

BoMM: MULTI-MODALITY LARGE-SMALL MODEL BIDIRECTIONAL COLLABORATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Different from existing single-modality large-small model collaborations, multi-modality large-small model collaboration is an under-explored paradigm where cloud-side multi-modality large model (MM-LM) collaborates with parties' small models (SMs) to achieve bidirectional domain-specific performance improvements. Nevertheless, this paradigm faces two key challenges. First, MM-LM inherently relies on abundant modality-aligned samples for training, but geographical and device diversity across parties inevitably lead to different collected samples and modalities. These differences significantly reduce overlapping sample entities across parties' multi-modality datasets, creating **modality alignment scarcity** challenge. Second, collected device failure and human annotation costs further lead to different modality missing problems in each party's dataset. Existing modality completion methods typically require enough modality-completed training samples to ensure generation quality, creating a **modality completeness gap** challenge. To address these challenges, we propose a multi-modality large-small model **bidirectional collaboration** framework, named **BoMM**, which consists of two key components. Specifically, *global prototype-guided alignment* strategy identifies potentially aligned samples through similarity distribution comparisons between unaligned data and established global prototypes, enabling knowledge transfer from SMs to MM-LM. With established prototypes, *preference-driven modality adaptive completion* method integrates direct preference optimization into generator training with real-time scheduling to dynamically complete missing modalities, enabling knowledge transfer from MM-LM to SMs. Theoretical analysis confirms BoMM's $O(1/\sqrt{T})$ convergence rate. Across three multi-modality scenarios, it outperforms state-of-the-art methods by up to 6.64% on two well-known datasets. Our code is available at <https://anonymous.4open.science/r/MultiLM-5D65>.

1 INTRODUCTION

In large language model (LLM) and small models (SMs) collaborative scenarios, cloud-side LLM leverages parties' domain data and SMs' expertise to address domain-specific weaknesses, while SMs of parties can also benefit from LLM' superior reasoning and generation capabilities Wang et al. (2024a); Liu et al. (2024). This collaborative relationship extends to federated learning scenarios Guo et al. (2024) when privacy becomes essential Zhang et al. (2025), enabling mutual enhancement while keeping all raw data local to satisfy privacy requirements.

Existing large-small model collaborations utilize three primary approaches: (i) distillation-based methods Zhou et al. (2021); Liu et al. (2021); Pham et al. (2021) leverage knowledge distillation to transfer knowledge between domain-specific SMs and LLM in both directions; (ii) generation-based methods Zhang et al. (2021); Li & Jin (2022); Nasser et al. (2024) leverage synthetic datasets to transfer knowledge from SMs to LLMs or use LLMs' generative capabilities to create task-specific datasets; and (iii) parameter-based methods Fan et al. (2024); Cheng et al. (2021); Yu et al. (2023) employ parameter-efficient fine-tuning (PEFT) techniques Han et al. (2024) to transmit knowledge bidirectionally or selectively transfer parametric knowledge from LLMs to SMs.

Despite significant achievements, these methods focus on single-modality large-small model collaboration, neglecting multi-modality scenarios. In reality, with multi-modality data becoming more

prevalent, multi-modality large model (MM-LM) are becoming an increasingly important research direction Tong et al. (2025). Since MM-LM often exhibits suboptimal performance on domain-specific tasks Ye et al. (2024), they require access to multi-modality domain data from parties. This necessitates collaboration between multi-modality large and small models. Although a recent study Chen et al. (2025) aligns representations between single-modality MM-LM and multi-modality SMs in medicine, it relies heavily on two idealistic assumptions: modality alignment and modality completeness for each sample. These rarely hold in real-world distributed settings. Additionally, this work also overlooks the need for mutual performance improvements in multi-modality collaborations, where the limitation of mutual performance improvement is also observed in single-modality collaborations.

Therefore, we explore an under-explored paradigm: multi-modality large-small model bidirectional collaboration, as shown in Fig. 1. Different from existing large-small model collaborations, due to geographical and device diversity across parties, this paradigm faces a **modality alignment scarcity** challenge where parties possess multi-modal data but lack aligned samples across modalities. These aligned samples are essential for training MM-LM’s alignment capability Wu et al. (2023); Song et al. (2023). Moreover, we identify an additional critical challenge: **modality completeness gap**, which occurs when parties have incomplete or missing modalities due to device failures or annotation expenses. Existing completion approaches Wu et al. (2024a) still struggle to generate high-quality modalities with severe or even complete modality missing in training datasets.

To address these challenges, we propose a *multi-modality large-small model bidirectional collaboration framework*, denoted as **BoMM**, where a *global prototype-guided alignment* method identifies potentially aligned samples by comparing similarity distributions to established global prototypes, enabling knowledge transfer from SMs to MM-LM. Building on these prototypes, *preference-driven modality adaptive completion* method optimizes the global generator to dynamically fill local missing modalities with a real-time sample scheduler, facilitating effective knowledge transfer from MM-LM to SMs. The main contributions are summarized as follows:

- (a) We present an under-explored paradigm: **multi-modality large-small model bidirectional collaboration**, and propose BoMM, where cloud-side MM-LM enhances domain performance by SMs, while parties’ SMs overcome modality missing using MM-LM’s generation capability.
- (b) We mitigate modality alignment scarcity challenge through a global prototype-guided alignment method, and address modality completeness gap challenge via a preference-driven modality adaptive completion method, where a real-time sample scheduler is designed to filter generated modalities, enabling effective bidirectional knowledge transfer between models.
- (c) BoMM can effectively handle different modality alignment situations, while supporting various modality missing rates of parties. Our theoretical analysis guarantees convergence at $O(1/\sqrt{T})$ rate, ensuring robust performance in multi-modality large-small model collaborations.
- (d) Extensive experiments on two real-world multi-modal datasets across three scenarios clearly demonstrate that BoMM achieves substantial improvements over state-of-the-art methods under varying rates of aligned samples and modality missing.

2 RELATED WORK

2.1 LARGE-SMALL MODEL COLLABORATION

large-small model collaboration leverages distributed private domain data to address public data scarcity and domain-specific poor performance challenges in large language models (LLMs). Con-

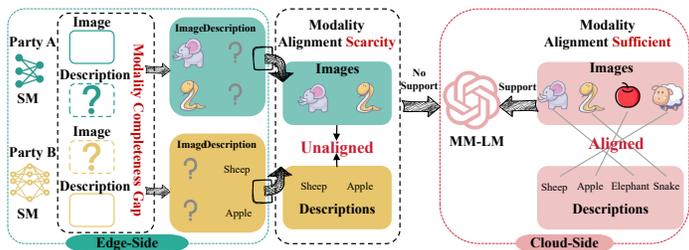


Figure 1: Multi-modality large-small model collaboration face to key challenges: modality completeness gaps from parties’ modality missing, and modality alignment scarcity from parties’ modality alignment problems for the same sample, hindering MM-LM training that relies on abundant modality alignments.

currently, the extensive knowledge within LLMs can improve complex reasoning capabilities and provide comprehensive background information for domain-specific small models (SMs). Three primary methods facilitate knowledge transfer between SMs and LLMs Liu et al. (2025): distillation-based Zhou et al. (2021); Liu et al. (2021); Pham et al. (2021), generation-based Zhang et al. (2021); Li & Jin (2022); Nasser et al. (2024), and parameter-based transfer methods Fan et al. (2024); Cheng et al. (2021); Yu et al. (2023). These approaches primarily concentrate on unidirectional knowledge transfer between large-small models, with few works Fan et al. (2024); Cheng et al. (2021) exploring bidirectional knowledge exchange. More importantly, existing works fail to cope with collaborative scenarios for a multi-modality large model (MM-LM) and SMs, where parties may possess heterogeneous data modalities, diverse model architectures, and face varying degrees of modality incompleteness in their local datasets.

2.2 LARGE MULTI-MODALITY MODEL

The rapid advancement of multi-modality large model has significantly demonstrated their remarkable ability to tackle various multi-modality challenges in real-world applications, including personal AI assistants Wang et al. (2024b), education Lee et al. (2025), medicine Wang et al. (2024c), autonomous driving Cui et al. (2024), and others. However, these works depend on training data where multiple modalities (such as images and text) are paired for each sample, a condition rarely met in distributed environments. In these distributed environments, each party typically holds only specific modalities for their local samples, while complementary modalities for these exact samples are often entirely absent from other parties’ collections. This modality misalignment generally emerges from regional differences, domain specializations, and other organizational factors, making complete multi-modality pairing exceedingly difficult to achieve. Leveraging such unaligned multi-modality data for MM-LM training can lead to weakened multi-modality understanding, increased hallucinations, and unreliable correlations between different modalities, ultimately limiting the model’s practical utility in complex real-world scenarios Song et al. (2023).

3 PRELIMINARY

3.1 MULTI-MODALITY LARGE-SMALL MODEL COLLABORATION

In our multi-modality large-small model collaboration scenarios with K parties and modality set M , each party $k \in \{1, 2, \dots, K\}$ has dataset D_k containing N_k samples and $M_k \subseteq M$ modalities. Parties may have identical or different modalities due to modality missing, with different sample alignment ratios. Within dataset D_k , each party has labeled samples $(x_i^l, y_i^l) \in D_k^l$ and unlabeled samples $(x_i^{ul}, y_i^{ul}) \in D_k^{ul}$. These data include aligned samples $(x_i^a, y_i^a) \in D_k^a$ (same sample ID across K parties) and unaligned samples $(x_i^{ua}, y_i^{ua}) \in D_k^{ua}$ (sample ID found only in party k).

3.2 OPTIMIZATION OBJECTIVE

Resource limitations lead parties to deploy lightweight small models (SMs) locally while collaboratively training a powerful domain-specialized multi-modality large model (MM-LM) on the cloud. Our framework establishes bi-directional knowledge transfer, simultaneously incorporating domain-specific expertise from SMs into MM-LM while enabling MM-LM to enhance SMs’ performance. Formally, we denote MM-LM with parameters ψ as f_ψ and SM of party k with parameters θ_k as f_{θ_k} . Our bi-directional optimization objective includes:

Global MM-LM: The cloud-side multi-modality large model f_ψ is optimized using parties’ multi-modality data for both domain-specific tasks and modality completion:

$$\min_{\psi} \mathcal{L}_{\text{global}}(\psi; \{D_k\}_{k=1}^K, \{\theta_k\}_{k=1}^K), \quad (1)$$

where $\psi = [\psi_1, \psi_2]$, with ψ_1 and ψ_2 representing classifier and generator parameters, respectively.

Local SMs: Each party’s small model f_{θ_k} is optimized using its local dataset D_k for domain-specific tasks, improved by enhanced domain data from MM-LM:

$$\min_{\theta_k} \mathcal{L}_{\text{local}}(\theta_k; D_k, \psi), \quad (2)$$

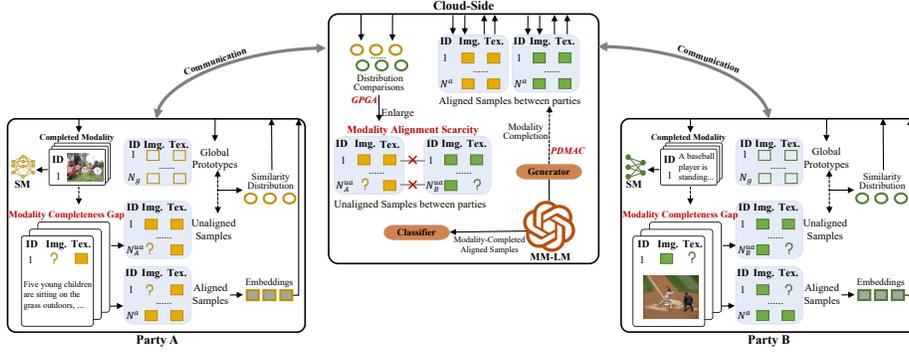


Figure 2: The overview of BoMM framework. The global prototype-guided alignment (GPGA) method identifies essential modality-aligned samples by analyzing similarity distributions between unaligned data and global prototypes, enhancing MM-LM’s domain performance. Meanwhile, parties leverage a preference-driven modality adaptive completion (PDMAC) method to synthesize missing modality using MM-LM’s generation ability, improving local domain performance.

where \mathcal{L}_{local} denotes the party-specific optimization loss for each local model, and \mathcal{L}_{global} represents the optimization loss for global MM-LM in the cloud.

4 OUR METHODOLOGY

In this section, we propose BoMM, a multi-modality large-small model bidirectional collaboration framework including a preference-driven modality adaptive completion method with a real-time sample scheduler and a global prototype-guided alignment strategy. The overall architecture is illustrated in Fig. 2.

4.1 PREFERENCE-DRIVEN MODALITY ADAPTIVE COMPLETION

The modality incompleteness of each party degrades local model performance on private domain tasks while simultaneously limiting knowledge contribution to cloud-side MM-LM training. Resource constraints prevent parties from training domain-specific MM-LM locally for modality completion. While general MM-LM is accessible via APIs, they generate generic outputs that fail to synthesize missing modalities for specialized domains accurately.

4.1.1 MODALITY ADAPTIVE COMPLETION

Inspired by direct preference optimization (DPO) Rafailov et al. (2023), which guides models to generate human-preferred outputs, we design a preference-driven modality adaptive completion (PDMAC), enabling MM-LM to generate domain-specific modality information tailored to local private data. Unlike existing DPO approaches Wu et al. (2024b); Karthik et al. (2024) for modality completion that require extensive modality-completed training datasets, our method effectively handles different modality incompleteness or entire modality missing. By leveraging MM-LM for modality completion, we substantially improve model performance on domain-specific tasks.

To be specific, for i -th sample with existing modality information $x_{i,k}^{\{v,t\}}$, we create a preference pair (g_i^p, g_i^q) . g_i^p is the preferred output, which is obtained from the closest global prototype o_i to the i -th sample, while g_i^q is the non-preferred output, which is the missing modality generated by the MM-LM’s generator \mathcal{G} . The optimization objective for \mathcal{G} is:

$$\mathcal{L}_{prefer} = -\mathbb{E}_{(x_{i,k}^{\{v,t\}}, g_i^p, g_i^q) \sim D_k} [\log \sigma(\beta(R(x_{i,k}^{\{v,t\}}, g_i^p)) - R(x_{i,k}^{\{v,t\}}, g_i^q))], \quad (3)$$

where $\{v, t\}$ indicates available image (v), text (t), or both modalities. σ is the logistic function, and β scales the reward function $R(\cdot)$. Since parties may have different missing modalities, $R(\cdot)$ is adapted based on the generated modality type, following Jiang et al. (2024) for text and Karthik et al. (2024) for images.

216 Additionally, to prevent generated modality from lacking personalization and becoming too similar
 217 to its global prototype, we introduce an additional loss term:

$$218 \mathcal{L}_{\text{special}} = \|g_i^q - g_i^p\|^2 (1 - \text{sim}(x_{i,k}^{\{v,t\}}, o_i^{\{v,t\}})), \quad (4)$$

219 where o_i and $x_{i,k}^{v,t}$ share modalities i -th sample already has, while g_i^q and g_i^p represent the target
 220 modality to generate. Our method needs just one complete-modality sample per class as a prototype,
 221 which is practical to obtain. If unavailable initially, existing generation methods Kwon et al. (2024);
 222 Laina et al. (2019) still can initialize these prototypes. Once established, samples grouped around
 223 these prototypes progressively receive high-quality completion results during training. For samples
 224 already obtaining the target modality, these real data serves directly as preferred output g_i^p . The
 225 generator \mathcal{G} is updated by $\mathcal{L}_{\mathcal{G}} = \mathcal{L}_{\text{prefer}} + \gamma \mathcal{L}_{\text{special}}$, where γ is a balancing hyperparameter. The
 226 MM-LM classification task loss term $\mathcal{L}_{\text{task}}$ is calculated by the standard cross-entropy loss Zhang
 227 & Sabuncu (2018) with completed modality data.

228 4.1.2 DYNAMIC SAMPLE SCHEDULER

229 Parties use the cloud-side MM-LM to generate missing modalities and expand local datasets. To
 230 handle potentially noisy generated data, we implement a scheduler that dynamic select high-reward
 231 samples in each batch based on sample reward. Low-quality samples are kept in the training set,
 232 with their current generations continuing to serve as non-preferred outputs for further optimization.

233 **Reward definition.** Since task loss can be applied to measure sample value Huang et al. (2025),
 234 we design a reward function that dynamically evaluates generator improvement after each batch
 235 update. Our reward function $r(\cdot)$ combines sample-level and batch-level assessments, with batch-
 236 level reward measuring the reduction in average generated loss before and after training. In practice,
 237 we use batch-average generated loss as an approximation of the expected generated loss across the
 238 entire data distribution. For the i -th sample in batch b , the reward is:

$$239 r(a_k^{i,b}, \theta_{C_k}^{b-1}, \theta_{C_k}^b) = \underbrace{-\log \mathcal{P}_{\theta_{C_k}^{b-1}}(y_k^i | \bar{x}_k^{i,b}) + \delta \mathcal{H}}_{\text{sample-level}} + \frac{\overbrace{\sum_{i=1}^{N_b} e^{\mathcal{L}_{\mathcal{G}}(a_k^{i,b}; \theta_{C_k}^{b-1})}}^A - \overbrace{\sum_{i=1}^{N_b} e^{\mathcal{L}_{\mathcal{G}}(a_k^{i,b+1}; \theta_{C_k}^b)}}^Q}{\underbrace{\max\left(\sum_{i=1}^{N_b} e^{\mathcal{L}_{\mathcal{G}}(a_k^{i,b}; \theta_{C_k}^{b-1})}, \sum_{i=1}^{N_b} e^{\mathcal{L}_{\mathcal{G}}(a_k^{i,b+1}; \theta_{C_k}^b)}\right)}_{\text{batch-level}}}, \quad (5)$$

240 where $a_k^{i,b} = (\bar{x}_k^{i,b}, y_k^i)$. $\theta_{C_k}^{b-1}$ represents classifier parameters in batch b , and $\mathcal{P}_{\theta_{C_k}^{b-1}}$ is the model's
 241 output probability distribution given input $a_k^{i,b}$. N_b is the batch size. Term A evaluates $\theta_{C_g}^{b-1}$ on batch
 242 $b-1$, while term Q evaluates $\theta_{C_g}^{b+1}$ on batch $b+1$. $\bar{x}_k^{i,b}$ combines intermediate representations across
 243 all layers to utilize the model's full representational capacity.

244 **Scheduler design.** The rewards of generated samples serve as supervisory signals for training the
 245 dynamic scheduler \mathcal{R} . This scheduling process addresses the exploration-exploitation trade-off
 246 Berger-Tal et al. (2014), requiring a balance between discovering new sample spaces and utiliz-
 247 ing known high-quality examples. We implement two networks after an encoder layer: \mathcal{R}^p for
 248 exploitation and \mathcal{R}^q for exploration.

249 Specifically, the exploit network \mathcal{R}^p predicts sample rewards by mapping inputs to observed re-
 250 wards, while the explore network \mathcal{R}^q estimates uncertainty and adds an exploration bonus. This de-
 251 sign balances exploitation and exploration during sample selection, following UCB Wen et al. (2015)
 252 and Thompson Sampling principles Zhang et al. (2020). For modality incomplete data $\bar{x}_k^{i,b} \in D_k^{\text{inc}}$
 253 in b -th batch, \mathcal{R}^p feedforward neural network with residual connections, denoted by $\mathcal{R}^p(\bar{x}_k^{i,b}; \theta_p^b)$.
 254 After obtaining the observed reward $r(a_k^{i,b}, \theta_p^{b-1}, \theta_p^b)$, the parameters θ_p^b are updated by:

$$255 \mathcal{L}^p(D_{k,b}^{\text{inc}}, \theta_p^b) = \frac{1}{2N_b} \sum_{i=1}^{N_b} [\mathcal{R}^p(\bar{x}_k^{i,b}; \theta_p^b) - r(a_k^{i,b}, \theta_p^{b-1}, \theta_p^b)]^2. \quad (6)$$

256 The explore network \mathcal{R}^q takes input $h_{i,b,k}^p$ created by concatenating intermediate hidden states of
 257 $\mathcal{R}^p(\bar{x}_k^{i,b}; \theta_p^{b-1})$ with the last dimension. This allows \mathcal{R}^q to consider the exploit network's internal
 258

states. Like \mathcal{R}^p , \mathcal{R}^q is a feedforward neural network with residual connections. Its parameters θ_q^b are updated by:

$$\mathcal{L}^q(D_{k,b}^{\text{inc}}, \theta_q^b) = \frac{1}{2N_b} \sum_{i=1}^{N_b} [\mathcal{R}^q(h_{i,b,k}^p; \theta_q^b) - (r(a_k^{i,b}, \theta_p^{b-1}, \theta_p^b) - \mathcal{R}^p(\bar{x}_k^{i,b}; \theta_p^b))]^2. \quad (7)$$

The final reward estimation combines both networks:

$$\hat{r}(\bar{x}_k^{i,b}; \theta_p^b, \theta_q^b) = \mathcal{R}^p(\bar{x}_k^{i,b}; \theta_p^b) + \lambda \mathcal{R}^q(h_{i,b,k}^p; \theta_q^b), \quad (8)$$

where λ controls exploration strength. Samples with predicted rewards $\hat{r} > \mu$ are selected as high-value samples for the current batch, which μ is a hyperparameter.

4.2 GLOBAL PROTOTYPE-GUIDED ALIGNMENT

Finding aligned data across parties is difficult due to geographic and domain differences. This creates a significant challenge for MM-LM training, which requires aligned multi-modality data. Current alignment techniques Feng (2022); Feng et al. (2022) typically need some pre-existing aligned data and can't handle scenarios with variable alignment quantities, especially with zero-aligned samples. While Guo et al. (2025) addresses zero-aligned scenarios, it assumes uniform modality distribution across parties and complete modalities within samples, failing in multi-modality collaboration settings. To address these limitations, we propose a global prototype-guided alignment (GPGA) strategy that works with any number of pre-existing aligned data, including zero. We use optimal transport (OT) distance to derive relationship vectors between unaligned samples and global prototypes, then identify new aligned data by analyzing these relationship vectors.

Specifically, even with zero pre-existing aligned data, we start by using local labeled data D_k^l to create class prototypes $\{p_k^c\}_{c=1}^C$ through weighted averaging Bien & Tibshirani (2011). We pair same-class prototypes across parties as initial aligned data to seed the global prototype set $\mathcal{O} = \{o^m | m = 1, \dots, N_{\text{global}}\}$. Using optimal transport Peyré et al. (2019), we calculate transport cost $M_{im}^k = 1 - \cos(x_i^{ua}, o^m)$ between unaligned samples and global prototypes. For each unaligned sample, we create a cost vector \mathbf{M}_i^k representing its relationship to all global prototypes. To find corresponding samples across parties, we compute:

$$j^* = \operatorname{argmax}_{j \in N_w^{ua}} \operatorname{sim}(x_k^{ua,i}, x_w^{ua,j}) = \operatorname{argmax}_{j \in N_w^{ua}} \frac{\mathbf{M}_i^k \cdot \mathbf{M}_j^w}{\|\mathbf{M}_i^k\|_2 \cdot \|\mathbf{M}_j^w\|_2}, \quad (9)$$

where N_w^{ua} is the size of the unaligned dataset D_w^{ua} in party w . This identifies cross-party alignment by comparing relationship patterns to shared prototypes rather than directly comparing incompatible modalities. The newly identified aligned sample pairs join the global prototype set until reaching N_{global} , after which these prototypes guide alignment of remaining samples. These enlarged aligned samples are utilized to improve the domain performance of local SM $\mathcal{L}_{\text{task}}^k$ by standard cross-entropy loss Zhang & Sabuncu (2018). These components support each other cyclically: alignment relies on completion for consistent representations, while completion depends on the global prototypes from alignment.

4.3 COMPLEXITY AND CONVERGENCE ANALYSIS

We analyze the convergence of our proposed BoMM. Let $\Psi := \{\psi\}$ and $\Theta := \{\theta_1, \theta_2, \dots, \theta_{K-1}, \theta_K\}$. The optimization problems in Eq. 1 and Eq. 2 thus form a bilevel optimization problem, which can be formulated as follows.

$$\begin{aligned} \min_{\Psi} F(\Psi) &:= \mathcal{L}_{\text{global}}(\Psi, \Theta^*(\Psi)) = \mathbb{E}_{\xi} [L_{\text{global}}(\Psi, \Theta^*(\Psi); \xi)] = \frac{1}{n} \sum_{i=1}^n L_{\text{global}}(\Psi, \Theta^*(\Psi); \xi_i), \\ \text{s.t. } \Theta^*(\Psi) &= \operatorname{arg min}_{\Theta} \mathcal{L}_{\text{local}}(\Psi, \Theta) = \mathbb{E}_{\zeta} [L_{\text{local}}(\Psi, \Theta; \zeta)] = \frac{1}{n} \sum_{i=1}^n L_{\text{local}}(\Psi, \Theta; \zeta_i), \end{aligned} \quad (10)$$

where $\mathcal{L}_{\text{global}}$ contains $\mathcal{L}_{\mathcal{G}}$ and $\mathcal{L}_{\text{task}}$. $\mathcal{L}_{\text{local}}$ consists only of $\mathcal{L}_{\text{task}}^k$. L_{global} and L_{local} are loss functions per sample for cloud and parties, where $L_{\text{local}}(\Psi, \Theta; \zeta) := l_{\text{local}}(\Psi, \Theta; \zeta) + \frac{\varphi}{2} \|\Theta\|_F^2$, with ζ representing a sample from K parties. To analyze convergence, we introduce the following definition.

Definition 1. A point $\hat{\Psi}$ is called an ϵ -accurate stationary point for the objective function $F(\Psi)$ if $\mathbb{E}\|\nabla F(\hat{\Psi})\|^2 \leq \epsilon$.

Proof Sketch. To explore essential insights, we first bound the tracking error $\|\Theta_t^{j-1} - \Theta^*(\Psi_t^0)\|$ between local parameters and optimal parameters at the j -th epoch. Then, using virtual updates technology (Yang et al., 2022), we establish an upper bound for the gradient approximation error $\left\| \frac{\partial \mathcal{L}_{\text{local}}(\Psi_t^j, \Theta_t^j)}{\partial \Psi_t^j} - \nabla F(\Psi_t^j) \right\|$.

Theorem 1. Under Assumptions 1-4 (detailed in Appendix B.1), define $\alpha := -L + \varphi$, choose step size η to be $\frac{2}{L_2 + \alpha}$, $N_b = \mathcal{O}(\frac{1}{\sqrt{\epsilon}})$, $\tau = \mathcal{O}(\log \frac{1}{\epsilon})$, $\eta' < \frac{1}{2L_0}$ and suppose $\alpha < L_2$, we have:

$$\frac{1}{\tau' T} \sum_{t=0}^{T-1} \sum_{j=0}^{\tau'-1} \|\nabla F(\Psi_t^j)\|^2 \leq \frac{F(\Psi_0^0) - \inf_{\Theta} F(\Psi)}{\tau' T (\frac{\eta'}{2} - L_0 \eta'^2)} + (\eta' + 2\eta'^2 L_0) L_2^2 \Delta^2 \frac{\tau' - 1}{\tau'} + \mathcal{O}(\epsilon), \quad (11)$$

where $\mathcal{O}(\cdot)$ indicates growth rate proportional to or slower than ϵ . With $L_0 := L_2 + \frac{2L_2^2 + L_1^2 L_3}{\alpha} + \frac{L_1 L_2 L_3 + L_1 L_2 L_4 + L_2^3}{\alpha^2} + \frac{L_1 L_2^2 L_4}{\alpha^3}$, setting $\tau' = 1$, yields a convergence rate of $\mathcal{O}(\frac{1}{\sqrt{T}})$.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Settings: We evaluate our proposed BoMM framework on two widely-used multi-modal datasets: MS COCO Lin et al. (2014) and CUB-200-2011 Wah et al. (2011). Specifically, we conduct experiments in a large-small collaboration scenario with heterogeneous multi-modality models, including two parties' SMs and one cloud-side MM-LM. Each party possesses both aligned and unaligned samples, where aligned samples share identical IDs but contain different textual descriptions, while unaligned samples are exclusively owned by individual parties. We evaluate three modality missing scenarios. In **scenario 1**, both parties and MM-LM suffer from text modality missing with probability Δ . In **scenario 2**, party 0 maintains complete image and text modalities while party 1 experiences text modality missing with probability Δ for both aligned and unaligned samples. In **scenario 3**, both parties and MM-LM suffer from image modality missing with probability Δ . Details of experimental datasets, settings and implementation details are provided in Appendix A.

Baselines: We compare our BoMM framework against five baseline methods: **1) LOCAL**, which trains local models using only local data without cloud-side collaboration; **2) MM-LOCAL**, which employs local multi-modal models to complete missing modalities before local training; **3) Vanilla** aggregates aligned data representations at the cloud-side MM-LM while parties train on their local data; **4) MM-AC**, which extends Vanilla by using the MM-LM to complete missing modalities in aligned data; and **5) MM-UC**, that further incorporates modality completion for both aligned and unaligned data at the cloud side.

5.2 OVERALL PERFORMANCE

Tab. 1 presents the experimental results under both scenario 1 and scenario 2 with modality missing probability Δ set to 0.5. Our findings reveal that text modality completion provides limited performance gains for party's models, as evidenced by comparing MM-AC and MM-UC results where the inclusion of unaligned data completion shows minimal improvement. While this factor also constrains the improvement potential of our BoMM framework on local performance, our method still achieves consistent gains. Most notably, BoMM demonstrates substantial improvements on the cloud-side MM-LM performance. Specifically, under scenario 1 with zero aligned samples, BoMM achieves 4.75% improvement over the best baseline on MS COCO. Under scenario 2, the cloud-side improvement reaches 6.64% on MS COCO, highlighting the effectiveness of our bidirectional collaboration framework in enhancing domain-specific multi-modality large model capabilities.

Moreover, as shown in Tab. 2, scenario 3 demonstrates that image modality completion yields significantly greater performance improvements compared to text completion scenarios. This is evidenced by the substantial gains observed when comparing MM-UC against Vanilla, where Party 1's accuracy dramatically increases from 38.96% to 50.08%, representing an 11.12% improvement.

Table 1: Model performance comparison of BoMM under scenario 1 and scenario 2. P1 and P2 represent party 1 and party 2’s performance, while C represents the cloud-side MM-LM performance.

Aligned Number	Method	Scenario 1						Scenario 2		
		MS COCO			CUB-200-2011			MS COCO		
		P1	P2	C	P1	P2	C	P1	P2	C
#0	LOCAL	38.67	38.28	-	56.34	29.98	-	50.48	34.96	-
	MM-LOCAL	36.13	36.52	-	57.62	31.84	-	50.68	36.52	-
	Vanilla	38.88	37.11	27.19	60.44	41.60	23.34	51.26	35.74	24.80
	MM-AC	39.06	37.40	27.64	59.45	46.38	23.27	50.20	35.84	26.26
	MM-UC	38.18	36.33	26.27	58.96	46.62	22.65	50.68	33.98	25.78
	BoMM	40.17	39.78	32.39	61.43	47.77	24.12	52.75	33.69	32.90
	Imp. (%)	1.99	2.38	4.75	0.99	1.15	0.78	0.49	-2.15	6.64
#200	LOCAL	39.47	37.41	-	60.85	37.25	-	50.98	37.41	-
	MM-LOCAL	39.14	36.10	-	60.36	35.58	-	51.97	36.10	-
	Vanilla	39.06	38.82	28.61	60.28	47.37	22.70	52.79	36.75	26.15
	MM-AC	38.15	39.80	28.45	60.03	43.25	23.11	51.94	36.84	25.65
	MM-UC	37.17	38.81	27.80	60.36	47.86	23.02	52.78	34.53	25.25
	BoMM	41.22	40.41	35.21	61.02	48.19	24.10	54.32	37.89	30.10
	Imp. (%)	1.75	0.61	6.60	0.17	0.82	0.99	1.53	1.05	3.95

Building upon this foundation, our proposed BoMM framework achieves even more pronounced improvements on local performance, with Party 1 and Party 2 achieving additional 2.27% and 2.28% accuracy gains respectively over the strongest baseline MM-UC. The cloud-side MM-LM also benefits significantly from this collaborative approach, improving from 44.90% to 51.18%, demonstrating a 6.28% enhancement. Furthermore, we evaluate the quality of generated missing modalities using CLIPScore metrics Hessel et al. (2021), where BoMM consistently outperforms all baselines with scores of 68.86 and 70.24 for both parties, indicating superior generation quality that contributes to the overall model performance improvements.

5.3 ABLATION STUDY

We conduct ablation studies on scenario 2 with zero aligned samples to validate the contribution of key components in our BoMM framework. Specifically, removing the preference loss corresponds to setting $\mathcal{L}_{\text{prefer}}$ to 0, while removing the special loss involves setting the balancing hyperparameter γ to zero. For the scheduler ablation, we set the reward threshold μ to 0, which disables the sample filtering mechanism and includes all generated samples in training. The results in Tab. 3 demonstrate that each component contributes meaningfully to the overall performance of BoMM. Removing these components both lead to consistent performance degradation across both parties and the cloud-side MM-LM.

5.4 CONVERGENCE EXPERIMENTS

Fig. 3(A) demonstrates the convergence behavior of our proposed BoMM framework compared with baseline methods across 5 training rounds. We evaluate convergence performance under the challenging zero-aligned samples scenario on both MS COCO and CUB-200-2011 datasets. Notably, BoMM exhibits faster convergence and reaches higher final accuracy compared to baseline methods, with both local parties and the MM-LM demonstrating steady performance improvements throughout the training process.

5.5 HYPERPARAMETER ANALYSIS

Fig. 3(B) presents the sensitivity analysis of two key hyperparameters in our BoMM framework. In Fig. 3(i), we evaluate the impact of modality missing rate Δ on performance across both datasets.

Table 2: Model performance under scenario 3.

Method	MS COCO				
	P1		P2		C
	Acc	CLIPScore	Acc	CLIPScore	Acc
LOCAL	36.13	-	34.47	-	-
MM-LOCAL	30.18	-	36.62	-	-
Vanilla	38.96	-	38.18	-	27.83
MM-AC	40.79	64.12	39.88	65.67	28.04
MM-UC	50.08	62.16	39.72	67.56	44.90
BoMM	52.35	68.86	42.16	70.24	51.18
Imp. (%)	2.27	4.74	2.28	2.68	6.28

Table 3: Ablation study results.

Model	MS COCO			CUB-200-2011		
	P1	P2	C	P1	P2	C
BoMM	40.17	39.78	32.39	61.43	47.77	24.12
w/o preference loss	37.51	38.13	30.82	60.48	45.18	20.34
w/o special loss	38.47	35.35	30.66	60.05	40.23	7.23
w/o scheduler	38.37	35.15	30.56	59.96	41.50	11.13

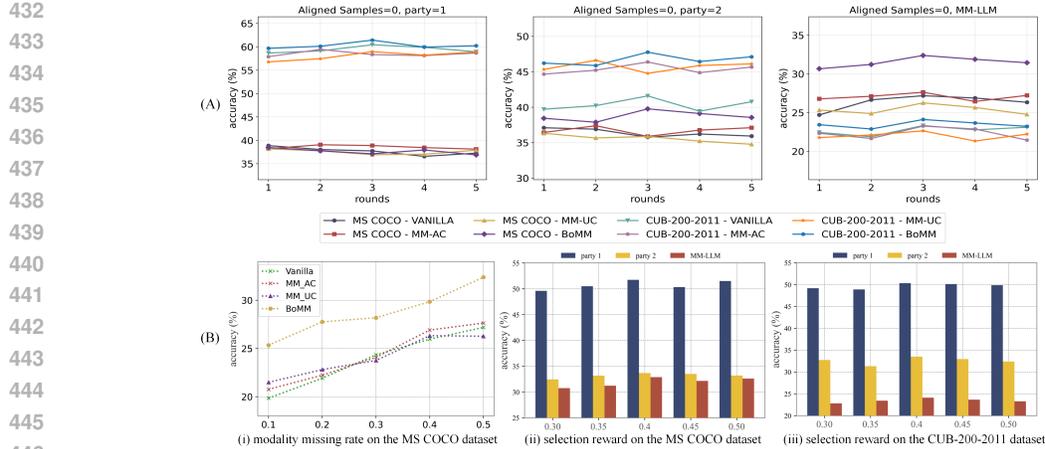


Figure 3: Convergence analysis and hyperparameter analysis of our proposed BoMM.



Figure 4: Case study results.

As expected, increasing modality missing rates lead to performance degradation, with our method maintaining stable performance even under severe missing conditions. Fig. 3(ii) and (iii) shows the effect of sample selection threshold μ on model performance. The results indicate that an optimal threshold around $\mu = 0.4$ achieves the best balance between sample quality and quantity, with too low thresholds including noisy samples and too high thresholds limiting data availability.

5.6 CASE STUDY

The case study results in Fig. 4 compares text descriptions generated by the original MM-LM against our proposed BoMM method across three representative bird images. The color coding in the descriptions indicates different types of attributes: yellow bold text represents attributes that accurately match the dataset’s standard descriptions, red bold text indicates attributes generated by the models that don’t align with the dataset descriptions, and green bold text highlights obviously incorrect attributes. Specifically, BoMM consistently produces more detailed and contextually rich descriptions. For example, while the original method provides a basic description of “a small bird perched on a fence” our BoMM method enhances this with additional details such as “displaying a pale belly” and “surrounded by greenery”. These improvements demonstrate that our preference-driven modality adaptive completion method successfully generates domain-specific, high-quality descriptions that better capture visual details.

6 CONCLUSION AND FUTURE DIRECTION

This paper explores an under-explored paradigm: multi-modality large-small model collaboration. We propose a global prototype-guided alignment method to identify potentially aligned samples to mitigate modality alignment scarcity, and design a preference-driven modality adaptive completion method to address modality completeness gap challenges, enabling bidirectional large-small model knowledge transfer. Although we avoid direct leakage of raw data by exchanging intermediate representations, enhancing privacy remains an important future direction.

REFERENCES

- 486
487
488 Oded Berger-Tal, Jonathan Nathan, Ehud Meron, and David Saltz. The exploration-exploitation
489 dilemma: a multidisciplinary framework. *PLoS one*, 9(4):e95693, 2014.
- 490 Jacob Bien and Robert Tibshirani. Prototype selection for interpretable classification. 2011.
491
- 492 Wanyi Chen, Zihua Zhao, Jiangchao Yao, Ya Zhang, Jiajun Bu, and Haishuai Wang. Multi-modal
493 medical diagnosis via large-small model collaboration. In *Proceedings of the Computer Vision
494 and Pattern Recognition Conference*, pp. 30763–30773, 2025.
- 495 Sijie Cheng, Jingwen Wu, Yanghua Xiao, and Yang Liu. Fedgems: Federated learning of larger
496 server models via selective knowledge fusion. *arXiv preprint arXiv:2110.11027*, 2021.
497
- 498 Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu
499 Lu, Zichong Yang, Kuei-Da Liao, et al. A survey on multimodal large language models for
500 autonomous driving. In *Proceedings of the IEEE/CVF winter conference on applications of com-
501 puter vision*, pp. 958–979, 2024.
- 502 Tao Fan, Guoqiang Ma, Yan Kang, Hanlin Gu, Yuanfeng Song, Lixin Fan, Kai Chen, and Qiang
503 Yang. Fedmkt: Federated mutual knowledge transfer for large and small language models. *arXiv
504 preprint arXiv:2406.02224*, 2024.
- 505 Siwei Feng. Vertical federated learning-based feature selection with non-overlapping sample uti-
506 lization. *Expert Systems with Applications*, 208:118097, 2022.
507
- 508 Siwei Feng, Boyang Li, Han Yu, Yang Liu, and Qiang Yang. Semi-supervised federated heteroge-
509 neous transfer learning. *Knowledge-Based Systems*, 252:109384, 2022.
- 510 Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv preprint
511 arXiv:1802.02246*, 2018.
512
- 513 Wei Guo, Fuzhen Zhuang, Xiao Zhang, Yiqi Tong, and Jin Dong. A comprehensive survey of
514 federated transfer learning: challenges, methods and applications. *Frontiers of Computer Science*,
515 18(6):186356, 2024.
- 516 Wei Guo, Yiqi Tong, Yiyang Duan, Fuzhen Zhuang, Xiao Zhang, Zhaojun Hu, and Jin Dong.
517 Feze: Alignment-flexible zero-shot vertical federated learning. In *Proceedings of the 31st ACM
518 SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pp. 745–754, 2025.
- 519 Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. Parameter-efficient fine-tuning
520 for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024.
521
- 522 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A
523 reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference
524 on Empirical Methods in Natural Language Processing*, pp. 7514–7528, 2021.
- 525 Zixuan Huang, Yikun Ban, Lean Fu, Xiaojie Li, Zhongxiang Dai, Jianxin Li, and Deqing Wang.
526 Adaptive sample scheduling for direct preference optimization. *arXiv preprint arXiv:2506.17252*,
527 2025.
- 528 Kaiyi Ji, Jason D Lee, Yingbin Liang, and H Vincent Poor. Convergence of meta-learning with task-
529 specific adaptation over partial parameters. *Advances in Neural Information Processing Systems*,
530 33:11490–11500, 2020.
531
- 532 Kaiyi Ji, Junjie Yang, and Yingbin Liang. Bilevel optimization: Convergence analysis and enhanced
533 design. In *International conference on machine learning*, pp. 4882–4892. PMLR, 2021.
- 534 Songtao Jiang, Yan Zhang, Ruizhe Chen, Tianxiang Hu, Yeying Jin, Qinglin He, Yang Feng, Jian
535 Wu, and Zuozhu Liu. Modality-fair preference optimization for trustworthy mllm alignment.
536 *arXiv preprint arXiv:2410.15334*, 2024.
537
- 538 Shyamgopal Karthik, Huseyin Coskun, Zeynep Akata, Sergey Tulyakov, Jian Ren, and Anil
539 Kag. Scalable ranked preference optimization for text-to-image generation. *arXiv preprint
arXiv:2410.18013*, 2024.

- 540 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
541 2014.
- 542
- 543 Soyeong Kwon, Taegyeong Lee, and Taehwan Kim. Zero-shot text-guided infinite image synthesis
544 with llm guidance. *arXiv preprint arXiv:2407.12642*, 2024.
- 545
- 546 Iro Laina, Christian Rupprecht, and Nassir Navab. Towards unsupervised image captioning with
547 shared multimodal embeddings. In *Proceedings of the IEEE/CVF International Conference on*
548 *Computer Vision*, pp. 7414–7424, 2019.
- 549 Gyeonggeon Lee, Lehong Shi, Ehsan Latif, Yizhu Gao, Arne Bewersdorff, Matthew Nyaaba,
550 Shuchen Guo, Zhengliang Liu, Gengchen Mai, Tianming Liu, et al. Multimodality of ai for
551 education: Towards artificial general intelligence. *IEEE Transactions on Learning Technologies*,
552 2025.
- 553
- 554 Lujun Li and Zhe Jin. Shadow knowledge distillation: Bridging offline and online knowledge trans-
555 fer. *Advances in Neural Information Processing Systems*, 35:635–649, 2022.
- 556
- 557 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
558 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European*
559 *conference on computer vision*, pp. 740–755. Springer, 2014.
- 560 Yang Liu, Bingjie Yan, Tianyuan Zou, Jianqing Zhang, Zixuan Gu, Jianbing Ding, Xidong Wang,
561 Jingyi Li, Xiaozhou Ye, Ye Ouyang, et al. Towards harnessing the collaborative power of large
562 and small models for domain tasks. *arXiv preprint arXiv:2504.17421*, 2025.
- 563
- 564 Yitao Liu, Tianxiang Sun, Xipeng Qiu, and Xuanjing Huang. Learning to teach with student feed-
565 back. *arXiv preprint arXiv:2109.04641*, 2021.
- 566
- 567 Zimo Liu, Kangjun Liu, Mingyue Guo, Shiliang Zhang, and Yaowei Wang. Cotuning: A large-small
568 model collaborating distillation framework for better model generalization. In *Proceedings of the*
569 *32nd ACM International Conference on Multimedia*, pp. 10487–10496, 2024.
- 570
- 571 Sahar Almahfouz Nasser, Nihar Gupte, and Amit Sethi. Reverse knowledge distillation: Training a
572 large model using a small one for retinal image matching on limited data. In *Proceedings of the*
IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 7778–7787, 2024.
- 573
- 574 Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data
575 science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- 576
- 577 Hieu Pham, Zihang Dai, Qizhe Xie, and Quoc V Le. Meta pseudo labels. In *Proceedings of the*
IEEE/CVF conference on computer vision and pattern recognition, pp. 11557–11568, 2021.
- 578
- 579 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
580 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
581 *in Neural Information Processing Systems*, 36:53728–53741, 2023.
- 582
- 583 Scott Reed, Zeynep Akata, Honglak Lee, and Bernt Schiele. Learning deep representations of
584 fine-grained visual descriptions. In *Proceedings of the IEEE conference on computer vision and*
pattern recognition, pp. 49–58, 2016.
- 585
- 586 Shezheng Song, Xiaopeng Li, Shasha Li, Shan Zhao, Jie Yu, Jun Ma, Xiaoguang Mao, and Weimin
587 Zhang. How to bridge the gap between modalities: A comprehensive survey on multimodal large
588 language model. *arXiv preprint arXiv:2311.07594*, 2023.
- 589
- 590 Yiqi Tong, Jiarui Zhang, Shaohang Wei, Wei Guo, Fuzhen Zhuang, Deqing Wang, Xi Yang, and
591 Richeng Xuan. Mindscore: quantifying human preference for text-to-image generation through
592 multi-view lens. *Science China Information Sciences*, 68(6):160105, 2025.
- 593
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd
birds-200-2011 dataset. 2011.

- 594 Fali Wang, Zhiwei Zhang, Xianren Zhang, Zongyu Wu, Tzuhao Mo, Qiuhaio Lu, Wanqing Wang,
595 Rui Li, Junjie Xu, Xianfeng Tang, et al. A comprehensive survey of small language models in the
596 era of large language models: Techniques, enhancements, applications, collaboration with llms,
597 and trustworthiness. *arXiv preprint arXiv:2411.03350*, 2024a.
- 598
599 Puyi Wang, Wei Sun, Zicheng Zhang, Jun Jia, Yanwei Jiang, Zhichao Zhang, Xionguo Min, and
600 Guangtao Zhai. Large multi-modality model assisted ai-generated image quality assessment. In
601 *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 7803–7812, 2024b.
- 602 Xiaosong Wang, Xiaofan Zhang, Guotai Wang, Junjun He, Zhongyu Li, Wentao Zhu, Yi Guo,
603 Qi Dou, Xiaoxiao Li, Dequan Wang, et al. Openmedlab: An open-source platform for multi-
604 modality foundation models in medicine. *arXiv preprint arXiv:2402.18028*, 2024c.
- 605 Zheng Wen, Branislav Kveton, and Azin Ashkan. Efficient learning in large-scale combinatorial
606 semi-bandits. In *International Conference on Machine Learning*, pp. 1113–1122. PMLR, 2015.
- 607
608 Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and Philip S Yu. Multimodal large
609 language models: A survey. In *2023 IEEE International Conference on Big Data (BigData)*, pp.
610 2247–2256. IEEE, 2023.
- 611 Renjie Wu, Hu Wang, Hsiang-Ting Chen, and Gustavo Carneiro. Deep multimodal learning with
612 missing modality: A survey. *arXiv preprint arXiv:2409.07825*, 2024a.
- 613
614 Xun Wu, Shaohan Huang, Guolong Wang, Jing Xiong, and Furu Wei. Multimodal large language
615 models make text-to-image generative models align better. *Advances in Neural Information Pro-
616 cessing Systems*, 37:81287–81323, 2024b.
- 617 Zhengjie Yang, Sen Fu, Wei Bao, Dong Yuan, and Albert Y Zomaya. Fastslowmo: Federated learn-
618 ing with combined worker and aggregator momenta. *IEEE Transactions on Artificial Intelligence*,
619 4(5):1041–1050, 2022.
- 620
621 Rui Ye, Wenhao Wang, Jingyi Chai, Dihan Li, Zexi Li, Yinda Xu, Yaxin Du, Yanfeng Wang, and
622 Siheng Chen. Openfedllm: Training large language models on decentralized private data via
623 federated learning. In *Proceedings of the 30th ACM SIGKDD conference on knowledge discovery
624 and data mining*, pp. 6137–6147, 2024.
- 625 Qiyang Yu, Yang Liu, Yimu Wang, Ke Xu, and Jingjing Liu. Multimodal federated learning via
626 contrastive representation ensemble. *arXiv preprint arXiv:2302.08888*, 2023.
- 627
628 Haozhe Zhang, Junjie Pang, Yan Huang, Zhenzhen Xie, and Zelei Liu. Fedbridgeicl: Federated
629 bridging of small and large models for in-context learning. In *International Conference on Wire-
630 less Artificial Intelligent Computing Systems and Applications*, pp. 1–10. Springer, 2025.
- 631 Huaao Zhang, Shigui Qiu, and Shilong Wu. Dual knowledge distillation for bidirectional neural
632 machine translation. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pp.
633 1–7. IEEE, 2021.
- 634 Weitong Zhang, Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural thompson sampling. *arXiv
635 preprint arXiv:2010.00827*, 2020.
- 636
637 Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks
638 with noisy labels. *Advances in neural information processing systems*, 31, 2018.
- 639 Wangchunshu Zhou, Canwen Xu, and Julian McAuley. Bert learns to teach: Knowledge distillation
640 with meta learning. *arXiv preprint arXiv:2106.04570*, 2021.
- 641
642

643 A IMPLEMENT DETAILS

644

645 We evaluate our proposed BoMM framework on two widely-used multi-modal datasets: MS COCO
646 Lin et al. (2014) and CUB-200-2011 Wah et al. (2011). Specifically, the MS COCO dataset contains
647 images across 80 object categories, where each image is paired with at least five textual descriptions.
For our image-text classification task, we assign each image to the category with the highest number

of target objects present. The CUB-200-2011 dataset comprises 11,788 images covering 200 fine-grained bird species, with 5,994 images for training and 5,794 for testing. Following with Reed et al. (2016), we utilize the extended version that incorporates fine-grained visual descriptions, where each image contains at least 10 descriptions, with each description comprising at least 10 words without species names, background details, or action information.

We conduct experiments in a large-small collaboration scenario with heterogeneous multi-modality models, including two parties’ SMs and one cloud-side MM-LM. The parties deploy small models including CLIP-base¹ and CLIP-large² locally, while the cloud hosts a large multi-modal model LLaVA-1.5-7B³. Each party possesses both aligned and unaligned samples, where aligned samples share identical IDs but contain different textual descriptions, while unaligned samples are exclusively owned by individual parties. We evaluate three modality missing scenarios. In scenario 1, both parties and MM-LM suffer from text modality missing with probability Δ . In scenario 2, party 0 maintains complete image and text modalities while party 1 experiences text modality missing with probability Δ for both aligned and unaligned samples. In scenario 3, both parties and MM-LM suffer from image modality missing with probability Δ . In addition, to assess the impact of data alignment, we vary the number of aligned samples across 0 and 200, while fixing unaligned samples at 1,000 for each party and MM-LM.

We set the collaboration process to 3 global rounds, with a batch size of 32. In each round, both parties and the MM-LM perform 20 training epochs. All models are optimized using the Adam optimizer Kingma (2014), with learning rates set to $5e-4$ for parties and $5e-5$ for the MM-LM. Finally, we use accuracy to measure classification performance, CLIPScore Hessel et al. (2021) to evaluate the quality of generated modalities, communication cost measured in MB to assess data transfer efficiency, and computational time in minutes to evaluate system efficiency.

B PROOF OF THEOREM 1

B.1 ASSUMPTIONS

Let $\mathcal{Z} := \{\Psi, \Phi\}$ denote all model parameters. We adopt standard assumptions for bi-level optimization problem (Ghadimi & Wang, 2018; Ji et al., 2020) equation 10 on L_{global} and L_{local} .

Assumption 1. (Lipschitz Condition) The loss function $L_{\text{global}}(\mathcal{Z}; \xi)$ is L_1 -Lipschitz for any given ξ .

Assumption 2. (Smoothness) The loss functions $L_{\text{global}}(\mathcal{Z}; \xi)$, $L_{\text{local}}(\mathcal{Z}; \xi)$ and $l_{\text{global}}(\mathcal{Z}; \xi)$ are L_2 -smooth, L_2 -smooth and L -smooth for any give ξ and ζ , respectively.

Assumption 3. (Lipschitz Condition for Second Derivatives) The second derivatives $\nabla_{\Psi} \nabla_{\Theta} L_{\text{local}}(\mathcal{Z}; \zeta)$ and $\nabla_{\Theta}^2 L_{\text{local}}(\mathcal{Z}; \zeta)$ are L_3 -Lipschitz and L_4 -Lipschitz for any give ζ , respectively.

Assumption 4. (Bounded Domain) The parameter Θ is in a bounded domain with a diameter Δ , i.e., for any Θ_1 and Θ_2 , we have:

$$\|\Theta_1 - \Theta_2\| \leq \Delta.$$

B.2 ABSTRACTION OF OUR ALGORITHM

For clarity in convergence analysis, we write our algorithm as Algorithm 1.

B.3 PROOF DETAILS

To give the convergence analysis of our algorithm, we need to give some useful lemmas first.

Firstly, with L -smoothness of $l_{\text{local}}(\Psi, \Theta; \zeta)$, we could give the strongly convexity property of $L_{\text{local}}(\Psi, \Theta; \zeta)$.

¹<https://huggingface.co/openai/clip-vit-base-patch32>

²<https://huggingface.co/openai/clip-vit-large-patch14>

³<https://huggingface.co/llava-hf/llava-1.5-7b-hf>

Algorithm 1 Algorithmic Abstraction of BoMM

Input: communication round T ; learning rate: η and η' ; epoch: τ, τ' ; initial parameters Θ_{-1}^τ and $\Psi_{-1}^{\tau'}$.

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

 Let $\Theta_t^0 = \Theta_{t-1}^\tau$ and $\Psi_t^0 = \Psi_{t-1}^{\tau'}$;

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

 Select S'_{j-1} , a subset of each party's dataset D_k ,

 update $\Psi_t^{j'} = \Psi_t^{j'-1} - \eta' \nabla_{\Psi} L_{\text{global}}(\Psi_t^{j'-1}, \Theta_{t-1}^\tau; S'_{j-1})$;

end for

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

 Update $\Theta_t^j = \Theta_t^{j-1} - \eta \nabla_{\Theta} L_{\text{local}}(\Psi_t^{\tau'}, \Theta_t^{j-1})$;

end for

end for

Lemma 1. Under Assumption 2, suppose $\alpha := -L + \varphi > 0$, $L_{\text{local}}(\Psi, \Theta; \zeta)$ is α -strongly convex w.r.t Θ .

Based on Assumptions 1 and 2, we could give Lemma 2 directly.

Lemma 2. Under Assumption 1 and 2, the derivatives $\nabla L_{\text{global}}(\mathcal{Z}; \xi)$, $\nabla_{\Psi} \nabla_{\Theta} L_{\text{local}}(\mathcal{Z}; \zeta)$ and $\nabla_{\Theta}^2 L_{\text{local}}(\mathcal{Z}; \zeta)$ have bounded variances, i.e., for any \mathcal{Z} , we have

$$\mathbb{E}_{\xi} \|\nabla L_{\text{global}}(\mathcal{Z}; \xi) - \nabla \mathcal{L}_{\text{global}}(\mathcal{Z})\|^2 \leq L_1^2, \quad (12)$$

$$\mathbb{E}_{\zeta} \|\nabla_{\Psi} \nabla_{\Theta} L_{\text{local}}(\mathcal{Z}; \zeta) - \nabla_{\Psi} \nabla_{\Theta} \mathcal{L}_{\text{local}}(\mathcal{Z})\|^2 \leq L_2^2, \quad (13)$$

$$\mathbb{E}_{\zeta} \|\nabla_{\Theta}^2 L_{\text{local}}(\mathcal{Z}; \zeta) - \nabla_{\Theta}^2 \mathcal{L}_{\text{local}}(\mathcal{Z})\|^2 \leq L_2^2. \quad (14)$$

We omit them since the proof of Lemma 1 and 2 is basic.

To show the smooth property of $F(\Psi)$, we introduce the following lemma which is proposed in Ghadimi & Wang (2018) firstly.

Lemma 3. (Lemma 2.2 in Ghadimi & Wang (2018)) Under Assumptions 1, 2 and 3, $F(\Psi)$ is L_0 -smooth where L_0 is given by

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

Tracking error $\mathbb{E} \|\Theta_t^{j-1} - \Theta^*(\Psi_t^0)\|$ is an important component in our convergence analysis. To give an upper bound on the tracking error, we utilized Lemma 9 in Ji et al. (2021). For simplicity, we assume equal batch sizes N_b for both cloud and local data.

Lemma 4. (Lemma 9 in Ji et al. (2021)) Under Assumptions 1 and 2, with stepsize η to be $\frac{2}{L_2 + \alpha}$, we have

$$\mathbb{E} \|\Theta_t^{j-1} - \Theta^*(\Psi_t^0)\|^2 \leq \left(\frac{L_2 - \alpha}{L_2 + \alpha} \right)^{2(j-1)} \mathbb{E} \|\Theta_t^0 - \Theta^*(\Psi_t^0)\|^2. \quad (15)$$

With the help of the above lemmas, we could give the estimation property of the $\frac{\partial L_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0}$ approximating $\nabla F(\Psi_t^0)$. The result is presented in the following Proposition 1.

Proposition 1. Under Assumptions 1-4, choose stepsize η to be $\frac{2}{L_2+\alpha}$ and suppose $\alpha < L_2$, we have

$$\begin{aligned} \mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^0, \Theta_t^\tau; \mathcal{S}'_0)}{\partial \Psi_t^0} - \nabla F(\Psi_t^0) \right\| &\leq (L_2 + \frac{L_2^2}{\alpha}) \left[\left(\frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \right. \\ &\quad \left. \sqrt{\Delta} \right] + L_1 \left[\frac{L_2(1 - \frac{2}{L_2+\alpha}\alpha)^\tau}{\alpha} \right. \\ &\quad \left. + \frac{2}{L_2 + \alpha} \left(\frac{L_2 L_4}{\alpha} + L_3 \right) \sqrt{\Delta} \right. \\ &\quad \left. \frac{(1 - \frac{2}{L_2+\alpha}\alpha)^\tau}{1 - \frac{2}{L_2+\alpha}\alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right] + \frac{L_1}{\sqrt{N_b}}. \end{aligned} \quad (16)$$

Proof. Using the triangle inequality, we have

$$\begin{aligned} &\mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^0, \Theta_t^\tau; \mathcal{S}'_0)}{\partial \Psi_t^0} - \nabla F(\Psi_t^0) \right\| \\ &= \mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^0, \Theta_t^\tau; \mathcal{S}'_0)}{\partial \Psi_t^0} - \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0} \right. \\ &\quad \left. + \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0} - \nabla F(\Psi_t^0) \right\| \\ &\leq \mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^0, \Theta_t^\tau; \mathcal{S}'_0)}{\partial \Psi_t^0} - \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0} \right\| \\ &\quad + \mathbb{E} \left\| \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0} - \nabla F(\Psi_t^0) \right\| \\ &\stackrel{(i)}{\leq} \frac{L_1}{\sqrt{N_b}} + \mathbb{E} \left\| \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0} - \nabla F(\Psi_t^0) \right\|. \end{aligned} \quad (17)$$

where (i) follows from equation 12.

Then, we need to give an upper bound for $\mathbb{E} \left\| \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0} - \nabla F(\Psi_t^0) \right\|$. Using

$$\nabla F(\Psi_t^0) = \nabla \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta^*(\Psi_t^0)) + \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \nabla_{\Theta} \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta^*(\Psi_t^0))$$

and

$$\frac{\partial L_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0} = \nabla_{\Psi} \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau) + \frac{\partial \Theta_t^\tau}{\partial \Psi_t^0} \nabla_{\Theta} \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau),$$

we have

$$\begin{aligned} &\mathbb{E} \left\| \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0} - \nabla F(\Psi_t^0) \right\| \\ &\leq L_2 \mathbb{E} \|\Theta_t^\tau - \Theta^*(\Psi_t^0)\| + L_1 \mathbb{E} \left\| \frac{\partial \Theta_t^\tau}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right\| \\ &\quad + L_2 \mathbb{E} \left(\left\| \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right\| \|\Theta_t^\tau - \Theta^*(\Psi_t^0)\| \right). \end{aligned} \quad (18)$$

Now we want to bound $\mathbb{E} \left\| \frac{\partial \Theta_t^\tau}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right\|$ first.

Recall the update method

$$\Theta_t^j = \Theta_t^{j-1} - \eta \nabla_{\Theta} \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta_t^{j-1})$$

and use the chain rule on it, we have

$$\begin{aligned} \frac{\partial \Theta_t^j}{\partial \Psi_t^0} &= \frac{\partial \Theta_t^{j-1}}{\partial \Psi_t^0} - \eta (\nabla_{\Psi} \nabla_{\Theta} \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta_t^{j-1})) \\ &\quad + \frac{\partial \Theta_t^{j-1}}{\partial \Psi_t^0} \nabla_{\Theta}^2 \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta_t^{j-1}). \end{aligned} \quad (19)$$

For $\Theta^*(\Psi_t^0)$ is the optimal solution of $\mathcal{L}_{\text{local}}(\Psi_t^0, \Theta)$, we have $\nabla_{\Theta} \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta^*(\Psi_t^0)) = 0$. Then, using the chain rule, we have

$$\nabla_{\Psi} \nabla_{\Theta} \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta^*(\Psi_t^0)) + \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \nabla_{\Theta}^2 \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta^*(\Psi_t^0)) = 0. \quad (20)$$

Combining equation 19 and equation 20, the following equation holds

$$\begin{aligned} \frac{\partial \Theta_t^j}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} &= \frac{\partial \Theta_t^{j-1}}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \\ &- \eta \left(\frac{\partial \Theta_t^{j-1}}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right) \nabla_{\Theta}^2 \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta_t^{j-1}) \\ &- \eta \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} (\nabla_{\Theta}^2 \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta_t^{j-1}) - \nabla_{\Theta}^2 \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta^*(\Psi_t^0))). \end{aligned} \quad (21)$$

Based on equation 20, we have

$$\left\| \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right\| \leq \frac{L_2}{\alpha}. \quad (22)$$

With the help of Assumption 3, Lemma 2, equation 21 and equation 22, the following bound holds

$$\begin{aligned} &\mathbb{E} \left\| \frac{\partial \Theta_t^j}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right\| \\ &\leq \mathbb{E} \left(\|I - \eta \nabla_{\Theta}^2 \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta_t^{j-1})\| \left\| \frac{\partial \Theta_t^{j-1}}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right\| \right) \\ &+ \eta \frac{L_2}{\alpha} (\mathbb{E} \|\nabla_{\Theta}^2 \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta_t^{j-1}) - \nabla_{\Theta}^2 \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta^*(\Psi_t^0))\|) \\ &+ \eta \mathbb{E} \|\nabla_{\Psi} \nabla_{\Theta} \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta_t^{j-1}) - \nabla_{\Psi} \nabla_{\Theta} \mathcal{L}_{\text{local}}(\Psi_t^0, \Theta^*(\Psi_t^0))\| \\ &\leq (1 - \eta\alpha) \mathbb{E} \left\| \frac{\partial \Theta_t^{j-1}}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right\| + \eta \left(\frac{L_2 L_4}{\alpha} + L_3 \right) \mathbb{E} \|\Theta_t^{j-1} - \Theta^*(\Psi_t^0)\|. \end{aligned} \quad (23)$$

For the choice of $\eta = \frac{2}{L_2 + \alpha}$, Lemma 4 holds. With the result in equation 15 and Assumption 4, we have

$$\begin{aligned} \mathbb{E} \|\Theta_t^{j-1} - \Theta^*(\Psi_t^0)\| &\leq \left(\frac{L_2 - \alpha}{L_2 + \alpha} \right)^{j-1} \sqrt{\mathbb{E} \|\Theta_t^0 - \Theta^*(\Psi_t^0)\|} \\ &\leq \left(\frac{L_2 - \alpha}{L_2 + \alpha} \right)^{j-1} \sqrt{\Delta}. \end{aligned} \quad (24)$$

Telescoping equation 23 over j from 0 to τ and combining the result with equation 24 yields

$$\begin{aligned} \mathbb{E} \left\| \frac{\partial \Theta_t^{\tau}}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right\| &\leq (1 - \eta\alpha)^{\tau} \mathbb{E} \left\| \frac{\partial \Theta_t^0}{\partial \Psi_t^0} - \frac{\partial \Theta^*(\Psi_t^0)}{\partial \Psi_t^0} \right\| \\ &+ \eta \left(\frac{L_2 L_4}{\alpha} + L_3 \right) \sum_{j=0}^{\tau-1} (1 - \eta\alpha)^{\tau-1-j} \\ &\left[\left(\frac{L_2 - \alpha}{L_2 + \alpha} \right)^j \sqrt{\Delta} \right] \\ &\leq \frac{L_2 (1 - \eta\alpha)^{\tau}}{\alpha} \\ &+ \eta \left(\frac{L_2 L_4}{\alpha} + L_3 \right) \sqrt{\Delta} \frac{(1 - \eta\alpha)^{\tau}}{1 - \eta\alpha - \frac{L_2 - \alpha}{L_2 + \alpha}}. \end{aligned} \quad (25)$$

Take expedition of both sides of equation 18, plugging equation 22, equation 24 and equation 25 into it yields

$$\begin{aligned}
\mathbb{E} \left\| \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^0, \Theta_t^\tau)}{\partial \Psi_t^0} - \nabla F(\Psi_t^0) \right\| &\leq (L_2 + \frac{L_2^2}{\alpha}) \left[\left(\frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \sqrt{\Delta} \right] \\
&+ L_1 \left[\frac{L_2(1 - \frac{2}{L_2 + \alpha} \alpha)^\tau}{\alpha} \right. \\
&+ \frac{2}{L_2 + \alpha} \left(\frac{L_2 L_4}{\alpha} + L_3 \right) \sqrt{\Delta} \\
&\left. \frac{(1 - \frac{2}{L_2 + \alpha} \alpha)^\tau}{1 - \frac{2}{L_2 + \alpha} \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right]. \tag{26}
\end{aligned}$$

Finally, plugging equation 26 into equation 17 yields equation 16. Thus, we complete the proof of Proposition 1. \square

To deal with multiple updates of the active party with a fixed t , we need to use a technology called virtual updates (Yang et al., 2022). Specifically, we introduce a virtual parameter $\Theta_{t,j'}^\tau$ which could be obtained by local updates when $\Psi_t^{j'}$ is given. With the help of Proposition 1, we could obtain the bound for $\mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_{t,j'}^\tau; S_{j'}')}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\|$. Now, we want to give a bound for $\mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_t^\tau; S_{j'}')}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\|$. The result is as follows

Proposition 2. Following the conditions in Proposition 1, we have

$$\begin{aligned}
\mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_t^\tau; S_{j'}')}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\| &\leq (L_2 + \frac{L_2^2}{\alpha}) \left[\left(\frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \right. \\
&\left. \sqrt{\Delta} \right] + L_1 \left[\frac{L_2(1 - \frac{2}{L_2 + \alpha} \alpha)^\tau}{\alpha} \right. \\
&+ \frac{2}{L_2 + \alpha} \left(\frac{L_2 L_4}{\alpha} + L_3 \right) \sqrt{\Delta} \\
&\left. \frac{(1 - \frac{2}{L_2 + \alpha} \alpha)^\tau}{1 - \frac{2}{L_2 + \alpha} \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right] + \frac{L_1}{\sqrt{N_b}} + L_2 \Delta. \tag{27}
\end{aligned}$$

Proof. Using the triangle inequality, we have

$$\begin{aligned}
&\left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_t^\tau; S_{j'}')}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\| \\
&\leq \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_t^\tau; S_{j'}')}{\partial \Psi_t^{j'}} - \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_{t,j'}^\tau; S_{j'}')}{\partial \Psi_t^{j'}} \right\| \\
&+ \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_{t,j'}^\tau; S_{j'}')}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\| \tag{28}
\end{aligned}$$

Take expectation of both sides of equation 28, with the help of Assumption 2 and 5, the following result holds

$$\begin{aligned}
& \mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_t^\tau; \mathcal{S}'_{j'})}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\| \\
& \leq \mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_t^\tau; \mathcal{S}'_{j'})}{\partial \Psi_t^{j'}} - \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_{t,j'}^\tau; \mathcal{S}'_{j'})}{\partial \Psi_t^{j'}} \right\| \\
& + \mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_{t,j'}^\tau; \mathcal{S}'_{j'})}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\| \tag{29} \\
& \leq L_2 \|\Theta_t^\tau - \Theta_{t,j'}^\tau\| + \mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_{t,j'}^\tau; \mathcal{S}'_{j'})}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\| \\
& \leq L_2 \Delta + \mathbb{E} \left\| \frac{\partial L_{\text{global}}(\Psi_t^{j'}, \Theta_{t,j'}^\tau; \mathcal{S}'_{j'})}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\|.
\end{aligned}$$

Plugging in the result in Proposition 1 into equation 29 yields equation 27. Thus, we complete the proof of Proposition 2. \square

Proof for Theorem 1. Based on the L_0 -smoothness of $F(\Psi)$ established in Lemma 3, we have

$$\begin{aligned}
F(\Psi_t^{j'+1}) & \leq F(\Psi_t^{j'}) + \langle \nabla F(\Psi_t^{j'}), \Psi_t^{j'+1} - \Psi_t^{j'} \rangle \\
& + \frac{L_0}{2} \|\Psi_t^{j'+1} - \Psi_t^{j'}\|^2 \\
& \leq F(\Psi_t^{j'}) - \eta' \langle \nabla F(\Psi_t^{j'}), \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^{j'}, \Theta_t^\tau)}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \rangle \\
& - \eta' \|\nabla F(\Psi_t^{j'})\|^2 + \eta'^2 L_0 \|\nabla F(\Psi_t^{j'})\|^2 \tag{30} \\
& + \eta'^2 L_0 \left\| \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^{j'}, \Theta_t^\tau)}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\|^2 \\
& \leq F(\Psi_t^{j'}) - \left(\frac{\eta'}{2} - \eta'^2 L_0 \right) \|\nabla F(\Psi_t^{j'})\|^2 \\
& + \left(\frac{\eta'}{2} + \eta'^2 L_0 \right) \left\| \frac{\partial \mathcal{L}_{\text{global}}(\Psi_t^{j'}, \Theta_t^\tau)}{\partial \Psi_t^{j'}} - \nabla F(\Psi_t^{j'}) \right\|^2.
\end{aligned}$$

Telescoping equation 30 over j' from 0 to $\tau' - 1$ and taking expectation of both sides yields

$$\begin{aligned}
\mathbb{E} F(\Psi_t^{\tau'}) & \leq \mathbb{E} F(\Psi_t^0) - \left(\frac{\eta'}{2} - \eta'^2 L_0 \right) \mathbb{E} \sum_{j'=0}^{\tau'-1} \|\nabla F(\Psi_t^{j'})\|^2 \\
& + (\eta' + 2\eta'^2 L_0) \tau' \left\{ \left(L_2 + \frac{L_2^2}{\alpha} \right) \left[\left(\frac{L_2 - \alpha}{L_2 + \alpha} \right)^\tau \right. \right. \\
& \left. \left. \sqrt{\Delta} \right] + L_1 \left[\frac{L_2 \left(1 - \frac{2}{L_2 + \alpha} \alpha \right)^\tau}{\alpha} \right. \right. \tag{31} \\
& \left. \left. + \frac{2}{L_2 + \alpha} \left(\frac{L_2 L_4}{\alpha} + L_3 \right) \sqrt{\Delta} \right. \right. \\
& \left. \left. \frac{\left(1 - \frac{2}{L_2 + \alpha} \alpha \right)^\tau}{1 - \frac{2}{L_2 + \alpha} \alpha - \frac{L_2 - \alpha}{L_2 + \alpha}} \right] + \frac{L_1}{\sqrt{N_b}} \right\}^2 \\
& + (\eta' + 2\eta'^2 L_0) L_2^2 \Delta^2 (\tau' - 1).
\end{aligned}$$

Telescoping equation 31 over t from 0 to $T - 1$ yields equation 11, we complete the proof of Theorem 1. \square