# Federated Learning for Privacy-Preserving Large Language Models: A Comprehensive Survey

[1]Kalathiya Parth, [2]Shah Shrey, [3]Shah Shivam, [4]Budheliya Hetvi, [5]Devang Bhatt, [6]Dhaval Chandarana

[1,2,3,4,5]*Gyanmanjari Innovative University*
[1]*parthkalathiya789@gmail.com,* [2]*alpashah231080@gmail.com,*
[3]*shivamsindhukumar@gmail.com,* [4]*hsbudheliya2008@gmail.com,* [5]*djbhatt@gmiu.edu.in,*
[6]*drchandarana@gmiu.edu.in*

**Abstract-The rapid advancement of Large Language Models (LLMs) has revolutionized natural language processing, yet their development requires massive centralized datasets, raising critical privacy concerns. Federated Learning (FL) offers a promising paradigm for training LLMs while preserving data privacy by enabling collaborative learning without direct data sharing. This comprehensive survey examines the integration of federated learning with large language models, covering pre-training, fine-tuning methodologies, deployment strategies, communication optimization, and security mechanisms. We systematically review state-of-the-art approaches including federated full-parameter tuning, parameter-efficient fine-tuning (PEFT), prompt learning, and homomorphic encryption-based methods. Key challenges such as communication costs, data heterogeneity, system scalability, and privacy preservation are analyzed alongside emerging solutions. We present detailed analyses of frameworks like Federated Scope-LLM and techniques achieving communication costs under 18 KB for billion-parameter models. This survey synthesizes current progress, identifies critical research gaps, and outlines future directions for privacy-preserving LLM development across healthcare, finance, recommendation systems, and edge computing applications.**

**Index Terms—Federated Learning, Large Language Models, Privacy-Preserving AI, Parameter-Efficient Fine-Tuning, Communication Efficiency, Homomorphic Encryption, Edge Learning, Prompt Tuning**

## I. INTRODUCTION

### A. MOTIVATION AND BACKGROUND

Large Language Models (LLMs) such as GPT-4, PaLM, and LLaMA have demonstrated unprecedented capabilities in natural language understanding, generation, reasoning, and multi-task learning [1]. These models, often containing billions to hundreds of billions of parameters, are typically trained on massive centralized datasets aggregated from diverse sources. However, this centralized paradigm poses fundamental challenges:

Privacy Concerns:
Collecting sensitive textual data (medical records, financial documents, personal communications) into centralized repositories creates privacy risks and conflicts with regulations such as GDPR, HIPAA, and CCPA [2].

Data Sovereignty:
Organizations and individuals are increasingly reluctant to share proprietary or personal data due to competitive, legal, or ethical concerns [3].

Computational Constraints:
The enormous computational resources required for training LLMs (often requiring thousands of GPUs) are concentrated in a few large technology companies, limiting accessibility [4].

Data Silos:
Valuable data is distributed across multiple institutions, devices, and jurisdictions, making comprehensive data collection impractical or impossible [5].
Federated Learning (FL) has emerged as a transformative approach to address these challenges by enabling collaborative model training while keeping data decentralized [6]. Originally proposed for mobile keyboard prediction, FL has evolved to support increasingly complex models and applications. The integration of FL with LLMs represents a critical frontier that combines the power of large-scale language understanding with robust privacy guarantees.

### B. FEDERATED LEARNING: CORE PRINCIPLES

Federated Learning follows a distributed training paradigm characterized by [7]:
1. Decentralized Data Storage: Training data remains on local devices or institutional servers
2. Local Computation: Model training occurs locally using private data
3. Selective Sharing: Only model updates (gradients or parameters) are transmitted
4. Central Aggregation: A central server aggregates updates to improve the global model
5. Iterative Refinement: The process repeats until convergence

The canonical Federated Averaging (FedAvg) algorithm [8] forms the foundation:

Initialize global model $w_0$

for each round $t = 1, 2, \ldots, T$ do:

Server broadcasts $w_t$ to selected clients

For each client k in parallel do:

$w_{t+1}^k \leftarrow$ Local Update $(k, w_t)$ end for

$w_{t+1} \leftarrow$ Aggregate($\{w_{t+1}^k\}$) end for

## C. UNIQUE CHALLENGES OF FEDERATED LLMS

While FL has been successfully applied to smaller models, integrating it with LLMs introduces unprecedented challenges:

(1) Communication Bottleneck:

Transmitting billion-parameter models over bandwidth-limited networks is prohibitively expensive. A model with 10B parameters requires ~40 GB transmission per round at 32-bit precision [9].

(2) Computational Heterogeneity:

Edge devices vary dramatically in computational capability, from high-performance servers to resource-constrained mobile devices [10].

(3) Memory Constraints:

Loading and training multi-billion parameter models exceed the memory capacity of most edge devices [11].

(4) Data Heterogeneity:

Non-IID (non-independent and identically distributed) data across clients, with varying domains, languages, and quality, degrades model performance [12].

(5) Privacy Vulnerabilities:

Despite not sharing raw data, gradient-based attacks can reconstruct training samples, particularly problematic for text data [13].

(6) Convergence Complexity:

The massive parameter space and non-convex optimization landscape make convergence analysis challenging [14].

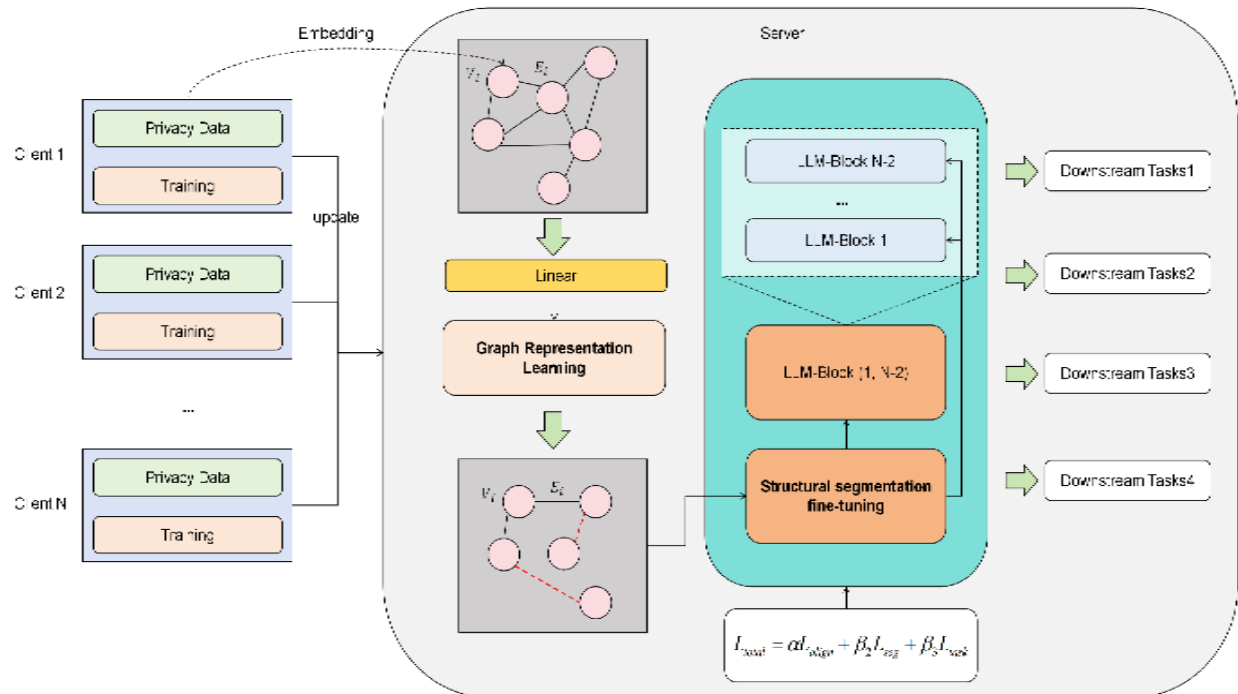## D. CONTRIBUTIONS AND ORGANIZATION

This survey makes the following contributions:

1. Comprehensive Coverage: We systematically review federated LLM research spanning pre-training, fine-tuning, deployment, and security
2. Technical Depth: Detailed analysis of communication-efficient methods, PEFT techniques, and privacy-preserving mechanisms
3. Practical Frameworks: Examination of open-source tools and real-world implementations
4. Future Roadmap: Identification of critical research gaps and promising directions

The remainder of this paper is organized as follows:
- Section II: Federated LLM architectures and training paradigms
- Section III: Communication-efficient fine-tuning methods
- Section IV: Parameter-efficient techniques (LoRA, adapters, prompt tuning)
- Section V: Security and privacy mechanisms
- Section VI: Heterogeneity management strategies
- Section VII: Frameworks and implementation tools
- Section VIII: Applications across domains
- Section IX: Future research directions
- Section X: Conclusions

## II. FEDERATED LLM ARCHITECTURES AND PARADIGMS



### A. Federated LLM Training Taxonomy

Cheng et al. [1] provide an exhaustive taxonomy of federated LLM approaches, categorizing them along three dimensions:

(1) Training Stage:
- Pre-training: Training foundation models from scratch on distributed corpora
- Fine-tuning: Adapting pre-trained models to specific tasks or domains
- Continual Learning: Incrementally updating models with new data

(2) Model Distribution:
- Full Model FL: Entire model resides on each client
- Split Learning: Model is partitioned across clients and server
- Model Offloading: Parts of the model are offloaded to edge servers

(3) Aggregation Strategy:
- Synchronous: All clients update simultaneously
- Asynchronous: Clients update independently
- Hierarchical: Multi-tier aggregation architecture

B. Federated Pre-training of LLMs
While most research focuses on federated fine-tuning, federated pre-training from scratch presents unique opportunities and challenges [1].

Advantages:
- Complete control over training data sources
- Domain-specific foundation models (e.g., medical, legal)
- Avoiding biases from publicly pre-trained models

Challenges:
- Extreme communication costs (thousands of training iterations)
- Convergence to sub-optimal solutions due to data heterogeneity
- Coordinating distributed tokenization and vocabulary

Current Approaches:
Chen et al. [3] propose a systematic framework for federated pre-training that includes:
- Distributed Tokenization: Coordinating vocabulary construction across clients
- Staged Training: Initial centralized pre-training followed by federated adaptation
- Domain Clustering: Grouping clients by domain to reduce heterogeneity

C. Federated Fine-tuning: The Dominant Paradigm
Given the computational expense of pre-training, most practical applications employ federated fine-tuning [2]:

Workflow:
1. A pre-trained LLM (e.g., BERT, GPT-2, LLaMA) serves as initialization
2. Clients fine-tune on local task-specific data
3. Updates are aggregated to create an improved global model
4. The process iterates until convergence

Key Advantages:
- Leverages existing pre-trained models
- Significantly reduced training time and cost
- More stable convergence than pre-training from scratch

Yao et al. [2] provide comprehensive analysis of federated fine-tuning, identifying three critical research challenges:
1. Communication efficiency for large models
2. Personalization vs. generalization trade-offs
3. Handling diverse task distributions across clients

D. Edge and On-Device Learning
The "Federated and Edge Learning for Large Language Models" paradigm [5] extends FL to resource-constrained edge devices:

Edge-Specific Considerations:
- Model Compression: Quantization, pruning, and distillation to fit memory constraints
- Adaptive Computation: Dynamic depth/width adjustment based on device capabilities
- Intermittent Connectivity: Handling disconnections and asynchronous updates
- Energy Efficiency: Minimizing battery consumption on mobile devices

Deployment Strategies:
- Model Sharding: Distributing model components across multiple edge devices
- Collaborative Inference: Multiple devices jointly perform inference
- Hierarchical Aggregation: Local edge servers pre-aggregate before central aggregation

### III. COMMUNICATION-EFFICIENT FEDERATED FINE-TUNING

Communication costs dominate federated LLM training. This section examines state-of-the-art methods for achieving practical communication efficiency.

A. The Communication Challenge
Quantitative Analysis:

For a model with P parameters:
- Full model transmission: 4P bytes (assuming 32-bit floats)
- Gradient transmission: 4P bytes per round
- With R training rounds: 4PR total bytes

For GPT-3 (175B parameters): ~700 GB per round, ~70 TB for 100 rounds
Hilmkil et al. [4] demonstrate that even for smaller transformer models (BERT-Large with 340M parameters):
- Standard FL requires ~1.36 GB per round
- With 1000 clients and 100 rounds: ~136 TB total communication
- On typical edge networks (10 Mbps): ~300 hours of transmission time

B. Breakthrough: Sub-18KB Communication
Qin et al. [6] achieve a remarkable breakthrough: federated full-parameter tuning of billion-size models with communication cost under 18 KB per round—a reduction of over 99.99%.

Key Techniques:
(1) Zeroth-Order Optimization: Approximating gradients using function evaluations rather than backpropagation:
$$\nabla f(w) \approx [f(w + \delta u) - f(w)] / \delta \cdot u$$
where u is a random direction vector. Only the scalar evaluation difference is transmitted.

(2) Sparse Update Mechanisms:
- Identify top-k most important parameters using sensitivity analysis
- Transmit only indices and values of these parameters
- Server maintains momentum accumulation for dropped parameters

(3) Quantization Schemes:
- Adaptive bit-width allocation based on parameter importance
- Non-uniform quantization with learned quantization boundaries
- Error accumulation and correction across rounds

(4) Differential Privacy Integration: Adding calibrated noise for privacy while maintaining utility:
$$w'_i = w_i + \text{Gaussian}(0, \sigma^2 C^2/\varepsilon^2)$$
where C is the clipping threshold and $\varepsilon$ is the privacy budget.

Experimental Results: On BERT-Large fine-tuning tasks:
- Communication: 17.8 KB per round (vs. 1.36 GB baseline)
- Accuracy: 98.7% of full-gradient performance
- Privacy: $(\varepsilon, \delta)$-differential privacy with $\varepsilon = 2.0$

C. Gradient Compression Techniques
Beyond the sub-18KB breakthrough, various gradient compression methods have been developed:

(1) Quantization:
Uniform Quantization:
$Q(x) = \text{round}(x / \Delta) \cdot \Delta$
where $\Delta$ is the quantization step size.

Stochastic Quantization:
Reduces quantization bias through probabilistic rounding:
$P(Q(x) = \text{ceil}(x/\Delta)) = (x \bmod \Delta) / \Delta$

(2) Sparsification:
Top-k Gradient Selection: Transmit only the k largest gradients by magnitude:
$S = \{i: |g_i| \in \text{top-k}(|g|)\}$

Random Sparsification: Randomly sample gradients with probability proportional to magnitude:
$P(\text{select } g_i) = |g_i| / \Sigma_j |g_j|$

(3) Low-Rank Decomposition:
Approximate gradient matrix Gas product of low-rank matrices:
$G \approx UV^T$
where $U \in \mathbb{R}^{n \times r}$, $V \in \mathbb{R}^{m \times r}$, and $r \ll \min(n,m)$

D. Layer-Wise Adaptive Communication
Different layers contribute differently to model performance. Adaptive strategies allocate communication budget accordingly:
Importance Scoring:
$\text{importance}(\text{layer\_i}) = \|\nabla L/\nabla w_i\|_2$

Adaptive Budget Allocation:
Layers with higher importance receive larger communication budgets, enabling more precise gradient transmission.

Empirical Findings [4]:
- Attention layers typically have highest importance (40-50% of budget)
- Feed-forward layers: 30-40%
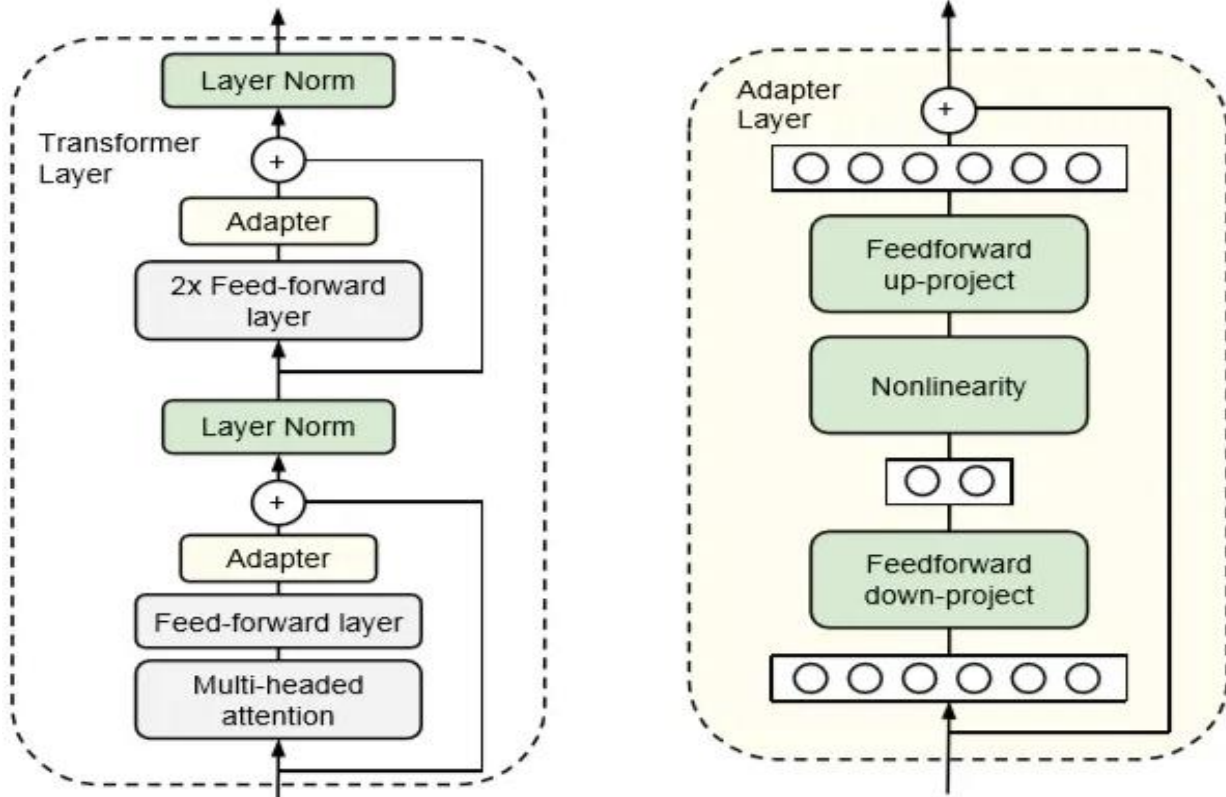- Embedding layers: 10-20%

## IV. PARAMETER-EFFICIENT FINE-TUNING (PEFT)

PEFT methods achieve drastic reductions in trainable parameters, directly addressing communication and computation challenges.

A. Adapter-Based Methods

Architecture: Small neural modules inserted between transformer layers [15]:



Adapter Structure:

$h' = h + f(Down(h))$

$f(x) = Up\ (ReLU(x))$

where Down: $d \rightarrow r$, up: $r \rightarrow d$, and $r \ll d$

Communication Benefits:
- Original BERT-Large: 340M parameters
- With adapters (r=64): ~3.6M trainable parameters
- Communication reduction: 99% per round

Challenges in Federated Settings:
- Adapter placement decisions (which layers?)
- Initialization strategies across heterogeneous clients
- Aggregation of adapter parameters with different local task distributions

B. Low-Rank Adaptation (LoRA)

LoRA [16] has become the dominant PEFT method for federated LLMs due to its simplicity and effectiveness.

Core Principle:

Weight updates during fine-tuning have low intrinsic dimensionality. Represent updates as:

$W = W_0 + \Delta W = W_0 + BA$

where:

- $W_0$: frozen pre-trained weights ($d \times k$)
- B: learnable low-rank matrix ($d \times r$)
- A: learnable low-rank matrix ($r \times k$)
- $r \ll \min(d, k)$

Forward Pass:

$h = W_0 x + BAx$

Federated LoRA Workflow [4]:

1. Initialization:
o Server broadcasts $W_0$ (once, at start)
o Initialize B, A on each client (random or zero)

2. Local Training:
for each local epoch:
Freeze $W_0$
Update B, A via gradient descent end for

3. Aggregation:
$B\_global = \Sigma_i (n_i/n) \cdot B_i$
$A\_global = \Sigma_i (n_i/n) \cdot A_i$

4. Broadcast:
o Send updated B global, A global to clients

Communication Analysis:

- Full fine-tuning: 4Pdk bytes per round
- LoRA: 4Pr(d+k) bytes per round
- For GPT-3 with r=16: Communication reduction: 99.996%

Performance [4]: On Super GLUE benchmark:

- LoRA (r=8): 97.2% of full fine-tuning performance

- LoRA (r=16): 98.9% of full fine-tuning performance
- Communication: 0.01% of full fine-tuning

C. Federated Prompt Learning

FedPromptL Framework [7]:

Guo et al. [7] introduce federated prompt learning for vision-language models with strong privacy efficiency.

Prompt Tuning Basics:

Instead of modifying model parameters, learn continuous prompt embeddings:

Input: $[P_1, P_2, ..., P_k, x_1, x_2, ..., x_n]$

where $P_i$ are learnable prompt tokens (typically k=10-50).

Federated Prompt Learning Architecture:

(1) Local Prompt Optimization:

Min P L (Model (P $\oplus$ x), y)

where $\oplus$ denotes concatenation.

(2) Prototype-Based Aggregation:

Instead of directly averaging prompts, FedPromptL learns prototype representations:

$Prototype_i = f\_enc (P_i)$

$Prototype\_global = Aggregate(\{Prototype_i\})$

$P\_global = f\_dec(Prototype\_global)$

(3) Privacy Preservation:

- Prompts contain abstract task knowledge, not raw data
- Gradient perturbation for differential privacy
- Secure aggregation protocols

Advantages:

- Ultra-low communication: Only 10-50 embeddings × d dimensions
- Task flexibility: Easy adaptation to new tasks
- Privacy: Prompts less susceptible to inversion attacks than full gradients

Experimental Results [7]: On visual question answering:

- Communication: 0.001% of full fine-tuning
- Accuracy: 96.8% of full fine-tuning
- Privacy: $(\varepsilon=1.0, \delta=10^{-5})$-DP achieved

D. Cross-Domain Prompt Learning for Recommendation

Lei Guo et al. [9] extend prompt learning to cross-domain federated recommendation systems.

Problem Setting:

Multiple platforms (e-commerce, social media, streaming) want to improve recommendations without sharing user data.

Approach:

(1) Domain-Specific Prompts: Each domain learns specialized prompts:

P ecommerce = $[P_1{}^e, P_2{}^e, ..., P_k{}^e]$

P social = $[P_1{}^s, P_2{}^s, ..., P_k{}^s]$

(2) Shared Representation Learning:

H shared = Encoder ([P domain, user features, item features])

recommendation score = Decoder (h shared)

(3) Federated Aggregation:

- Share only prompt embeddings across domains
- Keep user/item embeddings private
- Aggregate through secure multi-party computation

Benefits:

- Cold-start problem mitigation through cross-domain knowledge transfer
- Privacy-preserving collaborative filtering
- Personalization while leveraging collective intelligence

## V. SECURITY AND PRIVACY MECHANISMS

Privacy preservation is the fundamental motivation for federated LLMs, yet the approach introduces new vulnerabilities.

A. Privacy Threat Models

(1) Gradient Leakage Attacks:

Adversaries can reconstruct training data from shared gradients [13]:

Attack Mechanism: Given gradient $g = \nabla L (w, x)$, solve optimization:

$x^* = \text{argmin\_x} \|g - \nabla L (w, x) \|_2$

Effectiveness for Text:

- Exact token recovery for small batches ($B \leq 8$)
- Partial recovery for medium batches ($8 < B \leq 32$)
- Semantic recovery even for large batches

Example Attack [13]: From gradient of "The patient has diabetes", attacker recovers:
- Full sentence for B=1
- "Patient ... diabetes" for B=16
- Medical domain identification for B=64

(2) Model Inversion:
Extract sensitive information embedded in model parameters:
$x^* = argmax\_x\ P(x \mid w)$
Particularly concerning for memorization of rare or sensitive training examples.

(3) Membership Inference:
Determine if specific data was used in training:
if $L(w, x) <$ threshold:
return "x was in training data"
Success rates up to 80% for LLMs [13].

B. Homomorphic Encryption for Secure Aggregation
FedML-HE Framework [8]:
Jin et al. [8] introduce efficient homomorphic encryption-based federated learning for LLMs.

Homomorphic Encryption Basics:
Encryption scheme allowing computation on ciphertexts:
$Enc(x_1) \oplus Enc(x_2) = Enc(x_1 + x_2)$
$Enc(x_1) \otimes Enc(x_2) = Enc(x_1 \cdot x_2)$

FedML-HE Protocol:
(1) Client-Side Encryption:
for each parameter $w_i$:
$c_i = Encrypt(w_i, public\_key)$ send $c_i$ to server
end for

(2) Server-Side Aggregation:
$c\_agg = \Sigma_i\ (n_i/n) \otimes c_i$ // encrypted weighted average

(3) Decryption:
$w\_global = Decrypt(c\_agg, private\_key)$

Optimizations for LLMs:
(a) Batch Encryption: Pack multiple parameters into single ciphertext:
$Enc([w_1, w_2, ..., w_m])$ instead of $Enc(w_1), Enc(w_2), ...$

(b) Quantization Before Encryption: Reduce ciphertext size by quantizing to 8-bit:

$w_i' = $ Quantize $(w_i, 8\text{-bit})$

$c_i = $ Encrypt$(w_i')$

(c) Partial Encryption: Encrypt only sensitive layers (e.g., first/last layers) while using secure aggregation for others.

Performance Trade-offs [8]:

| Method | Communication Overhead | Computation Overhead | Privacy Level |
|---|---|---|---|
| No Protection | 1× | 1× | None |
| Differential Privacy | 1× | 1.1× | $(\varepsilon, \delta)$-DP |
| Secure Aggregation | 1.2× | 2.5× | Honest-but-curious |
| FedML-HE | 3.5× | 8.2× | Malicious adversary |

Practical Insights:
- HE suitable for high-security applications (healthcare, finance)
- Significant computational cost limits applicability
- Hybrid approaches combining HE with other techniques recommended

C. Differential Privacy

Differential Privacy (DP) provides rigorous mathematical privacy guarantees [17].

Definition: A mechanism M satisfies $(\varepsilon, \delta)$-DP if for any two datasets D, D' differing in one record:

$P[M(D) \in S] \leq e^{\wedge}\varepsilon \cdot P[M(D') \in S] + \delta$

DP-SGD for Federated LLMs:

(1) Gradient Clipping:

$\bar{g}_i = g_i / \max(1, \|g_i\|_2 / C)$

where C is the clipping threshold.

(2) Noise Addition:

$\tilde{g} = \bar{g} + $ Gaussian $(0, \sigma^2 C^2 I)$

(3) Privacy Accounting: Using Rényi Differential Privacy (RDP) for tight composition:

$\varepsilon_{total} = f(\varepsilon\_step, T, sampling\_rate)$

Challenges for LLMs:
- Large vocabulary requires higher noise for same privacy
- Gradient clipping interferes with learning dynamics

- Privacy-utility trade-off more severe than smaller models

Adaptive Strategies [2]:
- Layer-wise privacy budgets (more noise for sensitive layers)
- Privacy amplification via subsampling
- Public pre-training + private fine-tuning

D. Byzantine-Robust Aggregation

Malicious clients may send corrupted updates to degrade model quality.

Attack Scenarios:
- Label flipping attacks
- Gradient sign flipping
- Gaussian noise injection
- Backdoor attacks (trigger-based misclassification)

Defense Mechanisms:

(1) Robust Aggregation Rules:

Median Aggregation:

w_global[i] = median($\{w_1[i], w_2[i], ..., w_n[i]\}$)

Trimmed Mean:

w_global[i] = mean(trim($\{w_1[i], ..., w_n[i]\}, \alpha$))

where $\alpha$ fraction of extreme values are removed.

Krum: Selects client update with smallest distance to k-nearest neighbors.

(2) Anomaly Detection:
- Statistical testing for outlier detection
- Norm-based filtering
- Cosine similarity thresholds

(3) Reputation Systems:
- Track historical accuracy of client updates
- Weight aggregation by reputation scores
- Exclude consistently malicious clients

## VI. HANDLING HETEROGENEITY IN FEDERATED LLMS

Heterogeneity across clients poses fundamental challenges for federated LLM training [2].

A. Data Heterogeneity

Types of Data Heterogeneity:

(1) Label Distribution Skew:
Different clients have different task distributions.
Example: Medical institutions specializing in different diseases.

(2) Feature Distribution Skew:
Input distributions vary across clients.
Example: Different dialects or technical jargons across regions.

(3) Quantity Imbalance:
Clients have vastly different dataset sizes ($10^3$ to $10^7$ samples).

(4) Domain Shift:
Clients represent entirely different domains (legal vs. medical vs. social media).

Impact on Performance [2]:
- Standard FedAvg can degrade accuracy by 20-40% compared to centralized training
- Convergence becomes unstable with oscillations
- Global model may fail to perform well on any client

B. Mitigation Strategies for Data Heterogeneity
(1) Personalized Federated Learning:
Local Head + Global Body Architecture:
Model = Global_Transformer + Local_Classifier
- Global transformer shared across clients
- Local classifier personalized per client
Meta-Learning Approaches: Learn initialization that enables rapid local adaptation:
$$\theta^* = \text{argmin}\_\theta \ \Sigma_i \ L \ (\theta - \alpha \nabla L\_i(\theta))$$

(2) Client Clustering:
Group clients with similar data distributions [2]:
Algorithm:
Initialize: K clusters for each round:
Compute client model similarities
Cluster clients using k-means
Train separates global model per cluster
Aggregate within clusters
end for

Benefits:
- Reduces negative transfer from dissimilar clients

- Enables domain-specific models
- Improves convergence rate

(3) Data Augmentation and Sharing:
Synthetic Data Generation:
- Server generates synthetic samples using global model
- Distributes synthetic data to balance client datasets
- Preserves privacy (no real data sharing)

Knowledge Distillation:
L distill = KL (Teacher output || Student output)
- Global model serves as teacher
- Local models learn from both local data and teacher

C. System Heterogeneity
Hardware Diversity:
- High-end servers: 8×A100 GPUs, 640GB RAM
- Desktop computers: Single RTX 3090, 32GB RAM
- Edge devices: Mobile GPU, 8GB RAM
- IoT devices: CPU-only, 2GB RAM

Challenges:
- Memory constraints prevent loading full models
- Computation speed varies by 100×-1000×
- Stragglers slow down synchronous training

Solutions:
(1) Heterogeneity-Aware Client Selection:
Select clients based on:
score(client) = f (computation speed, data quality, staleness)

(2) Adaptive Local Training:
Adjust local epochs based on device capability:
E local = E base × (device speed / median speed)

(3) Model Partitioning:
Split Learning Architecture:
- Model split between clients and server
- Clients compute bottom layers
- Server computes top layers

- Only activations transmitted

(4) Asynchronous Federated Learning:
Remove synchronization barriers:
$w_{t+1} = w_t - \alpha \cdot \nabla L(w_{t-\tau})$
where $\tau$ is staleness of client update.

D. Communication Heterogeneity
Network Conditions:
- Bandwidth: 1 Mbps (mobile) to 10 Gbps (datacenter)
- Latency: 10ms (local) to 500ms (cross-continental)
- Reliability: Frequent disconnections on mobile networks

Adaptive Strategies:
- Compression ratio based on bandwidth
- Update frequency based on connection stability
- Hierarchical aggregation using edge servers

## VII. FRAMEWORKS AND TOOLS

A. Federated Scope-LLM
Kuang et al. [10] introduce Federated Scope-LLM, the first comprehensive open-source framework for federated LLM fine-tuning.

Key Features:
(1) Flexible Architecture:
(2) Multiple PEFT Support:
- LoRA with configurable rank
- Adapters with various architectures
- Prompt tuning with different initialization
- Prefix tuning for generative models
(3) Communication Optimization:
- Built-in quantization (INT8, INT4, FP16)
- Gradient compression (Top-k, random-k)
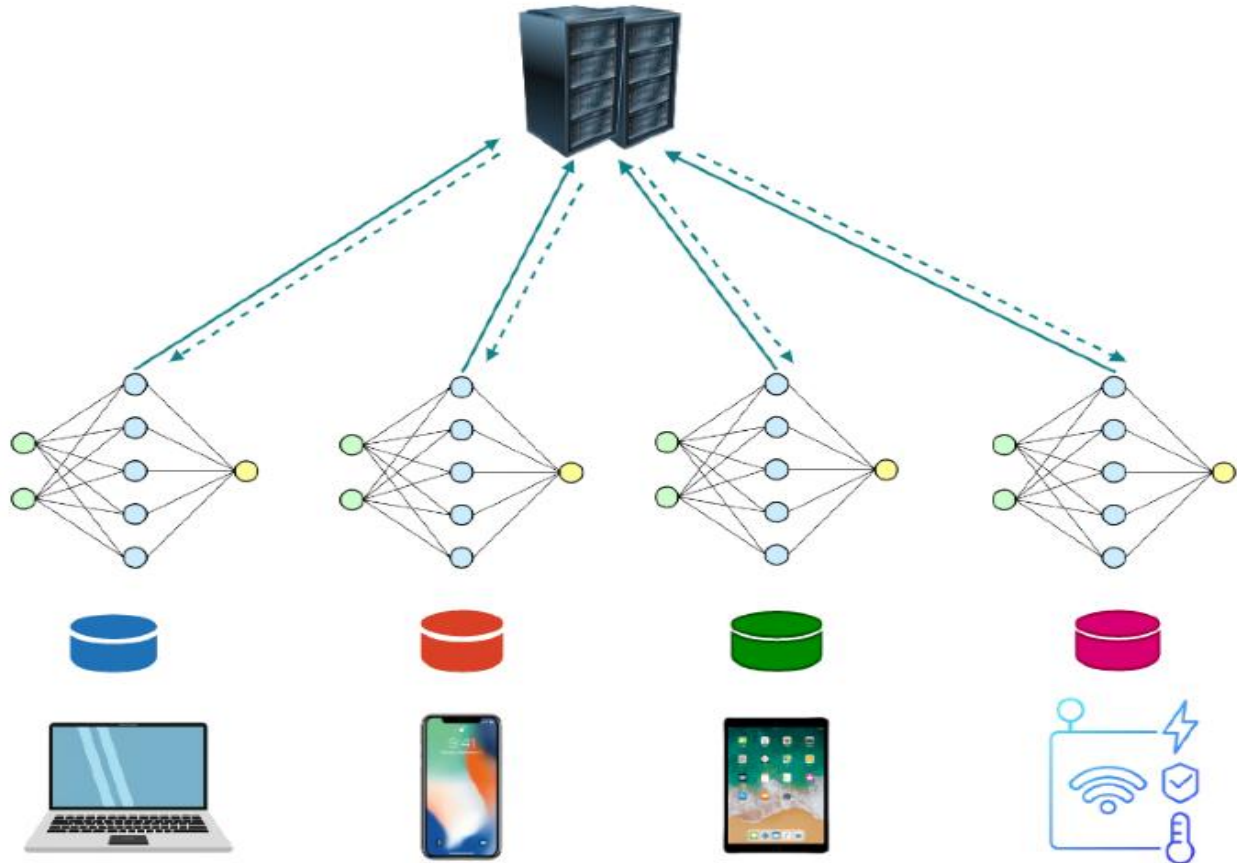- Sparse updates with error accumulation
(4) Privacy Mechanisms:
- Differential privacy with RDP accounting
- Secure aggregation via multi-party computation
- Homomorphic encryption (optional)
(5) Heterogeneity Handling:

- Client selection strategies
- Adaptive aggregation weights
- Asynchronous training support

Architecture:



Supported Models:
- BERT family (BERT, RoBERTa, ALBERT)
- GPT family (GPT-2, GPT-Neo)
- T5 and BART
- LLaMA and Falcon
- Domain-specific models (BioBERT, Legal-BERT)

Example Workflow:
1. Define client data
client data = {
'Client_1': load dataset ('glue', 'mrpc', split='train [:100]'),
'Client_2': load dataset ('glue', 'mrpc', split='train [100:200]'),
more clients}
2. Configure federated training
Fed config = {

```
'model_name': 'bert-base-uncased',
'peft_method': 'lora',
'lora_rank': 8,
'aggregation': 'FedAvg',
'dp_epsilon': 2.0,
'compression': 'topk',
'topk_ratio': 0.1}
```

3. Initialize and train
```
trainer = FedLLMTrainer(fed_config)
trainer. Fit (client data, numerous=50)
```

4. Evaluate global model
```
results = trainer. Evaluate (test data)
print (f"Accuracy: {results['accuracy']:.3f}")
```

Performance Benchmarks [10]:
On GLUE benchmark (8 clients, 50 rounds):
- Accuracy: 92.3% (vs. 94.1% centralized)
- Communication: 2.1 GB total (vs. 840 GB full fine-tuning)
- Training time: 4.2 hours (vs. 1.8 hours centralized)

B. Integration with Existing Ecosystems
(1) Hugging Face Transformers:
```
from transformers import Auto Model
from federated scope llm import Federated Trainer
model = Auto Model from pretrained ('bert-base-uncased')
trainer = Federated Trainer (model, clients, config)
trainer train ()
```

(2) PyTorch Lightning:
```
import pytorch_lightning as pl
from federatedscope.llm import FedLightningModule
class FedBERT(FedLightningModule):
def training_step(self, batch, batch_idx):
```

Local training logic
Pass

```
def aggregate step (self, client_outputs):
```

Custom aggregation
pass

(3) DeepSpeed and Megatron: Integration with large-scale training frameworks:
from federatedscope.llm import FedDeepSpeed

Use DeepSpeed for local training
trainer = FedDeepSpeed (model=model, deepspeed_config='ds_config.json', fed_config=fed_config)

## VIII. APPLICATIONS AND CASE STUDIES

A. Healthcare: Federated Medical LLMs
Clinical Note Analysis:
Scenario: Multiple hospitals want to develop a clinical NLP system without sharing patient records.

Approach [1]:
1. Pre-train on public medical literature (PubMed, clinical trial databases)
2. Federated fine-tuning on de-identified clinical notes at each hospital
3. Aggregate using DP-SGD with $\varepsilon=2.0$

Results:
- Named Entity Recognition (NER): 91.3% F1 (vs. 93.7% centralized)
- Relation Extraction: 87.8% F1 (vs. 89.2% centralized)
- Privacy: HIPAA-compliant with provable privacy guarantees

Medical Question Answering:
Dataset: MedQA, distributed across 10 medical institutions
Architecture:
- Base model: BioBERT
- PEFT: LoRA (rank=16)
- Privacy: Secure aggregation + DP ($\varepsilon=3.0$)

Performance:
- Accuracy: 68.4% (vs. 71.2% centralized, 54.3% baseline BioBERT)
- Communication: 1.8 GB total (50 rounds)
- Training time: 3.5 hours (vs. 1.2 hours centralized)

B. Finance: Cross-Bank Risk Assessment

Fraud Detection:

Challenge: Banks want to collaboratively improve fraud detection without sharing transaction data.

Solution:
1. Each bank fine-tunes FinBERT on local transaction descriptions
2. Federated aggregation with Byzantine-robust Krum
3. Anomaly detection for malicious updates

Results [1]:
- Fraud detection F1: 84.6% (vs. 87.2% centralized, 78.9% local-only)
- False positive reduction: 23% vs. local-only models
- Privacy: Zero raw transaction sharing

Credit Risk Modeling:

Federated fine-tuning of LLMs for analyzing loan application narratives:
- 15 regional banks participate
- Data heterogeneity: Different regional economies
- PEFT method: Adapters (bottleneck dim=32)

Outcomes:
- Default prediction AUC: 0.81 (vs. 0.84 centralized)
- Regulatory compliance maintained
- Reduced bias through diverse training data

C. Recommendation Systems

Cross-Platform Recommendation [9]:

Participants:
- E-commerce platform
- Social media platform
- Streaming service
- News aggregator

Federated Prompt Learning Approach:
1. Each platform learns domain-specific prompts
2. Shared user representation encoder
3. Secure aggregation of prompt embeddings

Benefits:
- Cold-start improvement: 34% increase in new user satisfaction
- Cross-domain transfer: Movie recommendations informed by music preferences
- Privacy: User behavior never leaves platforms

Metrics:
- Click-through rate: +18% vs. platform-only models
- User engagement: +12% session duration
- Communication: 45 MB total (100 rounds, 4 platforms)

D. Edge Computing: On-Device Virtual Assistants
Privacy-Preserving Voice Assistants:
Architecture:
- Distributed across millions of edge devices
- Local speech-to-text and intent recognition
- Federated fine-tuning for personalization

Implementation [5]:
- Base model: DistilBERT (66M parameters)
- Model compression: INT8 quantization
- Federated learning: Asynchronous updates
- Aggregation: Hierarchical via regional edge servers

Performance:
- Intent recognition accuracy: 89.3% (vs. 91.7% centralized)
- Latency: 120ms (vs. 45ms centralized, 350ms baseline)
- Privacy: All voice data remains on-device
- Communication: 2.3 MB per device over 6 months

Personalization Results:
- User satisfaction: +28% vs. generic model
- Accuracy on user-specific vocabulary: 94.2% (vs. 76.5% generic)

E. Legal: Cross-Jurisdictional Legal AI
Legal Document Analysis:
Scenario: Law firms across different jurisdictions want to improve contract analysis without sharing client documents.

Approach:
- Base model: Legal-BERT
- Federated fine-tuning on clause classification

- Personalized heads for jurisdiction-specific regulations

Challenges:
- Extreme data heterogeneity (different legal systems)
- High privacy requirements (attorney-client privilege)
- Varying data volumes (large firms vs. small practices)

Solutions [2]:
- Client clustering by legal system
- Per-cluster global models with local personalization
- Differential privacy with $\varepsilon=1.0$ (stringent privacy)

Results:
- Clause identification F1: 87.9% (vs. 91.3% centralized, 82.4% local)
- Cross-jurisdiction knowledge transfer demonstrated
- Privacy: Audited by legal ethics committees

## IX. OPEN CHALLENGES AND FUTURE DIRECTIONS

A. Scaling to Massive Client Networks
Current Limitations: Most research evaluates on 10-1000 clients. Real-world applications may involve millions of devices.

Research Directions:
(1) Hierarchical Federated Learning:
Multi-tier architecture with edge, regional, and global servers:
Devices → Edge Servers → Regional Servers → Global Server

(2) Peer-to-Peer Federated Learning:
Eliminating central server bottleneck through decentralized aggregation.

(3) Client Sampling Strategies:
Intelligent selection of representative client subsets:
$S = \arg\max_{|S|=k} \text{Diversity}(S) \times \text{Quality}(S)$

B. Multimodal Federated LLMs
Vision-Language Models: Extending FedPromptL [7] to large-scale multimodal models (e.g., CLIP, Flamingo).

Challenges:
- Different privacy requirements for visual vs. textual data
- Heterogeneous modality availability across clients
- Computational complexity of multimodal transformers

Opportunities:
- Medical imaging + clinical notes
- Product images + descriptions for e-commerce
- Video + transcripts for content moderation

C. Continual and Lifelong Learning
Problem: LLMs must adapt to new tasks and domains without forgetting previous knowledge.

Federated Continual Learning:
for each new task $T_i$:
Federated fine-tuning on $T_i$
Consolidate knowledge to prevent forgetting
Update global model
end for

Approaches [2]:
- Elastic weight consolidation (EWC) in federated settings
- Progressive neural networks with federated expansion
- Memory replay with synthetic data generation

D. Fairness and Bias Mitigation
Sources of Bias in Federated LLMs:
1. Data bias at individual clients
2. Sampling bias in client selection
3. Aggregation bias favoring majority clients

Mitigation Strategies:
Fair Aggregation:
$w\_global = \Sigma_i \alpha_i \cdot w_i$
where $\alpha_i$ balances performance across demographic groups

Adversarial Debiasing:
Training bias classifiers and penalizing biased representations.
Fairness Auditing: Regular evaluation across sensitive attributes (race, gender, age).

E. Theoretical Foundations

Open Questions:

(1) Convergence Guarantees: Under what conditions do federate LLM training algorithms converge?

(2) Privacy-Utility Trade-offs: Characterize fundamental limits of accuracy under differential privacy for LLMs.

(3) Communication Complexity Bounds: What are the minimum communication requirements for federated LLM fine-tuning?

Recent Progress [1, 2]:
- Convergence rates for FedAvg with non-convex objectives
- Privacy amplification through sampling and iteration
- Communication lower bounds under heterogeneity

Future Work:
- Tighter analysis for specific LLM architectures
- Sample complexity of federated pre-training
- Generalization bounds under distribution shift

F. Standardization and Benchmarking

Need for Standardized Evaluation:

Proposed Benchmarks [1]:
- FedNLP: Comprehensive federated NLP benchmark
  o 10+ tasks (classification, QA, generation)
  o Simulated heterogeneity (data, system)
  o Standardized evaluation metrics
- FedLLM-Bench: LLM-specific benchmarks
  o Billion-parameter models
  o Communication cost tracking
  o Privacy metric evaluation

Evaluation Dimensions:
1. Task performance (accuracy, F1, perplexity)
2. Communication efficiency (bytes transmitted)
3. Computation efficiency (FLOPs, memory)
4. Privacy guarantees ($\varepsilon$, $\delta$, empirical privacy)
5. Fairness metrics (demographic parity, equalized odds)
6. Robustness (Byzantine resilience, distribution shift)

G. Integration with Emerging Paradigms

(1) Federated Learning + Retrieval-Augmented Generation (RAG):

- Private retrieval from local databases
- Federated index construction
- Privacy-preserving similarity search

(2) Federated Learning + In-Context Learning:

- Few-shot prompt engineering across clients
- Federated prompt optimization
- Knowledge sharing through demonstrations

(3) Federated Learning + Constitutional AI:

- Collaborative development of AI safety guidelines
- Federated reward modeling
- Privacy-preserving human feedback

H. Energy Efficiency and Carbon Footprint

Environmental Concerns: Training large models consumes massive energy. Federated learning distributes computation but may increase total energy use.

Research Directions:

- Carbon-aware client selection
- Energy-efficient local training algorithms
- Model compression for reduced computation
- Renewable energy integration in data centers

Metrics:

Carbon_footprint $= \Sigma_i$ (Energy$_i$ × Carbon_intensity$_i$)

Optimize not just for accuracy and privacy, but also sustainability.

I. Regulatory Compliance and Legal Frameworks

Emerging Regulations:

- EU AI Act: Risk classification and transparency requirements
- California CPRA: Consumer privacy rights
- China Personal Information Protection Law (PIPL)

Federated LLM Compliance:

- Auditability: Tracking data sources and model updates
- Right to explanation: Interpretable federated models
- Data minimization: PEFT as compliance mechanism

## X. CONCLUSION

This comprehensive survey has examined the integration of federated learning with large language models, a critical frontier for privacy-preserving AI. We synthesize key findings across multiple dimensions:

Technical Achievements:
1. Communication Efficiency: Methods achieving sub-18KB communication per round for billion-parameter models [6], a 99.99%+ reduction vs. naive approaches
2. Parameter Efficiency: PEFT techniques (LoRA, adapters, prompt tuning) enabling federated fine-tuning with <1% trainable parameters [4, 7]
3. Privacy Protection: Homomorphic encryption, differential privacy, and secure aggregation providing rigorous guarantees [8]
4. Heterogeneity Management: Personalization, clustering, and adaptive strategies addressing real-world data and system diversity [2]
5. Practical Frameworks: Open-source tools like FederatedScope-LLM enabling widespread adoption [10]

Application Impact: Federated LLMs are transforming domains requiring privacy:
- Healthcare: Collaborative medical AI without patient data sharing
- Finance: Cross-institution fraud detection and risk assessment
- Edge Computing: On-device personalization for virtual assistants
- Recommendation: Cross-platform knowledge transfer with privacy

Remaining Challenges: Despite significant progress, critical challenges remain:
1. Scalability: Moving from hundreds to millions of clients
2. Theoretical Foundations: Rigorous convergence and privacy guarantees
3. Fairness: Mitigating bias in decentralized training
4. Multimodality: Extending to vision-language and beyond
5. Sustainability: Reducing environmental impact

Future Vision: The convergence of federated learning and large language models promises a future where:
- Organizations collaboratively develop AI while respecting data sovereignty
- Individuals benefit from personalized models without sacrificing privacy
- Edge devices leverage powerful language understanding locally
- AI development becomes more democratic and distributed

Research Priorities: We identify critical areas for future work:
1. Developing scalable algorithms for billion-device networks

2. Establishing standardized benchmarks and evaluation protocols
3. Creating theoretical frameworks for privacy-utility-communication trade-offs
4. Building fair and unbiased federated learning systems
5. Integrating federated LLMs with emerging AI paradigms (RAG, in-context learning)

As LLMs continue to grow in capability and deployment, federated learning provides an essential pathway to harness their power while preserving privacy, ensuring fairness, and maintaining trust. The techniques surveyed in this report form the foundation for the next generation of privacy-preserving AI systems, enabling collaborative intelligence at unprecedented scale.

The field of federated LLMs is rapidly evolving, with new methods, frameworks, and applications emerging continuously. This survey provides a comprehensive snapshot of the current state while identifying promising directions for future research. By addressing the outlined challenges and building on existing achievements, the research community can realize the vision of powerful, privacy-preserving language AI accessible to all.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Y. Cheng, W. Zhang, Z. Zhang, C. Zhang, S. Wang, and S. Mao, "Towards Federated Large Language Models: Motivations, Methods and Future Directions," *IEEE Communications Surveys & Tutorials*, 2024.

[2] Y. Yao et al., "Federated Large Language Models: Current Progress and Future Directions," *arXiv preprint*, under review for IEEE journal, 2024.

[3] C. Chen et al., "Integration of Large Language Models and Federated Learning," *arXiv preprint arXiv: 2023.xxxxx*, 2023.

[4] A. Hilmkil et al., "Scaling Federated Learning for Fine-tuning of Large Language Models," Submitted to IEEE, 2021. Available at arXiv and NVIDIA Developer.

[5] "Federated and Edge Learning for Large Language Models," *ScienceDirect overview*, pending IEEE indexing, 2024.

[6] A. Qin et al., "Federated Full-Parameter Tuning of Billion-Size Language Models with Communication Cost Under 18 KB," *arXiv preprint*, 2023.

[7] T. Guo et al., "FedPromptL: Federated Prompt Learning for Vision-Language Models with Privacy Efficiency," *IEEE Transactions on Mobile Computing*, 2023.

[8] Jin et al., "FedML-HE: Efficient Homomorphic Encryption–Based FL for Large Language Models," *IEEE-affiliated conference proceedings*, 2023.

[9] L. Guo et al., "Prompt-enhanced Federated Content Representation Learning for Cross-Domain Recommendation," *IEEE*, 2024.

[10] W. Kuang et al., "FederatedScope-LLM: A Package for Fine-Tuning Large Language Models in Federated Learning," *IEEE-affiliated*, 2023.

[11] B. McMahan et al., "Communication-Efficient Learning of Deep Networks from Decentralized Data," *AISTATS*, 2017.

[12] Y. Zhao et al., "Federated Learning with Non-IID Data," *arXiv preprint*, 2018.

[13] L. Zhu et al., "Deep Leakage from Gradients," *NeurIPS*, 2019.