

Synergizing Foundation Models and Federated Learning: A Survey

Anonymous ACL submission

Abstract

The recent development of Foundation Models (FMs), represented by large language models, vision transformers, and multimodal models, has been making a significant impact on both academia and industry. Compared with small-scale models, FMs have a much stronger demand for high-volume data during the pre-training phase. Although general FMs can be pre-trained on data collected from open sources such as the Internet, domain-specific FMs need proprietary data, posing a practical challenge regarding the amount of data available due to privacy concerns. Federated Learning (FL) is a collaborative learning paradigm that breaks the barrier of data availability from different participants. Therefore, it provides a promising solution to customize and adapt FMs to a wide range of domain-specific tasks using distributed datasets whilst preserving privacy. This survey paper discusses the potentials and challenges of synergizing FL and FMs and summarizes core techniques, future directions, and applications.

1 Introduction

The landscape of Artificial Intelligence (AI) has been revolutionized by the emergence of *Foundation Models (FMs)* (Bommasani et al., 2021), such as BERT (Devlin et al., 2019), GPT series (Brown et al., 2020; OpenAI, 2022, 2024), and LLaMA series (Touvron et al., 2023a,b) in Natural Language Processing (NLP); ViTs (Dosovitskiy et al., 2021) and SAM (Kirillov et al., 2023) in Computer Vision (CV); CLIP (Radford et al., 2021), DALL-E (Ramesh et al., 2021), Gemini (Google, 2023), and GPT-4o in multimodal applications. These FMs have become pivotal in a myriad of AI applications across diverse domains. Their superb capability to generalize across tasks and domains stems from their pre-training on extensive datasets (Gunasekar et al., 2023), which imbues them with a profound understanding of language, vision, and multimodal data.

While general-purpose FMs can leverage openly accessible data from the Internet, domain-specific FMs require proprietary data. It is, however, challenging to collect vast amounts of proprietary data and perform centralized pre-training or fine-tuning for domain-specific FMs, due to privacy restrictions (Jo and Gebru, 2020; GDPR, 2016; CCPA, 2023). Particularly in domains such as law, healthcare, and finance, where data is inherently privacy-sensitive, there is a pressing need for stringent privacy safeguards. Furthermore, given that data often constitutes a pivotal asset for enterprises, its widespread distribution is prohibitive. Consequently, there is an urgent need for novel strategies to handle data availability and facilitate model training, thereby unlocking the potential of domain-specific FMs whilst respecting data privacy.

To address the challenges associated with data privacy in model training, Federated Learning (FL) (McMahan et al., 2017) has emerged as a promising paradigm. FL facilitates collaborative model training across decentralized clients without the need for sharing raw data, thus ensuring privacy preservation. Concretely, FL encompasses periodic interactions between the server and decentralized clients for the exchange of trainable model parameters, without the requirement for private client data. Recognizing such a benefit, *integrating FMs with FL presents a compelling solution for domain-specific FMs* (Zhuang et al., 2023; Yu et al., 2023d).

Despite the potential synergies between FL and FMs, the field is still nascent, lacking a comprehensive understanding of challenges, methodologies, and directions. This survey aims to bridge this gap by providing a thorough exploration of the integration of FMs and FL. We delve into the motivations and challenges of combining these two paradigms, highlight representative techniques, and discuss future directions and applications. By elucidating the intersection of FL and FMs, we aim to catalyze further research and innovation in this burgeoning

083	ing area, ultimately advancing the development of	conversational agents, robotic control, and gaming.	133
084	privacy-aware, domain-specific FMs.		
085	The paper continues as follows: The next section	2.2 Federated Learning	134
086	introduces background on FMs and FL. Section 3	FL (McMahan et al., 2017) is a learning paradigm	135
087	presents the motivation and challenges for syner-	that enables a collection of clients to collabora-	136
088	gizing FMs and FL. Section 4 highlights represen-	tively learn a shared global model by leveraging	137
089	tative techniques. Before concluding, we discuss	their private datasets in a distributed manner, as-	138
090	representative future directions in Section 5.	sisted by the coordination of a central server. The	139
091		most representative algorithms in the FL literature	140
092	2 Background	are the FedAvg-family algorithms (McMahan et al.,	141
093		2017; Reddi et al., 2021). The standard FedAvg in-	142
094	2.1 Foundation Models	volves periodic interactions between the server and	143
095	An FM is a model that can be adapted to a wide	decentralized clients to exchange trainable model	144
096	array of tasks through fine-tuning after initial pre-	parameters. Many variants have been proposed to	145
097	training (Bommasani et al., 2021). The lifecycle	tackle issues such as convergence and local data	146
098	of FMs typically involves pre-training on extensive	heterogeneity (Diao et al., 2021). For example,	147
099	generic data to establish the basis of their abili-	FedProx (Li et al., 2020) and FedDyn (Acar et al.,	148
100	ties (Bubeck et al., 2023), followed by adaptation	2021) introduce regularizer terms to penalize client	149
101	to downstream tasks such as domain-specific ques-	updates that are far away from the server model. A	150
102	tion answering (Zhang et al., 2023d), and ultimately	general framework FedOpt (Reddi et al., 2021) uni-	151
103	application in various domains.	fies adaptive optimizers (<i>Adam, Yogi, etc.</i>) and	152
104	FMs have sparked a significant paradigm shift in	demonstrates superior convergence speed when	153
105	various fields of AI such as NLP, CV, speech and	compared to the naive FedAvg.	154
106	acoustics, and beyond. In the realm of NLP, the	FL offers an efficient privacy-preserving way	155
107	most prominent example is Large Language Mod-	to train models on large-scale and diverse data	156
108	els (LLMs) with substantial parameter sizes (Zhao	(Kairouz et al., 2021), leading to its application	157
109	et al., 2023). These models, such as ChatGPT and	across various domains such as healthcare (Lincy	158
110	GPT-4 (OpenAI, 2022, 2024), demonstrate excep-	and Kowshalya, 2020; Rieke et al., 2020; Joshi	159
111	tional abilities in natural language understanding	et al., 2022), finance (Chatterjee et al., 2023; Liu	160
112	and generation, enabling them to comprehend and	et al., 2023b), and smart cities (Ramu et al., 2022;	161
113	respond to user inputs with remarkable contextual	Pandya et al., 2023).	162
114	relevance. This capability proves invaluable in ap-	3 FM-FL: Motivation & Challenges	163
115	plications like customer service, virtual assistants,	In this section, we first motivate the synergy of	164
116	and chatbots, where effective communication is	FMs and FL (Section 3.1), then summarize the key	165
117	paramount. Moreover, LLMs eliminate the need	challenges (Section 3.2).	166
118	for training models from scratch for specific tasks,		
119	be it machine translation, document summarization,	3.1 Motivation	167
120	text generation, or other language-related tasks.	The integration of FMs and FL represents a com-	168
121	In the realm of CV and other modalities, FMs	pelling collaboration that leverages each other’s	169
122	have also made remarkable progress. Vision Trans-	strengths to address their respective limitations,	170
123	formers (ViTs) (Dosovitskiy et al., 2021) segment	embodying a complementary relationship (Zhuang	171
124	images into distinct patches, which serve as in-	et al., 2023; Li and Wang, 2024).	172
125	puts for transformer architectures. SAM (Kirillov	FL expands data availability for FMs. By lever-	173
126	et al., 2023) can segment anything in images ac-	aging data from a wide range of sources in a	174
127	cording to the input prompts. CLIP (Radford et al.,	privacy-preserving manner, FL makes it possible	175
128	2021) bridges the gap between text and images	to build models on sensitive data in specific do-	176
129	through contrastive learning. DALL-E, proposed	domains, such as healthcare (Lincy and Kowshalya,	177
130	by Ramesh et al. (2021), generates images from	2020; Joshi et al., 2022; Rieke et al., 2020) and	178
131	textual descriptions, expanding the possibilities of	finance (Chatterjee et al., 2023; Liu et al., 2023b).	179
132	creative image generation. Additionally, models	This enhances the diversity and volume of training	180
	like GAtO (Reed et al., 2022), exhibit versatility		
	by being applicable across various tasks such as		

181 data, improving model robustness and adaptability. 182 Moreover, FL enables the integration of personal 183 and task-specific data, allowing FMs to be cus- 184 tomized for personal applications. For instance, 185 Google has trained next-word-prediction language 186 models on mobile keyboard input data with FL to 187 improve user experience (Xu et al., 2023; Bonawitz 188 et al., 2021).

189 **FMs boost FL with feature representation and** 190 **few-shot learning capabilities.** By pre-training 191 on large-scale generic data, FMs acquire essential 192 knowledge and understanding capabilities (Brown 193 et al., 2020), providing multiple benefits to FL. 194 Firstly, they benefit FL systems by offering ad- 195 vanced feature representations and learning capa- 196 bilities from the outset. Secondly, leveraging the 197 pre-learned knowledge of FMs can accelerate the 198 FL process, enabling efficient and effective adapta- 199 tion to specific tasks with minimal additional train- 200 ing. Thirdly, FMs’ powerful generative capabilities 201 could help FL overcome the data heterogeneity 202 challenge by synthesizing extra data, thus acceler- 203 ating model convergence (Huang et al., 2024).

204 3.2 Core Challenges

205 In this part, we discuss challenges emerging from 206 the FM-FL marriage in three aspects: *efficiency*, 207 *adaptability*, as well as *trustworthiness*.

208 **Efficiency Challenges.** Efficiency challenges 209 stem from the mismatch between the significant 210 resource demands of FM training and the limited, 211 heterogeneous system resources (*e.g.*, mobile de- 212 vices) within FL systems, such as communication 213 bandwidth, computational power, and memory (Su 214 et al., 2023). The communication bottleneck of FL 215 is induced by frequently exchanging training infor- 216 mation between the server and clients over limited 217 bandwidth channels (Kairouz et al., 2021). The 218 substantial number of parameters in FMs further 219 exacerbates this burden, thus hindering the training 220 process.

221 **Adaptability Challenges.** Adaptability chal- 222 lenges arise from the adaptation of an FM to a 223 specific downstream task (*e.g.*, by fine-tuning) in 224 FL settings. Key challenges include *data hetero-* 225 *geneity* and *resource heterogeneity*. Performance 226 degradation in FL, attributed to heterogeneous data 227 distributions among clients, is a well-recognized is- 228 sue (Kairouz et al., 2021; Li et al., 2022). A recent 229 study (Babakniya et al., 2023a) has shown that such

230 performance penalty is even more substantial when 231 fine-tuning FMs. For NLP tasks, data heterogene- 232 ity can manifest as variations in language, style, 233 topic, or sentiment across datasets held by different 234 clients. In multi-modal scenarios, the challenge is 235 even more pronounced due to the inherent diversity 236 in data types (*e.g.*, text, images, and audio) (Yu 237 et al., 2023a). Addressing data heterogeneity in- 238 volves not just identifying and measuring it but also 239 developing algorithms that are robust to such diver- 240 sity, ensuring that the model can learn effectively 241 from varied data contributions without compromis- 242 ing on performance. In terms of resource hetero- 243 geneity, the memory and computational resources 244 of the devices for different participants may be di- 245 verse (Diao et al., 2021), which could cause delays 246 for model synchronization and inactivation of some 247 participants, *i.e.*, stragglers, making it challenging 248 to leverage the full potential of FMs in FL settings.

249 **Trustworthiness Challenges.** Trustworthiness 250 challenges emphasize the concerns regarding pri- 251 vacy, security, and ethical considerations in the life- 252 cycle of FM-FL, from the pre-training and model 253 adaptation to the application stages. We present 254 two representative challenges from this perspective: 255 (1) *intellectual property*: Intellectual Property (IP) 256 protection in FM-FL primarily involves attributing 257 ownership rights for both models and data. From 258 the server’s perspective, broadcasting a pre-trained 259 model to multiple nodes for fine-tuning poses IP 260 protection and security risks (*e.g.*, model theft), ne- 261 cessitating measures to safeguard IP rights and en- 262 sure model integrity (Kang et al., 2024); (2) *privacy* 263 *leakage*: Although FL does not immediately share 264 data, studies have shown that it may not always 265 guarantee sufficient privacy preservation (Geiping 266 et al., 2020), as model parameters (*e.g.*, weights or 267 gradients) may leak sensitive information to mali- 268 cious adversaries (Zhu et al., 2019).

269 4 Techniques

270 Recent work has begun to address challenges as- 271 sociated with adapting pre-trained FMs to specific 272 downstream tasks in FL settings. In this section, we 273 survey FM-FL techniques on three aspects, namely 274 *efficiency* (Section 4.1), *adaptability* (Section 4.2), 275 and *trustworthiness* (Section 4.3). As illustrated in 276 Figure 1, we further refine them according to the 277 key features of different methods.

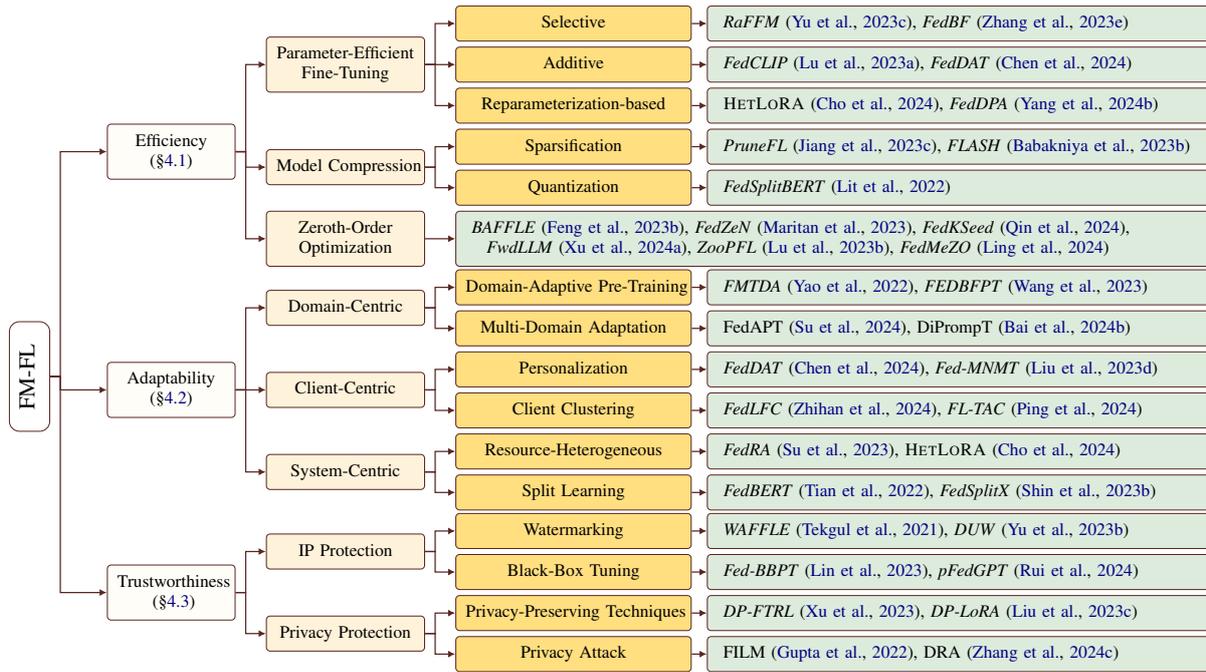


Figure 1: Taxonomy of research in foundation models with federated learning.

278 4.1 Efficiency

279 There has been a considerable focus on developing
 280 resource-efficient approaches. This part describes
 281 techniques that improve resource efficiency.

282 4.1.1 Parameter-Efficient Fine-Tuning

283 Federated Parameter-Efficient Fine-Tuning (Fed-
 284 PEFT), originating from the fine-tuning practices
 285 of FMs (Lester et al., 2021; Hu et al., 2022; Li and
 286 Liang, 2021), is a suite of techniques designed to
 287 reduce both the computational load and the asso-
 288 ciated communication overheads (Malaviya et al.,
 289 2023; Woisetschläger et al., 2024). In alignment
 290 with existing FM fine-tuning taxonomies (Lialin
 291 et al., 2023; Ding et al., 2023), we present FedPEFT
 292 methods in three categories: *selective methods*,
 293 *additive methods*, and *reparameterization-based*
 294 *methods*.

295 **Selective Methods.** Selective methods fine-tune
 296 a small subset of the parameters, leaving the ma-
 297 jority unchanged. In the field of LLMs, a promi-
 298 nent example of such methods is BitFit (Ben Za-
 299 ken et al., 2022), which only fine-tunes the bias
 300 terms. BitFit has inspired a series of studies
 301 in FedPEFT (Bu et al., 2022; Sun et al., 2022a;
 302 Zhang et al., 2023e), demonstrating the superior
 303 communication efficiency of only updating the
 304 bias terms while still achieving competitive per-
 305 formance. More sophisticated methods strive to find

306 sparse subnetworks for partial fine-tuning. Among
 307 them, various methods (Seo et al., 2021; Li et al.,
 308 2021a; Tamirisa et al., 2024) advocate for the Lot-
 309 tery Ticket Hypothesis (LTH) (Frankle and Carbin,
 310 2019), positing that a dense network contains many
 311 subnetworks whose inference capabilities are as ac-
 312 curate as that of the original network. FedSelect
 313 (Tamirisa et al., 2024) is a representative method
 314 that encourages clients to find optimal subnetworks
 315 based on LTH and continually fine-tunes these de-
 316 rived subnetworks to encapsulate local knowledge.
 317 As another important aspect, RaFFM (Yu et al.,
 318 2023c) proposes to prioritize specialized salient pa-
 319 rameters by ranking them using salience evaluation
 320 metrics such as the ℓ_1 and ℓ_2 norms.

321 **Additive Methods.** Instead of fine-tuning a sub-
 322 set of model parameters, additive methods incorpo-
 323 rate lightweight trainable blocks into frozen FMs
 324 and tune the additional parameters for model adap-
 325 tation. These methods not only enhance computa-
 326 tional and communicational efficiency but also in-
 327 troduce an extra benefit: personalization (Lu et al.,
 328 2023a), *i.e.*, the integration of these supplemen-
 329 tary parameters allows for the customization of het-
 330 erogeneous models tailored to specific local data
 331 characteristics or user preferences. Key branches
 332 within additive methods include *adapter tuning* and
 333 *prompt tuning*. Adapter tuning integrates small-
 334 scale neural networks (known as “adapters”) into

the pre-trained models (Houlsby et al., 2019; Hu et al., 2022). On the other hand, prompt tuning incorporates trainable task-specific continuous prompt vectors at the input layer (Liu et al., 2023a; Dong et al., 2023). More details on these methods are provided in Appendix A.1.

Reparameterization-based Methods. The hypothesis behind reparameterization-based methods is that fine-tuning adaptations can be reparameterized into optimization within low-rank subspaces (Aghajanyan et al., 2021). Low-Rank Adaptation (LoRA) (Hu et al., 2022), as a popular PEFT method from the area of LLMs, reduces the number of trainable parameters for downstream tasks by representing the weight updates with two smaller matrices (called update matrices) through low-rank decomposition (Ding et al., 2023). When optimizing a parameter matrix $\mathbf{W} \in \mathbb{R}^{m \times n}$, the update equation can be written as: $\mathbf{W} \leftarrow \mathbf{W} + \Delta\mathbf{W}$. The core idea of LoRA is to freeze the original matrix \mathbf{W} while approximating the parameter update $\Delta\mathbf{W}$ by low-rank decomposition matrices, *i.e.*, $\Delta\mathbf{W} = \mathbf{A} \cdot \mathbf{B}^\top$, where $\mathbf{A} \in \mathbb{R}^{m \times k}$ and $\mathbf{B} \in \mathbb{R}^{n \times k}$ are the trainable parameters for task adaptation and $k \ll \min(m, n)$ is the reduced rank. The trainable parameter size is then reduced from mn to $k(m + n)$. The major benefit of LoRA is that it can largely save memory and storage usage. A straightforward way to perform federated finetuning with LoRA is to train the LoRA modules \mathbf{A} and \mathbf{B} with homogeneous rank k across all clients with standard FL such as FedAvg (McMahan et al., 2017). Several studies have shown that this method can achieve an outstanding level of trade-off between performance and communication overhead for a wide range of FMs, including language models (Zhang et al., 2024b, 2023e), vision-language models (Nguyen et al., 2024), and speech-to-text models (Du et al., 2024).

To summarize, FedPEFT techniques have exhibited remarkable potential in effectively harnessing and adapting FMs within FL scenarios, leading to enhanced performance and efficiency across distributed systems. We refer to Appendix A.2 for a detailed comparison.

4.1.2 Model Compression

Model compression refers to the techniques used to reduce the size of models, thereby improving resource efficiency (Shah and Lau, 2023).

Sparsification. Model sparsification methods reduce communication burden by only transmitting a subset of FM parameters across the network (Jiang et al., 2023c). Typical methods focus on identifying and cultivating high-potential subnetworks (Frankle and Carbin, 2019; Tsouvalas et al., 2023).

Quantization. Quantization is well-established in both the FM and FL domains (Xu et al., 2024b; Reiszadeh et al., 2020), which involves decreasing the precision of floating-point parameters for mitigating the storage, computational, and communication demands. Quantization is orthogonal to other resource-efficient techniques, making it feasible to combine them for greater efficiency and flexibility (Lit et al., 2022).

4.1.3 Zeroth-Order Optimization

In contrast to the use of gradient descent in most FL optimization algorithms, a particular line of research advocates for the removal of BackPropagation (BP) (Malladi et al., 2023a) in favor of Zeroth-Order Optimization (ZOO) (Fang et al., 2022; Li and Chen, 2021). BP-free methods conserve memory needed for computing gradients and minimize communication overhead for model aggregation (Qin et al., 2024), making FMs more accessible for lower-end devices, thereby enhancing their applicability in diverse hardware environments.

ZOO methods primarily rely on perturbation methods to estimate gradients with forward propagation. Based on this, recent work, such as that by Xu et al. (2024a); Lu et al. (2023b), has initiated preliminary explorations into the deployment of both FedPEFT and full-model fine-tuning of billion-sized FMs, like LLaMA, on mobile devices. The naive ZOO methods remain impractical for training large FMs in standard FL frameworks such as FedAvg, as they still result in a significant communication burden for model aggregation. In light of this, FedKSeed (Qin et al., 2024) was proposed to further reduce communication overheads between the server and clients by using just a few random seeds and scalar gradients, requiring only a few thousand bytes for communication.

Although ZOO methods have shown promise in resource-efficient FL (Ling et al., 2024), they generally require many iterations to achieve strong performance (Malladi et al., 2023b). Compared to the well-established BP-based optimization, ZOO is still in the early stages of development, particularly for FM-FL settings, necessitating further research

434 and optimization.

435 **4.2 Adaptability**

436 Adaptation refers to the process of tailoring a pre-
437 trained FM to perform effectively across varying
438 FL settings and scenarios. This mainly includes the
439 capability to learn from different domains, cater to
440 individual user needs, and work across diverse de-
441 vices while retaining overall performance and effi-
442 ciency. We focus on three key aspects of adaptation,
443 namely *domain-centric adaptation*, *client-centric*
444 *adaptation*, and *system-centric adaptation*.

445 **4.2.1 Domain-Centric Adaptation**

446 Domain-centric adaptation focuses on adapting
447 FMs within specific domains by addressing the
448 domain diversity across client datasets.

449 **Domain-Adaptive Pre-Training.** Despite being
450 heavily reliant on large-scale and public datasets
451 for their initial training, FMs often require fur-
452 ther Domain-Adaptive Pre-Training (DAPT) with
453 domain-specific data for tasks that necessitate spe-
454 cialized knowledge (Gururangan et al., 2020; Guo
455 and Yu, 2022). In domains like healthcare, FL
456 allows for the continued pre-training of these mod-
457 els using sensitive, domain-specific data without
458 compromising privacy. Based on this idea, Jiang
459 et al. (2023b) proposed FFDAPT, a computational-
460 efficient further pre-training algorithm that freezes
461 a portion of consecutive layers while optimizing
462 the rest of the layers. Similarly, Wang et al. (2023)
463 proposed FEDBFPT that builds a local model for
464 each client, progressively training the shallower
465 layers of local models while sampling deeper lay-
466 ers, and aggregating trained parameters on a server
467 to create the final global model.

468 **Multi-Domain Adaptation.** Given that client
469 data may belong to various domains in real-world
470 FL scenarios, some efforts (Feng et al., 2023c;
471 Su et al., 2024) have been devoted to facilitating
472 multi-domain collaborative adaptation. Feng et al.
473 (2023c) applied a pre-trained CLIP to the multi-
474 domain scenario and proposed an adaptive prompt
475 tuning method that uses domain-specific keys to
476 generate prompts for each test sample. Further-
477 more, Su et al. (2024) employed knowledge distil-
478 lation to selectively distill global knowledge based
479 on an entropy measure, improving the generaliza-
480 tion across different domains.

481 **4.2.2 Client-Centric Adaptation**

482 Client-centric adaptation refers to the process of
483 tailoring an FM to meet the specific needs or pref-
484 erences of individual clients while leveraging the
485 decentralized and privacy-preserving nature of FL.
486 Particularly, we discuss two types of popular per-
487 sonalized methods as follows:

488 **Personalization.** Adapter-based methods intro-
489 duce small, trainable adapters into the frozen pre-
490 trained FMs, allowing for client-specific model
491 adaptation without altering the original FL. Fed-
492 DAT (Chen et al., 2024) leverages a dual-adapter
493 structure, with personalized adapters focusing on
494 client-specific knowledge and a global adapter
495 maintaining client-agnostic knowledge. FedDAT
496 executes bi-directional knowledge distillation be-
497 tween personalized adapters and the global adapter
498 to regularize the client’s updates and prevent over-
499 fitting. Prompt-based methods involve using client-
500 specific soft prompts to guide the model’s response.
501 pFedPG (Yang et al., 2023a) trains a prompt gener-
502 ator to exploit underlying client-specific character-
503 istics and produce personalized prompts for each
504 client, thereby enabling efficient and personalized
505 adaptation.

506 **Client Clustering.** This branch of study aims to
507 cluster clients based on the underlying relationships
508 and tailor FMs for the client group with similar data
509 distributions, thus reducing the negative impact of
510 data heterogeneity and improving accuracy. Zhihan
511 et al. (2024) proposed a FedPEFT-based frame-
512 work for multilingual modeling, which employs
513 language family clustering to alleviate parameter
514 conflicts of LoRA tuning.

515 **4.2.3 System-Centric Adaptation**

516 System-centric aims to improve adaptability at the
517 system level. This involves handling resource het-
518 erogeneity in the FL systems while ensuring train-
519 ing efficiency and model utility.

520 **Resource-Heterogeneous Methods.** Cross-
521 device FL systems may be composed of devices
522 equipped with heterogeneous resources, leading
523 to disparities where certain devices exhibit more
524 efficient model training than others (Chen et al.,
525 2024). To address this, several methods have
526 been developed to customize model architectures
527 for resource-heterogeneous FL systems. In FL
528 environments with heterogeneous resources,
529 LoRA-based FedPEFT exhibits distinctive flexi-

bility and adaptation in fine-tuning frozen FMs without overburdening client devices. [Su et al. \(2023\)](#) suggested assigning LoRA adapters to varying numbers of layers for heterogeneous clients according to a randomly generated mask matrix. An alternative and more targeted idea is to choose diverse LoRA ranks across clients based on their system capabilities. [Bai et al. \(2024a\)](#) proposed FlexLoRA to adjust local LoRA ranks dynamically. FlexLoRA reconstructs the uniform full-sized LoRA module $\Delta\mathbf{W}$ for server-side model aggregation followed by an SVD-based parameter redistribution. However, concurrent research by [Cho et al. \(2024\)](#) has empirically demonstrated that the reconstruct-redistribute method suffers from performance loss compared to homogeneous LoRA. Instead, they proposed HET-LoRA ([Cho et al., 2024](#)) that utilizes zero-padding to align module size before aggregation. It then truncates the global LoRA modules for the specific rank of the next selected clients.

Split Learning. Split learning addresses the resource heterogeneity between servers and clients by splitting a large model at a cut layer into client and server models ([Thapa et al., 2022](#)). For each training step, the output tensor, so-called smashed data, from the client model and the corresponding labels are transmitted over to the server. The server continues the forward propagation by processing the smashed data through its remaining layers; it then computes the loss using the transmitted label and performs backpropagation. The gradient generated at the first layer of the server model is then transmitted back to the client for further backpropagation. Along this line, FedBERT ([Tian et al., 2022](#)) proposes to leverage split learning for training the BERT model, showing the feasibility of training large FMs in FL settings. FedSplitX ([Shin et al., 2023b](#)) is a more fine-grained method that allows multiple partition points for model splitting, accommodating more diverse client capabilities. Compared to conventional FL, split learning scales better with the size of FMs as it communicates only small-sized smashed data instead of model parameters ([Singh et al., 2019](#)). Despite its merits, split learning is highly dependent on the network connection quality. Given that server-client interactions occur at every step of the optimization process ([Zheng et al., 2023](#)), communication delays cause a more significant impact on efficiency.

4.3 Trustworthiness

This line of work aims to enhance trustworthiness throughout the FM-FL lifecycle, covering a variety of key aspects including, but not limited to, *IP protection, privacy protection, and attack robustness*.

4.3.1 IP Protection

Existing IP protection involves safeguarding ownership of FMs from unauthorized use (*e.g.*, model theft) ([Tekgul et al., 2021](#)). We discuss the following two mainstream IP protection strategies: *watermarking* and *black-box tuning*.

Watermarking. Watermarking is a well-known deterrence technology for model IP protection by providing the identities of model owners to demonstrate ownership of their models ([Adi et al., 2018](#)). [Tekgul et al. \(2021\)](#) proposed WAFFLE, the first solution that addresses the ownership problem by injecting a watermark into the global model in FL environments. Recently, [Yu et al. \(2023b\)](#) proposed DUW that embeds a client-unique key into each client’s local model, aiming to identify the infringer of a leaked model while verifying the FL model’s ownership.

Black-Box Tuning. Black-Box Tuning (BBT) is a set of ZOO-based methods that fine-tune FMs without direct access to model parameters ([Sun et al., 2022c,b](#)). BBT methods are often additive, introducing additional parameters while keeping the original model frozen (See Section 4.1.1). FedBBPT ([Lin et al., 2023](#)) is a general prompt tuning framework that facilitates the joint training of a global lightweight prompt generator across multiple clients. FedBPT ([Sun et al., 2024a](#)) adopts a classic evolutionary-based ZOO method, CMA-ES ([Hansen and Ostermeier, 2001](#)), for training an optimal prompt that improves the performance of frozen FMs. ZooPFL ([Lu et al., 2023b](#)), on the other hand, applies coordinate-wise gradient estimate to learn input surgery that incorporates client-specific embeddings. BBT allows for local fine-tuning of FMs while not infringing IP constraints. However, current research in this line is limited to few-shot learning with small datasets for LLM fine-tuning ([Sun et al., 2022b](#)), while larger datasets and other modalities remain unexplored.

4.3.2 Privacy Protection

Protecting privacy in FM-FL requires both designing protective measures and studying privacy attack strategies.

Privacy-Preserving Techniques. Differential Privacy (DP) is a theoretical framework that governs privacy boundaries and manages the trade-off between privacy and model convergence (Wei et al., 2020; Xu et al., 2023). DP-based FL approaches often add artificial noise (e.g., Gaussian noise) to parameters at the clients’ side before aggregating to prevent information leakage (Xu et al., 2023). Besides, DP is compatible with most FedPEFT methods. For instance, Sun et al. (2024b) showed that DP noise can even be amplified by the locally “semi-quadratic” nature of LoRA-based methods, motivating the integration of LoRA with DP to improve resource efficiency while maintaining data privacy (Liu et al., 2023c). In addition to DP, Secure Multi-Party Computation (SMPC) (Munthan et al., 2019) and Homomorphic Encryption (HE) (Zhang et al., 2020) are also effective privacy-preserving mechanisms. However, they do not scale well enough for large-scale deployments in FM-FL.

Privacy Attack. Privacy attacks in FM-FL involve extracting sensitive information from the data used in training, even though the data itself is not directly shared. Major attacks include *membership inference attack* and *data reconstruction attack*, where the former aims to determine whether a specific data sample is in a victim client’s training set, and the latter strives to reconstruct original input data from the model parameters or gradients (Ren et al., 2024). Regarding membership inference attacks, Vu et al. (2024) revealed the vulnerabilities of popular LLMs, including BERT, DistilBERT, and OpenAI’s GPTs. In terms of data reconstruction attacks, Gupta et al. (2022) presented an attack *FILM*, which recovers private text data by extracting information from gradients transmitted during training despite employing a DP mechanism.

5 Future Directions

Although recent work has already begun to address the challenges discussed in Section 3.2, many critical open directions are yet to be explored. Here, we outline several representative ones.

Multimodal FM-FL. With the development of mobile technology and IoT infrastructures (Brunete et al., 2021), numerous edge devices produce data from a range of modalities, such as sensory, visual, and audio. In the era of FMs, the success of LLMs and their multimodal derivatives (Ramesh et al.,

2021; Google, 2023; OpenAI, 2024) have demonstrated the potential of multimodal FMs. The potential opportunities and challenges for multimodal FM-FL have yet to be explored.

Continual Learning. Continual learning enables models to adapt to new data over time, improving their performance and accuracy. By incorporating new data into the model training process, FL and FMs can continuously improve and adapt to changing environments and user needs (Yang et al., 2024a). Future directions may involve leveraging transfer learning techniques in continual learning for FL and FMs. Models can transfer knowledge from previous tasks or domains to new ones, enabling more efficient adaptation (Good et al., 2023).

Efficient Federated Black-Box Tuning. In scenarios where gradient access is unavailable, preliminary efforts have focused on federated fine-tuning black-box FMs (Lin et al., 2023; Sun et al., 2024a; Lu et al., 2023b; Rui et al., 2024) utilizing ZOO. However, ZOO’s noticeably slower convergence rates, especially in high-dimensional contexts compared to gradient-based methods (Golovin et al., 2020), indicate an important direction for further research. The impact of these slower convergence rates on overall efficiency and computational load within FL, particularly concerning large-scale FMs, has not been adequately investigated and understood.

FL with AI-Generated Content. AI-Generated Content (AIGC) denotes content produced via advanced generative FMs (Wu et al., 2023a). The strong generative capability of FMs offers the advantage of rapidly automating the creation of inexhaustible synthetic data. This capability positions AIGC as a valuable supplementary data source for model training and evaluation in many tasks (Xu et al., 2024c). Despite some efforts (Zhang et al., 2023b), more potential opportunities and challenges for AIGC-aided FL have yet to be explored.

6 Conclusions

In this survey, we have meticulously surveyed the intersection of FM and FL. We identified core challenges in efficiency, adaptability, and trustworthiness and proposed a comprehensive taxonomy of techniques in response to these challenges. In addition, we discussed future directions and applications in this research field, hoping to attract more breakthroughs in future research.

727 Limitations

728 FM and FL are very fast-moving fields. We have
729 put a lot of effort into including the latest research
730 efforts in the community in this survey. There-
731 fore, we believe that our survey will help to inspire
732 and push further research and innovation in these
733 important areas. Our survey does not focus on ex-
734 perimental evaluation of the available ideas and
735 systems. We believe that would be an important
736 next step that we are leaving for future work.

737 References

738 Durmus Alp Emre Acar, Yue Zhao, Ramon Matas,
739 Matthew Mattina, Paul Whatmough, and Venkatesh
740 Saligrama. 2021. [Federated learning based on dy-
741 namic regularization](#). In *International Conference on
742 Learning Representations*.

743 Yossi Adi, Carsten Baum, Moustapha Cisse, Benny
744 Pinkas, and Joseph Keshet. 2018. [Turning your weak-
745 ness into a strength: Watermarking deep neural net-
746 works by backdooring](#). In *27th USENIX Security
747 Symposium (USENIX Security 18)*, pages 1615–1631,
748 Baltimore, MD. USENIX Association.

749 Armen Aghajanyan, Sonal Gupta, and Luke Zettle-
750 moyer. 2021. [Intrinsic dimensionality explains the
751 effectiveness of language model fine-tuning](#). In *Pro-
752 ceedings of the 59th Annual Meeting of the Associa-
753 tion for Computational Linguistics and the 11th Inter-
754 national Joint Conference on Natural Language Pro-
755 cessing (Volume 1: Long Papers)*, pages 7319–7328,
756 Online. Association for Computational Linguistics.

757 Singh Akshay and Thakur Rahul. 2024. [Generalizable
758 multilingual hate speech detection on low resource
759 indian languages using fair selection in federated
760 learning](#). In *Proceedings of the 2024 Conference
761 of the North American Chapter of the Association
762 for Computational Linguistics: Human Language
763 Technologies*.

764 Muhammad Ammad-Ud-Din, Elena Ivannikova,
765 Suleiman A. Khan, Were Oyomno, Qiang Fu,
766 Kuan Eeik Tan, and Adrian Flanagan. 2019. [Feder-
767 ated collaborative filtering for privacy-preserving
768 personalized recommendation system](#). *arXiv
769 preprint arXiv:1901.09888*.

770 Sheikh Shams Azam, Tatiana Likhomanenko, Martin
771 Pelikan, and Jan “Honza” Silovsky. 2023a. [Import-
772 ance of smoothness induced by optimizers in fl4asr:
773 Towards understanding federated learning for end-to-
774 end asr](#). In *2023 IEEE Automatic Speech Recognition
775 and Understanding Workshop (ASRU)*, pages 1–8.

776 Sheikh Shams Azam, Martin Pelikan, Vitaly Feldman,
777 Kunal Talwar, Jan Silovsky, and Tatiana Likhoma-
778 nenko. 2023b. [Federated learning for speech recog-
779 nition: Revisiting current trends towards large-scale](#)

[ASR](#). In *International Workshop on Federated Learn-
ing in the Age of Foundation Models in Conjunction
with NeurIPS 2023*. 780
781
782

Sara Babakniya, Ahmed Elkordy, Yahya Ezzeldin,
Qingfeng Liu, Kee-Bong Song, MOSTAFA EL-
Khamy, and Salman Avestimehr. 2023a. [SLoRA:
Federated parameter efficient fine-tuning of language
models](#). In *International Workshop on Federated
Learning in the Age of Foundation Models in Con-
junction with NeurIPS 2023*. 783
784
785
786
787
788
789

Sara Babakniya, Souvik Kundu, Saurav Prakash, Yue
Niu, and Salman Avestimehr. 2023b. [Revisiting spar-
sity hunting in federated learning: Why does sparsity
consensus matter?](#) *Transactions on Machine Learn-
ing Research*. 790
791
792
793
794

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed,
and Michael Auli. 2020. [wav2vec 2.0: A framework
for self-supervised learning of speech representations](#).
In *Advances in Neural Information Processing Sys-
tems*, volume 33, pages 12449–12460. Curran Asso-
ciates, Inc. 795
796
797
798
799
800

Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deb-
orah Estrin, and Vitaly Shmatikov. 2020. [How to
backdoor federated learning](#). In *International Con-
ference on Artificial Intelligence and Statistics*, pages
2938–2948. PMLR. 801
802
803
804
805

Jiamu Bai, Daoyuan Chen, Bingchen Qian, Liuyi Yao,
and Yaliang Li. 2024a. [Federated fine-tuning of
large language models under heterogeneous lan-
guage tasks and client resources](#). *arXiv preprint
arXiv:2402.11505*. 806
807
808
809
810

Sikai Bai, Jie Zhang, Shuaicheng Li, Song Guo, Jingcai
Guo, Jun Hou, Tao Han, and Xiaocheng Lu. 2024b. [Diprompt:
Disentangled prompt tuning for multiple
latent domain generalization in federated learning](#).
arXiv preprint arXiv:2403.08506. 811
812
813
814
815

Elad Ben Zaken, Yoav Goldberg, and Shauli Ravfogel.
2022. [BitFit: Simple parameter-efficient fine-tuning
for transformer-based masked language-models](#). In
*Proceedings of the 60th Annual Meeting of the As-
sociation for Computational Linguistics (Volume 2:
Short Papers)*, pages 1–9, Dublin, Ireland. Associa-
tion for Computational Linguistics. 816
817
818
819
820
821
822

Peva Blanchard, El Mahdi El Mhamdi, Rachid Guer-
raoui, and Julien Stainer. 2017. [Machine learning
with adversaries: Byzantine tolerant gradient descent](#).
In *Proceedings of the 31st International Conference
on Neural Information Processing Systems*. 823
824
825
826
827

J. Bobadilla, F. Ortega, A. Hernando, and A. Gutiérrez.
2013. [Recommender systems survey](#). *Knowledge-
Based Systems*, 46:109–132. 828
829
830

Rishi Bommasani, Drew A Hudson, Ehsan Adeli,
Russ Altman, Simran Arora, Sydney von Arx,
Michael S Bernstein, Jeannette Bohg, Antoine Bosse-
lut, Emma Brunskill, et al. 2021. [On the opportuni-
ties and risks of foundation models](#). *arXiv preprint
arXiv:2108.07258*. 831
832
833
834
835
836

837	Kallista Bonawitz, Peter Kairouz, Brendan McMahan, and Daniel Ramage. 2021. Federated learning and privacy: Building privacy-preserving systems for machine learning and data science on decentralized data . <i>Queue</i> , 19(5):87–114.	891
838		892
839		893
840		894
841		895
842	Tom Brown et al. 2020. Language models are few-shot learners . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 1877–1901. Curran Associates, Inc.	896
843		897
844		898
845		899
846	Alberto Brunete, Ernesto Gambao, Miguel Hernando, and Raquel Cedazo. 2021. Smart assistive architecture for the integration of iot devices, robotic systems, and multimodal interfaces in healthcare environments . <i>Sensors</i> , 21(6):2212.	900
847		901
848		902
849		903
850		904
851	Zhiqi Bu, Yu-Xiang Wang, Sheng Zha, and George Karypis. 2022. Differentially private bias-term only fine-tuning of foundation models . In <i>Workshop on Trustworthy and Socially Responsible Machine Learning, NeurIPS 2022</i> .	905
852		906
853		907
854		908
855		909
856	Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4 . <i>arXiv preprint arXiv:2303.12712</i> .	910
857		911
858		912
859		913
860		914
861		915
862		916
863	Dongqi Cai, Yaozong Wu, Shangguang Wang, Felix Xiaozhu Lin, and Mengwei Xu. 2023. Efficient federated learning for modern nlp . In <i>Proceedings of the 29th Annual International Conference on Mobile Computing and Networking, ACM MobiCom '23</i> , New York, NY, USA. Association for Computing Machinery.	917
864		918
865		919
866		920
867		921
868		922
869		923
870	Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. 2021. Fltrust: Byzantine-robust federated learning via trust bootstrapping . In <i>ISOC Network and Distributed System Security Symposium (NDSS)</i> .	924
871		925
872		926
873		927
874		928
875	CCPA. 2023. California consumer privacy act (ccpa) .	929
876	Pushpita Chatterjee, Debashis Das, and Danda B Rawat. 2023. Use of federated learning and blockchain towards securing financial services . <i>arXiv preprint arXiv:2303.12944</i> .	930
877		931
878		932
879		933
880	Haokun Chen, Yao Zhang, Denis Krompass, Jindong Gu, and Volker Tresp. 2024. Feddat: An approach for foundation model finetuning in multi-modal heterogeneous federated learning . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(10):11285–11293.	934
881		935
882		936
883		937
884		938
885		939
886	Yudong Chen, Lili Su, and Jiaming Xu. 2017. Distributed statistical machine learning in adversarial settings: Byzantine gradient descent . <i>Proceedings of the ACM on Measurement and Analysis of Computing Systems</i> .	940
887		941
888		942
889		943
890		944
	Vijil Chenthamarakshan, Samuel C. Hoffman, C. David Owen, Petra Lukacik, Claire Strain-Damerell, Daren Fearon, Tika R. Malla, Anthony Tumber, Christopher J. Schofield, Helen M.E. Duyvesteyn, Wan-wisa Dejnirattisai, Loic Carrique, Thomas S. Walter, Gavin R. Sreaton, Tetiana Matviiuk, Aleksandra Mojsilovic, Jason Crain, Martin A. Walsh, David I. Stuart, and Payel Das. 2023. Accelerating drug target inhibitor discovery with a deep generative foundation model . <i>Science Advances</i> , 9(25):eadg7865.	945
		946
		947
	Yae Jee Cho, Luyang Liu, Zheng Xu, Aldi Fahrezi, and Gauri Joshi. 2024. Heterogeneous lora for federated fine-tuning of on-device foundation models . <i>arXiv preprint arXiv:2401.06432</i> .	948
		949
		950
	Yun-Wei Chu, Dong-Jun Han, and Christopher G. Brinton. 2024. Only send what you need: Learning to communicate efficiently in federated multilingual machine translation . In <i>Companion Proceedings of the ACM on Web Conference 2024</i> , page 1548–1557, New York, NY, USA. Association for Computing Machinery.	951
		952
		953
	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451, Online. Association for Computational Linguistics.	954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

948	2023, pages 14665–14675, Singapore. Association for Computational Linguistics.	Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haofen Wang, and Jiawei Zhang. 2023. Chat-rec: Towards interactive and explainable llms-augmented recommender system . <i>arXiv preprint arXiv:2303.14524</i> .	1000
949			1001
950	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An image is worth 16x16 words: Transformers for image recognition at scale . In <i>International Conference on Learning Representations</i> .		1002
951			1003
952			1004
953		GDPR. 2016. Regulation (eu) 2016/679 of the european parliament and of the council .	1005
954			1006
955		Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. 2020. Inverting gradients - how easy is it to break privacy in federated learning? In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 16937–16947. Curran Associates, Inc.	1007
956			1008
957			1009
958	Yichao Du, Zhirui Zhang, Linan Yue, Xu Huang, Yuqing Zhang, Tong Xu, Linli Xu, and Enhong Chen. 2024. Communication-efficient personalized federated learning for speech-to-text tasks . In <i>ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 10001–10005.		1010
959			1011
960			1012
961		Daniel Golovin, John Karro, Greg Kochanski, Chansoo Lee, Xingyou Song, and Qiuyi Zhang. 2020. Gradientless descent: High-dimensional zeroth-order optimization . In <i>International Conference on Learning Representations</i> .	1013
962			1014
963			1015
964			1016
965	Tao Fan, Yan Kang, Guoqiang Ma, Weijing Chen, Wenbin Wei, Lixin Fan, and Qiang Yang. 2023. Fate-llm: A industrial grade federated learning framework for large language models . <i>arXiv preprint arXiv:2310.10049</i> .		1017
966		Jack Good, Jimit Majmudar, Christophe Dupuy, Jixuan Wang, Charith Peris, Clement Chung, Richard Zemel, and Rahul Gupta. 2023. Coordinated replay sample selection for continual federated learning . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 331–342, Singapore. Association for Computational Linguistics.	1018
967			1019
968			1020
969			1021
970	Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. 2020. Local model poisoning attacks to byzantine-robust federated learning . In <i>29th USENIX Security Symposium (USENIX Security 20)</i> , pages 1605–1622.		1022
971			1023
972			1024
973			1025
974		Gemini Team Google. 2023. Gemini: a family of highly capable multimodal models . <i>arXiv preprint arXiv:2312.11805</i> .	1026
975	Wenzhi Fang, Ziyi Yu, Yuning Jiang, Yuanming Shi, Colin N. Jones, and Yong Zhou. 2022. Communication-efficient stochastic zeroth-order optimization for federated learning . <i>IEEE Transactions on Signal Processing</i> , 70:5058–5073.		1027
976			1028
977		Eduard Gorbunov, Samuel Horváth, Peter Richtárik, and Gauthier Gidel. 2023. Variance reduction is an antidote to byzantines: Better rates, weaker assumptions and communication compression as a cherry on the top . In <i>International Conference on Learning Representations</i> .	1029
978			1030
979			1031
980	FedML. 2023. Fedllm .		1032
981			1033
982	Chun-Mei Feng, Bangjun Li, Xinxing Xu, Yong Liu, Huazhu Fu, and Wangmeng Zuo. 2023a. Learning federated visual prompt in null space for mri reconstruction . In <i>2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 8064–8073.		1034
983		Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. 2023. Textbooks are all you need . <i>arXiv preprint arXiv:2306.11644</i> .	1035
984			1036
985			1037
986			1038
987	Haozhe Feng, Tianyu Pang, Chao Du, Wei Chen, Shuicheng Yan, and Min Lin. 2023b. Does federated learning really need backpropagation? <i>arXiv preprint arXiv:2301.12195</i> .		1039
988			1040
989			1041
990			1042
991	Xiachong Feng, Xiaocheng Feng, Xiyuan Du, Min-Yen Kan, and Bing Qin. 2023c. Adapter-based selective knowledge distillation for federated multi-domain meeting summarization . <i>arXiv preprint arXiv:2308.03275</i> .		1043
992		Lei Guo, Ziang Lu, Junliang Yu, Quoc Viet Hung Nguyen, and Hongzhi Yin. 2024. Prompt-enhanced federated content representation learning for cross-domain recommendation . In <i>Proceedings of the ACM on Web Conference 2024, WWW '24</i> , page 3139–3149, New York, NY, USA. Association for Computing Machinery.	1044
993			1045
994			1046
995			1047
996	Jonathan Frankle and Michael Carbin. 2019. The lottery ticket hypothesis: Finding sparse, trainable neural networks . In <i>International Conference on Learning Representations</i> .		1048
997			1049
998		Tao Guo, Song Guo, Junxiao Wang, Xueyang Tang, and Wenchao Xu. 2023. Promptfl: Let federated participants cooperatively learn prompts instead of models - federated learning in age of foundation model . <i>IEEE Transactions on Mobile Computing</i> , pages 1–15.	1050
999			1051
			1052
			1053
			1054

1055	Xu Guo and Han Yu. 2022. On the domain adaptation and generalization of pretrained language models: A survey . <i>arXiv preprint arXiv:2211.03154</i> .	1110
1056		1111
1057		1112
1058	Samyak Gupta, Yangsibo Huang, Zexuan Zhong, Tianyu Gao, Kai Li, and Danqi Chen. 2022. Recovering private text in federated learning of language models . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 8130–8143. Curran Associates, Inc.	1113
1059		1114
1060		1115
1061		1116
1062		1117
1063		1118
1064	Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8342–8360, Online. Association for Computational Linguistics.	1119
1065		1120
1066		1121
1067		1122
1068		1123
1069		
1070		
1071		
1072	David Ha, Andrew M. Dai, and Quoc V. Le. 2017. Hypernetworks . In <i>International Conference on Learning Representations</i> .	1124
1073		1125
1074		1126
1075	Nikolaus Hansen and Andreas Ostermeier. 2001. Completely derandomized self-adaptation in evolution strategies . <i>Evolutionary Computation</i> , 9(2):159–195.	1127
1076		1128
1077		1129
1078	Chaoyang He, Songze Li, Jinhyun So, Xiao Zeng, Mi Zhang, Hongyi Wang, Xiaoyang Wang, Praneeth Vepakomma, Abhishek Singh, Hang Qiu, Xinghua Zhu, Jianzong Wang, Li Shen, Peilin Zhao, Yan Kang, Yang Liu, Ramesh Raskar, Qiang Yang, Murali Annavaram, and Salman Avestimehr. 2020. Fedml: A research library and benchmark for federated machine learning . <i>arXiv preprint arXiv:2007.13518</i> .	1130
1079		1131
1080		1132
1081		1133
1082		1134
1083		1135
1084		1136
1085		1137
1086	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP . In <i>Proceedings of the 36th International Conference on Machine Learning</i> , volume 97 of <i>Proceedings of Machine Learning Research</i> , pages 2790–2799. PMLR.	1138
1087		1139
1088		1140
1089		1141
1090		1142
1091		1143
1092		1144
1093		1145
1094	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models . In <i>International Conference on Learning Representations</i> .	1146
1095		1147
1096		1148
1097		1149
1098		1150
1099	Xumin Huang, Peichun Li, Hongyang Du, Jiawen Kang, Dusit Niyato, Dong In Kim, and Yuan Wu. 2024. Federated learning-empowered ai-generated content in wireless networks . <i>IEEE Network</i> , pages 1–1.	1151
1100		1152
1101		1153
1102		1154
1103	Katharina Jeblick, Balthasar Schachtner, Jakob Dexl, Andreas Mittermeier, Anna Theresa Stüber, Johanna Topalis, Tobias Weber, Philipp Wesp, Bastian Oliver Sabel, Jens Rieke, et al. 2024. Chatgpt makes medicine easy to swallow: an exploratory case study on simplified radiology reports . <i>European radiology</i> , 34(5):2817–2825.	1155
1104		1156
1105		1157
1106		1158
1107		
1108		
1109		
	Malhar S Jere, Tyler Farnan, and Farinaz Koushanfar. 2020. A taxonomy of attacks on federated learning . <i>IEEE Security & Privacy</i> , 19(2).	1110
		1111
		1112
	Junteng Jia, Ke Li, Mani Malek, Kshitiz Malik, Jay Mahadeokar, Ozlem Kalinli, and Frank Seide. 2023. Joint federated learning and personalization for on-device asr . In <i>2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 1–8.	1113
		1114
		1115
		1116
		1117
		1118
	Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. 2022. Visual prompt tuning . In <i>Computer Vision – ECCV 2022</i> , pages 709–727, Cham. Springer Nature Switzerland.	1119
		1120
		1121
		1122
		1123
	Jingang Jiang, Xiangyang Liu, and Chenyou Fan. 2023a. Low-parameter federated learning with large language models . <i>arXiv preprint arXiv:2307.13896</i> .	1124
		1125
		1126
	Lekang Jiang, Filip Svoboda, and Nicholas Donald Lane. 2023b. FDAPT: Federated domain-adaptive pre-training for language models . In <i>International Workshop on Federated Learning in the Age of Foundation Models in Conjunction with NeurIPS 2023</i> .	1127
		1128
		1129
		1130
		1131
	Yuang Jiang, Shiqiang Wang, Víctor Valls, Bong Jun Ko, Wei-Han Lee, Kin K. Leung, and Leandros Tassiulas. 2023c. Model pruning enables efficient federated learning on edge devices . <i>IEEE Transactions on Neural Networks and Learning Systems</i> , 34(12):10374–10386.	1132
		1133
		1134
		1135
		1136
		1137
	Eun Seo Jo and Timnit Gebru. 2020. Lessons from archives: Strategies for collecting sociocultural data in machine learning . In <i>Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, FAccT '20</i> , pages 306–316, New York, NY, USA. Association for Computing Machinery.	1138
		1139
		1140
		1141
		1142
		1143
	Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation . <i>Transactions of the Association for Computational Linguistics</i> , 5:339–351.	1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
	Madhura Joshi, Ankit Pal, and Malaikannan Sankarabubbu. 2022. Federated learning for healthcare domain - pipeline, applications and challenges . <i>ACM Trans. Comput. Healthcare</i> , 3(4).	1152
		1153
		1154
		1155
	Peter Kairouz et al. 2021. Advances and open problems in federated learning . <i>Foundations and Trends in Machine Learning</i> , 14(1–2):1–210.	1156
		1157
		1158
	Yan Kang, Tao Fan, Hanlin Gu, Xiaojin Zhang, Lixin Fan, and Qiang Yang. 2024. Grounding foundation models through federated transfer learning: A general framework . <i>arXiv preprint arXiv:2311.17431</i> .	1159
		1160
		1161
		1162

1163	Yeachan Kim, Junho Kim, Wing-Lam Mok, Jun-Hyung Park, and SangKeun Lee. 2023. Client-customized adaptation for parameter-efficient federated learning . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1159–1172, Toronto, Canada. Association for Computational Linguistics.	1219
1164		1220
1165		1221
1166		1222
1167		
1168		
1169	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything . <i>arXiv preprint arXiv:2304.02643</i> .	
1170		
1171		
1172		
1173		
1174	Weirui Kuang, Bingchen Qian, Zitao Li, Daoyuan Chen, Dawei Gao, Xuchen Pan, Yuexiang Xie, Yaliang Li, Bolin Ding, and Jingren Zhou. 2023. Federatedscope-llm: A comprehensive package for fine-tuning large language models in federated learning . <i>arXiv preprint arXiv:2309.00363</i> .	
1175		
1176		
1177		
1178		
1179		
1180	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
1181		
1182		
1183		
1184		
1185		
1186		
1187	Ang Li, Jingwei Sun, Binghui Wang, Lin Duan, Sicheng Li, Yiran Chen, and Hai Li. 2021a. Lotteryfl: Empower edge intelligence with personalized and communication-efficient federated learning . In <i>2021 IEEE/ACM Symposium on Edge Computing (SEC)</i> , pages 68–79.	
1188		
1189		
1190		
1191		
1192		
1193	Guanghao Li, Wansen Wu, Yan Sun, Li Shen, Baoyuan Wu, and Dacheng Tao. 2024a. Visual prompt based personalized federated learning . <i>Transactions on Machine Learning Research</i> .	
1194		
1195		
1196		
1197	Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021b. Align before fuse: Vision and language representation learning with momentum distillation . In <i>Advances in Neural Information Processing Systems</i> , volume 34, pages 9694–9705. Curran Associates, Inc.	
1198		
1199		
1200		
1201		
1202		
1203		
1204	Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. 2022. Federated learning on non-iid data silos: An experimental study . In <i>2022 IEEE 38th International Conference on Data Engineering (ICDE)</i> , pages 965–978.	
1205		
1206		
1207		
1208		
1209	Shenghui Li, Edith Ngai, and Thiemo Voigt. 2023a. Byzantine-robust aggregation in federated learning empowered industrial iot . <i>IEEE Transactions on Industrial Informatics</i> , 19(2):1165–1175.	
1210		
1211		
1212		
1213	Shenghui Li, Edith Ngai, Fanghua Ye, Li Ju, Tianru Zhang, and Thiemo Voigt. 2024b. Blades: A unified benchmark suite for byzantine attacks and defenses in federated learning . In <i>2024 IEEE/ACM Ninth International Conference on Internet-of-Things Design and Implementation (IoTDI)</i> .	
1214		
1215		
1216		
1217		
1218		
	Shenghui Li, Edith C.-H. Ngai, and Thiemo Voigt. 2023b. An experimental study of byzantine-robust aggregation schemes in federated learning . <i>IEEE Transactions on Big Data</i> , pages 1–13.	1219
		1220
		1221
		1222
	Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. 2020. Federated optimization in heterogeneous networks . In <i>Proceedings of Machine Learning and Systems</i> , volume 2, pages 429–450.	1223
		1224
		1225
		1226
		1227
	Xi Li and Jiaqi Wang. 2024. Position paper: Assessing robustness, privacy, and fairness in federated learning integrated with foundation models . <i>arXiv preprint arXiv:2402.01857</i> .	1228
		1229
		1230
		1231
	Xi Li, Songhe Wang, Chen Wu, Hao Zhou, and Jiaqi Wang. 2023c. Backdoor threats from compromised foundation models to federated learning . In <i>International Workshop on Federated Learning in the Age of Foundation Models in Conjunction with NeurIPS 2023</i> .	1232
		1233
		1234
		1235
		1236
		1237
	Xi Li, Chen Wu, and Jiaqi Wang. 2023d. Unveiling backdoor risks brought by foundation models in heterogeneous federated learning . <i>arXiv preprint arXiv:2311.18350</i> .	1238
		1239
		1240
		1241
	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4582–4597. Association for Computational Linguistics.	1242
		1243
		1244
		1245
		1246
		1247
		1248
	Zan Li and Li Chen. 2021. Communication-efficient decentralized zeroth-order method on heterogeneous data . In <i>2021 13th International Conference on Wireless Communications and Signal Processing (WCSP)</i> , pages 1–6.	1249
		1250
		1251
		1252
		1253
	Vladislav Lialin, Vijeta Deshpande, and Anna Rumshisky. 2023. Scaling down to scale up: A guide to parameter-efficient fine-tuning . <i>arXiv preprint arXiv:2303.15647</i> .	1254
		1255
		1256
		1257
	Zihao Lin, Yan Sun, Yifan Shi, Xueqian Wang, Lifu Huang, Li Shen, and Dacheng Tao. 2023. Efficient federated prompt tuning for black-box large pre-trained models . <i>CoRR</i> , abs/2310.03123.	1258
		1259
		1260
		1261
	M Lincy and A Meena Kowshalya. 2020. Early detection of type-2 diabetes using federated learning . <i>International Journal of Scientific Research in Science, Engineering and Technology</i> , 12:257–267.	1262
		1263
		1264
		1265
	Zhenqing Ling, Daoyuan Chen, Liuyi Yao, Yaliang Li, and Ying Shen. 2024. On the convergence of zeroth-order federated tuning in large language models . <i>arXiv preprint arXiv:2402.05926</i> .	1266
		1267
		1268
		1269
	Zhengyang Lit, Shijing Sit, Jianzong Wang, and Jing Xiao. 2022. Federated split bert for heterogeneous text classification . In <i>2022 International Joint Conference on Neural Networks (IJCNN)</i> , pages 1–8.	1270
		1271
		1272
		1273

1274	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023a. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing . <i>ACM Comput. Surv.</i> , 55(9).	<i>Processing Systems</i> , volume 36, pages 53038–53075. Curran Associates, Inc.	1330 1331
1275			
1276			
1277			
1278			
1279	Tao Liu, Zhi Wang, Hui He, Wei Shi, Liangliang Lin, Ran An, and Chenhao Li. 2023b. Efficient and secure federated learning for financial applications . <i>Applied Sciences</i> , 13(10).	Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022. Peft: State-of-the-art parameter-efficient fine-tuning methods . https://github.com/huggingface/peft .	1332 1333 1334 1335 1336
1280			
1281			
1282			
1283	Xiao-Yang Liu, Rongyi Zhu, Daochen Zha, Jiechao Gao, Shan Zhong, and Meikang Qiu. 2023c. Differentially private low-rank adaptation of large language model using federated learning . <i>arXiv preprint arXiv:2312.17493</i> .	Andrea Manoel, Mirian del Carmen Hipolito Garcia, Tal Baumel, Shize Su, Jialei Chen, Robert Sim, Dan Miller, Danny Karmon, and Dimitrios Dimitriadis. 2023. Federated multilingual models for medical transcript analysis . In <i>Proceedings of the Conference on Health, Inference, and Learning</i> , volume 209 of <i>Proceedings of Machine Learning Research</i> , pages 147–162. PMLR.	1337 1338 1339 1340 1341 1342 1343 1344
1284			
1285			
1286			
1287			
1288	Yi Liu, Xiaohan Bi, Lei Li, Sishuo Chen, Wenkai Yang, and Xu Sun. 2023d. Communication efficient federated learning for multilingual neural machine translation with adapter . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5315–5328, Toronto, Canada. Association for Computational Linguistics.	Alessio Maritan, Subhrakanti Dey, and Luca Schenato. 2023. Fedzen: Towards superlinear zeroth-order federated learning via incremental hessian estimation . <i>arXiv preprint arXiv:2309.17174</i> .	1345 1346 1347 1348
1289			
1290			
1291			
1292			
1293			
1294			
1295	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach . <i>arXiv preprint arXiv:1907.11692</i> .	Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data . In <i>Artificial intelligence and statistics</i> , pages 1273–1282. PMLR.	1349 1350 1351 1352 1353
1296			
1297			
1298			
1299			
1300	Wang Lu, Xixu Hu, Jindong Wang, and Xing Xie. 2023a. Fedclip: Fast generalization and personalization for CLIP in federated learning . <i>IEEE Data Eng. Bull.</i> , 46(1):52–66.	Vaikkunth Mugunthan, Antigoni Polychroniadou, David Byrd, and Tucker Hybinette Balch. 2019. Smpai: Secure multi-party computation for federated learning . In <i>Proceedings of the NeurIPS 2019 Workshop on Robust AI in Financial Services</i> , volume 21. MIT Press Cambridge, MA, USA.	1354 1355 1356 1357 1358 1359
1301			
1302			
1303			
1304	Wang Lu, Hao Yu, Jindong Wang, Damien Teney, Hao-han Wang, Yiqiang Chen, Qiang Yang, Xing Xie, and Xiangyang Ji. 2023b. Zoopfl: Exploring black-box foundation models for personalized federated learning . <i>arXiv preprint arXiv:2310.05143</i> .	Duy Phuong Nguyen, J. Pablo Munoz, and Ali Jannesari. 2024. Flora: Enhancing vision-language models with parameter-efficient federated learning . <i>arXiv preprint arXiv:2404.15182</i> .	1360 1361 1362 1363
1305			
1306			
1307			
1308			
1309	Lingjuan Lyu, Han Yu, Xingjun Ma, Chen Chen, Lichao Sun, Jun Zhao, Qiang Yang, and Philip S. Yu. 2022. Privacy and robustness in federated learning: Attacks and defenses . <i>IEEE Transactions on Neural Networks and Learning Systems</i> , pages 1–21.	OpenAI. 2022. Chatgpt .	1364
1310			
1311			
1312			
1313			
1314	Shubham Malaviya, Manish Shukla, and Sachin Lodha. 2023. Reducing communication overhead in federated learning for pre-trained language models using parameter-efficient finetuning . In <i>Proceedings of The 2nd Conference on Lifelong Learning Agents</i> , volume 232 of <i>Proceedings of Machine Learning Research</i> , pages 456–469. PMLR.	OpenAI. 2024. Gpt-4 technical report . <i>arXiv preprint arXiv:2303.08774</i> .	1365 1366
1315			
1316			
1317			
1318			
1319			
1320			
1321	Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D. Lee, Danqi Chen, and Sanjeev Arora. 2023a. Fine-tuning language models with just forward passes . In <i>Workshop on Efficient Systems for Foundation Models @ ICML2023</i> .	Dimitrios P. Panagoulas, Maria Virvou, and George A. Tsihrantzis. 2024. Evaluating llm – generated multimodal diagnosis from medical images and symptom analysis . <i>arXiv preprint arXiv:2402.01730</i> .	1367 1368 1369 1370
1322			
1323			
1324			
1325			
1326	Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. 2023b. Fine-tuning language models with just forward passes . In <i>Advances in Neural Information</i>	Sharnil Pandya, Gautam Srivastava, Rutvij Jhaveri, M. Rajasekhara Babu, Sweta Bhattacharya, Praveen Kumar Reddy Maddikunta, Spyridon Mastorakis, Md. Jalil Piran, and Thippa Reddy Gadekallu. 2023. Federated learning for smart cities: A comprehensive survey . <i>Sustainable Energy Technologies and Assessments</i> , 55:102987.	1371 1372 1373 1374 1375 1376 1377
1327			
1328			
1329			
		Jungwuk Park, Dong-Jun Han, Minseok Choi, and Jaekyun Moon. 2021. Sageflow: Robust federated learning against both stragglers and adversaries . In <i>Advances in Neural Information Processing Systems</i> , volume 34, pages 840–851. Curran Associates, Inc.	1378 1379 1380 1381 1382

1383	Siqi Ping, Yuzhu Mao, Yang Liu, Xiao-Ping Zhang, and Wenbo Ding. 2024. FL-TAC: Enhanced fine-tuning in federated learning via low-rank, task-specific adapter clustering . In <i>ICLR 2024 Workshop on Large Language Model (LLM) Agents</i> .	1440
1384		1441
1385		1442
1386		1443
1387		1444
1388	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual bert? <i>arXiv preprint arXiv:1906.01502</i> .	1445
1389		1446
1390		1447
1391	Zhen Qin, Daoyuan Chen, Bingchen Qian, Bolin Ding, Yaliang Li, and Shuiguang Deng. 2024. Federated full-parameter tuning of billion-sized language models with communication cost under 18 kilobytes . In <i>Proceedings of the 41th International Conference on Machine Learning</i> .	1448
1392		1449
1393		1450
1394		1451
1395		1452
1396		
1397	Chen Qiu, Xingyu Li, Chaithanya Kumar Mummadi, Madan Ravi Ganesh, Zhenzhen Li, Lu Peng, and Wan-Yi Lin. 2024. Federated text-driven prompt generation for vision-language models . In <i>The Twelfth International Conference on Learning Representations</i> .	1453
1398		1454
1399		1455
1400		1456
1401		1457
1402		1458
1403	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision . In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	1459
1404		1460
1405		1461
1406		1462
1407		1463
1408		1464
1409	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine Mcleavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision . In <i>Proceedings of the 40th International Conference on Machine Learning</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 28492–28518. PMLR.	1465
1410		1466
1411		1467
1412		1468
1413		1469
1414		1470
1415		1471
1416	Alec Radford and Jeffrey Wu. 2019. Rewon child, david luan, dario amodei, and ilya sutskever. 2019. Language models are unsupervised multitask learners . <i>OpenAI blog</i> , 1(8):9.	1472
1417		1473
1418		1474
1419		1475
1420	Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation . In <i>International Conference on Machine Learning</i> , pages 8821–8831. PMLR.	1476
1421		1477
1422		1478
1423		1479
1424		
1425	Swarna Priya Ramu, Parimala Boopalan, Quoc-Viet Pham, Praveen Kumar Reddy Maddikunta, Thien Huynh-The, Mamoun Alazab, Thanh Thi Nguyen, and Thippa Reddy Gadekallu. 2022. Federated learning enabled digital twins for smart cities: Concepts, recent advances, and future directions . <i>Sustainable Cities and Society</i> , 79:103663.	1480
1426		1481
1427		1482
1428		1483
1429		
1430		
1431		
1432	Sashank J. Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and Hugh Brendan McMahan. 2021. Adaptive federated optimization . In <i>International Conference on Learning Representations</i> .	1484
1433		1485
1434		1486
1435		1487
1436		1488
1437	Scott Reed et al. 2022. A generalist agent . <i>Transactions on Machine Learning Research</i> . Featured Certification, Outstanding Certification.	1489
1438		
1439		
	Amirhossein Reiszadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramtin Pedarsani. 2020. Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization . In <i>Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics</i> , volume 108 of <i>Proceedings of Machine Learning Research</i> , pages 2021–2031. PMLR.	1490
		1491
		1492
		1493
	Chao Ren, Han Yu, Hongyi Peng, Xiaoli Tang, Anran Li, Yulan Gao, Alysia Ziyang Tan, Bo Zhao, Xiaoxiao Li, Zengxiang Li, and Qiang Yang. 2024. Advances and open challenges in federated learning with foundation models . <i>arXiv preprint arXiv:2404.15381</i> .	1494
		1495
	Nicola Rieke, Jonny Hancox, Wenqi Li, Fausto Milletari, Holger R Roth, Shadi Albarqouni, Spyridon Bakas, Mathieu N Galtier, Bennett A Landman, Klaus Maier-Hein, et al. 2020. The future of digital health with federated learning . <i>NPJ digital medicine</i> , 3(1):119.	
	Nuria Rodríguez-Barroso, Daniel Jiménez-López, M Victoria Luzón, Francisco Herrera, and Eugenio Martínez-Cámara. 2023. Survey on federated learning threats: Concepts, taxonomy on attacks and defenses, experimental study and challenges . <i>Information Fusion</i> , 90:148–173.	
	Holger R. Roth, Ziyue Xu, Yuan-Ting Hsieh, Adithya Renduchintala, Isaac Yang, Zhihong Zhang, Yuhong Wen, Sean Yang, Kevin Lu, Kristopher Kersten, Camir Ricketts, Daguang Xu, Chester Chen, Yan Cheng, and Andrew Feng. 2024. Empowering federated learning for massive models with nvidia flare . <i>arXiv preprint arXiv:2402.07792</i> .	
	Wang Rui, Yu Tong, Zhang Ruiyi, Kim Sungchul, Rossi Ryan A., Zhao Handong, Wu Junda, Mitra Subrata, Yao Lina, and Henao Ricardo. 2024. Personalized federated learning for text classification with gradient-free prompt tuning . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> .	
	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter . <i>arXiv preprint arXiv:1910.01108</i> .	
	Sejin Seo, Seung-Woo Ko, Jihong Park, Seong-Lyun Kim, and Mehdi Bennis. 2021. Communication-efficient and personalized federated lottery ticket learning . In <i>2021 IEEE 22nd International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)</i> , pages 581–585.	
	Suhail Mohamad Shah and Vincent K. N. Lau. 2023. Model compression for communication efficient federated learning . <i>IEEE Transactions on Neural Networks and Learning Systems</i> , 34(9):5937–5951.	
	Lujia Shen, Shouling Ji, Xuhong Zhang, Jinfeng Li, Jing Chen, Jie Shi, Chengfang Fang, Jianwei Yin,	

1496	and Ting Wang. 2021. Backdoor pre-trained models can transfer to all . In <i>Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security, CCS '21</i> . ACM.	1553
1497		1554
1498		1555
1499		1556
1500	Jaemin Shin, Hyungjun Yoon, Seungjoo Lee, Sungjoon Park, Yunxin Liu, Jinho Choi, and Sung-Ju Lee. 2023a. FedTherapist: Mental health monitoring with user-generated linguistic expressions on smartphones via federated learning . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 11971–11988, Singapore. Association for Computational Linguistics.	1557
1501		1558
1502		1559
1503		1560
1504		1561
1505		1562
1506		
1507		
1508	Jiyun Shin, Jinhyun Ahn, Honggu Kang, and Joonhyuk Kang. 2023b. Fedsplitx: Federated split learning for computationally-constrained heterogeneous clients . <i>arXiv preprint arXiv:2310.14579</i> .	1563
1509		1564
1510		1565
1511		1566
1512	Abhishek Singh, Praneeth Vepakomma, Otkrist Gupta, and Ramesh Raskar. 2019. Detailed comparison of communication efficiency of split learning and federated learning . <i>arXiv preprint arXiv:1909.09145</i> .	1567
1513		1568
1514		1569
1515		1570
1516	Shangchao Su, Bin Li, and Xiangyang Xue. 2023. Fedra: A random allocation strategy for federated tuning to unleash the power of heterogeneous clients . <i>arXiv preprint arXiv:2311.11227</i> .	1571
1517		1572
1518		1573
1519		1574
1520		1575
1521		1576
1522	Shangchao Su, Mingzhao Yang, Bin Li, and Xiangyang Xue. 2024. Federated adaptive prompt tuning for multi-domain collaborative learning . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(13):15117–15125.	1577
1523		1578
1524		1579
1525		1580
1526		1581
1527	Guangyu Sun, Matias Mendieta, Jun Luo, Shandong Wu, and Chen Chen. 2023. Fedperfix: Towards partial model personalization of vision transformers in federated learning . In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 4988–4998.	1582
1528		1583
1529		1584
1530		1585
1531	Guangyu Sun, Matias Mendieta, Taojiannan Yang, and Chen Chen. 2022a. Conquering the communication constraints to enable large pre-trained models in federated learning . <i>arXiv preprint arXiv:2210.01708</i> .	1586
1532		1587
1533		1588
1534		1589
1535	Jingwei Sun, Ziyue Xu, Hongxu Yin, Dong Yang, Daguang Xu, Yiran Chen, and Holger R Roth. 2024a. Fedbpt: Efficient federated black-box prompt tuning for large language models . In <i>Proceedings of the 41th International Conference on Machine Learning</i> .	1590
1536		1591
1537		
1538		
1539		
1540	Tianxiang Sun, Zhengfu He, Hong Qian, Yunhua Zhou, Xuanjing Huang, and Xipeng Qiu. 2022b. BBTv2: Towards a gradient-free future with large language models . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3916–3930, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1592
1541		1593
1542		1594
1543		1595
1544		1596
1545		1597
1546		
1547	Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. 2022c. Black-box tuning for language-model-as-a-service . In <i>Proceedings of the 39th International Conference on Machine Learning</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 20841–20855. PMLR.	1598
1548		1599
1549		1600
1550		1601
1551		1602
1552		1603
		1604
		1605
		1606
		1607
	Youbang Sun, Zitao Li, Yaliang Li, and Bolin Ding. 2024b. Improving loRA in privacy-preserving federated learning . In <i>The Twelfth International Conference on Learning Representations</i> .	1608
		1609
		1610
		1611
	Rishub Tamirisa, John Won, Chengjun Lu, Ron Arel, and Andy Zhou. 2024. Fedselect: Customized selection of parameters for fine-tuning during personalized federated learning . In <i>Federated Learning and Analytics in Practice: Algorithms, Systems, Applications, and Opportunities</i> .	1612
		1613
		1614
		1615
		1616
		1617
		1618
		1619
		1620
		1621
		1622
		1623
		1624
		1625
		1626
		1627
		1628
		1629
		1630
		1631
		1632
		1633
		1634
		1635
		1636
		1637
		1638
		1639
		1640
		1641
		1642
		1643
		1644
		1645
		1646
		1647
		1648
		1649
		1650
		1651
		1652
		1653
		1654
		1655
		1656
		1657
		1658
		1659
		1660
		1661
		1662
		1663
		1664
		1665
		1666
		1667
		1668
		1669
		1670
		1671
		1672
		1673
		1674
		1675
		1676
		1677
		1678
		1679
		1680
		1681
		1682
		1683
		1684
		1685
		1686
		1687
		1688
		1689
		1690
		1691
		1692
		1693
		1694
		1695
		1696
		1697
		1698
		1699
		1700
		1701
		1702
		1703
		1704
		1705
		1706
		1707
		1708
		1709
		1710
		1711
		1712
		1713
		1714
		1715
		1716
		1717
		1718
		1719
		1720
		1721
		1722
		1723
		1724
		1725
		1726
		1727
		1728
		1729
		1730
		1731
		1732
		1733
		1734
		1735
		1736
		1737
		1738
		1739
		1740
		1741
		1742
		1743
		1744
		1745
		1746
		1747
		1748
		1749
		1750
		1751
		1752
		1753
		1754
		1755
		1756
		1757
		1758
		1759
		1760
		1761
		1762
		1763
		1764
		1765
		1766
		1767
		1768
		1769
		1770
		1771
		1772
		1773
		1774
		1775
		1776
		1777
		1778
		1779
		1780
		1781
		1782
		1783
		1784
		1785
		1786
		1787
		1788
		1789
		1790
		1791
		1792
		1793
		1794
		1795
		1796
		1797
		1798
		1799
		1800
		1801
		1802
		1803
		1804
		1805
		1806
		1807
		1808
		1809
		1810
		1811
		1812
		1813
		1814
		1815
		1816
		1817
		1818
		1819
		1820
		1821
		1822
		1823
		1824
		1825
		1826
		1827
		1828
		1829
		1830
		1831
		1832
		1833
		1834
		1835
		1836
		1837
		1838
		1839
		1840
		1841
		1842
		1843
		1844
		1845
		1846
		1847
		1848
		1849
		1850
		1851
		1852
		1853
		1854
		1855
		1856
		1857
		1858
		1859
		1860
		1861
		1862
		1863
		1864
		1865
		1866
		1867
		1868
		1869
		1870
		1871
		1872
		1873
		1874
		1875
		1876
		1877
		1878
		1879
		1880
		1881
		1882
		1883
		1884
		1885
		1886
		1887
		1888
		1889
		1890
		1891
		1892
		1893
		1894
		1895
		1896
		1897
		1898
		1899
		1900
		1901
		1902
		1903
		1904
		1905
		1906
		1907
		1908
		1909
		1910
		1911
		1912
		1913
		1914
		1915
		1916
		1917
		1918
		1919
		1920
		1921
		1922
		1923
		1924
		1925
		1926
		1927
		1928
		1929
		1930
		1931
		1932
		1933
		1934
		1935
		1936
		1937
		1938
		1939
		1940
		1941
		1942
		1943
		1944
		1945
		1946
		1947
		1948
		1949
		1950
		1951
		1952
		1953
		1954
		1

1608	Eric Wang. 2023. Alpaca-lora .	
1609	Haoyu Wang, Handong Zhao, Yaqing Wang, Tong Yu,	
1610	Jiuxiang Gu, and Jing Gao. 2022. Fedkc: Federated	
1611	knowledge composition for multilingual natural	
1612	language understanding . In <i>Proceedings of the ACM</i>	
1613	<i>Web Conference 2022, WWW '22</i> , page 1839–1850,	
1614	New York, NY, USA. Association for Computing	
1615	Machinery.	
1616	Xin'ao Wang, Huan Li, Ke Chen, and Lidan Shou. 2023.	
1617	Fedbfpt: An efficient federated learning framework	
1618	for bert further pre-training . In <i>Proceedings of the</i>	
1619	<i>Thirty-Second International Joint Conference on Arti-</i>	
1620	<i>ficial Intelligence, IJCAI-23</i> , pages 4344–4352. Inter-	
1621	national Joint Conferences on Artificial Intelligence	
1622	Organization. Main Track.	
1623	Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H.	
1624	Yang, Farhad Farokhi, Shi Jin, Tony Q. S. Quek, and	
1625	H. Vincent Poor. 2020. Federated learning with dif-	
1626	ferential privacy: Algorithms and performance analy-	
1627	sis . <i>IEEE Transactions on Information Forensics and</i>	
1628	<i>Security</i> , 15:3454–3469.	
1629	Orion Weller, Marc Marone, Vladimir Braverman,	
1630	Dawn Lawrie, and Benjamin Van Durme. 2022. Pre-	
1631	trained models for multilingual federated learning .	
1632	In <i>Proceedings of the 2022 Conference of the North</i>	
1633	<i>American Chapter of the Association for Computa-</i>	
1634	<i>tional Linguistics: Human Language Technologies</i> ,	
1635	pages 1413–1421, Seattle, United States. Association	
1636	for Computational Linguistics.	
1637	Herbert Woisetschläger, Alexander Isenko, Shiqiang	
1638	Wang, Ruben Mayer, and Hans-Arno Jacobsen.	
1639	2024. A survey on efficient federated learning meth-	
1640	ods for foundation model training . <i>arXiv preprint</i>	
1641	<i>arXiv:2401.04472</i> .	
1642	Bichen Wu, Chenfeng Xu, Xiaoliang Dai, Alvin Wan,	
1643	Peizhao Zhang, Zhicheng Yan, Masayoshi Tomizuka,	
1644	Joseph Gonzalez, Kurt Keutzer, and Peter Vajda.	
1645	2020a. Visual transformers: Token-based image	
1646	representation and processing for computer vision .	
1647	<i>arXiv preprint arXiv:2006.03677</i> .	
1648	Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng	
1649	Wan, and Hong Lin. 2023a. Ai-generated content	
1650	(aigc): A survey . <i>arXiv preprint arXiv:2304.06632</i> .	
1651	Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang,	
1652	Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu,	
1653	Hengshu Zhu, Qi Liu, Hui Xiong, and Enhong Chen.	
1654	2023b. A survey on large language models for rec-	
1655	ommendation . <i>arXiv preprint arXiv:2305.19860</i> .	
1656	Panlong Wu, Kangshuo Li, Ting Wang, and Fangxin	
1657	Wang. 2023c. Fedms: Federated learning with mix-	
1658	ture of sparsely activated foundations models . <i>arXiv</i>	
1659	<i>preprint arXiv:2312.15926</i> .	
1660	Shijie Wu and Mark Dredze. 2020. Are all languages	
1661	created equal in multilingual BERT? In <i>Proceedings</i>	
1662	<i>of the 5th Workshop on Representation Learning for</i>	
1663	<i>NLP</i> , pages 120–130, Online. Association for Com-	
1664	putational Linguistics.	
	Zhaoxian Wu, Qing Ling, Tianyi Chen, and Geor-	1665
	gios B. Giannakis. 2020b. Federated variance-	1666
	reduced stochastic gradient descent with robustness	1667
	to byzantine attacks . <i>IEEE Transactions on Signal</i>	1668
	<i>Processing</i> , 68:4583–4596.	1669
	Chulin Xie, Keli Huang, Pin-Yu Chen, and Bo Li. 2020.	1670
	Dba: Distributed backdoor attacks against federated	1671
	learning . In <i>International Conference on Learning</i>	1672
	<i>Representations</i> .	1673
	Yuexiang Xie, Zhen Wang, Dawei Gao, Daoyuan Chen,	1674
	Liuyi Yao, Weirui Kuang, Yaliang Li, Bolin Ding,	1675
	and Jingren Zhou. 2023. Federatedscope: A flexible	1676
	federated learning platform for heterogeneity . <i>Proc.</i>	1677
	<i>VLDB Endow.</i> , 16(5):1059–1072.	1678
	Chang Xu, Yu Jia, Liehuang Zhu, Chuan Zhang, Guoxie	1679
	Jin, and Kashif Sharif. 2022. Tddl: Truth discovery	1680
	based byzantine robust federated learning . <i>IEEE</i>	1681
	<i>Transactions on Parallel and Distributed Systems</i> ,	1682
	33(12).	1683
	Mengwei Xu, Dongqi Cai, Yaozong Wu, Xiang Li,	1684
	and Shangguang Wang. 2024a. Fwdllm: Effi-	1685
	cient fedllm using forward gradient . <i>arXiv preprint</i>	1686
	<i>arXiv:2308.13894</i> .	1687
	Mengwei Xu, Wangsong Yin, Dongqi Cai, Rongjie	1688
	Yi, Daliang Xu, Qipeng Wang, Bingyang Wu,	1689
	Yihao Zhao, Chen Yang, Shihe Wang, Qiyang	1690
	Zhang, Zhenyan Lu, Li Zhang, Shangguang Wang,	1691
	Yuanchun Li, Yunxin Liu, Xin Jin, and Xuanzhe	1692
	Liu. 2024b. A survey of resource-efficient llm	1693
	and multimodal foundation models . <i>arXiv preprint</i>	1694
	<i>arXiv:2401.08092</i> .	1695
	Minrui Xu, Hongyang Du, Dusit Niyato, Jiawen Kang,	1696
	Zehui Xiong, Shiwen Mao, Zhu Han, Abbas Ja-	1697
	malipour, Dong In Kim, Xuemin Shen, Victor C. M.	1698
	Leung, and H. Vincent Poor. 2024c. Unleashing	1699
	the power of edge-cloud generative ai in mobile net-	1700
	works: A survey of aigc services . <i>IEEE Communica-</i>	1701
	<i>tions Surveys & Tutorials</i> , pages 1–1.	1702
	Zheng Xu, Yanxiang Zhang, Galen Andrew, Christopher	1703
	Choquette, Peter Kairouz, Brendan McMahan, Jesse	1704
	Rosenstock, and Yuanbo Zhang. 2023. Federated	1705
	learning of gboard language models with differential	1706
	privacy . In <i>Proceedings of the 61st Annual Meet-</i>	1707
	<i>ing of the Association for Computational Linguistics</i>	1708
	<i>(Volume 5: Industry Track)</i> , pages 629–639, Toronto,	1709
	Canada. Association for Computational Linguistics.	1710
	Fu-En Yang, Chien-Yi Wang, and Yu-Chiang Frank	1711
	Wang. 2023a. Efficient model personalization in fed-	1712
	erated learning via client-specific prompt generation .	1713
	In <i>2023 IEEE/CVF International Conference on Com-</i>	1714
	<i>puter Vision (ICCV)</i> , pages 19102–19111.	1715
	Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie,	1716
	Ziyan Kuang, and Sophia Ananiadou. 2023b. To-	1717
	wards interpretable mental health analysis with large	1718
	language models . <i>arXiv preprint arXiv:2304.03347</i> .	1719

1829 survey of large language models. *arXiv preprint*
1830 *arXiv:2303.18223*.

1831 Fei Zheng, Chaochao Chen, Lingjuan Lyu, and Binhui
1832 Yao. 2023. Reducing communication for split learn-
1833 ing by randomized top-k sparsification. In *Proceed-*
1834 *ings of the Thirty-Second International Joint Confer-*
1835 *ence on Artificial Intelligence, IJCAI '23*.

1836 Guo Zhihan, Zhang Yifei, Zhang Zhuo, Xu Zenglin,
1837 and King Irwin. 2024. Fedlfc: Towards efficient
1838 federated multilingual modeling with lora-based lan-
1839 guage family clustering. In *Proceedings of the 2024*
1840 *Conference of the North American Chapter of the*
1841 *Association for Computational Linguistics: Human*
1842 *Language Technologies*. Association for Computa-
1843 tional Linguistics.

1844 Ligeng Zhu, Zhijian Liu, and Song Han. 2019. Deep
1845 leakage from gradients. In *Advances in Neural In-*
1846 *formation Processing Systems*, volume 32. Curran
1847 Associates, Inc.

1848 Weiming Zhuang, Chen Chen, and Lingjuan Lyu. 2023.
1849 When foundation model meets federated learning:
1850 Motivations, challenges, and future directions. *arXiv*
1851 *preprint arXiv:2306.15546*.

1852 A Additional Details of FedPEFT

1853 A.1 Additive Methods

1854 A.1.1 Adapter Tuning

1855 Adapter tuning integrates small-scale neural net-
1856 works (known as “adapters”) into the pre-trained
1857 models (Houlsby et al., 2019; Hu et al., 2022). A
1858 straightforward implementation of adapter tuning
1859 is to collaboratively train a shared adapter among
1860 all clients in the FedAvg manner, as highlighted
1861 by Sun et al. (2022a). Based on FedAvg, Fed-
1862 CLIP (Lu et al., 2023a) incorporates an attention-
1863 based adapter for the image encoder in CLIP mod-
1864 els (Radford et al., 2021). In the domain of mul-
1865 tilingual machine translation, where different lan-
1866 guage pairs exhibit substantial discrepancies in data
1867 distributions, Fed-MNMT (Liu et al., 2023d) ex-
1868 plores clustering strategies that group adapter pa-
1869 rameters and makes inner-cluster parameters ag-
1870 gregation for alleviating the undesirable effect of
1871 data discrepancy. Another representative approach
1872 named C2A (Kim et al., 2023) employs hypernet-
1873 works (Ha et al., 2017) to generate client-specific
1874 adapters by conditioning on the client’s informa-
1875 tion, maximizing the utility of shared model param-
1876 eters while minimizing the divergence caused by
1877 data heterogeneity.

1878 A.1.2 Prompt Tuning

1879 Prompt tuning incorporates trainable task-specific
1880 continuous prompt vectors at the input layer (Liu
1881 et al., 2023a; Dong et al., 2023). Compared to full
1882 fine-tuning, it achieves comparable performance
1883 but with $1000\times$ less parameter storage and com-
1884 munication (Jia et al., 2022). A variation of prompt
1885 tuning, FedPerfix (Sun et al., 2023) uses a local
1886 adapter to generate the prefixes and aggregate the
1887 original self-attention layers.

1888 Depending on target modalities, prompt tuning
1889 in current literature can be further classified into
1890 three categories:

- 1891 • *Textual Prompt Tuning*. Task-specific prompt em-
1892 beddings are combined with the input text embed-
1893 dings, which are subsequently fed into language
1894 models. These soft prompts serve as instructive
1895 contexts to influence the generation process of
1896 LLMs by steering the probability distribution of
1897 the next token (Dong et al., 2023).
- 1898 • *Visual Prompt Tuning*. Taking inspiration from
1899 advances in efficiently tuning LLMs, prompts
1900 are also introduced in the input space of vision
1901 models (Jia et al., 2022). Naive implementa-
1902 tions introduce prompts at the pixel level, act-
1903 ing as a form of data augmentation (Li et al.,
1904 2024a). Alternatively, one could also insert the
1905 prompts as latent vectors for the first Transformer
1906 layer (Deng et al., 2024; Yang et al., 2023a). Nev-
1907 ertheless, an empirical study (Jia et al., 2022) has
1908 suggested that it is easier for visual prompts to
1909 learn condensed task-dependent signals in the
1910 latent input space of Transformers.
- 1911 • *Textual-Visual Prompt Tuning*. Unlike single-
1912 modal FMs, vision-language FMs can process
1913 and interpret both visual data and textual infor-
1914 mation, endowing them with powerful represen-
1915 tation ability and transferability (Radford et al.,
1916 2021). Based on vision-language FMs like CLIP,
1917 textual-visual prompt tuning shows promising ca-
1918 pabilities in FL (Guo et al., 2023), especially in
1919 cross-domain scenarios, where the model needs
1920 to generalize across varied domains and unseen
1921 classes (Qiu et al., 2024).

1922 A.2 Comparison of FedPEFT methods

1923 Figure 2 depicts the taxonomy of FedPEFT with
1924 representative methods. Note that some methods
1925 may belong to multiple overlapping categories. To

Table 1: Comparison of Federated Parameter-Efficient Fine-Tuning (FedPEFT) Methods.

Category	Representative Work	Modality	Model	# Full Params.	# Train. Params.	Training Accel.	Comm. Cost	
Selective	RaFFM (Yu et al., 2023c)	Txt.	BERT-Large (2019)	336M	100M	6.13×	29.8%	
	FedBF (Zhang et al., 2023e)	Txt.	Roberta-Base (2019)	125M	0.66M		1.6%	
Additive	Adapter	FedAP (Zhang et al., 2023e)	Txt.	Roberta-Base (2019)	125M	2M		1.6%
		FedCLIP (Lu et al., 2023a)	Vis.-Txt.	ViT-B/32 (2020a)	150M	0.53M		3.5%
		FedDAT (Chen et al., 2024)	Vis.-Txt.	ALBEF (2021b)	290M	2.86M		9.9%
		C2A (Kim et al., 2023)	Txt.	DistilBERT (2020)	66M	0.06M		0.1%
		Fed-MNMT (Liu et al., 2023d)	Txt.	mBART-50 (2020)	611M	8M		1.3%
		AdaFL (Cai et al., 2023)	Txt.	BERT (2019)	110M	0.61M	1.63×	0.6%
	Prompt	PromptFL (Guo et al., 2023)	Vis.-Txt.	ViT-B/16 (2021)	87M	0.87M	2.38×	0.9%
		MFPT (Zhao et al., 2024b)	Txt.	XLM-RoBERTa (2020)	270M	1.2M		0.4%
		FedAPT (Su et al., 2024)	Vis.-Txt.	ViT-B/32 (2020a)	88M	2.8M		3.2%
		FedSP (Dong et al., 2023)	Txt.	GPT2-XL (2019)	1.6B	111M		0.5%
Reparameterization-based Methods	SLoRA (Babakniya et al., 2023a)	Txt.	DistilBERT (2020)	67M	0.7M	13.47×	5.8%	
	LP-FL (Jiang et al., 2023a)	Txt.	BERT-Large (2019)	336M	100M		30%	
	FedMS (Wu et al., 2023c)	Vis.-Txt.	ViT-B/16 (2021)	87M	8.6M		10%	
	pFedS2T (Du et al., 2024)	Aud.	Whisper (2023)	254M	10.1M		4%	
	FFA-LoRA (Sun et al., 2024b)	Txt.	RoBERTa-Large (2019)	355M	0.39M		0.1%	

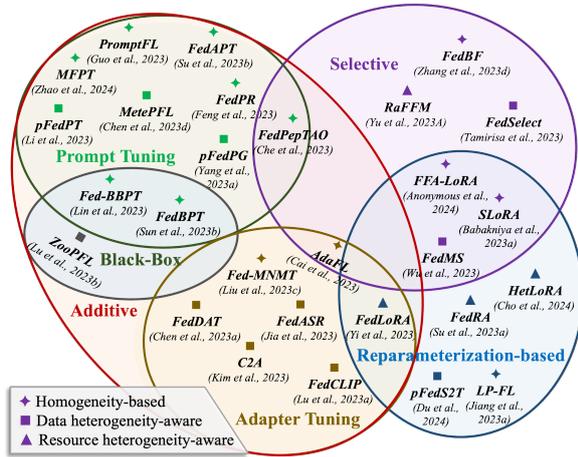


Figure 2: Taxonomy of Federated Parameter-Efficient Fine-Tuning (FedPEFT). Apart from efficiency, some methods also account for other considerations, such as data and resource heterogeneity challenges that are identified in Section 3.2 and black-box tuning (see Section 4.3).

compare the communication efficiency of different FedPEFT methods, Table 1 gives a brief overview of experimental evaluations from representative studies. Compared to full-model fine-tuning, FedPEFT methods only require 0.1%-30% communication overhead. We note that the differences can be attributed to several factors, including model complexity and implementation details.

B Additional Details of Trustworthiness

Due to the distributed characteristic of optimization, FL is vulnerable to poisoning attacks (Lyu

et al., 2022; Rodríguez-Barroso et al., 2023), wherein certain participants may deviate from the prescribed update protocol and upload arbitrary parameters to the central server.

Poisoning Attacks. Depending on the adversarial goals, poisoning attacks in FL can be classified as *targeted* and *untargeted* (Jere et al., 2020). Targeted attacks, like backdoor attacks, aim to manipulate the global model to generate attacker-desired misclassifications for some particular samples (Xie et al., 2020; Bagdasaryan et al., 2020), while untargeted attacks seek to degrade the overall performance of the model indiscriminately (Fang et al., 2020). In addition to the well-recognized attacks on conventional FL studies (Li et al., 2023b, 2024b), FM-FL also faces potential threats from compromised pre-trained FMs (Li et al., 2023c). Thus, The attacker can introduce backdoors to downstream tasks without prior knowledge (Shen et al., 2021). Specifically, Li et al. (2023d) proposed Fed-EBD that introduces a backdoor-compromised FM to generate a public, synthetic dataset for FL training. The clients’ models, pre-trained on this dataset, inherit the backdoor throughout the training.

Defense Techniques. As for defenses, robust aggregation rules are widely applied to make an attack-resilient estimation of the true updates and exclude the influence of malicious updates (Blanchard et al., 2017; Yin et al., 2018; Chen et al., 2017; Li et al., 2023a). Other research directions include trust-based strategies (Cao et al., 2021;

Xu et al., 2022; Park et al., 2021) and variance-reduced algorithms (Gorbunov et al., 2023; Wu et al., 2020b). Although these techniques have been widely examined in various FL settings, their effectiveness has yet to be explored in the FM-FL paradigm.

C Applications of FM-FL

In this part, we briefly review the recent progress on FM-FL applications. Table 2 lists representative work on specific applications and domains.

C.1 FM-FL for Multilingual NLP

Multilingual NLP refers to the techniques that handle multiple natural languages (Pires et al., 2019), often to perform equally well across them (Wu and Dredze, 2020). Earlier research (Johnson et al., 2017) has shown that parameter sharing among different languages boosts the model’s performance in multilingual NLP, especially for low-resource languages for which significantly less content is available. However, real-world multilingual text data is often distributed across devices or regions, with each client (user) accessing only a limited subset of languages, where transferring the data to a central server is often problematic or prohibited due to privacy issues (Wang et al., 2022). Thanks to its inherent privacy-preserving characteristic, FL holds promise in breaking the barriers of cross-lingual modeling and data isolation by allowing models to learn from decentralized datasets.

The pioneer work by Weller et al. (2022) has firstly demonstrated that fine-tuning pre-trained language models with FL can perform similarly to pre-trained models fine-tuned with the standard centralized method under multilingual NLP settings. Various subsequent studies have focused on adapting pre-trained FMs through FedPEFT techniques such as adapter tuning (Liu et al., 2023d), prompt tuning (Zhao et al., 2024b), and LoRA (Zhihan et al., 2024), aiming to enhance training efficiency.

Considering the adverse effect of conflicting parameters from diverse languages during federated fine-tuning, recent studies have exploited clustering strategies to alleviate this issue. For instance, Wang et al. (2022) applied k -means clustering on each client’s data to obtain representative knowledge, specifically the clustered data centroids. These centroids were then shared across clients for local training, enriching training data and addressing the challenges associated with data heterogeneity. Another

compelling strategy along this line is language family-based clustering. Liu et al. (2023d) explored various clustering strategies to group adapter parameters to mitigate the negative effects of multilingual data heterogeneity, showing that language family-based clustering significantly outperforms the other clustering strategies. Similarly, Zhihan et al. (2024) proposed fine-tuning FMs with LoRA and language family-based clustering to address the heterogeneity issue of multilingual modeling.

General downstream tasks include language modeling (Wang et al., 2022), machine translation (Liu et al., 2023d; Chu et al., 2024), and text classification (Weller et al., 2022). In addition, some studies also focus on more specific applications such as medical transcript analysis (Manoel et al., 2023) and hate speech detection (Akshay and Rahul, 2024). These advancements illustrate the applicability of FM-FL across a wide range of scenarios in multilingual NLP.

C.2 FM-FL for Speech

With the development of AI, researchers have also carried out many studies on speech-related FMs, *e.g.*, wav2vec 2.0 (Baevski et al., 2020) and Whisper (Radford et al., 2023). In this field, the adaptation of FMs often relies on FL to facilitate scenarios where the audio data is privacy-sensitive. Compared to other data modalities, speech-related FM-FL applications especially attract excessive attention to the aspects of **on-device training** and **personalization**, motivated by the following considerations: (1) Audio data is continually generated on end-devices such as mobile phones, and owned by individual users—thus it should be processed locally, rather than being transferred elsewhere; (2) Although FL takes advantage of all user data to collectively train one model that maximizes speaker-independent accuracy, such a one-model-fits-all solution can be sub-optimal for individual users (Jia et al., 2023). Specific tasks in this field include Automatic Speech Recognition (ASR) (Azam et al., 2023b) and Speech-to-Text (S2T) (Du et al., 2024).

C.3 FM-FL for Recommendation

Federated Recommendation (FR) strives to capture underlying user preferences and recommend appropriate information to users while safeguarding data privacy (Bobadilla et al., 2013; Zhang et al., 2023a). Typical FR systems consist of a server and multiple clients, where clients represent individual users or local data servers possessing

Table 2: A list of representative studies on the applications of FM-FL. Abbreviations: LoRA Tuning (LT), Adapter Tuning (AT), Full-Parameter Tuning (FT), Selective Tuning (ST), Prompt Tuning (PT).

Domain/Application	Task	Representative Work	On-Device Personalization	Modality	Backbone	Fine-Tuning
Multilingual NLP	Language Understanding	FedKC (Wang et al., 2022)	✗ ✗	Txt.	mBERT	FT
	Multi-Tasks	PMMFL (Weller et al., 2022)	✗ ✗	Txt.	mBERT	FT
	Machine Translation	Fed-MNMT (Liu et al., 2023d)	✗ ✗	Txt.	mBART-50	AT
	Machine Translation	FL-MetaSend (Chu et al., 2024)	✗ ✗	Txt.	M2M-100	ST
	Multi-Tasks	MFPT (Zhao et al., 2024b)	✓ ✗	Txt.	XLm-RoBERTa	PT
Speech	Speech-to-Text	pFedS2T (Du et al., 2024)	✗ ✓	Aud.	Conformer/Whisper	LT
	Speech Recognition	FedASR (Jia et al., 2023)	✓ ✓	Aud.	RNN-T	AT
	Speech Recognition	FedE2EASR (Azam et al., 2023a)	✗ ✓	Aud.	CTC-AED	FT
Recommendation	General	PPLR (Zhao et al., 2024a)	✗ ✓	Txt.	LLaMA-7B/LongFormer	FT
	General	TransFR (Zhang et al., 2024a)	✓ ✓	Txt.	DistBERT	AT
	General	GPT-FedRec (Zeng et al., 2024)	✗ ✗	Txt.	ChatGPT	NA
Healthcare	Mental Health Prediction	FedTherapist (Shin et al., 2023a)	✓ ✗	Txt.	BERT & LLaMa-7B	LT
	MRI Reconstruction	FedPR (Feng et al., 2023a)	✗ ✗	Vis.	Swin Transformers	PT

smaller datasets and retaining private user information (Ammad-Ud-Din et al., 2019). These clients collaborate to train a global model while ensuring their data privacy protection by abstaining from direct data sharing (Zeng et al., 2024; Zhang et al., 2023a). Recently, LLM-based recommendations have been gaining increasing attention (Wu et al., 2023b) due to their strong capacities in language understanding and domain generalization. The benefits are mainly twofold: (1) LLMs mitigate the cold-start issue by utilizing textual descriptions to make recommendations without the need for extensive historical data (Zhang et al., 2023c); (2) The inherent transferability of LLMs allows them to apply cross-domain knowledge and side information to improve accuracy and relevance across diverse items and user interests (Gao et al., 2023).

One straightforward way to adapt FMs for FR is by fine-tuning them with historical user-item data. More specifically, FedPEFT techniques such as adapter tuning (Zhang et al., 2024a) and split learning (Zhao et al., 2024a) can be employed to improve resource efficiency. Apart from parameter fine-tuning, LLMs can also be adapted to assist the recommendation in a zero-shot paradigm through prompt engineering (*i.e.*, without parameter tuning) (Gao et al., 2023). For example, Zeng et al. (2024) proposed GPT-FedRec, a two-stage FR framework that leverages ChatGPT for its powerful zero-shot generalization ability. Firstly, GPT-FedRec facilitates hybrid retrieval by collabora-

tively training ID and text retrievers, after which the retrieved results are transformed into text prompts and submitted to GPT for re-ranking in the second stage. Additionally, Guo et al. (2024) employed a pre-trained BERT to obtain the representation vectors of item descriptions, which are then fed into a recommender system as augmented input.

C.4 FM-FL for Healthcare

FMs, especially LLMs, have been found to excel in healthcare applications, showcasing impressive capabilities in tasks like mental health analysis (Yang et al., 2023b), disease diagnosis (Panagoulas et al., 2024), and drug discovery (Chenthamarakshan et al., 2023). However, it raises privacy concerns to upload the health information of patients (Tang et al., 2023) into a commercial server that supports the FMs. Meanwhile, FL has consistently received widespread attention in the healthcare domain (Lincy and Kowshalya, 2020; Rieke et al., 2020; Joshi et al., 2022), driven by the need for collaborative model training across different medical institutions without compromising patient data privacy. By breaking the barriers of private data availability, the FM-FL paradigm shows the potential to further harness the power of FMs in the healthcare domain.

A recent study (Shin et al., 2023a) presents a mobile mental health monitoring system, FedTherapist, which leverages user speech and keyboard input to fine-tune FMs with FL, demonstrating superior accuracy in mental health prediction tasks such

Table 3: A list of existing FM-FL libraries and benchmarks. **Missing or inapplicable details denoted by N/A.** ✓ denotes a strong focus or presence; ✗ indicates no focus or absence; ● signifies a moderate focus or partial inclusion.

Library/Benchmark	FL Backend	LLM Support	MultiModal FM Support	FedPEFT	On-Device Training	Distributed & Clustered	Differential Privacy	Description
FederatedScope-LLM (Kuang et al., 2023)	FederatedScope	✓	✗	✓	✗	✓	✗	An end-to-end benchmark for efficient fine-tuning LLMs with FL
NVIDIA FLARE (Roth et al., 2024)	NVFlare	✓	✗	✓	✗	✓	✓	Scalable and efficient fine-tuning LLMs with FL
FATE-LLM (Fan et al., 2023)	FATE	✓	✗	✓	✗	✓	✓	Focuses on IP and privacy protection in federated LLM
FedLLM (FedML, 2023)	FedML	✓	✗	●	✓	✓	✓	An MLOps-supported training pipeline based on FedML
OpenFedLLM (Ye et al., 2024)	N/A	✓	✗	●	✗	N/A	✗	An LLM framework focusing on FL instruction tuning/alignment
Shepherd (Zhang et al., 2024b)	N/A	✓	✗	✓	✗	✗	✗	Federated instruction tuning based on Huggingface
FedPETuning (Zhang et al., 2023e)	FedLab	✓	✗	✓	✗	✓	✗	A benchmark comprising four FedPEFT methods
FedLegal (Zhang et al., 2023d)	FedLab	✓	✗	✗	✗	✓	✗	A benchmark comprising six legal NLP tasks under FL settings

as depression, stress, and mood prediction. Another representative study (Feng et al., 2023a) focuses on Magnetic Resonance Imaging (MRI) reconstruction, which involves retrieving a complex-valued image from its under-sampled signal. The authors adopted an FM pre-trained on public datasets and trained visual prompts from decentralized clinical datasets via a personalized FL mechanism, thereby reducing communication costs and achieving competitive performance on limited local data.

Despite the efforts, it has been shown that FMs in healthcare risk generating misleading information due to their imperfect understanding of complex medical data (Jeblick et al., 2024).

D Libraries and Benchmarks

This part briefly introduces a series of available libraries and benchmarks for developing and examining FM-FL techniques. An overview is provided in Table 3.

- **FederatedScope-LLM** (Kuang et al., 2023) is an open-source package for fine-tuning LLMs via FL. Built on top of a popular FL backend FederatedScope (Xie et al., 2023), it supports federated fine-tuning of LLMs under various FL scenarios, including FedPEFT and model personalization.
- **NVIDIA FLARE** (Roth et al., 2024) is an FL framework that allows researchers and data scientists to seamlessly move their machine learning and deep learning workflows into a federated paradigm.
- **FATE-LLM** (Fan et al., 2023) is an industrial-grade FL framework for LLM. Apart from Fed-

PEFT, it provides a privacy hub integrating several IP protection and privacy-preserving mechanisms to protect model security and data privacy.

- **FedLLM** (FedML, 2023) is an MLOps-supported training pipeline built upon the FedML AI platform (He et al., 2020). FedLLM is compatible with popular LLM libraries such as HuggingFace and DeepSpeed to support a large range of FMs and datasets.
- **OpenFedLLM** (Ye et al., 2024) is a federated tuning framework for LLMs, which covers applications of instruction tuning and value alignment, diverse FL baselines, training datasets, and evaluation datasets.
- **Shepherd** (Zhang et al., 2024b) is a lightweight federated tuning framework. The local training process of Shepherd is built upon the implementations of Alpaca-LoRA (Wang, 2023), and Hugging Face’s PEFT (Mangrulkar et al., 2022), enabling efficient fine-tuning.
- **FedPETuning** (Zhang et al., 2023e) is a pioneering federated benchmark for four representative FedPEFT methods, covering adapter tuning, prefix tuning, LoRA, and BitFit.
- **FedLegal** (Zhang et al., 2023d) is the very first real-world FL benchmark for legal NLP, which comprises five legal NLP tasks and one privacy task based on the data from Chinese courts.