

# H<sup>2</sup>Tune: Federated Foundation Model Fine-Tuning with Hybrid Heterogeneity

Wei Guo<sup>†a</sup>, Siyuan Lu<sup>†b</sup>, Yiqi Tong<sup>c</sup>, Zhaojun Hu<sup>d</sup>, Fuzhen Zhuang<sup>\*a,e</sup>, Xiao Zhang<sup>f,\*</sup>, Tao Fan<sup>g</sup> and Jin Dong<sup>h</sup>

<sup>a</sup>School of Artificial Intelligence, Beihang University

<sup>b</sup>School of Computer Science and Technology, Heilongjiang University

<sup>c</sup>School of Computer Science and Engineering, Beihang University

<sup>d</sup>School of Statistics, Renmin University of China

<sup>e</sup>Zhongguancun Laboratory, China

<sup>f</sup>School of Computer Science and Technology, Shandong University

<sup>g</sup>WeBank Co., Ltd, Shenzhen, China

<sup>h</sup>Beijing Academy of Blockchain and Edge Computing

**Abstract.** Different from existing federated fine-tuning (FFT) methods for foundation models, hybrid heterogeneous federated fine-tuning (HHFFT) is an under-explored scenario where clients exhibit double heterogeneity in model architectures and downstream tasks. This hybrid heterogeneity introduces two significant challenges: *1) heterogeneous matrix aggregation*, where clients adopt different large-scale foundation models based on their task requirements and resource limitations, leading to dimensional mismatches during LoRA parameter aggregation; and *2) multi-task knowledge interference*, where local shared parameters, trained with both task-shared and task-specific knowledge, cannot ensure only task-shared knowledge is transferred between clients. To address these challenges, we propose *H<sup>2</sup>Tune*, a federated foundation model fine-tuning with hybrid heterogeneity. Our framework *H<sup>2</sup>Tune* consists of three key components: (i) *sparsified triple matrix decomposition* to align hidden dimensions across clients through constructing rank-consistent middle matrices, with adaptive sparsification based on client resources; (ii) *relation-guided matrix layer alignment* to handle heterogeneous layer structures and representation capabilities; and (iii) *alternating task-knowledge disentanglement* mechanism to decouple shared and specific knowledge of local model parameters through alternating optimization. Theoretical analysis proves a convergence rate of  $O(1/\sqrt{T})$ . Extensive experiments show our method achieves up to 15.4% accuracy improvement compared to state-of-the-art baselines. Our code is available at <https://anonymous.4open.science/r/H2Tune-1407>.

## 1 Introduction

Foundation models (FMs) integrated with federated learning (FL) [37] have emerged as a promising paradigm for optimizing client-specific applications [27, 28]. This integration enables the customization of FMs to address downstream tasks through federated fine-

**Table 1:** FFT framework comparisons with our proposed H<sup>2</sup>Tune.

Frameworks	Heterogeneity					
	Model			Task	Resource	
	Layer	Dimension	Architecture			
FFTs	×	×	×	×	×	×
HetLoRA [11]	×	×	×	✓	×	✓
FlexLoRA [5]	×	×	×	✓	✓	×
HeteroTune [24]	×	✓	×	×	×	×
pFedLoRA [39]	×	✓	×	×	×	×
<b>H<sup>2</sup>Tune</b>	✓	✓	✓	✓	✓	✓

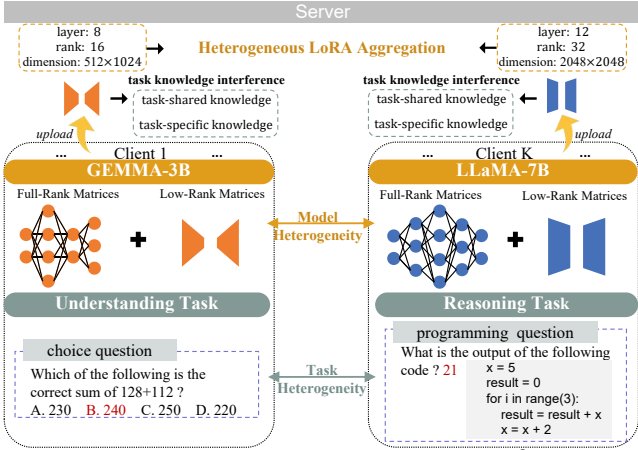
tuning (FFT). Current research presents two primary methodologies for FFT: full-parameter fine-tuning (FPFT) [17] and parameter-efficient fine-tuning (PEFT) [13]. The conventional FPFT approach updates all model parameters, which incurs prohibitive costs for FMs with billions of parameters on clients. In contrast, PEFT updates only a small subset of parameters while keeping others frozen, significantly reducing costs while maintaining comparable performance.

Despite these advances, effectively handling task heterogeneity across clients remains a core challenge in FFT. Previous works have explored task heterogeneity through diverse approaches. FedBone [9] addresses this challenge by separating general and task-specific models via server-client split learning, while FedMSplit [7] leverages dynamic graphs to capture inter-task correlations. Other approaches like FedLPS [25] addresses task heterogeneity by sharing a general encoder across tasks while using adaptive pruning and heterogeneous aggregation for different client resources. These approaches work effectively for small models by identifying and sharing common layers while allowing personalized remaining layers. However, these approaches aren't directly applicable to large-scale FMs where clients often aggregate PEFT parameters. This creates two problems: clients with different tasks may fine-tune different model components, making parameter alignment difficult; and even when the same components are tuned, the parameters contain entangled task-shared and task-private knowledge. Current federated FM frameworks cannot effectively align diverse fine-tuning structures or decouple this knowledge, limiting effective transfer across heterogeneous tasks.

Additionally, in practical FFT scenarios, clients may not only se-

\* Fuzhen Zhuang and Xiao Zhang are corresponding authors. Email: zhuang-fuzhen@buaa.edu.cn, xiaozhang@sdu.edu.cn.

† Equal contribution.



**Figure 1:** Challenges encountered in our proposed  $H^2Tune$  framework.  $H^2Tune$  focuses on hybrid heterogeneous FFT scenarios with large-scale foundation model heterogeneity and task heterogeneity. We address model heterogeneity across multiple levels, including layers, hidden dimensions, and model architectures, while solving the problem of inseparable task-shared and task-private knowledge in current federated fine-tuning frameworks.

lect varying model components for PEFT. Still, they might also select different FM scales or model families with different architectures based on their task or resource constraints. We term this scenario model heterogeneous federated fine-tuning. Traditional model heterogeneous approaches in FL, like knowledge distillation-based [43, 6, 10, 44, 1, 2] or mutual learning-based methods [32, 36], become infeasible due to the enormous computational and communication costs associated with FMs’ parameter counts. Although model split-based methods [21, 8, 3, 12, 18] achieve partial parameter sharing by splitting client models into shared and private parts, however, it strictly assume shared homogeneous models. Recent works in federated foundation models have begun addressing this heterogeneity. As shown in Table 1, HetLoRA [11] first attempted to address the challenge of heterogeneous LoRA aggregation with different ranks in FFT due to resource constraints through simple padding. FlexLoRA [5] advances this approach by distributing global LoRA parameters to client-specific ranks via SVD decomposition based on local requirements. These methods still require clients to apply LoRA to models from the same family with identical layer structures and dimensions. Although current works, HeteroTune [24] and pFedLoRA [39], try to support heterogeneous hidden state dimensions, they also require uniform layer counts and model architecture.

In this work, we identify an under-explored real-world scenario: hybrid heterogeneous federated fine-tuning (HHFFT), where clients simultaneously face task heterogeneity, resource heterogeneity, and model heterogeneity. For instance, in the healthcare field, large hospitals deploy LLaMA-7B [34] for radiology report generation and disease progression prediction, while smaller clinics use Qwen-1.8B [4] for medication recommendation and appointment scheduling. Similarly, in the financial sector, large investment banks leverage LLaMA-7B for trading strategy generation, while local credit unions adopt LLaMA-3B for credit scoring. In summary, we identify the following key challenges of HHFFT, as shown in Figure 1:

- **Heterogeneous matrix aggregation.** It refers that clients may fine-tune different model components or use different model scales or architectures based on their needs, creating fine-tuned weight matrix alignment challenges during aggregation due to mismatched layers, incompatible dimensions, and varied architectures.

- **Multi-task knowledge interference.** It refers to existing FFT methods that combine both shared and task-specific knowledge fail to prevent task-specific information from leaking into shared parameters. Given the task differences across clients, any knowledge specific to one task may interfere with other task performances.

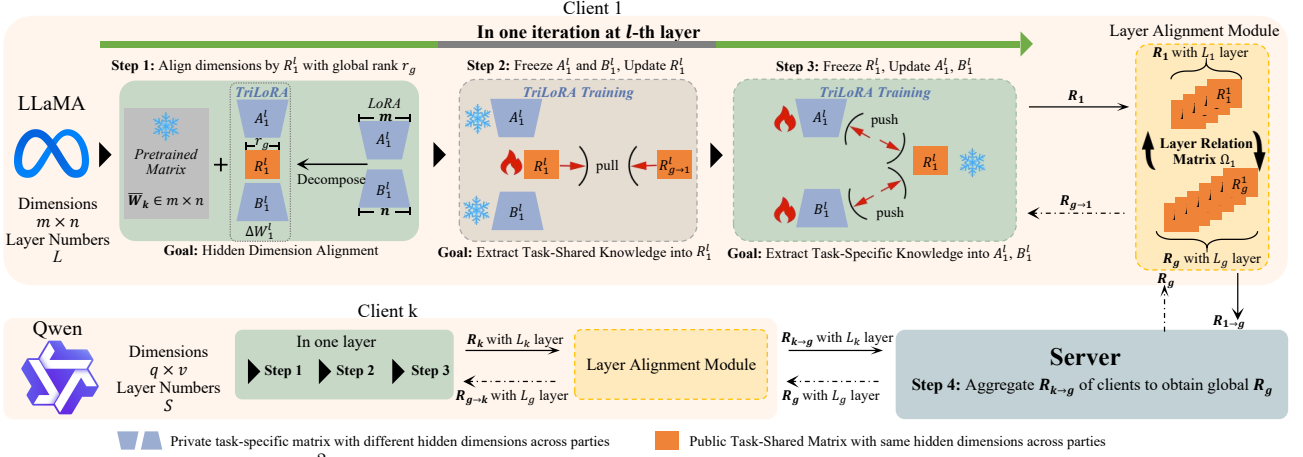
To address these challenges, we propose  $H^2Tune$ , a federated foundation model fine-tuning framework with hybrid heterogeneity.  $H^2Tune$  contains three key components: 1) *sparsified triple matrix decomposition* that aligns hidden dimensions (width) across clients by constructing consistent shared middle matrices, with adaptive sparsification based on each client’s resources; 2) *relation-guided matrix layer alignment* that enables heterogeneous layer (depth) aggregation by learning cross-layer matrix relationships via relation networks; and 3) *alternating task-knowledge disentanglement* that leverage alternating optimization to decouple local models into shared matrices and private matrices, and achieve federated knowledge transfer by shared matrices with task-shared knowledge. Our main contributions are summarized as follows:

- We identify an under-explored scenario in FFT named hybrid heterogeneous federated fine-tuning, where task, model, and resource heterogeneity coexist. This hybrid heterogeneity introduces unique challenges beyond current FFT methods.  $H^2Tune$  is developed to handle heterogeneity in task objectives, model structures, and resource constraints within FFT.
- $H^2Tune$  is carefully designed with two core components: 1) a sparse triple matrix decomposition with relation-guided alignment for heterogeneous matrix alignment across clients, and 2) an alternating disentanglement optimization mechanism that separates general and task-specific knowledge.
- Our approach accommodates both homogeneous and heterogeneous model settings and supports single or multi-task scenarios. And we theoretically establish the convergence guarantee of  $H^2Tune$ , proving that it achieves an  $O(1/\sqrt{T})$  convergence rate.
- Extensive experiments on two well-known benchmarks validate  $H^2Tune$  outperforms state-of-the-art baselines by up to 15.4% accuracy in both homogeneous and heterogeneous scenarios.

## 2 Related Work

Foundation models (FMs) have been widely adopted in federated learning, with fine-tuning approaches primarily divided into full-parameter fine-tuning (FPFT) [17] and parameter-efficient fine-tuning (PEFT) [13]. While FPFT offers superior performance, its computational and communication costs are substantial [29]. In contrast, PEFT methods like LoRA [19] provide efficient alternatives by updating only partial parameters, showing promising results in federated scenarios [26, 42]. However, multi-task federated fine-tuning for FMs remains challenging. Recent work [5, 7, 9] have made initial progress, but limitations persist. They either lack PEFT support for computational efficiency, or fail to effectively distill task-shared knowledge for cross-client transfer. More critically, they require identical foundation models across clients.

Model heterogeneity in FL has attracted substantial research attention. These works can be categorized into two main branches. Specifically, partially model heterogeneous methods allow clients to maintain different subnets of a global architecture that can be aggregated at the server, as demonstrated in FedRolex [3], HeteroFL [12], FjORD [18], and HFL [30]. On the other hand, the completely model-heterogeneous branch involves clients with entirely different



**Figure 2:** Overview of one client in  $H^2Tune$  in one communication round. Blue regions indicate task-specific knowledge, orange for task-shared knowledge, and white for sparsified parameters. Clients align matrices through triple factorization and relation matrices for global intermediate matrices, with sparsity adaptation to local resources. Alternating optimization decouples shared and specific knowledge, enabling cross-client fine-tuning through exchanging intermediate matrices containing task-shared knowledge.

model structures that cannot be directly aggregated, and is further divided into the following three categories: knowledge distillation-based methods that either use public datasets [6, 10] or employ alternative approaches like FedZKT [44] which generates synthetic datasets through generator training, and HFD [1, 2] which requires uploading local logits for server aggregation, all of which exacerbate substantial computational overhead for large-scale FMs in FL; model split-based methods that share either feature extractors [8] or classifiers [21] while experiencing performance bottlenecks from partial parameter sharing; and mutual learning-based methods that assign small homogeneous and large heterogeneous models to each client [32, 36], yet fail to optimize the relationship between model structure and parameter capacity.

However, in federated FM fine-tuning scenarios, the massive parameter space of FMs further amplifies the limitations of these existing approaches. Although PEFT techniques reduce update parameters for large FMs, they introduce a critical new challenge: clients may deploy PEFT methods like LoRA or Adaptor across different architectural layers or foundation model families, creating a novel form of model heterogeneity: model heterogeneous federated fine-tuning. Recent research has addressed LoRA heterogeneity within identical model frameworks: HetLoRA [11] mitigates rank heterogeneity through parameter alignment via simple padding, while FlexLoRA [5] employs SVD decomposition to distribute global parameters to client-specific ranks for personalized adaptation. Despite HeteroTune [24] and pFedLoRA [39] attempting to resolve hidden dimension heterogeneity, these approaches still remain constrained to uniform model families with same architectural configurations. Moreover, considering that clients in real-world scenarios often have diverse tasks and resource constraints, these approaches similarly overlook the challenges posed by such practical heterogeneity.

### 3 Preliminary

**Federated fine-tuning with hybrid heterogeneity.** A HHFFT framework coordinates  $K$  distributed clients with a central server, where each client  $k$  maintains a private dataset  $\mathcal{D}_k = \{(x_i, y_i), i = 1, \dots, N_k\}$  and a foundation model  $\mathbf{W}_k \subseteq \mathbb{R}^d$ . By sharing parameters  $\Delta\mathbf{W}_k \subseteq \mathbb{R}^q (q \ll d)$ , HHFFT enables collaborative learning across heterogeneous environments, allowing clients with different

foundation models, resource constraints, and tasks to enhance their individual model performance by  $\mathbf{W}_k = \bar{\mathbf{W}}_k + \Delta\mathbf{W}_k$ , where  $\bar{\mathbf{W}}_k$  denotes the frozen parameters of client  $k$ . The fine-tuning parameters  $\Delta\mathbf{W}_k$  can be represented as  $\Delta\mathbf{W}_k = \{W_k^l\}_{l=1}^{L_k}$ , where  $L_k$  denotes the number of layers and  $W_k^l \subseteq \mathbb{R}^{d_{in}^k \times d_{out}^k}$  represents the  $l$ -th layer parameters with input and output dimensions  $d_{in}^k$  and  $d_{out}^k$ . Each client  $k$  optimizes its task-specific function  $\mathcal{L}_k$  over local dataset  $\mathcal{D}_k$ .

Formally, we define  $\mathcal{M}$  as the foundation model,  $\mathcal{T}$  as the downstream task, and  $\mathcal{R}$  as computational resources. For any two clients  $k, j \in \{1, \dots, K\} (k \neq j)$ , the heterogeneous settings in HHFFT manifest: model heterogeneity ( $\mathcal{M}_k \neq \mathcal{M}_j$ ), task heterogeneity ( $\mathcal{T}_k \neq \mathcal{T}_j$ ), and resource heterogeneity ( $\mathcal{R}_k \neq \mathcal{R}_j$ ). Notably, model heterogeneity in FFT primarily manifests as heterogeneity in the  $\{\Delta\mathbf{W}_k\}_{k=1}^K$  being aggregated, which encompasses four key aspects:

- *Rank heterogeneity:*  $r_k \neq r_j$ ,
- *Dimensional heterogeneity:*  $d_{in}^k \neq d_{in}^j$  or  $d_{out}^k \neq d_{out}^j$ ,
- *Layer heterogeneity:*  $L_k \neq L_j$ ,
- *Architecture heterogeneity:*  $F_k \neq F_j$ ,

where  $r_k$  refers to the intermediate dimensions of low-rank fine-tuning matrix, with each layer-wise update  $\Delta W_k^l = A_k^l \times B_k^l$ , where  $A_k^l \in \mathbb{R}^{d_{in} \times r_k}$  and  $B_k^l \in \mathbb{R}^{r_k \times d_{out}}$  and ( $r_k \ll \min(d_{in}, d_{out})$ ).

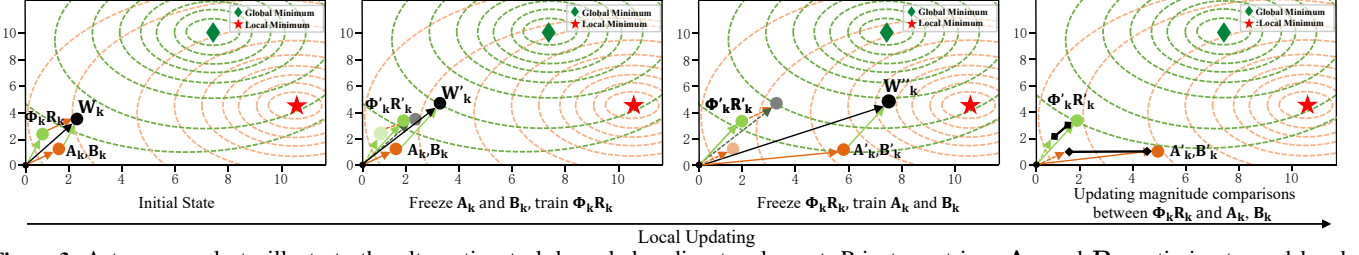
**Problem Formulation.** The optimization objective of  $H^2Tune$  is:

$$\begin{aligned} \min_{\Delta\mathbf{W}_1, \dots, \Delta\mathbf{W}_K} \sum_{k=1}^K \mathcal{L}_k(\Delta\mathbf{W}_k; \mathcal{M}_k, \mathcal{T}_k, \mathcal{R}_k), \\ \text{s.t. } \mathcal{M}_k \neq \mathcal{M}_j, \mathcal{T}_k \neq \mathcal{T}_j, \mathcal{R}_k \neq \mathcal{R}_j, \forall k \neq j. \end{aligned} \quad (1)$$

## 4 Methodology

### 4.1 Overview

We first provide an overview of  $H^2Tune$ . As illustrated in Figure 2, each layer parameter matrix  $W_k^l$  of client  $k$  is decomposed into a frozen pre-trained matrix  $\bar{W}_k^l$  plus a low-rank fine-tuning matrix  $\Delta W_k^l$ . Different from LoRA in PEFT which decomposes  $\Delta W_k^l$  into two components:  $A_k^l \subseteq \mathbb{R}^{a \times r_k}$  and  $B_k^l \subseteq \mathbb{R}^{r_k \times b}$ , we use three components:  $A_k^l \in \mathbb{R}^{a \times r_k}$ ,  $R_k^l \in \mathbb{R}^{r_k \times r_k}$ , and  $B_k^l \in \mathbb{R}^{r_k \times b}$ , where  $R_k^l$  is a dense real matrix rather than a diagonal intermediate matrix as in SVD. The training process in each communication



**Figure 3:** A toy example to illustrate the alternating task-knowledge disentanglement. Private matrices  $\mathbf{A}_k$  and  $\mathbf{B}_k$  optimize toward local optima with task-specific knowledge, while shared matrix  $\mathbf{R}_k$  optimizes toward the global optimum with shared knowledge. Due to their larger parameter number,  $\mathbf{A}_k$  and  $\mathbf{B}_k$  dominate the optimization direction toward local optima, while independent updates between private and shared matrices ensure optimization order does not affect the model optimization direction.

round involves: 1) each client  $k$  freezes its full-rank matrices  $\bar{\mathbf{W}}_k$  and task-shared matrix  $\mathbf{R}_k$ , while updating the task-specific matrices  $\mathbf{A}_k$  and  $\mathbf{B}_k$ . 2) each client  $k$  aligns the layers between received global task-shared matrices  $\mathbf{R}_g$  and local task-shared matrices using the layer relation matrix  $\Omega_k$ . Client then updates  $\mathbf{R}_k$  and  $\Omega_k$  while keeping task-specific matrices  $\mathbf{A}_k$  and  $\mathbf{B}_k$  frozen. After local updating, all clients only upload their task-shared matrices  $\mathbf{R}_k$  to server, maintaining privacy of their task-specific matrices. 3) server receives clients' task-shared matrices and aggregates them to generate global task-shared matrices  $\mathbf{R}_g$ . Once the aggregation is finished, server sends  $\mathbf{R}_g$  back to all clients. 4) Each client receives the global task-shared matrices and achieves local optimization. For notational simplicity, bold symbols denote collections across all  $L_k$  layers:  $\mathbf{W}_k = \{\mathbf{W}_k^l\}_{l=1}^{L_k}$ ,  $\bar{\mathbf{W}}_k = \{\bar{\mathbf{W}}_k^l\}_{l=1}^{L_k}$ ,  $\Delta\mathbf{W}_k = \{\Delta\mathbf{W}_k^l\}_{l=1}^{L_k}$ ,  $\mathbf{A}_k = \{\mathbf{A}_k^l\}_{l=1}^{L_k}$ ,  $\mathbf{B}_k = \{\mathbf{B}_k^l\}_{l=1}^{L_k}$ ,  $\mathbf{R}_k = \{\mathbf{R}_k^l\}_{l=1}^{L_k}$ ,  $\Phi_k = \{\Phi_k^l\}_{l=1}^{L_k}$  and  $\mathbf{R}_g = \{\mathbf{R}_g^l\}_{l=1}^{L_g}$ , where  $L_g$  represents the layer number of global task-shared matrix  $\mathbf{R}_g$ .

## 4.2 Sparsified Triple Matrix Decomposition

To address hidden dimension heterogeneity, we introduce TriLoRA, which decomposes each layer's parameters into private task-specific matrices  $\mathbf{A}_k^l$ ,  $\mathbf{B}_k^l$  and a public task-shared matrix  $\mathbf{R}_k^l$ . Though clients have varying computational resources that would normally lead to different local ranks  $r_k$ , we maintain a uniform global rank  $r_g$  across all task-shared matrices for effective knowledge sharing. We accommodate resource differences through client-specific sparsity rates  $\beta_k$  in sparse matrix  $\Phi_k^l$ , adjusting practical parameter update. Specifically, for each client  $k$ , we decompose local FM  $\mathbf{W}_k$  as:

$$\mathbf{W}_k = \bar{\mathbf{W}}_k + \Delta\mathbf{W}_k. \quad (2)$$

We decompose each layer's weight matrix  $\Delta\mathbf{W}_k^l$  of local model  $\Delta\mathbf{W}_k$  into two private task-specific matrices  $\mathbf{A}_k^l$ ,  $\mathbf{B}_k^l$  and one public task-shared matrix  $\mathbf{R}_k^l$ , combining them via residual connection:

$$\begin{aligned} \Delta\mathbf{W}_k^l &= (\mathbf{A}_k^l + \mathbf{A}_k^l \cdot (\Phi_k^l \cdot \mathbf{R}_k^l)) \mathbf{B}_k^l, \\ \text{s.t. } \mathbf{A}_k^l &\in \mathbb{R}^{a_k^l \times r_g}, \mathbf{R}_k^l, \Phi_k^l \in \mathbb{R}^{r_g \times r_g}, \mathbf{B}_k^l \in \mathbb{R}^{r_g \times b_k^l}, \end{aligned} \quad (3)$$

where  $a_k^l$  and  $b_k^l$  represent input and output dimensions at layer  $l$ . The global rank  $r_g$  satisfies  $0 < r \leq \min_{k \in [1, K]}(a_k^l, b_k^l)$ , ensuring compatibility across clients with heterogeneous model dimensions. The sparsification matrix  $\Phi_k^l$  has client-specific sparsity ratio  $\beta_k \in [0, 1]$ , controlling the proportion of non-zero elements and allowing clients to adjust updatable parameters based on available resources.

## 4.3 Relation-Guided Matrix Layer Alignment

Due to model scale or architectural heterogeneity across clients, clients have varying numbers of layers  $L_k$  in their shared matrices

$\mathbf{R}_k \in \mathbb{R}^{r_g \times r_g \times L_k}$ , preventing direct server aggregation. To solve this layer mismatch, we introduce a trainable layer relation alignment matrix  $\Omega_k \in \mathbb{R}^{L_k \times L_g}$  for each client  $k$ , where  $L_g = \max(L_k, k \in [1, K])$  represents the global layer count. The matrix  $\Omega_k$  transforms local task-shared matrices  $\mathbf{R}_k$  into uniform-sized representations for server aggregation:

$$\mathbf{R}_{k \rightarrow g} := \mathbf{R}_k \Omega_k, \Omega_k \in \mathbb{R}^{L_k \times L_g}. \quad (4)$$

The server averages these aligned matrices  $\{\mathbf{R}_{k \rightarrow g}\}_{k=1}^K$  to produce the global task-shared matrix  $\mathbf{R}_g \in \mathbb{R}^{r_g \times r_g \times L_g}$ . Upon receiving  $\mathbf{R}_g$ , each client use  $\Omega_k^\top$  to transform it into a locally compatible form:

$$\mathbf{R}_{g \rightarrow k} := \mathbf{R}_g \Omega_k^\top, \Omega_k^\top \in \mathbb{R}^{L_g \times L_k}. \quad (5)$$

These transformed matrices  $\mathbf{R}_{g \rightarrow k}$  and  $\Omega_k$  are further personalized to local distributions through loss  $\mathcal{L}_{\text{share}}^k$  in Equation 6.

## 4.4 Alternating Task-Knowledge Disentanglement

In existing task heterogeneous FFT methods [5, 7, 9], parameters are jointly trained with both task-specific and task-shared knowledge. This joint training causes the shared parameters to be contaminated by task-specific knowledge, hindering cross-client knowledge transfer. To address this interference, we propose an alternating optimization approach that separates task-shared knowledge into matrix  $\mathbf{R}_k$  and task-specific knowledge into matrices  $\mathbf{A}_k$  and  $\mathbf{B}_k$ . In each communication round  $t \in [1, T]$ , local iterations  $e \in [1, E]$  consist of two key optimization steps.

First, client  $k$  receives the global task-shared matrix  $\mathbf{R}_g$  and transforms it to local form  $\mathbf{R}_{g \rightarrow k}$  by  $\Omega_k^\top$ . Then, while freezing  $\mathbf{R}_{g \rightarrow k}$ ,  $\mathbf{A}_k$ , and  $\mathbf{B}_k$ , the client optimizes  $\mathbf{R}_k$ , sparse matrix  $\Phi_k$ , and layer relation matrix  $\Omega_k$  by minimizing:

$$\mathcal{L}_{\text{share}}^k = \text{CE}(y'_i, y_i) + \text{KL}(\mathbf{R}_k, \mathbf{R}_{g \rightarrow k}), \quad (6)$$

where  $y'_i$  is the predicted label and CE is the cross-entropy loss.

In the second step, to further disentangle shared and private parameters, each client freezes  $\mathbf{R}_k$ ,  $\Phi_k$ , and  $\Omega_k$ , and optimizes task-specific matrices  $\mathbf{A}_k$  and  $\mathbf{B}_k$  by:

$$\mathcal{L}_{\text{specific}}^k = \text{CE}(y''_i, y_i) - \text{KL}(y''_i, y'_i) + \frac{\mathcal{V}}{2}(\|\mathbf{A}_k\|_2 + \|\mathbf{B}_k\|_2), \quad (7)$$

where  $y''_i$  is the prediction based on optimized  $\mathbf{R}_k$ , and  $\mathcal{V}$  is the regularization coefficient. After local updating, each client uploads only  $\mathbf{R}_k$  to server while keeping  $\mathbf{A}_k$  and  $\mathbf{B}_k$  private.

**Toy example.** Figure 3 illustrates our local two-stage update process. The green  $\blacklozenge$  represents the optimal point for the task-shared model  $\mathbf{R}_k$  (trained on all clients' data), while the red  $\star$  represents



the optimal point for task-specific matrices  $\mathbf{A}_k$  and  $\mathbf{B}_k$  (trained on individual client data).  $\mathcal{L}_{\text{shared}}$  guides  $\mathbf{R}_k$  toward global knowledge, while  $\mathcal{L}_{\text{specific}}$  directs  $\mathbf{A}_k$  and  $\mathbf{B}_k$  toward client-specific knowledge. Since  $\mathbf{R}_k$  and  $\mathbf{A}_k, \mathbf{B}_k$  are updated independently, the optimization paths remain unbiased. However, as shown in the rightmost figure,  $\mathbf{R}_k$ 's lower dimensionality results in smaller update steps compared to  $\mathbf{A}_k$  and  $\mathbf{B}_k$ . This approach ensures that local parameters primarily serve the local task while still benefiting from global knowledge.

## 5 Convergence Analysis

In this section, we present theoretical guarantees for the convergence of our proposed  $H^2\text{Tune}$ . In contrast to traditional alternating optimization, our approach involves different optimization targets and models in its two-step client optimization, forming a bidirectional optimization process. Let  $\mathbf{R} := \{\mathbf{R}_1, \dots, \mathbf{R}_K\}$  and  $\mathbf{H} := \{(\mathbf{A}_1, \mathbf{B}_1), \dots, (\mathbf{A}_K, \mathbf{B}_K)\}$  represents task-shared parameters and task-specific parameters across all clients, respectively. We denote the complete set of parameters in the federated learning system as  $\Theta := \{\mathbf{R}, \mathbf{H}\}$ . Specifically, denote  $L(\mathbf{R}, \mathbf{H}) = \sum_{k=1}^K \mathcal{L}_{\text{share}}^k(\mathbf{R}_k, \mathbf{H}_k)$ ,  $h(\mathbf{R}) = \lambda \sum_{k=1}^K KL(\mathbf{R}_k, \mathbf{R}_g)$ ,  $G_0(\mathbf{R}, \mathbf{H}) = \sum_{k=1}^K \mathcal{L}_{\text{specific}}^k(\mathbf{R}_k, \mathbf{H}_k)$ , and  $G(\mathbf{R}, \mathbf{H}) = \sum_{k=1}^K \mathcal{L}_{\text{specific}}^k(\mathbf{R}_k, \mathbf{H}_k) + \frac{\mathcal{V}}{2} \sum_{k=1}^K \|\mathbf{H}_k\|^2$ . Then, our optimization problem could be written as follows:

$$\begin{aligned} \min_{\mathbf{R}} L(\mathbf{R}, \mathbf{H}^*(\mathbf{R})) + h(\mathbf{R}), \\ \text{s.t. } \mathbf{H}^*(\mathbf{R}) = \arg \min_{\mathbf{H}} G(\mathbf{R}, \mathbf{H}), \end{aligned} \quad (8)$$

where  $F(\mathbf{R}) := L(\mathbf{R}, \mathbf{H}^*(\mathbf{R}))$ . Let  $\Theta := \{\mathbf{R}, \mathbf{H}\}$  denote all parameters. Then, we make some standard assumptions ([15, 22, 31]) on  $L$  and  $G$  of our bi-level optimization problem.

**Assumption 1** (Lipschitz Condition). *The loss function  $\mathcal{L}(\Theta)$  is  $L_1$ -Lipschitz.*

**Assumption 2** (Smoothness). *The loss function  $\mathcal{L}(\Theta)$ ,  $G(\Theta)$  and  $G_0(\Theta)$  are  $L_2$ -smooth,  $L_2$ -smooth and  $L$ -smooth, respectively.*

**Assumption 3** (Lipschitz Condition for Second Derivatives). *The second derivatives  $\nabla_{\mathbf{R}} \nabla_{\mathbf{H}} G(\Theta)$  and  $\nabla_{\mathbf{H}}^2 G(\Theta)$  are  $L_3$ -Lipschitz and  $L_4$ -Lipschitz, respectively.*

**Assumption 4** (Bounded Domain). *The parameter  $\mathbf{H}$  is in a bounded domain with a diameter  $\Delta$ , i.e., for any  $\mathbf{H}$  and  $\mathbf{H}'$ , we have:*

$$\|\mathbf{H} - \mathbf{H}'\| \leq \Delta. \quad (9)$$

**Assumption 5.** *Function  $\mathcal{F} = F(\mathbf{R}) + h(\mathbf{R})$  is bounded below:*

$$\mathcal{F}^* = \inf_{\mathbf{R}} \mathcal{F}(\mathbf{R}) > -\infty. \quad (10)$$

**Theorem 1.** *Under Assumption 1 - 5, define  $\alpha = -L + \mathcal{V}$ , choose step size  $\eta$  to be  $\frac{2}{L_2 + \alpha}$ ,  $\eta'$  to be  $\frac{1}{6L_0}$ , and suppose  $\alpha < L_2$ ,  $h(\mathbf{R}_t^{\tau-1}) \leq h(\mathbf{R}_t^0)$ ,  $\forall t$ , we have:*

$$\frac{1}{\tau T} \sum_{t=0}^{T-1} \sum_{j=1}^{\tau-1} \|\mathcal{G}_t^j\|^2 \leq \mathcal{O}\left(\frac{1}{T}\right), \quad (11)$$

where  $\|\mathcal{G}_t^j\|$  is used to measure convergence properties of our algorithms. We define the generalized gradient at the  $j$ -th iteration of  $t$ -th round:  $\mathcal{G}_t^j = \frac{1}{\eta'} (\mathbf{R}_t^j - \mathbf{R}_t^{j+1})$ .  $\mathcal{V}$  refers to the regularization coefficient of  $\mathcal{L}_{\text{specific}}$ , and  $\tau$  is the number of local epochs.

**Table 2:** Experimental scenario settings for model heterogeneity.

Scenario	Model	Scale	Layer	Dimension
Scenario 1	Gemma	2B	18	2,048
	Llama3.2	3B	28	3,072
	SmolLM	1.7B	24	2,048
Scenario 2	Llama3	8B	32	4,096
	Llama3.2	3B	28	3,072
	Llama3.2	1B	16	2,048
Scenario 3	Gemma	7B	28	3,072
	Llama3	8B	32	4,096
	Yi	6B	32	4,096

## 6 Experiments

### 6.1 Settings

**Dataset.** We conduct experiments on two well-known benchmarks: MATHInstruct [41] and GLUE [35]. Specifically, MATHInstruct contains various mathematical tasks including fill-in-the-blank (FIBQ), multiple choice (MCQ), and programming questions (PQ). And GLUE is a collection of natural language understanding tasks such as single sentence classification (SSC), sentence pair classification (SPC), and natural language inference (NLI). To ensure balanced comparison across different tasks, we randomly sample 5,000 training examples and 500 test examples for each task.

**Baselines.** As shown in Table 2, we design three model heterogeneous scenarios based on seven FMs, including Gemma [33], Llama [14], SmolLM, and Yi [40], with parameter scales ranging from 1B to 8B. These models feature diverse architectural configurations with varying layer depths and hidden dimensions. To evaluate our proposed  $H^2\text{Tune}$  comprehensively, we compare it against four representative baselines and an upper boundary. Specifically, *LOCAL* trains models independently on each task without FL, serving as the lower performance bound. *FLLM* represents the conventional federated setup where clients share identical foundation models and LoRA ranks, enabling direct parameter aggregation. *HetLoRA* [11] allows clients to adopt different LoRA ranks while sharing the same base foundation models. *FLTLA* aggregates only the tail layers across clients with different foundation models. As an upper boundary, we fine-tune local models on combined datasets of all clients without federation to indicate the maximum achievable performance.

**Implementation Details.** We conduct experiments using two GPU configurations with four L20 GPUs for models exceeding 5B parameters and four RTX 4090 GPUs for smaller models. The system involves 3 client nodes participating in the training process. We set the learning rates ranging from  $2e-7$  to  $2e-3$  with 5 communication rounds. Finally, we evaluate model performance through accuracy and quantify communication cost in MB and minutes.

### 6.2 Main Result

We evaluate  $H^2\text{Tune}$  in both homogeneous and heterogeneous task scenarios. Table 3 presents the comprehensive results, demonstrating  $H^2\text{Tune}$ 's superior performance across all scenarios. Our experimental results show that  $H^2\text{Tune}$  achieves consistent improvements across all scenarios compared to baseline methods. In Scenario 1,  $H^2\text{Tune}$  demonstrates performance gains of 4.3% on average for MATHInstruct tasks and 3.7% for GLUE tasks. Similar trends are observed in Scenario 2, where the model achieves improvements of 2.8% and 3.0% respectively. In Scenario 3,  $H^2\text{Tune}$  shows average gains of 3.0% on MATHInstruct and 3.3% on GLUE tasks. While  $H^2\text{Tune}$  does not always reach the upper boundary performance, it maintains a reasonable gap while providing the significant advantage of using a single model for multiple tasks. Overall, these results

**Table 3:** The overall performance comparison results in different scenarios and task settings, where each client handles a distinct task and @ $r$  denotes the rank of low-rank updates supported by the client’s resource constraints.

Scenario	Model	Homogeneous Task Scenario						Heterogeneous Task Scenario					
		MATHInstruct			GLUE			MATHInstruct			GLUE		
		MCQ@16	MCQ@64	MCQ@128	SSC@16	SSC@64	SSC@128	FIBQ@16	MCQ@64	PQ@128	SSC@16	SPC@64	NLI@128
Scenario 1	LOCAL	26.0	30.6	17.2	93.8	93.6	92.2	28.8	32.0	12.2	90.0	83.8	54.2
	FLLM	26.4	32.2	18.0	94.2	93.6	92.4	29.4	35.6	14.4	90.6	86.2	77.6
	HetLoRA	27.0	33.6	18.2	94.4	93.8	92.4	29.4	29.0	12.2	95.6	87.0	86.8
	FLTLA	24.2	21.0	16.8	93.6	92.8	91.8	23.2	27.8	7.8	94.0	84.4	64.0
	$H^2Tune$	<b>29.2</b>	<b>34.4</b>	<b>20.6</b>	<b>94.4</b>	<b>94.2</b>	<b>92.6</b>	<b>32.4</b>	<b>35.0</b>	<b>17.2</b>	<b>95.8</b>	<b>87.0</b>	<b>86.0</b>
	avg. Imp	3.3 $\uparrow$	5.1 $\uparrow$	3.1 $\uparrow$	0.4 $\uparrow$	0.8 $\uparrow$	0.4 $\uparrow$	4.7 $\uparrow$	3.9 $\uparrow$	5.6 $\uparrow$	3.3 $\uparrow$	1.7 $\uparrow$	15.4 $\uparrow$
Scenario 2	LOCAL	36.4	30.6	23.6	81.8	93.6	72.2	53.2	32.0	3.2	52.2	83.8	81.4
	FLLM	37	31.2	24.4	83.6	93.6	79.0	56.4	32.2	7.2	52.4	88.2	85.0
	HetLoRA	40.8	32.0	24.4	84.8	93.6	80.2	57.0	33.0	8.8	53.0	88.4	85.2
	FLTLA	28.8	29.0	21.2	81.4	85.4	72.2	53.4	31.0	6.6	52.2	83.0	81.4
	$H^2Tune$	<b>42.2</b>	<b>33.2</b>	<b>25.2</b>	<b>87.8</b>	<b>93.8</b>	<b>81.4</b>	<b>57.2</b>	<b>33.0</b>	<b>9.0</b>	<b>53.0</b>	<b>88.4</b>	<b>85.2</b>
	avg. Imp	6.5 $\uparrow$	2.5 $\uparrow$	1.8 $\uparrow$	4.9 $\uparrow$	2.3 $\uparrow$	5.5 $\uparrow$	2.3 $\uparrow$	1.0 $\uparrow$	2.6 $\uparrow$	0.6 $\uparrow$	2.6 $\uparrow$	2.0 $\uparrow$
Scenario 3	LOCAL	30.8	41.4	32.4	87.8	83.6	94.4	57.0	45.0	5.6	94.2	83.6	89.8
	FLLM	32.0	41.4	33.0	87.8	84.8	94.6	57.6	45.4	10.4	95.2	86.2	89.8
	HetLoRA	35.8	43.2	36.0	92.4	87.0	95.4	57.6	45.8	9.6	96.0	87.6	90.4
	FLTLA	30.4	40.2	31.0	82.2	83.0	93.2	55.8	44.0	6.0	94.8	84.6	88.8
	$H^2Tune$	<b>36.0</b>	<b>43.6</b>	<b>37.0</b>	<b>92.8</b>	<b>89.8</b>	<b>95.6</b>	<b>61.4</b>	<b>46.8</b>	<b>10.6</b>	<b>96.4</b>	<b>88.2</b>	<b>91.4</b>
	avg. Imp	3.6 $\uparrow$	2.1 $\uparrow$	3.9 $\uparrow$	7.8 $\uparrow$	5.2 $\uparrow$	1.2 $\uparrow$	4.4 $\uparrow$	1.5 $\uparrow$	2.7 $\uparrow$	1.4 $\uparrow$	2.7 $\uparrow$	1.7 $\uparrow$
Scenario 3	Upper Boundary	39.2	47.0	37.6	95.0	95.2	96.8	65.0	48.6	12.6	97.2	88.4	92.4

**Table 4:** The influence of global ranks under heterogeneous and homogeneous scenarios.

Settings	Client 1	Client 2	Client 3	Client 1	Client 2	Client 3
Support max rank	16	64	64	16	64	128
Heterogeneous LLM	LLAMA-3B	GEMMA-2B	SmolLM-1.7B	LLAMA-3B	GEMMA-2B	SmolLM-1.7B
Heterogeneous Task	FIBQ	MCQ	PQ	FIBQ	MCQ	PQ
Global rank $r = 64$	30.6	33.4	10.4	-	-	-
Global rank $r = 128$	31.0	33.5	10.7	31.6	34.8	17.2
Global rank $r = 192$	31.1	33.5	10.9	31.8	34.8	17.0
Heterogeneous Task	SSC	SPC	NLI	SSC	SPC	NLI
Global rank $r = 64$	95.4	86.8	84.8	-	-	-
Global rank $r = 128$	95.8	87.1	85.5	95.8	87.0	86.0
Global rank $r = 192$	96.4	87.5	85.8	96.2	87.6	85.2
Homogeneous Task	MCQ					
Global rank $r = 64$	27.0	33.2	17.3	27.5	33.7	18.2
Global rank $r = 128$	28.7	34.0	19.9	29.8	34.6	21.3
Global rank $r = 192$	28.6	34.2	18.5	29.4	35.0	20.1

**Table 5:** Results of ablation experiments, where TSM, ATKD, and SM represent Task-shared matrix, Alternating Task-Knowledge Disentanglement, and Sparsification Matrix, respectively.

	FIBQ	MCQ	PQ	SSC	SPC	NLI
$H^2Tune$	<b>32.4</b>	<b>35.0</b>	<b>17.2</b>	<b>95.8</b>	<b>87.0</b>	<b>86.0</b>
w/o TSM	28.8	27.2	10.8	90.4	66.6	54.4
w/o ATKD	29.0	27.8	11.2	95.4	86.2	86.8
w/o SM	30.2	29.4	19.8	95.8	86.4	88.2

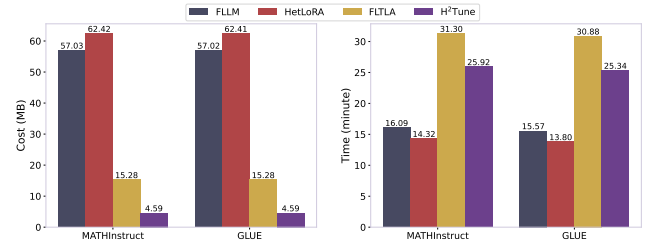
clearly demonstrate that  $H^2Tune$  can effectively improve model performance in both homogeneous and heterogeneous task scenarios.

### 6.3 Ablation Study

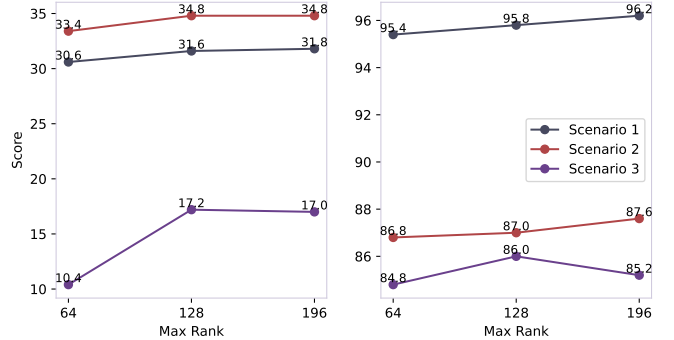
To investigate the effectiveness of proposed components in  $H^2Tune$ , we conduct ablation experiments as shown in Table 5. Specifically, removing TSM causes the most performance drop, with average accuracy decreasing by 5.9% on MATHInstruct and 19.1% on GLUE, indicating its crucial role in capturing shared knowledge. Without ATKD, we observe an average performance reduction of 5.5% on MATHInstruct and minimal impact on GLUE. And the removal of SM shows moderate performance degradation on MATHInstruct with a 1.7% average drop, while maintaining comparable performance on GLUE tasks. This highlights that sparsification plays a task-dependent role in performance.

### 6.4 Communication Cost

We evaluate the communication costs and training efficiency using MATHInstruct and GLUE benchmarks, as shown in Figure 4. While



**Figure 4:** Comparison of local and strategy communication costs.



**Figure 5:** Performance of heterogeneous models with different maximum ranks on MATHInstruct (left) and GLUE (right).

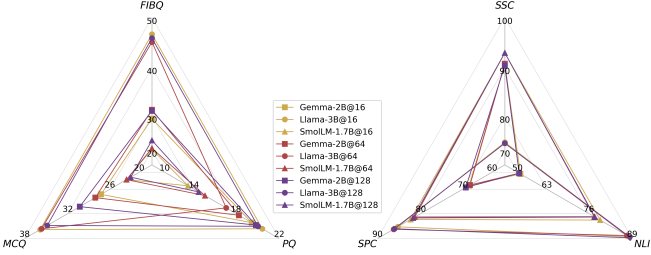
FLLM and HetLoRA require substantial communication overhead, our proposed  $H^2Tune$  achieves efficient communication at only 4.59 MB, significantly outperforming FLTLA. In terms of training time,  $H^2Tune$  shows similar computational overhead to FLTLA, which is moderately higher than FLLM and HetLoRA due to the task-knowledge disentanglement process. Notably, this slight increase in training time is well justified by the significant performance improvements demonstrated in our main results, while maintaining better communication efficiency compared with FLTLA.

### 6.5 Influence of Global Ranks

**Homogeneous task scenario.** Table 4 presents the impact of increasing global rank on model performance in homogeneous task scenario where all clients perform MCQ tasks. Results indicate that increasing the global rank from 64 to 128 yields consistent performance improvements across all clients. Further increasing the global rank to 192 produces minimal additional gains in some cases, suggest-

**Table 6:** The influence of global layers under heterogeneous and homogeneous scenarios.

Settings	Client 1	Client 2	Client 3	Client 1	Client 2	Client 3	Client 1	Client 2	Client 3
Heterogeneous LLM	LLaMA-3B	GEMMA-2B	SmolLM-1.7B	LLaMA-3B	GEMMA-2B	SmolLM-1.7B	LLaMA-3B	LLaMA-3.2B	LLaMA-3.2B
Heterogeneous Task	FIBQ	MCQ	PQ	SSC	SPC	NLI	FIBQ	MCQ	PQ
Real layer	28	18	24	28	18	24	32	28	16
Global layer	L=18	30.4	29.2	15.4	49.7	82.3	78.9	60.4	40.6
	L=24	32.4	29.2	16.6	51.3	84.7	82.1	62.8	42.9
	L=28	34.8	31.6	17.2	52.8	87.4	83.7	63.9	44.1
	L=32	31.4	31.6	17.2	53.0	88.4	85.2	64.1	46.8
	L=40	31.6	31.8	20.8	52.2	86.9	84.8	63.2	45.5
	L=48	29.6	32.4	15.4	47.9	85.6	84.0	59.4	45.3
Homogeneous Task	MCQ								
Real layer	28	18	24	28	18	24	32	28	16
Global layer	L=18	26.8	31.2	18.7	38.4	29.7	24.7	33.5	41.8
	L=24	27.7	33.1	20.0	40.7	32.5	24.3	34.2	42.9
	L=28	28.4	33.8	19.8	41.3	32.9	24.7	35.1	41.8
	L=32	29.2	34.4	20.6	42.2	33.2	25.2	36.0	43.6
	L=40	28.6	33.7	20.3	41.8	32.4	24.5	35.5	40.7
	L=48	27.9	33.9	19.2	42.1	32.0	24.1	35.8	42.9



**Figure 6:** Performance of homogeneous models with different ranks. @ $r$  denotes the resource-constrained rank for client updates.

ing that the benefits of higher rank representations plateau after a certain threshold. Notably, the configuration with support max rank of 16/64/128 consistently outperforms the 16/64/64 configuration across all global ranks, with client 3 showing the most significant improvement when given a higher maximum rank capacity.

**Heterogeneous task scenario.** For heterogeneous task settings, we observe nuanced performance patterns across different client configurations. In the first task set, lifting the global rank from 64 to 128 yields gains of 0.4, 0.1, and 0.3 percentage points for clients 1, 2, and 3, respectively. For the second task set involving classification tasks, performance improvements are more pronounced as the global rank increases, with consistent gains across SSC, SPC, and NLI tasks.

**Different model heterogeneity scenarios.** Figure 5 shows the results of Gemma-2B, Llama-3B, and SmolLM-1.7B at different maximum ranks using the  $H^2Tune$  strategy, where the maximum ranks for the three clients varying from 64 to 196. The results demonstrate that increasing the maximum rank from 64 to 128 improves performance. However, further increasing to 196 yields marginal gains, indicating that an appropriate rank setting helps balance between feature expressiveness and potential noise introduction.

## 6.6 Influence of Global Layers

**Homogeneous task scenario.** Table 6 presents the impact of global layer count on model performance in homogeneous task scenario. As global layer count increases from 18 to 32, LLaMA-3B, GEMMA-2B, and SmolLM-1.7B show accuracy gains of 2.4%, 3.2%, and 1.9% in scenario 1. Similar trends appear in Scenarios 2 and 3. However, further increasing to 40 or 48 layers yields diminishing returns or even slight performance decreases.

**Heterogeneous task scenario.** For heterogeneous tasks, we observe more varied responses to global layer count increases. For example, FIBQ performance peaks at  $L = 28$  with 34.8% accuracy, while MCQ and PQ achieve optimal results at  $L = 48$  and  $L = 40$  with 32.4% and 20.8% accuracy in scenario 1. The results indicate hetero-

**Table 7:** Hyper-parameter analysis results for  $\beta_k$ .

Model	$\beta_k$	Task	MATHInstruct	Task	GLUE
LLaMA-3B	0.125	FIBQ	31.6	SSC	95.3
	0.500		32.4		95.8
	1.000		32.2		95.4
GEMMA-2B	0.125	MCQ	32.4	SPC	86.2
	0.500		34.8		87.0
	1.000		35.0		86.8
SmolLM-1.7B	0.125	PQ	11.8	NLI	86.0
	0.500		17.0		85.1
	1.000		17.2		85.9

geneous scenarios generally require higher global layer, with most optimal performance observed between  $L = 28$  and  $L = 32$ .

## 6.7 Hyper-parameter Analysis

We investigate how the local parsity ratio  $\beta_k$  affects model performance in both homogeneous and heterogeneous model settings.

**Homogeneous models with different local ranks.** Figure 6 illustrates the performance of Gemma-2B, Llama-3B, and SmolLM-1.7B across different ranks in a homogeneous model setting, where all clients use identical models and ranks. The performance is visualized as radar charts, where a larger triangle area indicates better overall performance, and the angle bias reflects performance imbalance across tasks. We observe that increasing rank values consistently leads to better performance, suggesting that larger LoRA parameters can capture more task-relevant information.

**Heterogeneous models with different local ranks.** As shown in Table 7, we observe that increasing sparsity ratio  $\beta_k$  generally improves performance across models and tasks, with optimal values varying by task type. These results indicate that mathematical reasoning tasks benefit from higher sparsity ratios, while language understanding tasks often perform optimally at moderate values around  $\beta_k = 0.500$ , suggesting task-specific adaptation of sparsity is beneficial for heterogeneous model deployments.

## 7 Conclusion

This paper explores an under-explored direction: federated fine-tuning with heterogeneous models, tasks, and resources. We propose triple matrix decomposition and layer relation matrices to handle heterogeneous matrix fusion. Our alternating optimization separates shared and private knowledge to mitigate multi-task interference, enabling cross-client knowledge transfer. Extensive experiments show our method improves accuracy by up to 15.4% over baselines.

## References

- [1] J.-H. Ahn, O. Simeone, and J. Kang. Wireless federated distillation for distributed edge learning with heterogeneous data. In *2019 IEEE 30th Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, pages 1–6. IEEE, 2019.
- [2] J.-H. Ahn, O. Simeone, and J. Kang. Cooperative learning via federated distillation over fading channels. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8856–8860. IEEE, 2020.
- [3] S. Alam, L. Liu, M. Yan, and M. Zhang. Fedrolex: Model-heterogeneous federated learning with rolling sub-model extraction. *Advances in neural information processing systems*, 35:29677–29690, 2022.
- [4] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [5] J. Bai, D. Chen, B. Qian, L. Yao, and Y. Li. Federated fine-tuning of large language models under heterogeneous tasks and client resources. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [6] H. Chang, V. Shejwalkar, R. Shokri, and A. Houmansadr. Cronus: Robust and heterogeneous collaborative learning with black-box knowledge transfer. *arXiv preprint arXiv:1912.11279*, 2019.
- [7] J. Chen and A. Zhang. Fedmsplit: Correlation-adaptive federated multi-task learning across multimodal split networks. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 87–96, 2022.
- [8] J. Chen, R. Zhang, J. Guo, Y. Fan, and X. Cheng. Fedmatch: Federated learning over heterogeneous question answering data. In *Proceedings of the 30th ACM international conference on information & knowledge management*, pages 181–190, 2021.
- [9] Y.-Q. Chen, T. Zhang, X.-L. Jiang, Q. Chen, C.-L. Gao, and W.-L. Huang. Fedbone: Towards large-scale federated multi-task learning. *Journal of Computer Science and Technology*, 39(5):1040–1057, 2024.
- [10] S. Cheng, J. Wu, Y. Xiao, and Y. Liu. Fedgems: Federated learning of larger server models via selective knowledge fusion. *arXiv preprint arXiv:2110.11027*, 2021.
- [11] Y. J. Cho, L. Liu, Z. Xu, A. Fahrezi, and G. Joshi. Heterogeneous lora for federated fine-tuning of on-device foundation models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 12903–12913, 2024.
- [12] E. Diao, J. Ding, and V. Tarokh. Heterofit: Computation and communication efficient federated learning for heterogeneous clients. *arXiv preprint arXiv:2010.01264*, 2020.
- [13] N. Ding, Y. Qin, G. Yang, F. Wei, Z. Yang, Y. Su, S. Hu, Y. Chen, C.-M. Chan, W. Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.
- [14] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [15] S. Ghadimi and M. Wang. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018.
- [16] S. Ghadimi, G. Lan, and H. Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1):267–305, 2016.
- [17] Z. Han, C. Gao, J. Liu, J. Zhang, and S. Q. Zhang. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024.
- [18] S. Horvath, S. Laskaridis, M. Almeida, I. Leontiadis, S. Venieris, and N. Lane. Fjord: Fair and accurate federated learning under heterogeneous targets with ordered dropout. *Advances in Neural Information Processing Systems*, 34:12876–12889, 2021.
- [19] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [20] F. Huang, J. Li, S. Gao, and H. Huang. Enhanced bilevel optimization via bregman distance. *Advances in Neural Information Processing Systems*, 35:28928–28939, 2022.
- [21] J. Jang, H. Ha, D. Jung, and S. Yoon. Fedclassavg: Local representation learning for personalized federated learning on heterogeneous neural networks. In *Proceedings of the 51st international conference on parallel processing*, pages 1–10, 2022.
- [22] K. Ji, J. D. Lee, Y. Liang, and H. V. Poor. Convergence of meta-learning with task-specific adaptation over partial parameters. *Advances in Neural Information Processing Systems*, 33:11490–11500, 2020.
- [23] K. Ji, J. Yang, and Y. Liang. Bilevel optimization: Convergence analysis and enhanced design. In *International conference on machine learning*, pages 4882–4892. PMLR, 2021.
- [24] R. Jia, W. Xie, J. Lei, H. Qin, J. Ma, and L. Fang. Towards efficient model-heterogeneity federated learning for large models. *arXiv preprint arXiv:2411.16796*, 2024.
- [25] Y. Jia, X. Zhang, A. Beheshti, and W. Dou. Fedlps: heterogeneous federated learning for multiple tasks with local parameter sharing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 12848–12856, 2024.
- [26] J. Jiang, H. Jiang, Y. Ma, X. Liu, and C. Fan. Low-parameter federated learning with large language models. In *International Conference on Web Information Systems and Applications*, pages 319–330. Springer, 2024.
- [27] W. Kuang, B. Qian, Z. Li, D. Chen, D. Gao, X. Pan, Y. Xie, Y. Li, B. Ding, and J. Zhou. Federatedscope-llm: A comprehensive package for fine-tuning large language models in federated learning. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5260–5271, 2024.
- [28] X.-Y. Liu, R. Zhu, D. Zha, J. Gao, S. Zhong, M. White, and M. Qiu. Differentially private low-rank adaptation of large language model using federated learning. *ACM Transactions on Management Information Systems*, 2023.
- [29] Y. Liu, Y. Zhang, Q. Li, T. Liu, S. Feng, D. Wang, Y. Zhang, and H. Schütze. Hift: A hierarchical full parameter fine-tuning strategy. *arXiv preprint arXiv:2401.15207*, 2024.
- [30] X. Lu, Y. Liao, C. Liu, P. Lio, and P. Hui. Heterogeneous model fusion federated learning mechanism based on model mapping. *IEEE Internet of Things Journal*, 9(8):6058–6068, 2021.
- [31] S. Rajput, A. Gupta, and D. Papailiopoulos. Closing the convergence gap of sgd without replacement. In *International Conference on Machine Learning*, pages 7964–7973. PMLR, 2020.
- [32] T. Shen, J. Zhang, X. Jia, F. Zhang, G. Huang, P. Zhou, K. Kuang, F. Wu, and C. Wu. Federated mutual learning. *arXiv preprint arXiv:2006.16765*, 2020.
- [33] G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Riviere, M. S. Kale, J. Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- [34] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [35] A. Wang. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- [36] C. Wu, F. Wu, L. Lyu, Y. Huang, and X. Xie. Communication-efficient federated learning via knowledge distillation. *Nature communications*, 13(1):2032, 2022.
- [37] Q. Yang, Y. Liu, T. Chen, and Y. Tong. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):1–19, 2019.
- [38] Z. Yang, S. Fu, W. Bao, D. Yuan, and A. Y. Zomaya. Fastslowmo: Federated learning with combined worker and aggregator momenta. *IEEE Transactions on Artificial Intelligence*, 4(5):1041–1050, 2022.
- [39] L. Yi, H. Yu, G. Wang, and X. Liu. Fedlora: Model-heterogeneous personalized federated learning with lora tuning. *arXiv preprint arXiv:2310.13283*, 2023.
- [40] A. Young, B. Chen, C. Li, C. Huang, G. Zhang, G. Zhang, H. Li, J. Zhu, J. Chen, J. Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024.
- [41] X. Yue, X. Qu, G. Zhang, Y. Fu, W. Huang, H. Sun, Y. Su, and W. Chen. Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*, 2023.
- [42] J. Zhang, S. Vahidian, M. Kuo, C. Li, R. Zhang, T. Yu, G. Wang, and Y. Chen. Towards building the federatedgpt: Federated instruction tuning. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6915–6919. IEEE, 2024.
- [43] L. Zhang, L. Shen, L. Ding, D. Tao, and L.-Y. Duan. Fine-tuning global model via data-free knowledge distillation for non-iid federated learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10174–10183, 2022.
- [44] L. Zhang, D. Wu, and X. Yuan. Fedzkt: Zero-shot knowledge transfer towards resource-constrained federated learning with heterogeneous on-device models. In *2022 IEEE 42nd International Conference on Distributed Computing Systems (ICDCS)*, pages 928–938. IEEE, 2022.



## A Algorithm Pseudocode

---

### Algorithm 1 $H^2$ Tune

---

**Input:** communication round  $T$ ; learning rate:  $\eta$  and  $\eta'$ ; epoch:  $\tau$ ; training dataset  $(x_i, y_i) \in D_k$  of client  $k$ .

**for**  $t = 0, 1, \dots, T - 1$  **do**

parties receive global task-shared matrix  $\mathbf{R}_g^t$ ;

**for** party  $k = 1, \dots, K$  **in parallel do**

**for**  $j = 1, \dots, \tau$  **do**

**for**  $(x_i, y_i) \in D_k$  **do**

**Freeze** task-specific matrices  $\mathbf{A}_k$  and  $\mathbf{B}_k$ ;

update  $\Phi_k^{t,j} \mathbf{R}_k^{t,j} = \Phi_k^{t,j} \mathbf{R}_k^{t,j-1} - \eta' \nabla_{\mathbf{R}} \mathcal{L}_{share}(x_i, y_i)$ ;

update  $\Omega_k^{t,j} = \Omega_k^{t,j-1} - \eta' \nabla_{\Omega} \mathcal{L}_{share}(x_i, y_i)$ ;

**Freeze** task-shared matrix  $\Phi_k \mathbf{R}_k$  and relation matrix  $\Omega_k$ ;

update  $\mathbf{A}_k^{t,j} = \mathbf{A}_k^{t,j} - \eta \nabla_{\mathbf{A}} \mathcal{L}_{specific}(x_i, y_i)$ ;

update  $\mathbf{B}_k^{t,j} = \mathbf{B}_k^{t,j} - \eta \nabla_{\mathbf{B}} \mathcal{L}_{specific}(x_i, y_i)$ ;

**end for**

**end for**

upload task-shared matrix  $\mathbf{R}_k^t$  to server;

**end for**

server aggregates  $\mathbf{R}_k^t$  from party  $k = 1, \dots, K$  by  $\mathbf{R}_g^t = \frac{1}{K} \sum_{k=1}^K \mathbf{R}_k^t$ ;

**end for**

---

## B Convergence Analysis

We begin with introducing the convergence metric  $\|\mathcal{G}_t^j\|$  to measure convergence properties of our algorithms. We define the generalized gradient at the  $j$ -th iteration of  $t$ -th round:

$$\mathcal{G}_t^j = \frac{1}{\eta'} (\mathbf{R}_t^j - \mathbf{R}_t^{j+1}).$$

When  $j \leq \tau - 2$ , the  $j$ -th iteration corresponds to the gradient descent step, we have  $\|\mathcal{G}_t^j\| = \|\nabla F(\mathbf{R}_t^j)\|$  which is a common convergence metric used in [15, 23]. When  $j = \tau - 1$ , the  $j$ -th iteration corresponds to the proximal step, where  $\|\mathcal{G}_t^j\|$  was also used as an important metric ([20]). To simplify the convergence analysis, we can rewrite Algorithm 1 as Algorithm 2, and present the convergence analysis under Algorithm 2.

---

### Algorithm 2 Proximity Version of $H^2$ Tune

---

**Input:** communication round  $T$ ; learning rate:  $\eta$  and  $\eta'$ ; epoch:  $\tau$ .

**for**  $t = 0, 1, \dots, T - 1$  **do**

Let  $\mathbf{H}_t^0 = \mathbf{H}_{t-1}^\tau$  and  $\mathbf{R}_t^0 = \mathbf{R}_{t-1}^\tau$ ;

**for**  $j = 1, 2, \dots, \tau - 1$  **do**

update  $\mathbf{R}_t^j = \mathbf{R}_t^{j-1} - \eta' \nabla_{\mathbf{R}} \mathcal{L}(\mathbf{R}_t^{j-1}, \mathbf{H}_{t-1}^\tau)$ ;

**end for**

update  $\mathbf{R}_t^\tau = \arg \min_{\mathbf{R}} \{ \langle \frac{\partial \mathcal{L}(\mathbf{R}_t^{\tau-1}, \mathbf{H}_{t-1}^\tau)}{\partial \mathbf{R}_t^{\tau-1}}, \mathbf{R} \rangle + h(\mathbf{R}) + \frac{1}{\eta'} \|\mathbf{R} - \mathbf{R}_t^{\tau-1}\|^2 \}$

**for**  $j = 1, 2, \dots, \tau$  **do**

update  $\mathbf{H}_t^j = \mathbf{H}_t^{j-1} - \eta \nabla_{\mathbf{H}} \mathcal{G}(\mathbf{R}_t^j, \mathbf{H}_t^{j-1})$ ;

**end for**

**end for**

---

In Algorithm 2,  $\mathbf{R}_t^\tau = \arg \min_{\mathbf{R}} \{ \langle \frac{\partial \mathcal{L}(\mathbf{R}_t^{\tau-1}, \mathbf{H}_{t-1}^\tau)}{\partial \mathbf{R}_t^{\tau-1}}, \mathbf{R} \rangle + h(\mathbf{R}) + \frac{1}{\eta'} \|\mathbf{R} - \mathbf{R}_t^{\tau-1}\|^2 \}$  is equivalent to doing the last step in Algorithm 1 first, and then assigns a personalized parameter  $\mathbf{R}_k$  for each client.

### B.1 Proof Sketch

**Proof Sketch.** To explore the essential insights, we first bound the difference between the local parameter and the local optimal parameter at the  $t$ -th round (i.e. tracking error  $\|\mathbf{H}_t^{j-1} - \mathbf{H}^*(\mathbf{R}_t^0)\|$ ). Then, the upper bound of the difference between the approximate gradient and the exact gradient (i.e.  $\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_t^j)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \|$ ) is provided by a proof technology called virtual updates ([38]). Lastly, we make use of the property of the proximal operator in our bi-level optimization problem to deal with the proximal step in our algorithm.

### B.2 Proof Details

To provide the convergence analysis of our algorithm, we need to give some auxiliary lemmas first. Firstly, with  $L$ -smoothness of  $G_0(\mathbf{R}, \mathbf{H})$ , we could obtain the strongly-convexity property of  $G(\mathbf{R}, \mathbf{H})$ .

**Lemma 2.** Under Assumption 2, suppose  $\alpha := -L + \mathcal{V} > 0$ ,  $G(\mathbf{R}, \mathbf{H})$  is  $\alpha$ -strongly convex w.r.t  $\mathbf{H}$ .

The proof of this lemma is trivial. Thus, we omit it.

**Lemma 3** (Lemma 2.2 in [15]). Under Assumptions 1, 2, and 3,  $F(\mathbf{R})$  is  $L_0$ -smooth, where  $L_0$  is given by

$$L_0 := L_2 + \frac{2L_2^2 + L_1^2 L_3}{\alpha} + \frac{L_1 L_2 L_3 + L_1 L_2 L_4 + L_2^3}{\alpha^2} + \frac{L_1 L_2^2 L_4}{\alpha^3}. \quad (12)$$

Tracking error  $\|\mathbf{H}_t^{j-1} - \mathbf{H}^*(\mathbf{R}_t^0)\|$  is an important component in our convergence analysis. To give an upper bound on the tracking error, we utilized Lemma 9 in [23].

**Lemma 4** (Lemma 9 in [23]). Under Assumptions 1 and 2, with step size  $\eta$  to be  $\frac{2}{L_2 + \alpha}$ , we have

$$\|\mathbf{H}_t^{(j-1)} - \mathbf{H}^*(\mathbf{R}_t^0)\| \leq \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^{(j-1)} \|\mathbf{H}_t^0 - \mathbf{H}^*(\mathbf{R}_t^0)\|. \quad (13)$$

With the help of the above lemmas, we could give the estimation property of the  $\frac{\partial \mathcal{L}(\mathbf{R}_t^0, \mathbf{H}_t^0)}{\partial \mathbf{R}_t^0}$  approximating  $\nabla F(\mathbf{R}_t^0)$ . The result is presented in the following Proposition 5.

**Proposition 5.** Under Assumptions 1-5, choose step size  $\eta$  to be  $\frac{2}{L_2 + \alpha}$ , we have

$$\begin{aligned} \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^0, \mathbf{H}_t^0)}{\partial \mathbf{R}_t^0} - \nabla F(\mathbf{R}_t^0) \right\| &\leq \left( L_2 + \frac{L_2^2}{\alpha} \right) \left[ \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \cdot \Delta \right] \\ &+ L_1 \left[ \frac{L_2 \left( 1 - \frac{2}{L_2 + \alpha} \cdot \alpha \right)^\tau}{\alpha} + \frac{2}{L_2 + \alpha} \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \cdot \Delta \right. \\ &\quad \left. \left( 1 - \frac{2}{L_2 + \alpha} \cdot \alpha \right)^\tau \right] \cdot \frac{1}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}}. \end{aligned} \quad (14)$$

**Proof.** Using:

$$\begin{aligned} \nabla F(\mathbf{R}_t^0) &= \nabla_{\mathbf{R}} \mathcal{L}(\mathbf{R}_t^0, \mathbf{H}^*(\mathbf{R}_t^0)) \\ &+ \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \nabla_{\mathbf{H}} \mathcal{L}(\mathbf{R}_t^0, \mathbf{H}^*(\mathbf{R}_t^0)), \end{aligned} \quad (15)$$

and:

$$\frac{\partial \mathcal{L}(\mathbf{R}_t^0, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^0} = \nabla_{\mathbf{R}} \mathcal{L}(\mathbf{R}_t^0, \mathbf{H}_t^\tau) + \frac{\partial \mathbf{H}_t^\tau}{\partial \mathbf{R}_t^0} \nabla_{\mathbf{H}} \mathcal{L}(\mathbf{R}_t^0, \mathbf{H}_t^\tau), \quad (16)$$

we have:

$$\begin{aligned} \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^0, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^0} - \nabla F(\mathbf{R}_t^0) \right\| &\leq L_2 \|\mathbf{H}_t^\tau - \mathbf{H}^*(\mathbf{R}_t^0)\| + L_1 \left\| \frac{\partial \mathbf{H}_t^\tau}{\partial \mathbf{R}_t^0} \right. \\ &\quad \left. - \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \right\| + L_2 \left\| \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \right\| \cdot \|\mathbf{H}_t^\tau - \mathbf{H}^*(\mathbf{R}_t^0)\|. \end{aligned} \quad (17)$$

Now we first bound:  $\left\| \frac{\partial \mathbf{H}_t^\tau}{\partial \mathbf{R}_t^0} - \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \right\|$ . Recall the update method:

$$\mathbf{H}_t^j = \mathbf{H}_t^{j-1} - \eta \nabla_{\mathbf{H}} G(\mathbf{R}_t^0, \mathbf{H}_t^{j-1}), \quad (18)$$

and we apply the chain rule on it, we have:

$$\begin{aligned} \frac{\partial \mathbf{H}_t^j}{\partial \mathbf{R}_t^0} &= \frac{\partial \mathbf{H}_t^{j-1}}{\partial \mathbf{R}_t^0} - \eta \left( \nabla_{\mathbf{R}} \nabla_{\mathbf{H}} G(\mathbf{R}_t^0, \mathbf{H}_t^{j-1}) \right. \\ &\quad \left. + \frac{\partial \mathbf{H}_t^{j-1}}{\partial \mathbf{R}_t^0} \nabla_{\mathbf{H}}^2 G(\mathbf{R}_t^0, \mathbf{H}_t^{j-1}) \right). \end{aligned} \quad (19)$$

Since  $\mathbf{H}^*(\mathbf{R}_t^0)$  is the optimal solution of  $G(\mathbf{R}_t^0, \mathbf{H})$ , we have  $\nabla_{\mathbf{H}} G(\mathbf{R}_t^0, \mathbf{H}^*(\mathbf{R}_t^0)) = 0$ . Then, using the chain rule, we have:

$$\nabla_{\mathbf{R}} \nabla_{\mathbf{H}} G(\mathbf{R}_t^0, \mathbf{H}^*(\mathbf{R}_t^0)) + \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \nabla_{\mathbf{H}}^2 G(\mathbf{R}_t^0, \mathbf{H}^*(\mathbf{R}_t^0)) = 0. \quad (20)$$

Combining (19) and (20), the following equation holds:

$$\begin{aligned} \frac{\partial \mathbf{H}_t^j}{\partial \mathbf{R}_t^0} - \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} &= \frac{\partial \mathbf{H}_t^{j-1}}{\partial \mathbf{R}_t^0} - \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \\ &\quad - \eta \left( \nabla_{\mathbf{R}} \nabla_{\mathbf{H}} G(\mathbf{R}_t^0, \mathbf{H}_t^{j-1}) - \nabla_{\mathbf{R}} \nabla_{\mathbf{H}} G(\mathbf{R}_t^0, \mathbf{H}^*(\mathbf{R}_t^0)) \right) \\ &\quad - \eta \left( \frac{\partial \mathbf{H}_t^{j-1}}{\partial \mathbf{R}_t^0} - \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \right) \nabla_{\mathbf{H}}^2 G(\mathbf{R}_t^0, \mathbf{H}_t^{j-1}) \\ &\quad + \eta \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \left( \nabla_{\mathbf{H}}^2 G(\mathbf{R}_t^0, \mathbf{H}^*(\mathbf{R}_t^0)) - \nabla_{\mathbf{H}}^2 G(\mathbf{R}_t^0, \mathbf{H}_t^{j-1}) \right). \end{aligned} \quad (21)$$

Combining (20) and Assumption 2 yields:

$$\left\| \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \right\| \leq \frac{L_2}{\alpha}. \quad (22)$$

With the help of Lemma 3, (21) and (22), the following bound holds:

$$\begin{aligned} \left\| \frac{\partial \mathbf{H}_t^j}{\partial \mathbf{R}_t^0} - \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \right\| &\leq (1 - \eta\alpha) \left\| \frac{\partial \mathbf{H}_t^{j-1}}{\partial \mathbf{R}_t^0} - \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \right\| \\ &\quad + \eta \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \left\| \mathbf{H}_t^{j-1} - \mathbf{H}^*(\mathbf{R}_t^0) \right\|. \end{aligned} \quad (23)$$

For the choice of  $\eta = \frac{2}{L_2 + \alpha}$ , Lemma 4 holds. With Lemma 4 and Assumption 4, we have:

$$\begin{aligned} \left\| \mathbf{H}_t^{j-1} - \mathbf{H}^*(\mathbf{R}_t^0) \right\| &\leq \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^{j-1} \left\| \mathbf{H}_t^0 - \mathbf{H}^*(\mathbf{R}_t^0) \right\| \\ &\leq \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^{j-1} \cdot \Delta. \end{aligned} \quad (24)$$

Telescoping (23) over  $j$  from 0 to  $\tau$  and combining (24) yields:

$$\begin{aligned} \left\| \frac{\partial \mathbf{H}_t^\tau}{\partial \mathbf{R}_t^0} - \frac{\partial \mathbf{H}^*(\mathbf{R}_t^0)}{\partial \mathbf{R}_t^0} \right\| &\leq \frac{L_2(1 - \eta\alpha)^\tau}{\alpha} + \eta \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \\ &\quad \Delta \cdot \frac{(1 - \eta\alpha)^\tau}{1 - \eta\alpha - \frac{L_2 - \alpha}{L_2 + \alpha}}. \end{aligned} \quad (25)$$

Plugging (22) (24) (25) into (17) yields:

$$\begin{aligned} \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^0, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^0} - \nabla F(\mathbf{R}_t^0) \right\| &\leq \left( L_2 + \frac{L_2^2}{\alpha} \right) \left[ \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \cdot \Delta \right] \\ &\quad + L_1 \left[ \frac{L_2(1 - \frac{2}{L_2 + \alpha} \cdot \alpha)^\tau}{\alpha} + \frac{2}{L_2 + \alpha} \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \Delta \right. \\ &\quad \left. \cdot \frac{\left( 1 - \frac{2}{L_2 + \alpha} \cdot \alpha \right)^\tau}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right]. \end{aligned} \quad (26)$$

Hence, we complete the proof of Proposition 5.

To give an upper bound on  $\left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \right\|$ , we make use of a technology called virtual updates ([38]). Specifically, we introduce a virtual parameter  $\mathbf{H}_{t,j}^\tau$  obtained by  $\tau$  updates when  $\mathbf{R}_t^j$  is given. Using Proposition 5, We can bound  $\left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_{t,j}^\tau)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \right\|$ . Now, we aim to bound  $\left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \right\|$ . The result is shown in Proposition 6.

**Proposition 6.** *Following the conditions in Proposition 5, we have:*

$$\begin{aligned} \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \right\| &\leq \left( L_2 + \frac{L_2^2}{\alpha} \right) \left[ \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \cdot \Delta \right] \\ &\quad + L_1 \left[ \frac{L_2 \left( 1 - \frac{2}{L_2 + \alpha} \right)^\tau}{\alpha} + \frac{2}{L_2 + \alpha} \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \Delta \right. \\ &\quad \left. \frac{\left( 1 - \frac{2}{L_2 + \alpha} \right)^\tau}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right] + L_2 \Delta. \end{aligned} \quad (27)$$

**Proof.** Using the triangle inequality, we have:

$$\begin{aligned} \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \right\| &\leq \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^j} - \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_{t,j}^\tau)}{\partial \mathbf{R}_t^j} \right\| \\ &\quad + \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_{t,j}^\tau)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \right\|. \end{aligned} \quad (28)$$

With the help of Assumptions 2 and 4, the following result holds:

$$\left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \right\| \leq L_2 \Delta + \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^j, \mathbf{H}_{t,j}^\tau)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \right\|. \quad (29)$$

Plugging in (14) into (29) yields (27). Thus, we complete the proof of Proposition 6.

Then, we restate some useful lemmas of [16] to deal with the proximal operator.

**Lemma 7** (Lemma 1 in [16]). *Let  $\mathbf{R}_t^\tau = \arg \min_{\mathbf{R}} \left\{ \left\langle \frac{\partial \mathcal{L}(\mathbf{R}_t^{\tau-1}, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^{\tau-1}}, \mathbf{R} \right\rangle + h(\mathbf{R}) + \frac{1}{\eta^\tau} \cdot \frac{1}{2} \|\mathbf{R} - \mathbf{R}_t^{\tau-1}\|^2 \right\}$*

and  $\tilde{\mathcal{G}}_t^{\tau-1} = \frac{1}{\eta'} (\mathbf{R}_t^{\tau-1} - \mathbf{R}_t^\tau)$ . For all  $\tau \geq 1$ , we have

$$\left\langle \frac{\partial \mathcal{L}(\mathbf{R}_t^{\tau-1}, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^{\tau-1}}, \tilde{\mathcal{G}}_t^{\tau-1} \right\rangle \geq \|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 + \frac{1}{\eta'} (h(\mathbf{R}_t^\tau) - h(\mathbf{R}_t^{\tau-1})). \quad (30)$$

**Lemma 8** (Lemma 2 in [16]). Define  $(\mathbf{R}_t^\tau)^+ = \arg \min_{\mathbf{R}} \left\{ \langle \nabla F(\mathbf{R}_t^{\tau-1}), \mathbf{R} \rangle + h(\mathbf{R}) + \frac{1}{\eta'} \cdot \frac{1}{2} \|\mathbf{R} - \mathbf{R}_t^{\tau-1}\|^2 \right\}$ , and let  $\mathcal{G}_t^{\tau-1} = \frac{1}{\eta'} (\mathbf{R}_t^{\tau-1} - (\mathbf{R}_t^\tau)^+)$ ,  $\tilde{\mathcal{G}}_t^{\tau-1} = \frac{1}{\eta'} (\mathbf{R}_t^{\tau-1} - \mathbf{R}_t^\tau)$ , we have

$$\|\mathcal{G}_t^{\tau-1} - \tilde{\mathcal{G}}_t^{\tau-1}\| \leq \|\nabla F(\mathbf{R}_t^{\tau-1}) - \frac{\partial L(\mathbf{R}_t^{\tau-1}, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^{\tau-1}}\|. \quad (31)$$

**Theorem 9** (Restatement of Theorem 1). Under Assumptions 1–5, define  $\alpha := -L + \mathcal{V}$ , and  $\mathcal{G}_t^j = \frac{1}{\eta'} (\mathbf{R}_t^j - \mathbf{R}_t^{j+1})$ ,  $0 \leq j \leq \tau-2$ , choose step size  $\eta$  to be  $\frac{2}{L_2 + \alpha}$ ,  $\eta'$  to be  $\frac{1}{6L_0}$ , and suppose  $\alpha < L_2$ ,  $h(\mathbf{R}_t^{\tau-1}) \leq h(\mathbf{R}_t^0)$ ,  $\forall t$ , we have:

$$\frac{1}{\tau T} \sum_{t=0}^{T-1} \sum_{j=1}^{\tau-1} \|\mathcal{G}_t^j\|^2 \leq \mathcal{O}\left(\frac{1}{T}\right). \quad (32)$$

**Proof.** According to the above Lemma 3, the function  $\nabla F(R)$  is  $L_0$ -Lipschitz. Let  $\tilde{\mathcal{G}}_t^\tau = \frac{1}{\eta'} (\mathbf{R}_t^{\tau-1} - \mathbf{R}_t^\tau)$ , we have:

$$\begin{aligned} F(\mathbf{R}_t^\tau) &\leq F(\mathbf{R}_t^{\tau-1}) + \langle \nabla F(\mathbf{R}_t^{\tau-1}), \mathbf{R}_t^\tau - \mathbf{R}_t^{\tau-1} \rangle \\ &+ \frac{L_0}{2} \|\mathbf{R}_t^\tau - \mathbf{R}_t^{\tau-1}\|^2 = F(\mathbf{R}_t^{\tau-1}) - \eta' \left\langle \frac{\partial \mathcal{L}(\mathbf{R}_t^{\tau-1}, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^{\tau-1}}, \tilde{\mathcal{G}}_t^{\tau-1} \right\rangle \\ &+ \eta' \left\langle \frac{\partial \mathcal{L}(\mathbf{R}_t^{\tau-1}, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^{\tau-1}} - \nabla F(\mathbf{R}_t^{\tau-1}), \tilde{\mathcal{G}}_t^{\tau-1} \right\rangle + \frac{\eta'^2 L_0}{2} \|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 \end{aligned} \quad (33)$$

$$\begin{aligned} &\leq F(\mathbf{R}_t^{\tau-1}) - \eta' \|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 - h(\mathbf{R}_t^\tau) + h(\mathbf{R}_t^{\tau-1}) \\ &+ \eta' \left\langle \frac{\partial \mathcal{L}(\mathbf{R}_t^{\tau-1}, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^{\tau-1}} - \nabla F(\mathbf{R}_t^{\tau-1}), \tilde{\mathcal{G}}_t^{\tau-1} \right\rangle \\ &+ \frac{\eta'^2 L_0}{2} \|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 \text{ (i)} \leq F(\mathbf{R}_t^{\tau-1}) + \left( \frac{\eta'^2 L_0}{2} - \frac{3\eta'}{4} \right) \|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 \\ &- h(\mathbf{R}_t^\tau) + h(\mathbf{R}_t^{\tau-1}) + \eta' \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^{\tau-1}, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^{\tau-1}} - \nabla F(\mathbf{R}_t^{\tau-1}) \right\| \text{ (ii)}, \end{aligned} \quad (34)$$

where (i) holds by the above Lemma 7, and (ii) holds by the Cauchy inequality and the basic inequality.

Let  $\mathcal{F}(\mathbf{R}) = F(\mathbf{R}) + h(\mathbf{R})$ , plugging the result in Proposition 6

into (31), we have:

$$\begin{aligned} \mathcal{F}(\mathbf{R}_t^\tau) &\leq F(\mathbf{R}_t^{\tau-1}) + h(\mathbf{R}_t^{\tau-1}) + \left( \frac{\eta'^2 L_0}{2} - \frac{3\eta'}{4} \right) \|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 \\ &+ \eta' \left\| \frac{\partial \mathcal{L}(\mathbf{R}_t^{\tau-1}, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^{\tau-1}} - \nabla F(\mathbf{R}_t^{\tau-1}) \right\| \\ &\leq F(\mathbf{R}_t^{\tau-1}) + h(\mathbf{R}_t^{\tau-1}) + \left( \frac{\eta'^2 L_0}{2} - \frac{3\eta'}{4} \right) \|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 \\ &+ \eta' (L_2 + \frac{L_2^2}{\alpha}) \left[ \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \Delta \right] \\ &+ \eta' L_1 \left[ \frac{L_2 \left( 1 - \frac{2}{L_2 + \alpha} \right)^\tau}{\alpha} + \frac{2}{L_2 + \alpha} \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \right. \\ &\quad \left. \Delta \frac{\left( 1 - \frac{2}{L_2 + \alpha} \cdot \alpha \right)^\tau}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right] + \eta' L_2 \Delta. \end{aligned} \quad (35)$$

According to Lemma 8, the difference between  $\tilde{\mathcal{G}}_t^{\tau-1}$  and  $\mathcal{G}_t^{\tau-1}$  are bounded, we have:

$$\begin{aligned} \|\mathcal{G}_t^{\tau-1}\|^2 &\leq 2\|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 + 2\|\tilde{\mathcal{G}}_t^{\tau-1} - \mathcal{G}_t^{\tau-1}\|^2 \\ &\leq 2\|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 + 2\left\| \frac{\partial L(\mathbf{R}_t^{\tau-1}, \mathbf{H}_t^\tau)}{\partial \mathbf{R}_t^{\tau-1}} - \nabla F(\mathbf{R}_t^{\tau-1}) \right\|^2 \\ &\leq 2\|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 + 4(L_2 + \frac{L_2^2}{\alpha})^2 \left[ \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \Delta \right]^2 \\ &+ 4L_1^2 \left[ L_2 \left( 1 - \frac{2}{L_2 + \alpha} \right)^\tau \frac{1}{\alpha} + \frac{2}{L_2 + \alpha} \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \right. \\ &\quad \left. \Delta \frac{\left( 1 - \frac{2}{L_2 + \alpha} \cdot \alpha \right)^\tau}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right]^2 + 4L_2^2 \Delta^2. \end{aligned} \quad (36)$$

Thus, we have:

$$\begin{aligned} -\|\mathcal{G}_t^{\tau-1}\|^2 &\leq -\frac{1}{2} \|\tilde{\mathcal{G}}_t^{\tau-1}\|^2 + 2(L_2 + \frac{L_2^2}{\alpha})^2 \left[ \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \Delta \right]^2 \\ &+ 2L_1^2 \left[ L_2 \left( 1 - \frac{2}{L_2 + \alpha} \right)^\tau \frac{1}{\alpha} + \frac{2}{L_2 + \alpha} \right. \\ &\quad \left. \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \Delta \frac{\left( 1 - \frac{2}{L_2 + \alpha} \cdot \alpha \right)^\tau}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right]^2 + 2L_2^2 \Delta^2. \end{aligned} \quad (37)$$

Plugging (37) into (35), we have:

$$\begin{aligned} \mathcal{F}(\mathbf{R}_t^\tau) &\leq F(\mathbf{R}_t^{\tau-1}) + h(\mathbf{R}_t^{\tau-1}) + \left( \frac{\eta'^2 L_0}{4} - \frac{3\eta'}{8} \right) \|\mathcal{G}_t^{\tau-1}\|^2 \\ &+ 2 \left( \frac{7\eta'}{4} - \frac{\eta'^2 L_0}{2} \right) (L_2 + \frac{L_2^2}{\alpha})^2 \left[ \left( \frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \Delta \right]^2 \\ &+ 2 \left( \frac{7\eta'}{4} - \frac{\eta'^2 L_0}{2} \right) L_1^2 \left[ \frac{L_2 \left( 1 - \frac{2}{L_2 + \alpha} \right)^\tau}{\alpha} \right. \\ &\quad \left. + \frac{2}{L_2 + \alpha} \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \Delta \frac{\left( 1 - \frac{2}{L_2 + \alpha} \cdot \alpha \right)^\tau}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right]^2 \\ &+ 2 \left( \frac{7\eta'}{4} - \frac{\eta'^2 L_0}{2} \right) L_2^2 \Delta^2. \end{aligned}$$

Based on the  $L_0$ -smoothness of  $F(\mathbf{R})$  established in Lemma 3, for

$j \geq 1$ , we have:

$$\begin{aligned}
F(\mathbf{R}_t^{j+1}) &\leq F(\mathbf{R}_t^j) - \left(\frac{\eta'}{2} - \eta'^2 L_0\right) \|\nabla F(\mathbf{R}_t^j)\|^2 \\
&\quad + \left\| \frac{\partial L(\mathbf{R}_t^j, \mathbf{H}_t^j)}{\partial \mathbf{R}_t^j} - \nabla F(\mathbf{R}_t^j) \right\|^2 \cdot \left(\frac{\eta'}{2} + \eta'^2 L_0\right) \\
&\leq F(\mathbf{R}_t^j) - \left(\frac{\eta'}{2} - \eta'^2 L_0\right) \|\nabla F(\mathbf{R}_t^j)\|^2 \\
&\quad + (\eta' + 2\eta'^2 L_0) \left(L_2 + \frac{L_2^2}{\alpha}\right)^2 \left(\frac{L_2 - \alpha}{L_2 + \alpha}\right)^{2\tau} \Delta^2 \\
&\quad + (\eta' + 2\eta'^2 L_0) L_1^2 \left[ \frac{L_2 \left(1 - \frac{2}{L_2 + \alpha}\right)^\tau}{\alpha} + \frac{2}{L_2 + \alpha} \right. \\
&\quad \left. \left( \frac{L_2 L_4}{\alpha} + L_3 \right) \Delta \frac{\left(1 - \frac{2}{L_2 + \alpha} \cdot \alpha\right)^\tau}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right]^2 \\
&\quad + (\eta' + 2\eta'^2 L_0) L_2^2 \Delta^2.
\end{aligned} \tag{38}$$

Then, with  $\eta' = \frac{1}{6L_0}$ , telescoping (38) over  $j$  from 0 to  $\tau - 2$  yields:

$$\begin{aligned}
F(\mathbf{R}_t^{\tau-1}) &\leq F(\mathbf{R}_t^0) - \frac{1}{18L_0} \sum_{j=0}^{\tau-2} \|\nabla F(\mathbf{R}_t^j)\|^2 + \frac{2}{9L_0} (\tau - 1) \\
&\quad \left\{ \left(L_2 + \frac{L_2^2}{\alpha}\right)^2 \left(\frac{L_2 - \alpha}{L_2 + \alpha}\right)^{2\tau} \Delta^2 + L_1^2 \right. \\
&\quad \left[ \frac{L_2 \left(1 - \frac{2}{L_2 + \alpha}\right)^\tau}{\alpha} + \frac{2}{L_2 + \alpha} \left(\frac{L_2 L_4}{\alpha} + L_3\right) \right. \\
&\quad \left. \left. \Delta \frac{\left(1 - \frac{2}{L_2 + \alpha} \alpha\right)^\tau}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right]^2 \right\} + \frac{2}{9L_0} (\tau - 2) L_2^2 \Delta^2.
\end{aligned} \tag{39}$$

Plugging (39) into (38), and  $h(\mathbf{R}_t^{\tau-1}) \leq h(\mathbf{R}_t^0)$  yields:

$$\begin{aligned}
\mathcal{F}(\mathbf{R}_t^\tau) &\leq F(\mathbf{R}_t^0) + h(\mathbf{R}_t^0) - \frac{1}{18L_0} \|\mathcal{G}_t^{\tau-1}\|^2 - \frac{1}{18L_0} \sum_{j=0}^{\tau-2} \|\mathcal{G}_t^j\|^2 \\
&\quad + \left[ \frac{5}{9L_0} + \frac{2}{9L_0} (\tau - 1) \right] \left\{ \left(L_2 + \frac{L_2^2}{\alpha}\right)^2 \left(\frac{L_2 - \alpha}{L_2 + \alpha}\right)^{2\tau} \Delta^2 \right. \\
&\quad \left. + L_1^2 \left[ \frac{L_2 \left(1 - \frac{2}{L_2 + \alpha}\right)^\tau}{\alpha} + \frac{2}{L_2 + \alpha} \left(\frac{L_2 L_4}{\alpha} + L_3\right) \Delta \right. \right. \\
&\quad \left. \left. \frac{\left(1 - \frac{2}{L_2 + \alpha} \alpha\right)^\tau}{1 - \frac{2}{L_2 + \alpha} \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right]^2 \right\} + \left( \frac{5}{9L_0} + \frac{2}{9L_0} (\tau - 2) \right) L_2^2 \Delta^2.
\end{aligned} \tag{40}$$

Telescoping (41) over  $t$  from 0 to  $T - 1$  yields:

$$\begin{aligned}
\frac{1}{\tau T} \sum_{t=0}^{T-1} \sum_{j=0}^{\tau-1} \|\mathcal{G}_t^j\|^2 &\leq 18L_0 \frac{\mathcal{F}(\mathbf{R}_0^0) - \inf_{\mathbf{R}} \mathcal{F}(\mathbf{R})}{\tau \cdot T} \\
&\quad + \frac{1}{\tau} \left[ \frac{5}{9L_0} + \frac{2}{9L_0} (\tau - 1) \right] \left\{ \left(L_2 + \frac{L_2^2}{\alpha}\right)^2 \left(\frac{L_2 - \alpha}{L_2 + \alpha}\right)^{2\tau} \Delta^2 \right. \\
&\quad \left. + L_1^2 \left[ \frac{L_2 \left(1 - \frac{2}{L_2 + \alpha}\right)^\tau}{\alpha} + \frac{2}{L_2 + \alpha} \left(\frac{L_2 L_4}{\alpha} + L_3\right) \Delta \right. \right. \\
&\quad \left. \left. \frac{\left(1 - \frac{2}{L_2 + \alpha} \cdot \alpha\right)^\tau}{1 - \frac{2}{L_2 + \alpha} \cdot \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right]^2 \right\} + \frac{1}{\tau} \left[ \frac{5}{9L_0} + \frac{2}{9L_0} (\tau - 2) \right] L_2^2 \Delta^2.
\end{aligned} \tag{41}$$

$$\tag{42}$$

$$\tag{43}$$

Thus, we complete the proof of Theorem 1.