properly used by current software development lifecycle approaches. This study proposes a unique methodology for integrating timely engineering and management into the creation of GenAI applications. Organizations may benefit from using "PromptOps" to create GenAI applications more quickly, effectively, and securely. It offers a technique to lower the danger of bias, increase the accuracy and dependability of GenAI systems, and decrease the cost of development and implementation. Our platform facilitates the seamless integration of several automated technologies in software development by performing prompt operations (PromptOps). These include Continuous Integration/Continuous Deployment (CI/CD) pipelines, workflows, APIs, and more. Our approach enables developers to easily include automated technologies, leading to a more simplified and efficient process. Furthermore, this study indicates that the framework may enable all stakeholders, including non-engineering units, to convert prompts into services, expanding their use in the building of applications. This study emphasizes the critical significance of prompts in GenAI and shows how their incorporation may improve AI application development, eventually stimulating creativity and driving the adoption of Generative AI technology.

1 Introduction

The advent of Large Language Models (LLMs) represents a monumental milestone in the field of Natural Language Processing (NLP). Attention-based architectures have revolutionized the field by replacing complex recurrent or convolutional neural networks with a network architecture solely based on attention mechanisms, resulting in models that excel in both quality and efficiency, significantly influencing subsequent LLM advancements (Kaplan et al., 2020; Vaswani et al., 2017). These powerful LLMs, such as GPT-3 (Generative Pre-trained Transformer 3), Palm2 (Pathways Language Model 2) and their successors, have demonstrated remarkable capabilities as versatile computational engines. Their ability to process and generate natural language text has found applications across a wide range of domains, from language translation to content generation. However, harnessing the full potential of LLMs and effectively utilizing them in real-world applications necessitates a nuanced understanding of the critical role that prompts play in steering these models.

Emergence of Large Language Models (LLMs): The emergence of LLMs has reshaped the landscape of AI-driven applications. These models have the capacity to understand, interpret, and generate human language at an unprecedented scale and complexity. They have demonstrated the capability to perform tasks that range from text completion to question-answering and text summarization. This transformative technology has laid the foundation for a new era of AI-driven applications, where human-machine interactions are facilitated through natural language.

The Trend Towards LLM-Based Application Development: In recent years, there has been a noticeable trend towards building applications that leverage LLMs as the underlying

intelligence. This trend can be attributed to the versatility and adaptability of LLMs, which allow developers to create a wide array of applications without the need for custom-built machine learning models. From simple text-based chatbots to sophisticated content generators, LLMs have become the go-to choice for developers seeking to integrate natural language understanding and generation capabilities into their applications. This shift is evident in domains like customer support, content creation, and data analysis.

The Crucial Role of Prompts in LLM Communication: To harness the power of LLMs effectively, communication with these models is achieved primarily through prompts. A prompt is a natural language input or instruction that is provided to the LLM to elicit a desired response or behavior. It acts as the bridge between human intent and machine execution. For example, a prompt might instruct an LLM to translate a sentence from English to French or generate a summary of a news article. The quality, clarity, and specificity of prompts are paramount in obtaining the desired output from the LLM.

The Impact of Prompts and Model Variants on Output: However, it is essential to recognize that prompts alone are not the sole determinants of LLM behavior. The choice of prompt and the specific LLM variant being used can significantly influence the output generated. Different prompts may yield varying results, even when targeting the same task. Moreover, LLMs often exist in multiple versions or variants, each with its own characteristics and performance nuances. The interaction between prompts and model versions is complex and requires careful consideration. The significance of prompts and their management has become increasingly evident as more organizations and developers integrate LLMs into their workflows and applications. The quality of prompts directly impacts the utility, reliability, and accuracy of LLM-based systems. Hence, the effective control and optimization of prompts have emerged as critical areas of research and development in the field of NLP.

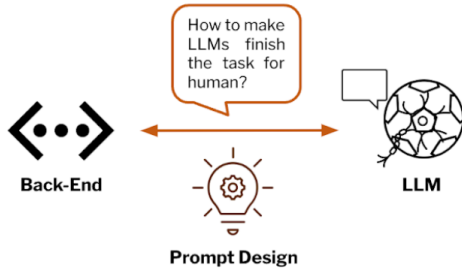
As LLMs continue to evolve and grow in power, the importance of prompt engineering becomes a prominent concern (Reynolds & McDonell, 2021). There has been some prior research of identifying suitable prompts (Sanh et al., 2021; Wei et al., 2021; Min et al., 2021; Mishra et al., 2021) along with the proposal of relevant tools (Bach et al., 2022; Zhou

et al., 2023). However, there is a lack of research concerning prompt management and operation. To address this gap, we propose the "GenFlow" framework, aimed at streamlining prompt management and operations. Our main contributions are:

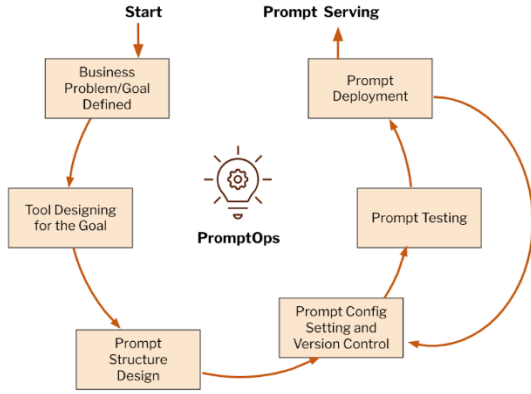
- We introduce the concept of incorporating prompt management into the DevOps (Development and Operations) flow. By aligning PromptOps with established software development practices, we facilitate the seamless integration of prompt operations into development pipelines, workflows, and APIs. This integration ensures that prompts are managed consistently, enhancing the reliability and efficiency of GenAI application development.
- Our proposed method, "GenFlow", empowers both developers and non-coders. This tool democratizes Prompt usage by providing an accessible interface for creating, modifying, and optimizing prompts. Its user-friendly design allows stakeholders from diverse backgrounds to harness the potential of prompts, thereby broadening their utility in application development.
- Within our framework, we introduce the concept of Prompt as a Service (PaaS). This extends the reach of prompts beyond development teams, enabling various stakeholders to utilize prompts as integral components in application building. This extension aligns with the growing recognition that prompt engineering will be a pivotal focus in the future of GenAI technology.

2 Related works

Prompt Engineering: The utilization of prompts provides an intuitive and natural way for human interaction with generative models, such as text-to-image model and LLMs. (Brown et al., 2020; Schick & Schütze, 2021; Sanh et al., 2022; Rombach et al., 2022). However, the effectiveness of LLMs requires accurate prompt design, either through manual intervention (Reynolds & McDonell, 2021) or automated methods (Gao et al., 2021; Shin et al., 2020). It is primarily since LLMs



(a) Prompt is the bridge between application and LLM



(b) The PromptOps workflow

Figure 1:(a) The diagram depicts the role of prompt designing in software development. (b) The diagram which describes each step in the workflow of prompt management.

do not interpret prompts in the same manner as humans (Webson & Pavlick, 2021; Lu et al., 2021). While numerous successful prompt tuning methods have leveraged gradient-based optimization over continuous spaces (Liu et al., 2021; Qin & Eisner, 2021; Lester et al., 2021), the computational cost of computing gradients escalates, especially when access to models shifts to APIs that may not offer gradient access. This raises a practical challenge: How can we empower users to create, design, and refine prompts effectively? This process, known as prompt engineering, is crucial for successful deployment due to the significant impact downstream predictions caused by prompt choices, especially in zero-shot settings (Perez et al., 2021; Zhao et al., 2021; Webson and Pavlick, 2021).

200

Prompt Management: Prompt management is a vital aspect of the emerging field of Natural Language Processing (NLP) where prompts serve as a bridge for interaction with Large Language Models (LLMs). Researchers have developed various techniques to search suitable prompts, such

as prompt generation (Gao et al., 2021; Ben-David et al., 2021), prompt scoring (Davison et al., 2019), and prompt paraphrasing (Jiang et al., 2020; Yuan et al., 2021). However, the journey doesn't end with well-engineered prompts, there arises a need for prompt management. Prompt management encompasses the organization, storage, and retrieval of these well-designed prompts to facilitate their efficient use in different contexts. The most notable tools in prompt management are OpenPrompt (Ding et al., 2022) and PromptSource (Bach et al., 2022).

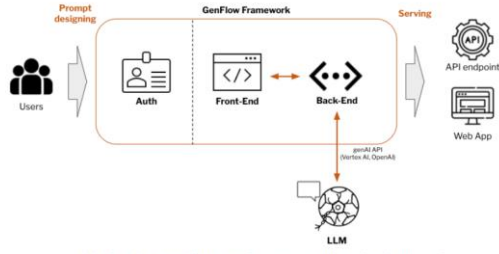
OpenPrompt simplifies prompt management within the world of pre-trained language models (PLMs). With the rise of prompt learning in natural language processing, there emerged a pressing need for a standardized framework. OpenPrompt fills this void by offering a user-friendly, research-focused toolkit. It's designed to cater to diverse PLMs, task formats, and prompting modules, providing a unified platform for effective prompt management.

Alongside OpenPrompt, PromptSource steps in as an essential tool for prompt development and sharing within the NLP landscape. As a prompt repository, its templating language for crafting data-linked prompts, an intuitive iterative development interface, and a vibrant community that contributes and collaborates, PromptSource empowers users with a wealth of over 2,000 prompts spanning around 170 datasets. This invaluable resource fosters seamless prompt management and utilization in language model training and querying, ultimately enhancing your NLP endeavors.

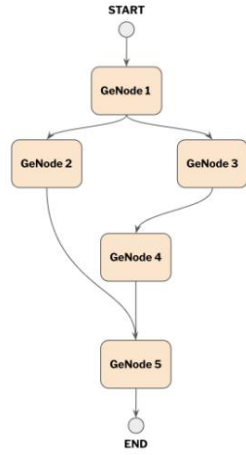
Although these tools have been introduced, the increasing prevalence of incorporating large language models (LLMs) into diverse applications highlights the imperative requirement for establishing **operational workflows for prompts** within this context. In pursuit of enhanced operational capabilities for prompts, encompassing Continuous Delivery (CD), version control, and more, we introduce a novel framework known as "GenFlow."

3 System design and workflow

The conventional approach to leveraging Large Language Models (LLMs) often involves the use of numerous prompts and processes to construct an application, particularly when each prompt requires fine-tuning. This method, while effective



(a) This diagram illustrates how users utilize this framework



(b) The diagram of a "FLOW"

Figure 2: The diagram of *GenFlow* framework

in achieving specific results, can lead to complications in maintenance, deployment, and management.

3.1 *GenFlow* framework

Through the abstraction of prompts into a configurable format, we simplify the development "process" significantly. Once prompts are transformed into configurations, they consist of parameters (variables) and prompt text. The prompt's parameters can be set through a user-friendly front-end interface. Prompt configurations become a resource introduced to the backend, allowing the backend or core application to focus on communication with the generative AI model, reducing its coupling with backend or core applications.

PromptConfig can be engineered by AI consumers from various domains and can be configured accordingly. In other words, the same backend/core application can generate different applications simply by employing different PromptConfigs.

With PromptConfigs in place, we can establish independent version control processes for prompts. This includes version tracking, diff comparisons,

and version control. Prompts become independently manageable, effectively transforming prompts into PromptOps that seamlessly integrate into DevOps practices (Figure 1b).

However, this architecture has its limitations, as it does not facilitate a sequence of operations, even though many applications are composed of a series of prompts. To address this limitation, we elevate the concept of promptConfig to that of a web API. Web APIs are constructed using URLs (Endpoints) and parameters, where each Endpoint can represent a functional module. This implies that every promptConfig can become an individual API. The structure of a PromptConfig as a web API is represented as follows:

PromptConfig:

```
{
  "param1": "value1",
  "param2": "value2",
  ...
}
```

PromptAPI(GET):

```
https://domain.name/{promptName}?param1=value1&param2=value2...
```

PromptAPI (POST):

```
https://domain.name/{promptName}
```

POST data:

```
{
  "param1": "value1",
  "param2": "value2",
  ...
}
```

When each prompt becomes an API, it becomes straightforward to manage and integrate through workflow management tools. Consequently, diverse prompts can be seamlessly combined to create extensive applications, and even the composition of prompts can be exposed as an API. promptCombosAPI = (prompt 1 API -> prompt 2 API -> prompt 3 API)

This approach allows us to regard Prompts as a Service. We term this entire process a "FLOW," where each PromptConfig/PromptAPI serves as a "NODE." This sequence of operations forms a "FLOW", as Figure 2b, where PromptConfigs are

used as "NODES" that can be orchestrated into a coherent process.

Title	Social Media Post Assistant		
Set Input Fields			
Variable	socialPlatform	Name	RelatedTopic
Placeholder	Please enter the name of the social platform	Your nickname on the platform	The topic of this post
Prompt	You are a super editor known as {{Name}} on the social platform {{socialPlatform}}. You will write an article about {{RelatedTopic}} on {{socialPlatform}}. Please generate a specific post in the json format: { "socialPlatform": "The name of the social platform", "Name": "Your nickname", "RelatedTopic": "The generated content" }		

Table 1: An example of *GeNode* design.

3.2 *GeNode*

Within the *GenFlow* framework, there exists a fundamental component known as "GeNode", which assumes the role of facilitating prompt editing and serves as a fundamental unit for prompt configuration. Within the user interface of *GeNode*, users are provided with the flexibility to tailor various facets of *GeNode* in accordance with their specific requirements. This entails the ability to designate a user-defined title, specify personalized variables, and construct the prompt content according to their preferences. For instance, users have the autonomy to create a *GeNode* themselves, as exemplified in Table 1. Once these configurations have been meticulously defined and subsequently saved, the *GeNode* is preserved as a prompt template, effectively culminating in the creation of a *GeNode* aptly named "Social Media Post Generator".

Through the innovative power of *GeNode*, users gain the remarkable ability to effortlessly create a multitude of applications, all without the need for advanced coding skills. By simply configuring various parameters within the *GeNode* interface, users can easily craft a diverse range of utilities. These encompass a broad spectrum of functions, including translation services, computational tools, code generation aids, and much more.

By wholeheartedly embracing *GeNode's* capabilities, users not only unlock the untapped potential of prompts for comprehending natural

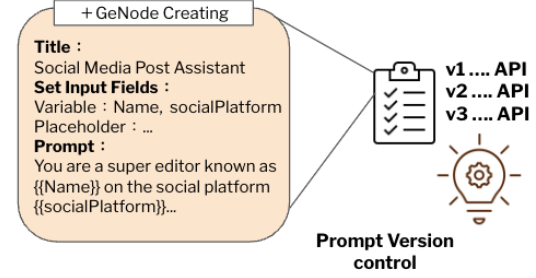


Figure 3: The diagram of *GenFlow* framework.

language but also effectively harness this resource to build an impressive array of real-world applications. Each of these applications is meticulously tailored to meet their unique requirements. This adaptability and flexibility undeniably position *GeNode* as an indispensable and cardinal component within the continually evolving landscape of prompt-driven tool development.

3.3 *GeNode* version list

Within the *GenFlow* framework, we have incorporated the concept of version control. When creating a *GeNode*, users can establish new versions for the same *GeNode*. These new versions can be directly modified from the original *GeNode* prompt. Consequently, a single application (*GeNode*) can have multiple distinct versions, with the option to select one as the published version. Users can access this published version through a user-friendly URL link that leads to the application's web user interface. Alternatively, they can serve the application directly by copying its API URL link.

Through this meticulously designed system, we have seamlessly integrated prompt configuration, version control, and serving as API endpoints and web UIs within our framework. This integration greatly simplifies the processes of application development and deployment, rendering them more accessible, manageable, and user-friendly.

3.4 *GeNode* List

Once a *GeNode* has been successfully created, it will be listed in the *GeNode* List page. Each *GeNode* will be accompanied by essential details such as its title, published version, version Create time, and actions. Specifically, under the "Actions"

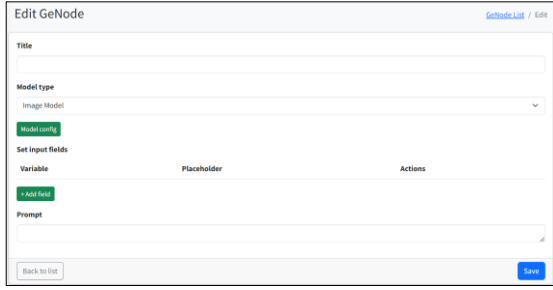
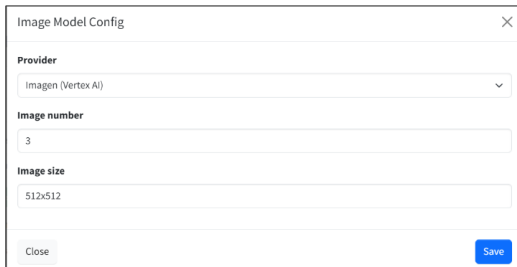


Figure 4: With the *GeNode* view, users can build their own application by designing the prompt and choosing proper models.



(a) The detail of LLM model configuration setting interface



(b) The detail of image model configuration setting interface

Figure 5: With the Config view, users can choose proper models for different applications.

column, users will find convenient links to perform actions such as creating a new version, accessing the UI via URL link, and serving the application using the API URL link. In the *GenFlow* UI section, visual representations will provide further insights and guidance on these features.

4 Implement Based on our framework

To facilitate user interaction and streamline the process further, we have developed the *GenFlow* web interface, as illustrated in Figure 4. Within this interface, we have integrated various models, including options for users to choose from within the interface, encompassing both LLMs and image generation models. Currently available language

models include Google PALM2-related models and OpenAI GPT (Figure 5a); image generation models include Google Imagen (Vertex AI) and DALL-E (Figure 5b). Users can choose to generate text or images based on their needs. Through this interface, it becomes possible to utilize different sources and types of models on a single platform, enabling the creation of customized applications.

Furthermore, in the *GenFlow* web interface, users can seamlessly engage in prompt design and editing on the *GeNode* edit page. This flexibility allows a single *GeNode* to be edited into multiple versions, with users having the prerogative to select which version to publish. Figure 6 showcases a user-friendly version list UI, which enhances the management of these different *GeNode* versions. This framework empowers users to rapidly create customized applications. Through the utilization of UI links and API links, these applications can be served effortlessly. Notably, this entire workflow is devoid of coding complexities, as users can solely focus on prompt editing. This concept aligns harmoniously with our previously introduced notion of "Prompt as a service," revolutionizing the application development paradigm.

By using nodes (*GeNode*) strung together as a flow can be achieved through various methods. Currently, there are similar services available. For instance, GCP offers Workflow (as shown in Figure 6), and AWS provides Step Functions.

5 Conclusion

In this study, we have introduced a comprehensive framework, *GenFlow*, to streamline and enhance the development of applications powered by Large Language Models (LLMs) through prompt management and operation. As the utilization of LLMs becomes increasingly prevalent in diverse applications, the role of prompts in shaping the behavior of these models has become paramount.

Our framework, *GenFlow*, addresses the critical need for prompt management within the context of LLM-driven application development. By seamlessly integrating prompt operations (PromptOps) into established DevOps practices, we have demonstrated the potential to enhance the reliability and efficiency of GenAI application development. This integration extends to Continuous Integration/Continuous Deployment

GeNode Version List				
Title	Version	Status	Version Create Time	Actions
Social Media Post Assistant	1.0.0.274	Draft	2023-09-28 08:16:50	Delete Copy Link Copy Link Copy Link
Social Media Post Assistant	1.0.0.403	Published	2023-09-28 08:14:57	Delete Copy Link Copy Link Copy Link

Figure 6: *GeNode Version List* view.

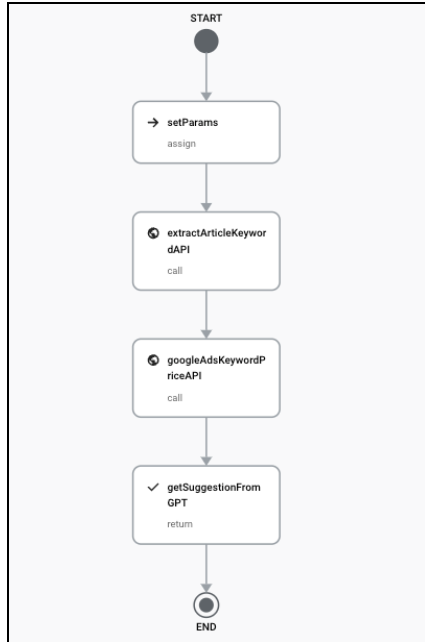


Figure 7: An example of GCP *Workflow* service.

(CI/CD) pipelines, workflows, APIs, and more, ensuring that prompts are consistently managed.

One of the key contributions of our approach is the democratization of prompt usage. *Genflow* empowers not only developers but also non-engineering units to create, modify, and optimize prompts through its accessible interface. This democratization broadens the utility of prompts in application development, facilitating collaboration and creativity across diverse teams.

Furthermore, our framework introduces the concept of Prompt as a Service (PaaS), allowing various stakeholders to leverage prompts as integral components in application building. This extension aligns with the growing recognition that prompt engineering will play a pivotal role in the future of GenAI technology.

Overall, the *GenFlow* framework offers a holistic solution to the challenges of prompt management and operation in GenAI application development. It empowers users to harness the full potential of LLMs while simplifying the process and promoting collaboration across disciplines. As the field of GenAI continues to evolve, prompt management will remain a critical aspect, and

GenFlow is poised to play a pivotal role in driving innovation and efficiency in this domain.

In conclusion, the effective control and optimization of prompts through *GenFlow* hold the promise of transforming the GenAI landscape, ultimately fostering creativity, reliability, and widespread adoption of Generative AI technology. We look forward to further research and development in this exciting field, with the aim of advancing the state of the art in GenAI application development.

Limitations

Limited Validation: The paper primarily focuses on proposing the *GenFlow* framework without providing extensive empirical validation or case studies to demonstrate its effectiveness in real-world scenarios. While the theoretical foundation is well-articulated, empirical evidence is crucial to validate the practical applicability and performance of the proposed framework.

Scalability Concerns: regarding the *GenFlow* framework are particularly pertinent, especially considering its intended deployment on the Google Cloud Platform (GCP). While the paper outlines the potential benefits of *GenFlow* in streamlining prompt management and operation, it's essential to consider how well the framework can scale within the infrastructure provided by GCP. Additionally, the scalability of network infrastructure plays a vital role in facilitating communication between different components of the *GenFlow* framework deployed across distributed environments. GCP provides scalable networking solutions, such as Virtual Private Clouds (VPCs) and load balancers, to ensure low-latency communication and high availability of applications.

References

- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. *Scaling laws for neural language models*. *Computing Research Repository*, arXiv:2001.08361.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser and Illia Polosukhin. 2017. *Attention Is All You Need*. *Computing Research Repository*, arXiv:1706.03762.
- Laria Reynolds and Kyle McDonell. 2021. *Prompt programming for large language models: Beyond the few-shot paradigm*. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in*

- 544 *Computing Systems*, pages 1–7. 599 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
545 <https://doi.org/10.1145/3411763.3451760>. 600 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
- 546 Victor Sanh, Albert Webson, Colin Raffel, Stephen H. 601 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
547 Bach, Lintang Sutawika, Zaid Alyafeai, Antoine 602 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
548 Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, 603 Gretchen Krueger, Tom Henighan, Rewon Child,
549 Manan Dey, M Saiful Bari, Canwen Xu, Urmish 604 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
550 Thakker, Shanya Sharma Sharma, Eliza Szczechla, 605 Clemens Winter, Christopher Hesse, Mark Chen,
551 Taewoon Kim, Gunjan Chhablani, Nihal Nayak, 606 Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin
552 Debajyoti Datta, Jonathan Chang, Mike Tian-Jian 607 Chess, Jack Clark, Christopher Berner, Sam
553 Jiang, Han Wang, Matteo Manica, Sheng Shen, 608 McCandlish, Alec Radford, Ilya Sutskever and
554 Zheng Xin Yong, Harshit Pandey, Rachel Bawden, 609 Dario Amodei. 2020. *Language models are few-shot*
555 Thomas Wang, Trishala Neeraj, Jos Rozen, 610 *learners*. *Computing Research Repository*,
556 Abheesht Sharma, Andrea Santilli, Thibault Fevry, 611 arXiv:2005.14165.
- 557 Jason Alan Fries, Ryan Teehan, Tali Bers, Stella 612 Timo Schick and Hinrich Schütze. 2021. *Exploiting*
558 Biderman, Leo Gao, Thomas Wolf and Alexander M. 613 *cloze-questions for few-shot text classification and*
559 Rush. 2022. *Multitask Prompted Training Enables* 614 *natural language inference*. In *Proceedings of the*
560 *Zero-Shot Task Generalization*. *Computing* 615 *16th Conference of the European Chapter of the*
561 *Research Repository*, arXiv:2110.08207. 616 *Association for Computational Linguistics: Main*
617 *Volume*, pages 255–269.
618 <https://doi.org/10.48550/arXiv.2001.07676>.
- 562 Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin 619 Victor Sanh, Albert Webson, Colin Raffel, Stephen H.
563 Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew 620 Bach, Lintang Sutawika, Zaid Alyafeai, Antoine
564 M. Dai and Quoc V. Le. 2021. *Finetuned language* 621 Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja,
565 *models are zero-shot learners*. *Computing Research* 622 Manan Dey, M Saiful Bari, Canwen Xu, Urmish
566 *Repository*, arXiv:2109.01652. 623 Thakker, Shanya Sharma Sharma, Eliza Szczechla,
624 Taewoon Kim, Gunjan Chhablani, Nihal Nayak,
567 Sewon Min, Mike Lewis, Luke Zettlemoyer and 625 Debajyoti Datta, Jonathan Chang, Mike Tian-Jian
568 Hannaneh Hajishirzi. 2021. *MetalCL: Learning to* 626 Jiang, Han Wang, Matteo Manica, Sheng Shen,
569 *Learn In Context*. *Computing Research Repository*, 627 Zheng Xin Yong, Harshit Pandey, Rachel Bawden,
570 arXiv:2110.15943. 628 Thomas Wang, Trishala Neeraj, Jos Rozen,
629 Abheesht Sharma, Andrea Santilli, Thibault Fevry,
571 Swaroop Mishra, Daniel Khashabi, Chitta Baral and 630 Jason Alan Fries, Ryan Teehan, Tali Bers, Stella
572 Hannaneh Hajishirzi. 2021. *Cross-Task* 631 Biderman, Leo Gao, Thomas Wolf and Alexander M.
573 *Generalization via Natural Language* 632 Rush. 2022. *Multitask Prompted Training Enables*
574 *Crowdsourcing Instructions*. *Computing Research* 633 *Zero-Shot Task Generalization*. *Computing*
575 *Repository*, arXiv:2104.08773. 634 *Research Repository*, arXiv:2110.08207.
- 576 Stephen H. Bach, Victor Sanh, Zheng-Xin Yong, Albert 635 Yong chao Zhou, Andrei Ioan Muresanu, Ziwen Han,
577 Webson, Colin Raffel, Nihal V. Nayak, Abheesht 636 Keiran Paster, Silviu Pitis, Harris Chan and Jimmy
578 Sharma, Taewoon Kim, M Saiful Bari, Thibault 637 Ba. 2023. *Large Language Models Are Human-*
579 Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, 638 *Level Prompt Engineers*. arXiv:2211.01910.
580 Zhiqing Sun, Srulik Ben-David, Canwen Xu, 639 Version 2
581 Gunjan Chhablani, Han Wang, Jason Alan Fries, 640 Robin Rombach, Andreas Blattmann, Dominik Lorenz,
582 Maged S. Al-shaibani, Shanya Sharma, Urmish 641 Patrick Esser and Björn Ommer. 2022. *High-*
583 Thakker, Khalid Almubarak, Xiangru Tang, 642 *Resolution Image Synthesis with Latent Diffusion*
584 Dragomir Radev, Mike Tian-Jian Jiang and 643 *Models*. *Computing Research Repository*,
585 Alexander M. Rush. 2022. *PromptSource: An* 644 arXiv:2112.10752. Version 2
586 *Integrated Development Environment and*
587 *Repository for Natural Language Prompts*. In
588 *Proceedings of the 60th Annual Meeting of the*
589 *Association for Computational Linguistics: System*
590 *Demonstrations*, pages 93–104.
591 <https://aclanthology.org/2022.acl-demo.9>.
- 592 Ning Ding, Shengding Hu, Weilin Zhao, Yulin Chen, 645 Tianyu Gao, Adam Fisch and Danqi Chen. 2021.
593 Zhiyuan Liu, Haitao Zheng and Maosong Sun. 2022. 646 *Making pre-trained language models better few-*
594 *Openprompt: An open-source framework for* 647 *shot learners*. *Computing Research Repository*,
595 *prompt-learning*. In *Proceedings of the 60th Annual* 648 arXiv:2012.15723. Version 2
596 *Meeting of the Association for Computational*
597 *Linguistics: System Demonstrations*, pages 105–113.
598 <https://aclanthology.org/2022.acl-demo.10>. 650 Taylor Shin, Yasaman Razeghi, Robert L. Logan IV,
651 Eric Wallace and Sameer Singh. 2020. *AutoPrompt:*
652 *Eliciting knowledge from language models with*
653 *automatically generated prompts*. *Computing*
Research Repository, arXiv:2010.15980. Version 2

A Appendices

- Albert Webson and Ellie Pavlick. 2022. *Do prompt-based models really understand the meaning of their prompts?*. *Computing Research Repository*, arXiv:2109.01247. Version 2
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel and Pontus Stenetorp. 2022. *Fantastically Ordered Prompts and Where to Find Them: Overcoming Few-Shot Prompt Order Sensitivity*. *Computing Research Repository*, arXiv:2104.08786. Version 2
- Guanghui Qin and Jason Eisner. 2021. *Learning how to ask: Querying lms with mixtures of soft prompts*. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5203–5212. <https://aclanthology.org/2021.naacl-main.410>.
- Brian Lester, Rami Al-Rfou and Noah Constant. 2021. *The power of scale for parameter-efficient prompt tuning*. In *Proceedings of the 2021 Conference on Empirical Methods Natural Language Processing*, pages 3045–3059. <https://aclanthology.org/2021.emnlp-main.243>.
- Ethan Perez, Douwe Kiela and Kyunghyun Cho. 2021. *True few-shot learning with language models*. *Advances in Neural Information Processing Systems*, 34:11054–11070.
- Tony Z. Zhao, Eric Wallace, ShiFeng, Dan Klein and Sameer Singh. 2021. *Calibrate Before Use: Improving Few-Shot Performance of Language Models*. *Computing Research Repository*, arXiv:2102.09690. Version 2
- Eyal Ben-David, Nadav Oved and Roi Reichart. 2022. *PADA: Example-based Prompt Learning for on-the-fly Adaptation to Unseen Domains*. *Computing Research Repository*, arXiv:2102.12206. Version 4
- Joe Davison, Joshua Feldman and Alexander M Rush. 2019. *Common sense knowledge mining from pretrained models*. In *Proceedings of the 2019 Conference on Empirical Methods Natural Language Processing and the 9th International Joint Conference on natural language processing*, pages 1173–1178. <https://aclanthology.org/D19-1109>.
- Zhengbao Jiang, Frank F. Xu, Jun Araki and Graham Neubig. 2020. *How Can We Know What Language Models Know?*. *Transactions of the Association for Computational Linguistics*, 8:423–438. <https://aclanthology.org/2020.tacl-1.28>.
- Weizhe Yuan, Graham Neubig and Pengfei Liu. 2021. *BartScore: Evaluating generated text as text generation*. *Advances in Neural Information Processing Systems*, 34:27263–27277.