

Fine Tuning LLM using Federated Learning

^[1] Chetan More, ^[2] Dhananjay Singh, ^[3] Tanishq Mohite, ^[4] Aditya Bahiram, ^[5] Shital Girme, ^[6] Sachin Gupta

^[1] ^[2] ^[3] ^[4] ^[5] Pune Institute of Computer Technology

^[6] Veritas, Pune Institute of Computer Technology

Corresponding Author Email: ^[1] chetanmore8605@gmail.com, ^[2] adityab16@gmail.com, ^[3] sdhananjay0187@gmail.com, ^[4] sngirme@pict.edu, ^[5] mtanishq@yahoo.com, ^[6] sachin.gupta@veritas.com

Abstract— Large language models (LLMs) have revolutionized natural language processing, pushing the boundaries of what machines can understand and generate text. These complex models, trained on massive datasets, excel at various tasks like summarizing factual topics, creating all kinds of creative content, and translating languages. However, their strength lies in their generality, and to truly shine on specific tasks, they require fine-tuning. This fine-tuning process tailors the LLM to a particular domain or application, significantly boosting its performance. Traditionally, this fine-tuning relies on centralized data storage, where vast amounts of user data are aggregated in one location. This approach raises significant challenges regarding data privacy and security. Users often hesitate to share their data, and regulations regarding data ownership and transfer can create roadblocks. Federated learning offers a promising solution to these challenges. It enables collaborative training on decentralized datasets stored on individual devices or local servers. This approach protects user privacy by keeping the data local while allowing the LLM to leverage the collective knowledge from these distributed sources. By exploring federated learning for fine-tuning LLMs, this research aims to bridge the gap between generic capabilities and domain-specific expertise, boosting the development of efficient and privacy-preserving language models.

Index Terms— Federated Learning, Large Language Models, Fine Tuning, Security, Encryption, OpenSSL, Flower, Natural Language Processing.

I. INTRODUCTION

The field of natural language processing [1], [19]–[21], [29], [30] has witnessed a paradigm shift with the emergence of large language models (LLMs) [4]–[7], [13], [19], [20], [22]. These behemoths, trained on colossal datasets of text and code, have demonstrated remarkable capabilities in tasks ranging from generating creative text formats to translating languages and answering complex questions. However, their very strength – their generality – can be a double-edged sword. To unlock their full potential for specific applications, LLMs require fine-tuning [4], [23], [29], a process that customizes the model for a particular domain. Traditionally, fine-tuning has relied on centralized data repositories, where user data is aggregated in a single location. This approach, while effective, raises significant concerns about data privacy and security. Users are increasingly wary of sharing their data, and regulations around data ownership and transfer pose additional hurdles.

Our research delves into federated learning [1], [5], [7]–[13], [17], [25]–[27] as a powerful alternative for fine-tuning LLMs. Federated learning upends the traditional paradigm by enabling collaborative training on decentralized datasets. Here, data resides on individual devices or local servers, fostering a privacy-preserving approach. The core principle lies in training a lightweight model copy on each local dataset. These local models then communicate updates to a central server, which aggregates the knowledge to improve the global LLM without ever requiring the raw data to leave its source. This approach offers a compelling

solution: it empowers LLMs to leverage the collective knowledge from distributed datasets while ensuring user privacy remains paramount.

This paper explores the potential of federated learning for fine-tuning LLMs. We investigate the challenges and opportunities associated with this approach, aiming to bridge the gap between the generic capabilities of LLMs and the need for domain-specific expertise. By successfully integrating federated learning with LLM fine-tuning, we pave the way for developing powerful and privacy-preserving language models that can excel in diverse real-world applications.

Leveraging federated learning, we present a fine-tuned Large Language Model (LLM) using the Flower library in a horizontal federated setup [10], [14]–[16] for classifying movie reviews into positive or negative sentiments. Our research utilizes the renowned IMDB dataset, a benchmark in sentiment analysis tasks, to train and validate the performance of our model. The Language Model employed in our study is GPT-2 (Generative Pre-trained Transformer 2) [5], [21], [24], a state-of-the-art natural language processing model renowned for its capability to understand and generate human-like text. By fine-tuning GPT-2 on the IMDB dataset, we aimed to enhance its ability to discern sentiment nuances within movie reviews, facilitating more accurate classification into positive or negative categories.

II. RELATED WORK

A. Federated Learning

In recent years, machine learning has become increasingly centralized, with large amounts of data being collected and processed by a single entity. While this centralized approach has led to significant advancements in the field, it also raises several concerns, including privacy issues, scalability limitations, and lack of representation. To address these challenges, a new paradigm called federated learning has emerged.

Federated learning is a distributed machine learning approach that enables multiple devices or organizations to collaboratively train models without sharing their data. By leveraging the collective computing power of these devices, federated learning can achieve better performance than traditional centralized methods while maintaining data privacy and security. In this conference paper, we explore the challenges associated with centralized machine learning and present federated learning as a viable alternative.

We begin by discussing the limitations of centralized machine learning and highlighting the need for a decentralized approach. We then introduce the fundamental concepts and techniques of federated learning, including the various approaches used to implement it. Finally, we provide case studies and experimental results demonstrating the effectiveness of federated learning in various applications. Our goal is to encourage researchers, policymakers, and industry leaders to embrace federated learning and work towards developing solutions that benefit society as a whole.

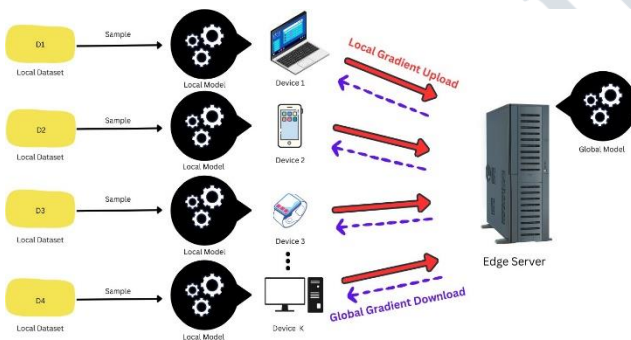


Fig. 1. Federated Learning Architecture

The rapid growth of federated learning (FL) has been accompanied by a surge in the development of specialized frameworks, each addressing the diverse needs and challenges that arise in various FL applications. Unlike monolithic solutions, each framework offers unique strengths and caters to specific use cases. For instance, TensorFlow Federated (TFF) [1]–[3], [8] and PySyft [2] prioritize general-purpose FL, while LEAF [2], [8], [15] focuses on mobile device optimization. Frameworks like FATE [3] and Substra prioritize secure enclaves for enhanced data privacy in medical fields. This heterogeneity within the FL landscape, with frameworks like XayNet, Nvidia Clara, PyVertical, MindSpore, Fedn, Synergos, and others, reflects the ongoing innovation and specialization required to address

the multifaceted demands of federated machine learning.

The existence of numerous FL frameworks also highlights the ongoing exploration of different algorithmic approaches. Frameworks like FedTree [3] and SecureBoost [28] explore applying decision tree and gradient boosting algorithms in a federated setting, respectively. Domain-specific frameworks like FedBioMed for healthcare and Nvidia Clara for medical imaging further demonstrate this targeted development. Other notable frameworks include FedML [15], PaddleFL, FederatedScope [6], Fedlearner, FLUTE, Flower [25], FedScale, CrypTen, OpenFL [18], PyFederate, IBM Federated Learning, LocalFed, BlockFL, FedAWS, ML-agents (Unity), NVFlare, FBLearner, PyTorch Lightning FL, PySyft + PyGrid, FLSim, OpenFed [12]. This variety allows researchers and practitioners to choose the most suitable tool for their needs, considering data privacy, communication efficiency, and model architecture. The ongoing development of these frameworks fosters continued advancement in the field of federated learning.

B. Flower

The Flower framework is an open-source platform designed for building, deploying, and managing machine learning workflows. It tackles the challenges of federated learning, like scale and device variations, by offering a customizable and scalable approach. Flower is unique because it allows researchers to seamlessly switch between simulated environments and real-world deployments with minimal code changes. With its modular architecture, Flower enables users to build custom components or integrate existing ones into their workflow. Additionally, it offers features such as distributed training, model versioning, and experiment tracking, making it well-suited for large-scale machine-learning projects. Furthermore, Flower supports various machine learning libraries, including TensorFlow, PyTorch, and Scikit-learn, providing flexibility in choosing the right tool for the job. Overall, Flower streamlines the machine learning development process by automating repetitive tasks, facilitating collaboration between team members, and ensuring the reproducibility of results.

C. Large Language Model and Fine Tuning

Language Model Learning (LLM) refers to the process of teaching machines how to understand and generate human language. This involves training deep neural networks on vast amounts of text data to learn patterns and relationships within language. Once a base language model has been trained, fine-tuning can be used to adapt the model to specific downstream applications, such as sentiment analysis, question answering, or translation.

Fine-tuning involves continuing the training process on a smaller, task-specific dataset, allowing the model to specialize in the desired application. During fine-tuning, only a small fraction of the weights in the network are updated, while the majority remain fixed, preserving the knowledge acquired during pretraining. By using transfer learning

through fine-tuning, models can achieve high performance even with limited labeled data, reducing the amount of time and resources required for training. Moreover, fine-tuned models often exhibit better generalization capabilities compared to models trained from scratch, resulting in more robust and accurate predictions. Therefore, fine-tuning plays a crucial role in developing practical NLP applications that deliver real-world value.

D. Generative Pretrained Transformer 2

The GPT-2 (Generative Pre-trained Transformer 2) model is considered the most advanced model in NLP, known for its superior ability to understand and generate human-like text. At its core, GPT-2 uses the Transformer architecture. It revolutionizes the field through attention mechanisms, allowing models to capture wide range of dependencies within text sequences. This architecture consists of multiple layers of self-attention mechanisms and feedforward neural networks that enable GPT-2 to process large amounts of text data and generate consistent responses.

The working principle of GPT-2 is pre-training on a diverse corpus of text data, where the model learns to predict the next word in a sequence based on the previous context. GPT-2 develops a deep understanding of language structure, semantics, and syntax through this unsupervised pre-training process. This pre-trained model can be tuned for specific downstream tasks such as text classification, sentiment analysis, language translation, and text generation by providing task-specific labeled data.

Through fine-tuning, GPT-2 can adapt its learned representations to the nuances of the target task, delivering state-of-the-art performance on various natural language processing tasks. GPT-2's applications span broad range of fields, reflecting its versatility and effectiveness in tackling various text-related tasks.

Regarding content generation, GPT-2 excels at producing consistent and context-relevant text, making it extremely useful for content summarization, dialog generation, and story writing. Additionally, GPT-2 is widely used in sentiment analysis to detect sentiment polarity in text input to help understand public opinion, analyze customer feedback, and monitor social media. Additionally, GPT-2 supports chatbots and virtual assistants, enabling natural and engaging interactions with users on many different platforms. Its applications have extended to areas such as machine translation, question answering, and text completion, and it has proven to have widespread impact on natural language understanding and natural language generation tasks.

E. Open Secure Socket Layer

The OpenSSL library serves as a basic tool for secure communication and cryptographic operations in software applications. The extensive feature set allows developers to implement encryption, decryption, digital signatures, and secure communication protocols such as SSL/TLS into their software systems. OpenSSL provides a robust

implementation of cryptographic algorithms and protocols to ensure data

confidentiality, integrity, and authenticity in network communications, data storage, and other sensitive operations. Due to its versatility and widespread adoption, OpenSSL has become the foundation for building secure software applications in various domains, including web servers, network devices, and embedded systems.

III. EXPERIMENT

1) Linear Regression using OpenFL (Synthetic Dataset):

In this study, a federated learning project was conducted using the OpenFL package in Python, focusing on executing a linear regression algorithm through a horizontal federated structure on a synthesized dataset. By leveraging the federated dataset library along with PyTorch, the investigation enabled joint model training across several entities while maintaining data security by confining raw information locally. To simulate realistic situations emphasizing privacy concerns, the fabricated dataset dispersed characteristics amongst involved participants. Utilizing federated learning tactics, model adjustments were traded between the groups instead of sharing raw data, resulting in creating a unified model without compromising private datasets. Overall, this research highlights the practicality and potential of horizontal federated learning approaches in cooperative machine-learning initiatives, especially when prioritizing data protection.

Experiment Results- The experimental results showcased an impressive accuracy of 0.73 achieved through the implementation of a linear regression algorithm on the synthetic dataset using the OpenFL package for horizontal federated learning. This notable accuracy underscores the effectiveness of the federated learning approach in collaborative model training while maintaining data privacy.

2) Linear Regression using Flower (Synthetic Dataset):

In this research endeavor, an experiment in federated learning was undertaken employing the Flower package in Python. The focus centered on utilizing a horizontal federated setup to implement a linear regression algorithm on a synthetic dataset. Leveraging the federated dataset library alongside TensorFlow, the experiment facilitated collaborative model training across multiple parties while preserving data privacy by keeping raw data localized. The synthetic dataset was strategically crafted to distribute features among the participating parties, mirroring real-world scenarios where privacy is paramount. Through the federated learning approach, model updates were exchanged among the parties rather than raw data, ensuring the development of a global model while safeguarding the confidentiality of individual datasets. This experiment serves to demonstrate the efficacy and applicability of horizontal federated

learning techniques in collaborative machine learning endeavors, particularly in contexts where data privacy is of paramount importance.

Experiment Results- The experimental results showed an impressive accuracy of 0.81 achieved through the implementation of a linear regression algorithm on the synthetic dataset using the Flower package for horizontal federated learning. This notable accuracy underscores the effectiveness of the federated learning approach in collaborative model training while maintaining data privacy.

- 3) Linear Regression using Flower (Boston Housing Dataset):** In this research study, an experiment in federated learning was conducted using the Flower package in Python. The focus was on implementing a linear regression algorithm on the Boston Housing dataset within a horizontal federated setup. Leveraging the federated dataset library and TensorFlow, the experiment facilitated collaborative model training across multiple parties while ensuring data privacy by keeping raw data localized. The Boston Housing dataset, a well-known benchmark dataset in machine learning, provided real-world housing-related features distributed among the participating parties. Through the federated learning approach, model updates were exchanged among the parties instead of raw data, allowing for the development of a global model while preserving the confidentiality of individual datasets. This experiment demonstrates the practical application of horizontal federated learning techniques in real-world scenarios, particularly in domains such as housing prediction, where data privacy and collaborative model training are paramount concerns.

Experiment Results- The experimental results revealed a notable accuracy of 0.83 achieved through the application of a linear regression algorithm on the Boston Housing dataset using the Flower package for horizontal federated learning. The obtained accuracy serves as compelling evidence of the feasibility and efficacy of horizontal federated learning methodologies, particularly in real-world applications like housing prediction, where accurate models are crucial for decision-making processes.

- 4) Multiclass Classification using Flower (MNIST Dataset):** In this research endeavor, an experiment in federated learning was conducted utilizing the Flower package in Python. The focal point of the study was the application of a CNN [15] model for classification tasks on the MNIST dataset within a horizontal federated setup. Leveraging the federated dataset library and TensorFlow, the experiment facilitated collaborative model training across multiple parties while ensuring the privacy of raw data. The MNIST dataset, renowned for its digit images, served as the benchmark for this classification task, with data distributed horizontally among participating parties. Through the federated learning approach, model updates

were exchanged among the parties rather than raw data, enabling the development of a global model while safeguarding the confidentiality of individual datasets. This experiment highlights the efficacy and practicality of horizontal federated learning methodologies, particularly in image classification tasks, where data privacy and collaborative model training are essential considerations.

Experiment Results- The experimental findings revealed an accuracy of 0.78 achieved through the application of CNN model for classification tasks on the MNIST dataset within a horizontal federated setup. This accuracy metric serves as a tangible measure of the effectiveness of the federated learning for classification purposes. The obtained accuracy underscores the feasibility and potential of horizontal federated learning methodologies in real-world scenarios, particularly in image classification tasks like those encountered in the MNIST dataset, where accurate models are imperative for reliable classification outcomes.

- 5) Logistic Regression using Flower (Iris Dataset):** In this research endeavor, a comprehensive experiment in federated learning was conducted utilizing the Flower package in Python. The central focus of the study was the application of a logistic regression algorithm for classification tasks on the Iris dataset within a horizontal federated setup. Leveraging the federated dataset library and TensorFlow, the experiment facilitated collaborative model training across multiple parties while ensuring the privacy of raw data. The Iris dataset, a well-known benchmark in machine learning, provided botanical features distributed horizontally among participating parties. Through the federated learning approach, model updates were exchanged among the parties rather than raw data, enabling the development of a global model while preserving the confidentiality of individual datasets. This experiment underscores the practical application of horizontal federated learning methodologies, particularly in classification tasks, where data privacy and collaborative model training are essential considerations. **Experiment Results-** The experimental findings revealed an accuracy of 0.79 achieved through the application of logistic regression for classification tasks on the Iris dataset within a horizontal federated setup. This accuracy metric serves as a tangible measure of the effectiveness of the federated learning approach in collaborative model training while preserving data privacy. The obtained accuracy underscores the feasibility and potential of horizontal federated learning methodologies in real-world scenarios, particularly in classification tasks like those encountered in the Iris dataset, where accurate models are crucial for reliable classification outcomes.

6) LLM Fine Tuning using Flower: During our exploration, we undertook an experimental evaluation of Large Language Model (LLM) fine-tuning utilizing the versatile Flower platform within a Python environment.

Focused on honing LLMs for domain-specific applications like sentiment analysis, we employed the IMDB movie review dataset to assess performance improvements gained via fine-tuning. With the aid of the adaptable Flower interface, we designed bespoke workflows targeting LLM fine-tuning efficiently. Distributing the fine-tuning procedure across multiple nodes expedited parallel processing and considerably shrank convergence times. Moreover, Flower's decentralized nature eliminated dependencies on central servers, subsequently minimizing latency and boosting scalability. Reflecting on our experience, the successful integration of Flower showcased remarkable enhancements in LLM fine-tuning processes, yielding actionable insights destined for cultivating advanced NLP architectures readily adapted to varied application domains. Most notably, these discoveries emphasize the indispensable role cutting-edge platforms like Flower play in propelling AI technology forward while nurturing thriving interoperable ecosystems.

Experiment Results- The experimental results revealed a notable accuracy of 0.78 achieved through the application of a LLM fine tuning for sentiment analysis using the Flower package for horizontal federated learning. Our experiment faced reduced accuracy linked predominantly to restricted computational power. Insufficient memory capacity impeded the simultaneous execution of essential tasks necessary for intricate procedures, requiring extra iterations and negatively influencing output quality. As a result, the model couldn't sufficiently explore broad search spaces, eventually opting for sub-optimal solutions, causing disappointing accuracy rates. Therefore, ample computational resources prove crucial for obtaining favorable outcomes, encouraging efforts to secure better facilities or innovate ways around current restrictions.

TABLE I: Experiment Results

Framework	Model	Dataset	Accuracy
OpenFL	Linear Regression	Synthetic Dataset	0.73
Flower	Linear Regression	Synthetic Dataset	0.81
Flower	Linear Regression	Boston Housing	0.83
Flower	CNN	MNIST	0.78
Flower	Logistic Regression	Iris	0.79
Flower	GPT-2	IMDB movie review	0.78

IV. CONCLUSION & FUTURE SCOPE

In summary, our study attempts to exploit the potential of federated learning to fine-tune large language models (LLMs) for specific tasks, exemplified by classifying reviews of movies into positive or negative emotions. Through applying federated learning techniques, especially using the Flower library and the FedAvg aggregation algorithm, we have fine-tuned the GPT-2 model on the IMDb data set, achieving a remarkable accuracy of 0.78.

Our research demonstrates the feasibility and effectiveness

of horizontal federated learning in the domain of natural language processing, paving the way for decentralized model training without compromising data privacy. By distributing the training process across multiple clients while updating the model as a group, federated learning alleviates privacy concerns & also enables the use of data residing on edge devices, expanding the scope of AI applications in real-life situations.

Looking ahead, the scope of our future research includes several promising avenues. First, exploring alternative aggregation algorithms beyond FedAvg could improve model performance and convergence. Additionally, investigating the integration of differentiated security mechanisms can further enhance data privacy protection in federated learning settings. Additionally, expanding our research to include more datasets and LLM architectures will facilitate a deeper understanding of the applicability of federated learning in different domains. Essentially, our study contributes to the evolving landscape of federated learning by demonstrating its feasibility and effectiveness in fine-tuning LLM for specific tasks, while also highlighting enabling future avenues of discovery and refinement. As the field continues to develop, the combination of federated learning and natural language processing holds great promise for democratizing AI and driving innovation in various fields.

REFERENCES

- [1] Khan M, Glavin FG, Nickles M. Federated learning as a privacy solution-an overview. *Procedia Computer Science*. 2023 Jan 1;217:316-25.
- [2] Riviera W, Galazzo IB, Menegaz G. FeLebrities: a user-centric assessment of Federated Learning frameworks. *IEEE Access*. 2023 Sep 6.
- [3] Liu X, Shi T, Xie C, Li Q, Hu K, Kim H, Xu X, Li B, Song D. Unifed: A benchmark for federated learning frameworks. *arXiv preprint arXiv:2207.10308*. 2022 Jul 21.
- [4] Ye R, Wang W, Chai J, Li D, Li Z, Xu Y, Du Y, Wang Y, Chen S. OpenFedLLM: Training Large Language Models on Decentralized Private Data via Federated Learning. *arXiv preprint arXiv:2402.06954*. 2024 Feb 10.
- [5] Fan T, Kang Y, Ma G, Chen W, Wei W, Fan L, Yang Q. Fate-llm: A industrial grade federated learning framework for large language models. *arXiv preprint arXiv:2310.10049*. 2023 Oct 16.
- [6] Kuang W, Qian B, Li Z, Chen D, Gao D, Pan X, Xie Y, Li Y, Ding B, Zhou J. Federatedscope-llm: A comprehensive package for fine-tuning large language models in federated learning. *arXiv preprint arXiv:2309.00363*. 2023 Sep 1.
- [7] Hard A, Rao K, Mathews R, Ramaswamy S, Beaufays F, Augenstein S, Eichner H, Kiddon C, Ramage D. Federated learning for mobile keyboard prediction. *arXiv preprint arXiv:1811.03604*. 2018 Nov 8.
- [8] Li T, Sahu AK, Talwalkar A, Smith V. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine*. 2020 May 1;37(3):50-60.
- [9] Nilsson A, Smith S, Ulm G, Gustavsson E, Jirstrand M. A performance evaluation of federated learning algorithms.

- InProceedings of the second workshop on distributed infrastructures for deep learning 2018 Dec 10 (pp. 1-8).
- [10] Zhang C, Xie Y, Bai H, Yu B, Li W, Gao Y. A survey on federated learning. *Knowledge-Based Systems*. 2021 Mar 15;216:106775.
- [11] Wei K, Li J, Ding M, Ma C, Yang HH, Farokhi F, Jin S, Quek TQ, Poor HV. Federated learning with differential privacy: Algorithms and performance analysis. *IEEE transactions on information forensics and security*. 2020 Apr 17;15:3454-69.
- [12] Chen D, Tan VJ, Lu Z, Wu E, Hu J. OpenFed: A comprehensive and versatile open-source federated learning framework. InProceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition 2023 (pp. 5017-5025).
- [13] Zhang Y, Ramage D, Xu Z, Zhang Y, Zhai S, Kairouz P. Private Federated Learning in Gboard. *arXiv preprint arXiv:2306.14793*. 2023 Jun 26.
- [14] Danish Z, Khan IR. A review of Federated Learning. InICIDSSD 2022: Proceedings of the 3rd International Conference on ICT for Digital, Smart, and Sustainable Development, ICIDSSD 2022, 24-25 March 2022, New Delhi, India 2023 May 16 (p. 146). European Alliance for Innovation.
- [15] Li Q, Wen Z, Wu Z, Hu S, Wang N, Li Y, Liu X, He B. A survey on federated learning systems: Vision, hype and reality for data privacy and protection. *IEEE Transactions on Knowledge and Data Engineering*. 2021 Nov 2;35(4):3347-66.
- [16] Banabilah S, Aloqaily M, Alsayed E, Malik N, Jararweh Y. Federated learning review: Fundamentals, enabling technologies, and future applications. *Information processing & management*. 2022 Nov 1;59(6):103061.
- [17] Kairouz P, McMahan HB, Avent B, Bellet A, Bennis M, Bhagoji AN, Bonawitz K, Charles Z, Cormode G, Cummings R, D'Oliveira RG. Advances and open problems in federated learning. *Foundations and trends® in machine learning*. 2021 Jun 22;14(1-2):1-210.
- [18] Foley P, Sheller MJ, Edwards B, Pati S, Riviera W, Sharma M, Moorthy PN, Wang SH, Martin J, Mirhaji P, Shah P. OpenFL: the open federated learning library. *Physics in Medicine & Biology*. 2022 Oct 19;67(21):214001.
- [19] Ouyang L, Wu J, Jiang X, Almeida D, Wainwright C, Mishkin P, Zhang C, Agarwal S, Slama K, Ray A, Schulman J. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*. 2022 Dec 6;35:27730-44.
- [20] Zhao J, Wang W, Xu C, Ren Z, Ng SK, Chua TS. LLM-based Federated Recommendation. *arXiv preprint arXiv:2402.09959*. 2024 Feb 15.
- [21] Bai Y, Jones A, Ndousse K, Askell A, Chen A, DasSarma N, Drain D, Fort S, Ganguli D, Henighan T, Joseph N. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*. 2022 Apr 12.
- [22] Touvron H, Martin L, Stone K, Albert P, Almahairi A, Babaei Y, Bashlykov N, Batra S, Bhargava P, Bhosale S, Bikel D. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*. 2023 Jul 18.
- [23] Chen C, Feng X, Zhou J, Yin J, Zheng X. Federated large language model: A position paper. *arXiv preprint arXiv:2307.08925*. 2023 Jul 18.
- [24] Webb T, Holyoak KJ, Lu H. Emergent analogical reasoning in large language models. *Nature Human Behaviour*. 2023 Sep;7(9):1526-41.
- [25] Beutel DJ, Topal T, Mathur A, Qiu X, Fernandez-Marques J, Gao Y, Sani L, Li KH, Parcollet T, de Gusmao PP, Lane ND. Flower: A friendly federated learning research framework. *arXiv preprint arXiv:2007.14390*. 2020 Jul 28.
- [26] Lyu L, Yu J, Nandakumar K, Li Y, Ma X, Jin J, Yu H, Ng KS. Towards fair and privacy-preserving federated deep models. *IEEE Transactions on Parallel and Distributed Systems*. 2020 May 21;31(11):2524-41.
- [27] AbdulRahman S, Tout H, Ould-Slimane H, Mourad A, Talhi C, Guizani M. A survey on federated learning: The journey from centralized to distributed on-site learning and beyond. *IEEE Internet of Things Journal*. 2020 Oct 12;8(7):5476-97.
- [28] Cheng K, Fan T, Jin Y, Liu Y, Chen T, Papadopoulos D, Yang Q. Secureboost: A lossless federated learning framework. *IEEE Intelligent Systems*. 2021 May 25;36(6):87-98.
- [29] Brown T, Mann B, Ryder N, Subbiah M, Kaplan JD, Dhariwal P, Neelakantan A, Shyam P, Sastry G, Askell A, Agarwal S. Language models are few-shot learners. *Advances in neural information processing systems*. 2020;33:1877-901.
- [30] Du Z, Qian Y, Liu X, Ding M, Qiu J, Yang Z, Tang J. Glm: General language model pretraining with autoregressive blank infilling. *arXiv preprint arXiv:2103.10360*. 2021 Mar 18.