
Position Paper: Assessing Robustness, Privacy, and Fairness in Federated Learning Integrated with Foundation Models

Xi Li^{*1} Jiaqi Wang^{*1}

Abstract

Federated Learning (FL), while a breakthrough in decentralized machine learning, contends with significant challenges such as limited data availability and the variability of computational resources, which can stifle the performance and scalability of the models. The integration of Foundation Models (FMs) into FL presents a compelling solution to these issues, with the potential to enhance data richness and reduce computational demands through pre-training and data augmentation. However, this incorporation introduces novel issues in terms of robustness, privacy, and fairness, which have not been sufficiently addressed in the existing research. We make a preliminary investigation into this field by systematically evaluating the implications of FM-FL integration across these dimensions. We analyze the trade-offs involved, uncover the threats and issues introduced by this integration, and propose a set of criteria and strategies for navigating these challenges. Furthermore, we identify potential research directions for advancing this field, laying a foundation for future development in creating reliable, secure, and equitable FL systems.

1. Introduction

In the evolving landscape of machine learning, Federated Learning (FL) has emerged as a pivotal framework, enabling collaborative model training across multiple devices while preserving data privacy. This decentralized approach, however, suffers from inherent challenges such as ineffective training and resource limitations. Ineffective training arises as clients often hold limited datasets, insufficient for training robust models. Coupled with the privacy-preserving nature of FL, where data remains localized, this scarcity hampers the development of effective models. Additionally, the dis-

tributed nature of FL introduces resource limitations, with client devices exhibiting wide variations in computational and communication capabilities, impacting the efficiency and scalability of model training.

Foundation Models (FMs), characterized by their vast knowledge and versatility, offer a promising solution to these challenges. FMs, particularly when integrated into FL, can mitigate issues of ineffective training by, *e.g.*, local data augmentation and model pre-training. By producing data that mirrors real-world distributions, FMs improve FL models' performance and generalization. Furthermore, FMs facilitate a reduction in computational burden through transfer learning, allowing clients to fine-tune pre-trained models with their local data, thereby requiring less computational power. Additionally, FMs can minimize communication overhead by serving as efficient encoders, reducing the need for extensive data transmission during the model aggregation phase.

However, the integration of FMs into FL introduces complex challenges concerning robustness, privacy, and fairness. The interaction between FMs and FL systems can amplify vulnerabilities, leading to new issues. For instance, the reliance on synthetic data generation and model pre-training with FMs may inadvertently introduce biases or facilitate evasion attacks, compromising the integrity and robustness of the federated models. Moreover, the privacy-preserving promises of FL could be undermined if the integration with FMs is not carefully managed, risking unintended data leakage or exploitation through sophisticated attack vectors.

Despite the potential of FM-FL integration to address the pressing challenges of federated learning, there is a noticeable gap in research exploring the responsibility of FM-FL. Our work stands at the forefront of this emerging field, systematically assessing the implications of integrating FMs with FL systems. We aim to shed light on the intricate dynamics between these two technologies, exploring their potential to revolutionize collaborative learning environments while investigating the complexities they introduce. This position paper delves into the challenges brought about by the integration of FMs with FL. It presents a detailed examination of key issues related to robustness, privacy, and fairness, thereby laying the groundwork for future research

^{*}Equal contribution ¹The Pennsylvania State University. Correspondence to: Xi Li <XiLi@psu.edu>, Jiaqi Wang <jqwang@psu.edu>.

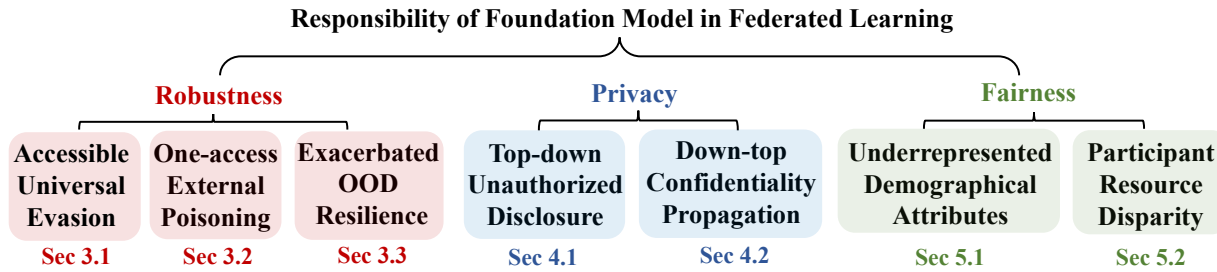


Figure 1. Overview of the FM-FL responsibility: robustness, privacy, and fairness.

and development in this burgeoning interdisciplinary field. Additionally, we propose potential research directions aimed at both understanding and addressing these concerns, further contributing to the advancement of FM-FL integration.

The responsibility issues explored in this paper are summarized in Fig. 1. Section 2 lays the groundwork by discussing the incorporation of FMs into FL. Section 3 assesses the robustness from the perspectives of accessible universal evasion, one-access external poisoning, and exacerbated OOD resilience. Section 4 assesses the privacy concern from the perspectives of top-down unauthorized disclosure and down-top confidentiality propagation. In Section 5, we address issues of fairness from underrepresented demographical attributes and participant resource disparity. Finally, Section 6 proposes prospective research directions that could fortify and refine this integration.

2. Foundation Models in FL

2.1. Selected Challenges in FL

FL faces several challenges, with ineffective training and constraints on resources being particularly significant.

Ineffective Training. In FL, clients often possess limited datasets inadequate for training robust models (Imteaj et al., 2021; Tuor et al., 2021). This scarcity is compounded by the privacy-preserving nature of FL, where data remains localized, preventing the pooling of resources to overcome individual data limitations (Kairouz et al., 2021; Bouacida & Mohapatra, 2021b). This leads to a trade-off between maintaining user privacy and obtaining sufficient data for effective model training.

Resource Limitation. In FL, the training of models is distributed across various client devices, which may vary greatly in their computational capabilities and communication abilities, leading to issues with efficiency and scalability (Wang et al., 2019; Shi et al., 2020; Wang et al., 2023b). This diversity in computational power significantly affects the overall efficiency of the FL process and influences the consistency and reliability of the aggregated global model (Khan et al., 2021; Zhang et al., 2022). Clients with limited capabilities might contribute less efficiently, while

simpler models, which are more feasible for such clients, might not significantly enhance the global model’s performance (Liu et al., 2022; Zeng et al., 2021).

2.2. Benefits Brought by FMs to FL

FMs can be leveraged to address the challenges of data scarcity and resource limitations in FL as follows:

Effective Training with Synthetic Data. One of the primary applications of FMs in FL is to generate synthetic data that closely mirrors real-world distributions (Eigenschink et al., 2021; Torres, 2018; Assefa et al., 2020). This capability of FMs to produce diverse and comprehensive datasets helps overcome data scarcity in FL by augmenting the limited data of individual clients, thereby enhancing model performance and generalization capabilities.

Reducing Computational Burden through Transfer Learning. FMs can also be used to initialize the model in FL with a strong base performance level (Bommasani et al., 2021a; Zhuang et al., 2023). Clients can fine-tune these models with their local data, which requires significantly less computational power than training a model from scratch. This approach not only speeds up the training process but also ensures that clients with lower computational capabilities can still effectively participate in the FL process.

Minimizing Communication Overhead with FM Feature Extraction. FMs, pre-trained on extensive datasets, offer a comprehensive understanding of data patterns (Zhou et al., 2023). FL clients can leverage these FMs as encoders, attaching simpler models that are more suited to their limited local data and resources (Yi et al., 2023). By uploading only the parameters of these simpler models during global communication phases, the communication load is significantly reduced, thereby enhancing the overall efficiency of the FL process (Kalra et al., 2023).

Facilitating Performance via Knowledge Distillation/Imitation Training. FMs can significantly enhance the performance of FL systems by serving as teachers (Xing et al., 2022; Zhu et al., 2021; Li & Wang, 2019; Yang et al., 2023; Liang et al., 2021). Through the process of knowl-

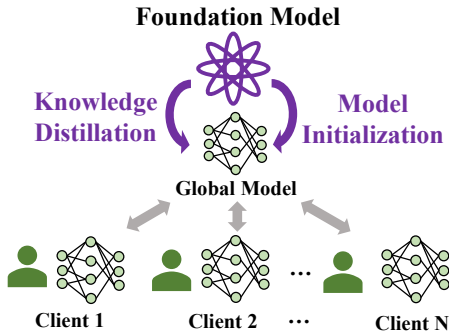


Figure 2. Foundation models at the server side.

edge distillation, FMs transfer their extensive and diverse understanding, acquired from training on large datasets to the FL models to solve their performance dilemma. This approach effectively bridges the gap between the advanced capabilities of FMs and the collaborative, distributed nature of FL, leading to more robust and capable FL systems.

2.3. Integration of FMs in FL

Integration on the Server. As shown in Fig. 2, the server utilizes FMs for synthetic data generation, driven by prompts collected from clients. This approach ensures privacy as the data generated is similar to the clients’ real data but does not directly expose it. The synthetic data serves dual purposes – initializing the global model and facilitating knowledge distillation from the FM to the aggregated global model. This approach addresses the scarcity of publicly available data that matches the local data utilized in (Lin et al., 2020; Li & Wang, 2019), closely resembling real client data, and is used to impart the broad knowledge of FMs to the global model.

Deployment at the Clients. Given the constraints of available FMs and resources, several clients may query the same FM to generate synthetic data that resembles their real data. As demonstrated in Fig. 3, this approach enables clients to produce data useful for initializing their individual models and facilitates imitation training, where the client models mimic the behaviors of FMs on the synthetic data. Besides, the clients could use FM parameters as a starting point for the following transfer learning on their local data, accelerating the convergence. Furthermore, clients can employ FMs as feature extractors, using the produced embeddings to enhance simpler, less demanding models. This method alleviates the computational load on client devices and ensures efficient training with their limited data.

3. Robustness

In this paper, we focus on the robustness challenges in machine learning including evasion attacks that deceive models

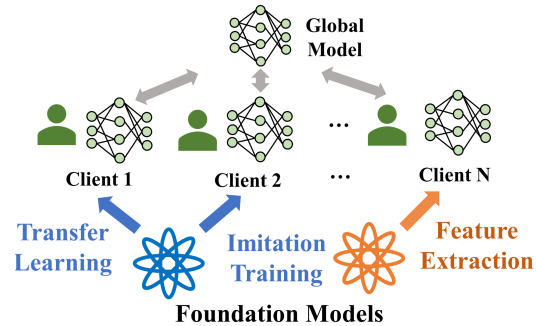


Figure 3. Foundation models at the client side.

during inference, poisoning attacks that corrupt training data, and the need for Out of Distribution (OOD) resilience to ensure models perform well on novel data, collectively posing significant threats to model integrity and performance.

Evasion attacks in FL and FM. Evasion attacks manipulate inputs to cause misclassification, notably during the inference phase. These can occur under full model knowledge (white-box) (Biggio et al., 2013; Carlini & Wagner, 2017) or without it (black-box) through surrogate models (Demontis et al., 2019) or gradient approximation (Chen et al., 2017). In Federated Learning, evasion attacks exploit the distributed nature of the system (Lyu et al., 2020b; Bouacida & Mohapatra, 2021a). In a black-box attack scenario, an attacker can compromise one client, gaining partial white-box access to craft adversarial inputs that deceive models at other clients. In the context of LLMs, these attacks, known as “jailbreaks,” involve manipulating prompts to exploit model biases, leading to outputs diverging from intended human values (Liu et al., 2023; Wei et al., 2023; Lapid et al., 2023; Chao et al., 2023). Advances in automatically engineering jailbreak prompts pose new challenges to LLM robustness (Chao et al., 2023; Lapid et al., 2023).

Poisoning attacks in FL and FM. Poisoning attacks, targeting the training phase, embed misbehavior into the victim model by inserting malicious instances into its training dataset, classified into mislabeled (Biggio et al., 2012; Gu et al., 2017) and clean-labeled (Shafahi et al., 2018) attacks. In FL, the distributed nature and privacy-preserving mechanisms introduce unique vulnerabilities where the attacker compromises a few clients and infuses malicious updates into the global model, gradually misleading its decision-making (Fang et al., 2020; Tolpegin et al., 2020). FMs, especially LLMs with in-context learning (ICL) capabilities (Dong et al., 2022), face risks from poisoning at inference time, where attackers embed malicious behavior in prompts, circumventing the need for direct manipulation of training data (Dong et al., 2022; Kandpal et al., 2023; Wang et al., 2023a).

OOD resilience in FL and FM. OOD data, not being an

adversarial attack, challenges the assumption of identical training and testing distributions in machine learning, leading to performance issues on unseen data (Hendrycks & Gimpel, 2017; Hsu et al., 2020). In FL, the non-IID nature of client data further complicates OOD robustness, affecting generalization and potentially reducing the global model’s effectiveness on non-participating clients (Reisizadeh et al., 2020; Guo et al., 2023). FMs, also struggle with OOD robustness due to the vast variety of real-world scenarios, making their ability to handle or recognize OOD queries critical for reliability (Li et al., 2023a). Specifically, LLMs might “hallucinate” or generate made-up responses to OOD queries (Wang et al., 2023a; Bubeck et al., 2023).

3.1. Accessible Universal Evasion

Accessible Universal Evasion Attacks leverage the intelligence of FMs in FM-FL, enabling attackers to initiate attacks that are effective across various domains, from images to text, without requiring specialized knowledge like optimization techniques. These attacks universally compromise FL systems by utilizing simple malicious prompts.

In the context of FM-FL, specifically large language models (LLMs), the interplay between client-server interactions introduces a novel vulnerability to evasion attacks. In this paradigm, clients contribute by providing prompts that describe their data needs for synthetic data generation by the server’s LLM, aiming to enhance the FL training process without compromising data privacy. However, this mechanism inherently exposes the system to potential evasion attacks through the injection of malicious prompts. Such attacks, often orchestrated by compromised clients, involve crafting prompts designed to manipulate the LLM into generating synthetic data that deviates from legitimate data distributions. This tactic, known as “jailbreaking” the LLM, can surreptitiously introduce biases, misinformation, and other detrimental elements into the FL system. The resulting corrupted synthetic data, once disseminated across the federated network, can compromise the integrity of the model training process, leading to a degradation of the system’s overall robustness. This vulnerability underscores the critical need for robust detection and mitigation strategies to safeguard against the subtle yet significant threat of accessible universal evasion, ensuring the preservation of the system’s integrity, reliability, and trustworthiness in the face of adversarial challenges.

3.2. One-Access External Poisoning

Integrating FMs to FL has given rise to innovative poisoning attack strategies, termed One-Access External Poisoning Attacks. These attacks allow external attackers to effectively inject poison into the FL ecosystem without persistent involvement in FL, thereby compromising the integrity of the system with minimal effort. These attacks are primarily

executed through two distinct strategies.

Strategy 1: Attack by a Third Party. In this strategy, the attacker is external to the FL system, targeting neither as a client nor as a part of the server. They target the FM before its integration into the FL system. For instance, they could fine-tune the FM with poisoning instructions (Xu et al., 2023) or insert malicious system prompts (Wang et al., 2023a). The compromised FM, obtained from an open source by the server or clients, is then used for *e.g.*, synthetic data generation and knowledge transfer. The synthetic data produced by the FM, carrying the poison, is utilized in the FL model initialization and in the mutual information-sharing phase, either on the server (Lin et al., 2020) or between clients (Li & Wang, 2019). Consequently, the malicious behaviors embedded in the poisoned data are transmitted to the client models during initialization and further reinforced through model aggregation on the server (Li et al., 2023b; 2024). On the other hand, the poison embedded within the FM could persist in FL models even after the local fine-tuning, starting with poisoned FM parameter initialization, potentially compromising the integrity of the FL models (Shen et al., 2021; Wang et al., 2022).

Strategy 2: Attacker Compromising a Client. This strategy diverges from the classic FL poisoning attacks, which generally require compromising numerous clients and maintaining persistent participation. In the FM-FL context, the attacker only needs to compromise a single client. By exploiting the shared FM, particularly harnessing its in-context learning abilities like those in LLMs, the attacker can compromise the FM. This is done by injecting malicious prompts into the LLM during queries, leading to the generation of malicious instances alongside benign ones (Wang et al., 2023a; Kandpal et al., 2023; Shi et al., 2023). This method allows the attacker to efficiently disseminate the poison to other clients querying the same FM, thereby compromising the overall integrity of the global model more effectively, due to the shared usage of the compromised FM.

3.3. Exacerbated Out of Distribution Resilience

Exacerbated OOD Resilience describes the intensified challenges in OOD resilience resulting from integrating FMs into FL. This occurs when the downstream tasks or requirements of FL extend beyond the FM’s expertise, leading FMs to generate outputs misaligned with FL’s diverse data distributions, thereby exacerbating the OOD resilience challenges within FL. Our preliminary investigations indicate that in the context of FM-FL, where the FM can be deployed either on the server or at clients, the exacerbated OOD resilience manifests in two areas: synthetic data generation and knowledge distillation.

With synthetic data generation. The central server or a single client typically lacks direct access to diverse client

data, which restricts its ability to fine-tune the FM for the specific domain of the FL’s downstream task. This limitation becomes crucial when the downstream task lies outside the scope of the FM’s pre-existing knowledge. In such scenarios, the FM may generate synthetic data that does not accurately reflect the overall distribution of the clients’ data or even “nonsense” data. Due to their pivotal role in FL training, misaligned synthetic data can hinder the FL training process, leading to suboptimal training, slower convergence, and reduced model performance.

Through knowledge distillation. Another critical aspect is the transfer of knowledge from the FM to the FL model through knowledge distillation. If the FM’s knowledge base does not encompass the specifics of FL’s downstream task, the distilled knowledge could be based on inaccuracies or fabrications. This transferred “made-up” knowledge can significantly impact the FL model’s performance.

3.4. Discussion

Accessible universal evasion attacks exploit LLM vulnerabilities with jailbreak prompts during the FL training phase, differing from traditional FL attacks that occur post-deployment. Attackers can repurpose jailbreak prompts, such as those used in popular LLMs like ChatGPT, bypassing the need for individual perturbations per adversarial input. This method stands out because traditional text-based adversarial inputs are less effective, significantly differentiating from conventional FL attack approaches.

The one-access external poisoning attack significantly increases FM-FL’s vulnerability by embedding threats directly in FMs, making subsequent threat transmission during FL independent of the attacker’s continuous involvement. This strategy can potentially compromise all clients in large-scale FL scenarios, where traditional methods requiring the compromise of many clients are impractical. Furthermore, it can evade existing FL defenses designed to counter conventional attacks by filtering outliers, as it leverages clean local datasets for client updates, thus avoiding detection.

The integration of FM into FL introduces unique OOD robustness challenges beyond the non-IID data variance seen in FL alone. In FM-FL, issues extend to the creation and use of synthetic data and knowledge transfer from FMs, complicating initial training and knowledge distillation. These added complexities can result in suboptimal training, slower model convergence, and decreased FL model performance.

4. Privacy

Privacy Leakage in FL and FMs. FL aims to enhance privacy through local training but faces privacy risks as model updates to a central server might reveal sensitive information. These updates are vulnerable to attacks like membership inference (Shokri et al., 2017), data pattern

reconstruction (Hitaj et al., 2017), and private data reverse-engineering (Zhu et al., 2019). Similarly, FMs, especially LLMs with ICL abilities, risk memorizing and leaking sensitive data during inference, exemplified by instances of revealing personal information, such as ChatGPT exposing email addresses from training materials (Wang et al., 2023a; Carlini et al., 2021; Yu et al., 2023).

Integrating FMs into FL presents novel and complex privacy leakage challenges. Within the FM-FL framework, we identify two primary pathways for potential privacy leakage: (1) Top-down unauthorized disclosure, where sensitive information embedded in the FM may be inadvertently passed on to FL clients, and (2) Down-top confidentiality propagation, wherein the FM may memorize sensitive data from local client queries.

4.1. Top-down unauthorized disclosure

In scenarios where a central FM (at the server) is utilized to enhance or initialize local models within an FL framework, there is a risk of transmitting sensitive information embedded within the FM to these local models. This could occur in several ways. Through Synthetic Data Generation: triggered by certain prompts, sensitive information might be embedded in the synthetic data generated by the FM, which is then transmitted to FL clients. Via Knowledge Distillation: FM could be used as a teacher to assist the learning of FL models. Knowledge of normal data, as well as the sensitive information memorized by the FM, could be transferred to the FL models through techniques like knowledge distillation. Furthermore, FL models that are initialized with FM parameters and subsequently fine-tuned on local datasets may also retain and memorize sensitive information.

4.2. Down-top confidentiality propagation

In scenarios where FMs are deployed to local clients to assist local training, multiple clients might access the same FM, due to resource limitations or data availability. In this setup, clients may request synthetic data from the FM, such as an LLM, aiming to bolster their local training efforts. These requests typically embed sensitive details reflective of the clients’ own training datasets, leading to the FM inadvertently memorizing this sensitive information during its inference processes. As a result, when another client later accesses the same FM for synthetic data or to facilitate knowledge transfer, there exists a risk that the sensitive information memorized from previous queries could be unintentionally exposed. This scenario underscores a significant privacy concern, where the FM’s capacity to retain information from individual queries can facilitate unintended information leakage among clients. Consequently, the practice of sharing FMs among multiple clients for local

training enhancement introduces a potential vector for the inadvertent dissemination of sensitive information across the federated network.

4.3. Discussion

In FL, privacy leakage predominantly occurs from local clients to the central server. Sensitive information, such as membership, data distribution, and even specific training data, are reverse-engineered from local updates through inference attacks. However, for FM-FL, the leakage is more complex, involving not just the transfer of model information, but also the potential embedding and transmission of sensitive data within synthetic data generated by FMs or through knowledge distillation processes. FM-FL introduces a two-way leakage risk: from FM to FL (where sensitive information within the FM can be transferred to FL clients) and from FL to FM (where sensitive local client data can be memorized by the FM during queries for synthetic data or knowledge). The advanced capabilities of FMs and the interaction between FM and FL makes it a more dynamic and multifaceted privacy challenge.

5. Fairness

Fairness in FL and FMs. The decentralized nature and data heterogeneity inherent in FL present challenges in achieving group fairness (Dwork et al., 2012; Ezzeldin et al., 2023). The non-IID data distribution across clients can lead to biases and affect model generalization (He & Garcia, 2009; Krawczyk, 2016). Additionally, difficulty in representing all demographics in local training data leads to inequities regarding sensitive attributes. client participation variability further contributes to potential biases, particularly favoring data characteristics of consistently involved demographics (Lyu et al., 2020a). FMs trained on internet-sourced data risk inheriting biases, as this data may not accurately represent diverse human language and behavior (Bommasani et al., 2021b; Si et al., 2023), potentially amplifying societal biases. The universal application of FMs may fail to reflect the diverse ethical, cultural, and linguistic nuances of global communities, leading to less equitable models in certain contexts. Additionally, the complexity of FMs hinders transparency in their decision-making, impeding efforts to identify and mitigate inherent biases.

5.1. Underrepresented Demographical Attributes

Demographical Attributes, such as gender, race, and age, define diverse demographic groups. Underrepresented Demographical Attributes, referring to insufficient representation of demographical attributes in datasets and studies, can result in biased and stereotypical outputs from models. This underrepresentation leads to models delivering inequitable outputs for these marginalized groups. The concern is heightened in the era of FMs, given their broad

societal impact and the inherent biases they often learn from extensive, internet-sourced datasets.

Integrating FMs into FL can significantly impact fairness, especially regarding underrepresented demographical attributes. **Challenge 1: Transmission of Inherent Unfairness from FMs to FL.** The use of FMs for synthetic data generation and knowledge distillation in FL, whether on servers or at client sites, risks transmitting any existing biases from FMs to FL models. This transmission occurs mainly in two ways. Firstly, biases within FMs, including those against underrepresented demographics, can be ingrained in the synthetic data they generate, affecting the initialization of local FL models. Secondly, if knowledge distillation involves data with underrepresented attributes, it can trigger and propagate FM’s unfairness to local models. **Challenge 2: Unfairness Injected into FM from Compromised FL Clients.** When FMs are deployed at client locations for local training assistance, new risks arise. Adversaries might introduce demographically unfair content through prompts to FMs, leading to unfair model outputs. This is particularly concerning with models capable of in-context learning, like Large Language Models, as the induced bias can quickly spread to other clients querying the same FM, amplifying unfairness across the network.

5.2. Participate Resource Disparity

Participation Resource Disparity refers to the unequal distribution of resources among clients (participants) in federated learning, including local data volume, computational power, memory, network connectivity, and funding. This disparity divides clients into two categories: **resource-abundant clients**, with abundant resources, and **resource-constrained clients**, with limited resources. The participation resource disparity influences each participant’s capacity to process data, develop sophisticated models, and contribute efficiently to the collaborative system, thereby creating imbalances in terms of influence and outcomes within FL. This disparity is further emphasized with the integration of FMs, as efficiently running these models requires significant resources, particularly in terms of memory and computing power. Consequently, in the FM-FL environment, this resource disparity amplifies differences in local data size, the capacity for effective local training, and communication frequency. As a result, it leads to considerable unfairness in both the training process and the distribution of benefits among participants, with resource-abundant clients potentially dominating the learning outcomes while resource-constrained clients may not benefit equally.

In FM-FL scenarios where FMs are integrated with clients, resource-abundant clients, with the capability to build complex models and run FMs locally, gain a significant advantage. They have unrestricted access to FMs, allowing them to greatly enhance their local training sets, train their lo-

cal models effectively, acquire extensive knowledge from the FM, and participate more actively in FL. Consequently, these resource-abundant clients tend to dominate the FL training process. The substantial contributions heavily influence the aggregated global model, potentially leading to a bias that favors the data characteristics and preferences of these larger entities. Conversely, resource-constrained clients, despite contributing to the learning process, might not be able to reap the same benefits. They often lack the resources to utilize the weights and insights derived from resource-abundant clients effectively. This asymmetry in learning and contribution means that resource-constrained clients can aid the learning of resource-abundant clients but cannot benefit equally from the reverse. Plus, clients lacking the resources to run FMs locally resort to using API-based services, sharing FMs with other resource-constrained clients. However, funding limitations and restricted FM availability present similar obstacles to direct FM integration. The constraints hinder a client's ability to equally contribute to and benefit from the FL system, perpetuating the disparity in participation and advantages within FL.

5.3. Dissusion

The difference in fairness issues between FL and FM-FL lies in the source and complexity of biases. In FL, biases arise from data heterogeneity and client participation variability, leading to models potentially favoring certain client data patterns. In FM-FL, the integration of FMs introduces additional layers of bias, stemming from the FM's inherent biases in training and in-context learning abilities, which can further propagate to FL models through mechanisms like synthetic data generation and knowledge distillation. This integration results in more complex fairness challenges in FM-FL compared to FL alone.

Besides, incorporating FM introduces a new dimension of unfairness, resulting in aggregation bias and undermining the fairness of participation in FL. The disproportionate influence of resource-abundant clients skews the learning process, potentially marginalizing resource-constrained clients and creating a feedback loop that further entrenches the dominance of larger entities in the FL ecosystem. Addressing this disparity is crucial for maintaining the collaborative and equitable spirit of Federated Learning.

6. Future Direction

Considering the integration of FMs into FL impacts robustness, privacy, and fairness, future research should focus on evaluating these effects and devising solutions. Interdisciplinary efforts combining machine learning, cybersecurity, and ethics are vital for developing robust, transparent, and ethical FL systems, necessitating ongoing innovation and societal impact assessments.

6.1. Robustness

For research investigation the robustness of FM-FL, we propose several future research directions to address the threats of FM-FL discussed in Sec. 3:

1. Assessing the Susceptibility of FM-FL under the novel threats. A comprehensive assessment of the susceptibility of FM-FL under the novel threats is essential. This evaluation should address the effects of evasion attacks, particularly through "jailbreak" prompts in LLMs, on FL model performance, convergence, and misinformation spread, including societal impacts and attacker cost-effectiveness. It must also investigate the influence of compromised FMs, like those in backdoor attacks, on FL integrity, focusing on vulnerabilities introduced by synthetic data. Assessing the attacker's effort in poisoning FL systems, including FM fine-tuning and exploiting ICL, is critical. Lastly, measuring the discrepancy between FM-generated and actual client data, alongside the FM's performance on FL-specific tasks, is vital for identifying enhancements in data generation and knowledge distillation to boost FL robustness and reliability.

2. Effectiveness of Current FL Defenses. A comprehensive evaluation of the current defense mechanisms in FL against emerging threats is essential. This assessment should encompass the effectiveness of robust aggregation strategies and post-training detection methods in countering these novel threats. Such evaluations will provide insights into the adequacy of existing defense approaches and the necessity for the development of robust defensive strategies tailored to the FM-FL context.

3. Strengthening the Robustness and Security of FM-FL. Explore prompt validation techniques to identify and eliminate malicious inputs before they're processed by FMs, such as using anomaly detection to spot harmful prompts, and preventing adversarial data manipulation. Address the limitations of current FL defenses designed for decentralized threats, recognizing that new attacks target FMs directly and do not appear as statistical outliers. Develop defenses against centralized attacks and those compromising numerous clients without relying on anomaly detection. Investigate dynamic knowledge distillation methods tailored to FL's varied tasks, focusing on enhancing knowledge transfer where FMs are confident and leveraging learner models for areas where FMs are less certain, ensuring effective learning across FL scenarios.

6.2. Privacy

Addressing the complex privacy challenges in FM-FL requires a comprehensive approach that encompasses both impact assessment and solution development, including:

1. Privacy Exposure Assessment. This involves assessing privacy leakage by evaluating the presence of sensitive data

in FM-generated synthetic datasets, transfer mechanisms of sensitive data during the knowledge distillation, and assessing the persistence of privacy risks in FL models fine-tuned with local data. It also includes investigating vulnerabilities to inference attacks to understand if leaked sensitive information can be reverse-engineered.

2. Privacy Enhancement Techniques. This involves crafting methods that not only detect and anonymize sensitive information effectively but also preserve the utility of the data for FL models. The focus should be on striking a balance between privacy preservation and maintaining the richness of data that supports the development of robust and accurate models. This subsection also explores the need for designing algorithms capable of unlearning or discarding sensitive data post-training, ensuring that FL can adaptively protect privacy without compromising on performance.

3. Privacy-Preserving Knowledge Transfer and Effective Unlearning Algorithms. Develop knowledge distillation algorithms that are adaptive to privacy requirements, specifically obscuring sensitive data details. Design algorithms for effective unlearning, enabling FL systems to forget or discard sensitive data during/post-training. These unlearning algorithms must be optimized for efficiency, minimizing their impact on the computational resources and overall performance of the FL system, thus maintaining operational integrity while safeguarding privacy.

4. Ethical Guidelines and Policy Development for FM-FL. Formulate ethical guidelines specifically for FM-FL integration, addressing privacy and data protection challenges. Develop comprehensive policy frameworks to govern data collection, use, and sharing within FM-FL systems, ensuring a balance between innovation and privacy.

6.3. Fairness

To foster a more equitable and bias-aware FM-FL ecosystem, we propose a series of strategic research directions aimed at addressing fairness issues in the FM-FL:

1. Bias Transmission Assessment from FMs to FL This involves a thorough examination of how biases present in synthetic data and during the knowledge distillation processes might be carried over into FL model outputs, with a particular focus on the impact these biases have on underrepresented groups. Understanding how biases transfer from FMs to FL enables the creation of strategies to mitigate these biases, making FL models fairer and more representative.

2. Fair Data Creation and Knowledge Distillation By developing methods that meticulously scrutinize and neutralize biases in data produced by FMs, we ensure that all demographic groups are equitably represented. This effort is complemented by the innovation of knowledge distillation techniques that are specifically designed to identify and

correct biases, thereby facilitating a fairer transfer of knowledge to FL models. This dual strategy enhances FL system integrity and inclusivity by reducing bias at both the data generation and training phases.

3. Impact Analysis of Resource-Abundant Clients in FL This includes understanding how the contributions from these well-resourced clients might disproportionately shape the FL model’s development and outcomes. For instance, whether the dominance of resource-rich participants skews model performance and fairness, particularly affecting the less-resourced, or resource-limited clients.

4. Equitable Participation and Collaborative Learning We propose the creation of equitable participation mechanisms and collaborative learning models that bridge the gap between resource-diverse clients. This involves implementing strategies such as round-robin participation and fair access policies to FMs, ensuring all clients, irrespective of their resource capacity, contribute to and benefit from the federated learning process equally. Furthermore, we propose the exploration of collaborative learning frameworks that allow resource-constrained clients to leverage the advanced capabilities and insights of resource-abundant counterparts. Such models could facilitate the efficient transfer of knowledge to smaller, more communication-efficient models suitable for the FMFL environment, promoting a more inclusive and effective learning ecosystem. This dual strategy aims to democratize access to FMs in the FL context, enhancing the robustness, fairness, and overall performance of FMFL systems.

5. Formulation of Fairness-Driven Policies in FM-FL Integration. Policies should mitigate biases and ensure equitable participation in FM-FL integrations, fostering an inclusive digital ecosystem. They must prioritize resource distribution, fair data representation, and fairness metrics evaluation, steering FM-FL systems towards technological excellence that benefits a diverse society and promotes an accessible, fair digital future for all.

7. Conclusion

In conclusion, this paper presents an in-depth analysis of the robustness, privacy, and fairness challenges arising from the integration of FMs with FL, which remains underexplored in current research. Additionally, we shed lights on potential future directions in this field, advocating for the development of a responsible FM-FL ecosystem. This ecosystem would not only leverage the strengths of FMs to enhance FL but also prioritize the establishment of secure, reliable, and ethical practices to safeguard against any adverse implications of their deployment. Moving forward, future research must further explore the FM-FL relationship, ensuring advancements are driven by a commitment to societal welfare and justice.

References

- Assefa, S. A., Dervovic, D., Mahfouz, M., Tillman, R. E., Reddy, P., and Veloso, M. Generating synthetic data in finance: opportunities, challenges and pitfalls. In *Proceedings of the First ACM International Conference on AI in Finance*, pp. 1–8, 2020.
- Biggio, B., Nelson, B., and Laskov, P. Poisoning attacks against support vector machines. In *ICML*, 2012.
- Biggio, B., Corona, I., Maiorca, D., Nelson, B., Srndic, N., Laskov, P., Giacinto, G., and Roli, F. Evasion attacks against machine learning at test time. In *ECML PKDD*, 2013.
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosse-lut, A., Brunskill, E., et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021a.
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R. B., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosse-lut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N. S., Chen, A. S., Creel, K., Davis, J. Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K., Goodman, N. D., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D. E., Hong, J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P. W., Krass, M. S., Krishna, R., Kuditipudi, R., and et al. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258, 2021b.
- Bouacida, N. and Mohapatra, P. Vulnerabilities in federated learning. *IEEE Access*, 2021a.
- Bouacida, N. and Mohapatra, P. Vulnerabilities in federated learning. *IEEE Access*, 9:63229–63249, 2021b.
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S. M., Nori, H., Palangi, H., Ribeiro, M. T., and Zhang, Y. Sparks of artificial general intelligence: Early experiments with GPT-4. *CoRR*, abs/2303.12712, 2023.
- Carlini, N. and Wagner, D. A. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy*, 2017.
- Carlini, N., Tramèr, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T. B., Song, D., Erlingsson, Ú., Oprea, A., and Raffel, C. Extracting training data from large language models. In *30th USENIX Security Symposium, USENIX Security*, 2021.
- Chao, P., Robey, A., Dobriban, E., Hassani, H., Pappas, G. J., and Wong, E. Jailbreaking black box large language models in twenty queries. *CoRR*, abs/2310.08419, 2023.
- Chen, P., Zhang, H., Sharma, Y., Yi, J., and Hsieh, C. ZOO: zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *AISec@CCS*, 2017.
- Demontis, A., Melis, M., Pintor, M., Jagielski, M., Biggio, B., Oprea, A., Nita-Rotaru, C., and Roli, F. Why do adversarial attacks transfer? explaining transferability of evasion and poisoning attacks. In *28th USENIX Security Symposium, USENIX Security 2019, Santa Clara, CA, USA, August 14-16, 2019*, 2019.
- Dong, Q., Li, L., Dai, D., Zheng, C., Wu, Z., Chang, B., Sun, X., Xu, J., and Sui, Z. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. S. Fairness through awareness. In *Innovations in Theoretical Computer Science*, 2012.
- Eigenschink, P., Vamosi, S., Vamosi, R., Sun, C., Reutterer, T., and Kalcher, K. Deep generative models for synthetic data. 2021.
- Ezzeldin, Y. H., Yan, S., He, C., Ferrara, E., and Avestimehr, A. S. Fairfed: Enabling group fairness in federated learning. In *AAAI*, 2023.
- Fang, M., Cao, X., Jia, J., and Gong, N. Z. Local model poisoning attacks to byzantine-robust federated learning. In *29th USENIX Security Symposium, USENIX Security*, 2020.
- Gu, T., Dolan-Gavitt, B., and Garg, S. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *CoRR*, abs/1708.06733, 2017.
- Guo, Y., Guo, K., Cao, X., Wu, T., and Chang, Y. Out-of-distribution generalization of federated learning via implicit invariant relationships. In *ICML*, 2023.
- He, H. and Garcia, E. A. Learning from imbalanced data. *IEEE Trans. Knowl. Data Eng.*, 2009.
- Hendrycks, D. and Gimpel, K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- Hitaj, B., Ateniese, G., and Pérez-Cruz, F. Deep models under the GAN: information leakage from collaborative deep learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS*, 2017.

- Hsu, Y.-C., Shen, Y., Jin, H., and Kira, Z. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- Imteaj, A., Thakker, U., Wang, S., Li, J., and Amini, M. H. A survey on federated learning for resource-constrained iot devices. *IEEE Internet of Things Journal*, 9(1):1–24, 2021.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- Kalra, S., Wen, J., Cresswell, J. C., Volkovs, M., and Tizhoosh, H. Decentralized federated learning through proxy model sharing. *Nature communications*, 14(1):2899, 2023.
- Kandpal, N., Jagielski, M., Tramèr, F., and Carlini, N. Backdoor attacks for in-context learning with language models. *CoRR*, abs/2307.14692, 2023.
- Khan, L. U., Saad, W., Han, Z., Hossain, E., and Hong, C. S. Federated learning for internet of things: Recent advances, taxonomy, and open challenges. *IEEE Communications Surveys & Tutorials*, 23(3):1759–1799, 2021.
- Krawczyk, B. Learning from imbalanced data: open challenges and future directions. *Prog. Artif. Intell.*, 2016.
- Lapid, R., Langberg, R., and Sipper, M. Open sesame! universal black box jailbreaking of large language models. *CoRR*, abs/2309.01446, 2023.
- Li, D. and Wang, J. Fedmd: Heterogenous federated learning via model distillation. *CoRR*, abs/1910.03581, 2019. URL <http://arxiv.org/abs/1910.03581>.
- Li, X., Fang, Y., Liu, M., Ling, Z., Tu, Z., and Su, H. Distilling large vision-language model with out-of-distribution generalizability. In *ICCV*, 2023a.
- Li, X., Wang, S., Wu, C., Zhou, H., and Wang, J. Backdoor threats from compromised foundation models to federated learning. *FL@FM-NeurIPS 23*, 2023b.
- Li, X., Wu, C., and Wang, J. Unveiling backdoor risks brought by foundation models in heterogeneous federated learning. In *PAKDD*, 2024.
- Liang, K. J., Hao, W., Shen, D., Zhou, Y., Chen, W., Chen, C., and Carin, L. Mixkd: Towards efficient distillation of large-scale language models. In *9th International Conference on Learning Representations, ICLR, 2021*.
- Lin, T., Kong, L., Stich, S. U., and Jaggi, M. Ensemble distillation for robust model fusion in federated learning. In *NeurIPS*, 2020.
- Liu, Y., Deng, G., Xu, Z., Li, Y., Zheng, Y., Zhang, Y., Zhao, L., Zhang, T., and Liu, Y. Jailbreaking chatgpt via prompt engineering: An empirical study. *CoRR*, abs/2305.13860, 2023.
- Liu, Z., Chen, Y., Zhao, Y., Yu, H., Liu, Y., Bao, R., Jiang, J., Nie, Z., Xu, Q., and Yang, Q. Contribution-aware federated learning for smart healthcare. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 12396–12404, 2022.
- Lyu, L., Xu, X., Wang, Q., and Yu, H. Collaborative fairness in federated learning. In *Federated Learning - Privacy and Incentive*. 2020a.
- Lyu, L., Yu, H., and Yang, Q. Threats to federated learning: A survey. *CoRR*, abs/2003.02133, 2020b.
- Reisizadeh, A., Farnia, F., Pedarsani, R., and Jadbabaie, A. Robust federated learning: The case of affine distribution shifts. In *NeurIPS*, 2020.
- Shafahi, A., Huang, W. R., Najibi, M., Suci, O., Studer, C., Dumitras, T., and Goldstein, T. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *NeurIPS*, 2018.
- Shen, L., Ji, S., Zhang, X., Li, J., Chen, J., Shi, J., Fang, C., Yin, J., and Wang, T. Backdoor pre-trained models can transfer to all. In *CCS '21: 2021 ACM SIGSAC Conference on Computer and Communications Security, Virtual Event, Republic of Korea, November 15 - 19, 2021*, 2021.
- Shi, J., Liu, Y., Zhou, P., and Sun, L. Badgpt: Exploring security vulnerabilities of chatgpt via backdoor attacks to instructgpt. *CoRR*, abs/2304.12298, 2023.
- Shi, W., Zhou, S., Niu, Z., Jiang, M., and Geng, L. Joint device scheduling and resource allocation for latency constrained wireless federated learning. *IEEE Transactions on Wireless Communications*, 20(1):453–467, 2020.
- Shokri, R., Stronati, M., Song, C., and Shmatikov, V. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy, SP*, 2017.
- Si, C., Gan, Z., Yang, Z., Wang, S., Wang, J., Boyd-Graber, J. L., and Wang, L. Prompting GPT-3 to be reliable. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*, 2023.

- Tolpegin, V., Truex, S., Gursoy, M. E., and Liu, L. Data poisoning attacks against federated learning systems. In *Computer Security - ESORICS 2020 - 25th European Symposium on Research in Computer Security, ESORICS 2020, Guildford, UK, September 14-18, 2020, Proceedings, Part I*, 2020.
- Torres, D. G. *Generation of synthetic data with generative adversarial networks*. PhD thesis, Royal Institute of Technology Stockholm, Sweden, 2018.
- Tuor, T., Wang, S., Ko, B. J., Liu, C., and Leung, K. K. Overcoming noisy and irrelevant data in federated learning. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 5020–5027. IEEE, 2021.
- Wang, B., Chen, W., Pei, H., Xie, C., Kang, M., Zhang, C., Xu, C., Xiong, Z., Dutta, R., Schaeffer, R., Truong, S. T., Arora, S., Mazeika, M., Hendrycks, D., Lin, Z., Cheng, Y., Koyejo, S., Song, D., and Li, B. Decodingtrust: A comprehensive assessment of trustworthiness in GPT models. *CoRR*, abs/2306.11698, 2023a.
- Wang, J., Zeng, S., Long, Z., Wang, Y., Xiao, H., and Ma, F. Knowledge-enhanced semi-supervised federated learning for aggregating heterogeneous lightweight clients in iot. In *Proceedings of the 2023 SIAM International Conference on Data Mining (SDM)*, pp. 496–504. SIAM, 2023b.
- Wang, S., Tuor, T., Salonidis, T., Leung, K. K., Makaya, C., He, T., and Chan, K. Adaptive federated learning in resource constrained edge computing systems. *IEEE journal on selected areas in communications*, 37(6):1205–1221, 2019.
- Wang, S., Nepal, S., Rudolph, C., Grobler, M., Chen, S., and Chen, T. Backdoor attacks against transfer learning with pre-trained deep learning models. *IEEE Trans. Serv. Comput.*, 2022.
- Wei, A., Haghtalab, N., and Steinhardt, J. Jailbroken: How does LLM safety training fail? *CoRR*, abs/2307.02483, 2023.
- Xing, H., Xiao, Z., Qu, R., Zhu, Z., and Zhao, B. An efficient federated distillation learning system for multitask time series classification. *IEEE Trans. Instrum. Meas.*, 2022.
- Xu, J., Ma, M. D., Wang, F., Xiao, C., and Chen, M. Instructions as backdoors: Backdoor vulnerabilities of instruction tuning for large language models. *CoRR*, abs/2305.14710, 2023.
- Yang, X., Li, Q., Zhang, C., and Woodland, P. C. Knowledge distillation from multiple foundation models for end-to-end speech recognition. *CoRR*, abs/2303.10917, 2023.
- Yi, L., Wang, G., Liu, X., Shi, Z., and Yu, H. Fedgh: Heterogeneous federated learning with generalized global header. *arXiv preprint arXiv:2303.13137*, 2023.
- Yu, W., Pang, T., Liu, Q., Du, C., Kang, B., Huang, Y., Lin, M., and Yan, S. Bag of tricks for training data extraction from language models. In *ICML*, 2023.
- Zeng, H., Zhou, T., Guo, Y., Cai, Z., and Liu, F. Fedcav: contribution-aware model aggregation on distributed heterogeneous data in federated learning. In *Proceedings of the 50th International Conference on Parallel Processing*, pp. 1–10, 2021.
- Zhang, T., Gao, L., He, C., Zhang, M., Krishnamachari, B., and Avestimehr, A. S. Federated learning for the internet of things: Applications, challenges, and opportunities. *IEEE Internet of Things Magazine*, 5(1):24–29, 2022.
- Zhou, C., Li, Q., Li, C., Yu, J., Liu, Y., Wang, G., Zhang, K., Ji, C., Yan, Q., He, L., Peng, H., Li, J., Wu, J., Liu, Z., Xie, P., Xiong, C., Pei, J., Yu, P. S., and Sun, L. A comprehensive survey on pretrained foundation models: A history from BERT to chatgpt. *CoRR*, 2023.
- Zhu, L., Liu, Z., and Han, S. Deep leakage from gradients. In *NeurIPS*, 2019.
- Zhu, Z., Hong, J., and Zhou, J. Data-free knowledge distillation for heterogeneous federated learning. In *Proceedings of the 38th International Conference on Machine Learning, ICML*, 2021.
- Zhuang, W., Chen, C., and Lyu, L. When foundation model meets federated learning: Motivations, challenges, and future directions. *arXiv preprint arXiv:2306.15546*, 2023.