

DaMo: Data Mixing Optimizer in Fine-tuning Multimodal LLMs for Mobile Phone Agents

Anonymous ACL submission

Abstract

Mobile Phone Agents (MPAs) have emerged as a promising research direction due to their broad applicability across diverse scenarios. While Multimodal Large Language Models (MLLMs) serve as the foundation for MPAs, their effectiveness in handling multiple mobile phone tasks simultaneously remains limited. Although multitask supervised fine-tuning (SFT) is widely adopted for multitask learning, existing approaches struggle to determine optimal training data compositions for peak performance. To address this challenge, we propose DaMo (Data Mixture Optimizer) – a novel solution employing a trainable network that predicts optimal data mixtures by forecasting downstream task performance for any given dataset ratio. To support comprehensive evaluation, we introduce PhoneAgentBench, the first specialized benchmark to evaluate MLLMs on multimodal mobile phone tasks, comprising 1,235 QA pairs spanning diverse real-world industrial mobile application scenarios. Demonstrating strong predictive capability ($R^2=0.81$) in small-scale pilot experiments, DaMo efficiently extrapolates optimal data mixing configurations. Our results show DaMo achieves 3.06% average score improvement on PhoneAgentBench and open-source benchmarks, including BFCL-v3, MME-Reasoning, MME-Perception, and OCRBench, compared to alternative methods. Through predicting optimal data mixture only on open-source benchmarks, DaMo outperforms other approaches by 6.70% in terms of average score. Moreover, DaMo improves the metrics by 12.74% than other methods when used solely for MLLM optimization on the BFCL-v3 task. Notably, DaMo maintains robust scalability, preserving its effectiveness when applied to other model architectures.

1 Introduction

Mobile phone agents (MPAs) have attracted huge attention due to their practicability in a multitude of

scenarios. An ideal MPA has to master multiple capabilities, such as environment perception (Zhang et al., 2024; Ingold, 2021), task planning (Song et al., 2023; Liu et al., 2024c), multimodal reasoning (Lu et al., 2022; Wang et al., 2024), function call (Chen et al., 2024a; Basu, 2024), and personalized memory (Li et al., 2024a; Yuan et al., 2023).

The advent of multimodal large language models (MLLMs) provides a promising solution for the ideal agent. However, existing MLLMs encounter significant challenges in effectively integrating these diverse capabilities. Consequently, developing a versatile model capable of handling multiple tasks is critical for creating an advanced phone agent.

Multitask supervised fine-tuning (SFT) is the predominant approach utilized to empower MLLMs in addressing multiple tasks. Nevertheless, in light of numerous training datasets and downstream tasks, identifying optimal data blending strategies to maximize model performance remains a significant research challenge. The existing works on data mixture optimization (Xie et al., 2023a; Ge et al., 2024; Albalak et al., 2023; Liu et al., 2024a) focus on the pretraining phase by predicting validation loss of LLM. However, these methods are inadequate to determine the optimal data mixture for fine-tuning MLLMs, as they fail to directly correlate with model performance on downstream tasks. The sampling strategies commonly used in industrial scenarios, e.g., uniform sampling, natural sampling, and random sampling, can not yield optimal mixture due to fixed or random data mixing ratios. Grid search (Liashchynskyi, 2019) needs huge cost to find optimal mixture when the number of the training datasets becomes large.

We investigate whether downstream task performance can be reliably predicted for any given data mixture prior to actual model training, including identifying the optimal mixture that would yield optimal performance. To this end, we propose the

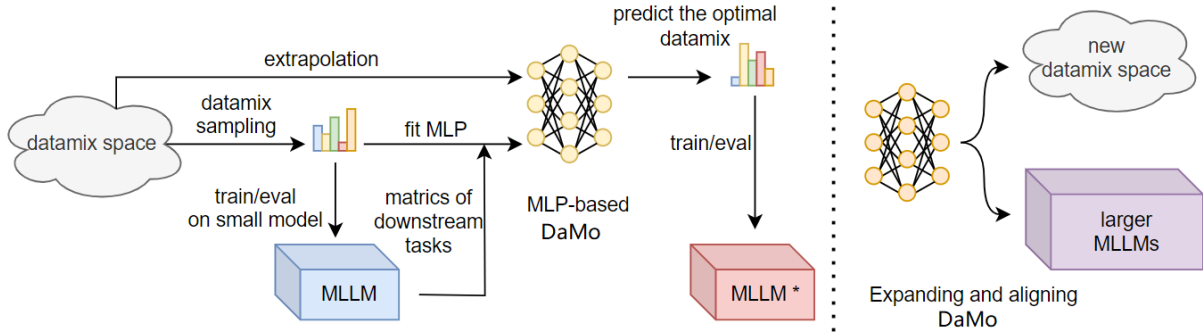


Figure 1: Illustration of our pipeline for obtaining the optimal data mixture. Left: Given m training sets with a batch size of b , all possible mixture combinations constitute the data mixing space. We sample a small number of data mixture from this space, train them on a small MLLM, and then evaluate downstream task performance. Using the data mixture as inputs and the metrics as outputs, we fit a MLP to establish the DaMo. By extrapolating from the data mixing space, we predict the optimal data mixture to train the MLLM. Right: Demonstrates the extension and alignment of DaMo to other MLLMs and new data mixing spaces.

downstream task performance prediction (DaPP) method to build **Data Mixing Optimizer (DaMo)**. DaPP leverages a function to straightly predict model performance at downstream tasks. Considering that exponential functions used in (Xie et al., 2023a; Ge et al., 2024; Albalak et al., 2023) are not well-aligned with SFT performance trajectories for specific downstream applications (Huang et al., 2019; Xie et al., 2024; Isik et al., 2024), we propose to utilize a trainable neural network for target fitting. The optimal data mixture is obtained through extrapolation via DaMo.

Another obstacle in developing an ideal mobile phone agent is the absence of comprehensive real-world industrial benchmarks for evaluating MPA performance. Current benchmarks (Gao et al., 2024; Cheng et al., 2024; Li et al., 2025a; Wang et al., 2025) in this domain predominantly focus on Graphical User Interface (GUI) tasks, which fail to capture the full spectrum of practical application scenarios. To address this critical gap, we introduce PhoneAgentBench - a thorough benchmark encompassing four fundamental capabilities: 1) complex task planning, 2) device-native tool usage, 3) multimodal memory, and 4) screen context understanding. Our benchmark comprises 1235 meticulously validated test cases that simulate real-world phone interactions.

Our proposed DaMo demonstrates three key advantages. First, it achieves 3.06% average score gain on PhoneAgentBench and general benchmarks, including BFCL-V3 (Patil et al., 2025a), MME-perception (Fu et al., 2023), MME-reasoning (Yuan et al., 2025), and OCRBench (Liu et al., 2024d), compared to state-of-the-art method,

DML (Ye et al., 2024). Second, when only optimized on general benchmarks, DaMo surpasses state-of-the-art method by 6.70% in terms of average score on open-source benchmarks. Third, DaMo exhibits robust scalability across other models, while introducing significant gains on downstream tasks over other methods.

Our core contributions are as follows.

- We propose Downstream Task Performance Prediction method to establish a Data Mixing Optimizer, which directly estimates model performance on downstream tasks for optimal data mixing.
- We construct PhoneAgentBench, a benchmark spanning four critical dimensions: complex task planning, device-native tool usage, multimodal memory, and screen context understanding, mirroring real-world mobile interaction scenarios.
- Through systematic experiments, our method demonstrates exceptional generalization and scalability, outperforming other methods on PhoneAgentBench, and achieving state-of-the-art performance on the BFCL-V3 leaderboard among 4B-scale models, while also maintaining stable prediction accuracy with efficient adaptation to other models.

2 Related work

Data Mixing Early heuristic approaches like uniform sampling (Michel et al., 2021) gave way to learnable solutions; DoReMi (Xie et al., 2023b) uses Group DRO (Sagawa et al., 2020) for domain

weights; ODM (Albalak et al., 2023) frames selection as a bandit problem; BiMix (Ge et al., 2024) jointly optimizes domain proportions and data scaling. These approaches are designed for pre-training stage of LLM, which can not be directly applicable to SFT of MLLM since the intricate interactions among downstream tasks. SFTMix (Xiao et al., 2024) optimizes intra-dataset ratios but cannot handle multi-source data. Data from transfer is leveraged to estimate validation loss for LLM (Li et al., 2025b). A dynamic sampling strategy (Zhu et al., 2024) is proposed to fine-tune MOE LLM via recording the routing tokens and calculating the corresponding gate load. Key gaps remain in developing general multi-source mixing schemes for fine-tuning MLLMs.

Agent Benchmark PlanBench (Valmeekam et al., 2023) and REALM-Bench (Geng and Chang, 2025) assess planning capabilities. ToolBench (Qin et al., 2023), BFCL (Patil et al., 2025a), and API-Bank (Li et al., 2023) evaluate tool invocation and ReflectionBench (Li et al., 2024b) measures self-reflection. LTM Benchmark (Castillo-Bolado et al., 2024) tests memory retention. These benchmarks are limited to single-dimensional evaluations, lacking holistic assessment. GAIA (Mialon et al., 2023) uses end-to-end evaluation to assess general agents, but lacks granularity. AgentBench (Liu et al., 2023) and KAgentBench (Pan et al., 2023) are unimodal, ignoring multimodal interaction. ScreenSpot-Pro (Cheng et al., 2024), MobileViews (Gao et al., 2024), VisualAgentBench (Liu et al., 2024b), ScreenSpot-Pro (Li et al., 2025a), and MMBench-GUI L2 (Wang et al., 2025) can evaluate phone agents, but they are designed mainly for GUI tasks. A critical gap remains: the absence of a comprehensive benchmark supporting multimodal interaction while systematically evaluating mobile phone agents across planning, tool usage, memory, and other dimensions.

3 PhoneAgentBench

To develop a mobile phone agent benchmark tailored to real-world industrial application scenarios, we design a novel benchmark supporting systematic evaluation across key dimension such as multimodal interaction, planning, tool use, and memory. This benchmark encompasses six carefully curated datasets focusing on key mobile phone application tasks. We use Multimodal Task Planning task (MT-Plan) as a case to describe data construction.

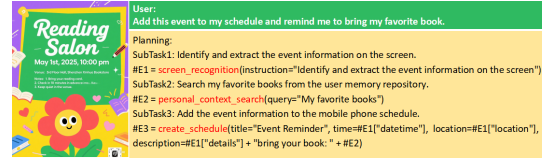


Figure 2: MT-Plan example.

MT-Plan MT-Plan is designed to evaluate multimodal task planning capabilities. Unlike T-Eval planning (Chen et al., 2023), it focuses on multimodal complex task interactions in phone agent scenarios. As shown in Figure 2, MT-Plan takes `<image + query>` as input and outputs a planning structured as a directed acyclic graph (DAG). Images are sourced from real photos or mobile screenshots, while tools are derived from APIs provided by operating systems or app ecosystems. Queries and plannings are carefully constructed by annotators based on the images. Queries are required to be concise, colloquial, and aligned with real users' daily needs. Meanwhile, tasks must be sufficiently complex to require plannings to invoke at least 2 tools. To ensure data accuracy, three annotators were invited to conduct cross-validation, and the data with inconsistent annotations was removed. Additionally, to evaluate the dataset's complexity and diversity, we compared the metrics of MT-Plan and T-Eval planning, as presented in Table 7.

The evaluator adopts the T-Eval planning evaluator (Chen et al., 2023): it compares the predicted plan with the golden plan, and calculates the score based on the length of the longest ordered action sequence derived from similarity-matched pairs.

The construction methods of the remaining five datasets and the statistics of all the tasks are presented in Appendix A.

4 Methodology

This section formalizes multitask fine-tuning optimization as identifying the optimal data mixture to maximize downstream task metrics. We propose predicting unseen mixture performance by fitting the performance of downstream tasks with limited training configurations. The process is shown in Figure 1.

4.1 Problem Formulation

Consider fine-tuning a MLLM using a mixture of m heterogeneous training datasets, denoted as $\mathcal{D} = \cup_{i=1}^m \mathcal{D}_i$. Each \mathcal{D}_i contains n_i labeled samples with the total number of samples being

245 $N = \sum_{i=1}^m n_i$. We fine-tune the MLLM starting
 246 from initial parameters θ_0 , using a batch size b , for
 247 a maximum of $T = \lceil N/b \rceil$ training steps.

248 We define the **data mixture proportion** as
 249 $\mathbf{p} = [p_1, p_2, \dots, p_m]$, where p_i represents the pro-
 250 portion of samples drawn from dataset \mathcal{D}_i . The
 251 data mixture proportion \mathbf{p} satisfies $\sum_{i=1}^m p_i = 1$.

252 Similarly, we consider k downstream test
 253 datasets, denoted as $\mathcal{D}^{test} = \cup_{j=1}^k \mathcal{D}_j^{test}$. Let $\mathbf{s} =$
 254 $[s_1, \dots, s_k] \in [0, 1]^k$ represent the score of each test
 255 dataset. The overall average score of the MLLM
 256 with parameters θ is given by $S_\theta = \frac{1}{k} \sum_{j=1}^k s_j$.

257 We aim to find the optimal data mixture propor-
 258 tion $\mathbf{p}^* \in \mathcal{P}$ (where \mathcal{P} denotes the complete data
 259 mixing space, $\mathbf{p} \in \mathbb{R}^m$) that maximizes the overall
 260 average score of downstream tasks:

$$261 \quad \mathbf{p}^* = \arg \max_{\mathbf{p} \in \mathcal{P}, t \leq T} \mathbb{E}_{\theta \sim \mathcal{A}(\mathbf{p}, t, \theta_0)} S_\theta, \quad (1)$$

262 where \mathcal{A} denotes the fine-tuning process that pro-
 263 duces the MLLM’s parameters θ based on the ini-
 264 tial parameters θ_0 for t steps using the data mixture
 265 strategy \mathbf{p} .

266 Without any constraints, the size of the set \mathcal{P}
 267 that represents batch-wise permutations is given by
 268 $|\mathcal{P}| = \frac{N!}{(b!)^T}$, which is computationally intractable.
 269 Therefore, we introduce some necessary assump-
 270 tions to prune the space \mathcal{P} . By disregarding the or-
 271 der of samples within the same dataset and keeping
 272 the data mixture fixed throughout the entire number
 273 of training steps T , we obtain a smaller data mixing
 274 space \mathcal{P}_{fix} . According to the principle of combi-
 275 nation with repetition, the size of this fixed data
 276 mixing space \mathcal{P}_{fix} is given by $|\mathcal{P}_{fix}| = C_{m+b-1}^{m-1}$.

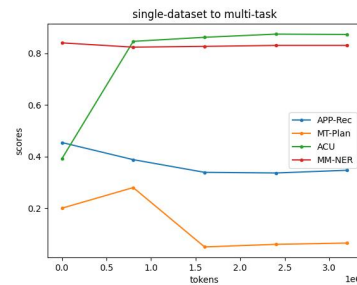
277 4.2 Performance Prediction of Tasks

278 We aim to find the optimal mixture $\mathbf{p}^* \in \mathcal{P}$. Given
 279 the high training cost of MLLMs, an exhaustive
 280 brute-force search is clearly impractical. To ad-
 281 dress this problem, we propose DaMo which is
 282 able to estimate model performance at downstream
 283 tasks without training, given any mixture propor-
 284 tions of training data. Towards this target, we fit
 285 a function f to predict performance based on data
 286 mixtures. To obtain accurate f , a efficient sampling
 287 approach is proposed to generate training samples.
 288 The sampling process is detailed as: 1) Randomly
 289 select a small set of m -dimensional mixing ratios
 290 from \mathcal{P}_{fix} . 2) Train MLLM while saving check-
 291 points at every τ steps. 3) Evaluate each check-
 292 point to obtain the performance of downstream

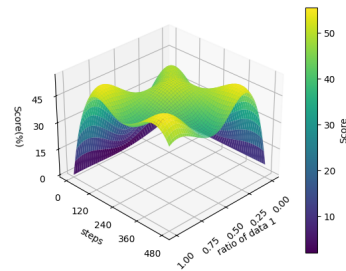
293 tasks. This process yields the mapping: (data mix-
 294 ture, training steps) \rightarrow performance of downstream
 295 tasks. Based on these samples, we fit f to predict
 296 the performance trajectory of unseen mixture:

$$297 \quad \hat{\mathbf{s}} = f(\mathbf{p}, t; \theta_0), \quad (2)$$

298 where θ_0 is initial model state and $t = \tau * i$ is train
 299 steps of the i -th checkpoint. With an accurate fitting
 300 of f , we can extrapolate performance estimates
 301 across the entire \mathcal{P}_{fix} space, dramatically reducing
 302 the model training costs required to identify the
 303 optimal data mixture.



(a) Performance on single-dataset training.



(b) Performance under dual-dataset mixtures

Figure 3: Training dynamics on downstream tasks.

304 The critical challenge lies in selecting an ap-
 305 propriate function f . While conventional expo-
 306 nential or power-law functions (Achiam et al.,
 307 2023; Grattafiori et al., 2024) are widely adopted
 308 for pretraining loss convergence, we hypothesize
 309 their inadequacy in multi-task fine-tuning scenar-
 310 ios involving interacting datasets. To validate this,
 311 we systematically analyze training dynamics un-
 312 der two configurations: (1) single-dataset training
 313 (MultiModal-Understanding, MMU) and (2) dual-
 314 dataset mixtures (APP Recognition (APP-Rec) +
 315 MMU, see Section 5.1).

316 We trained a MLLM on the MMU dataset and
 317 evaluated its performance on PhoneAgentBench.
 318 As shown in Figure 3(a), the results reveal distinct
 319 task-specific patterns: (1) **Enhancement**: MMU

320 significantly improves ACU performance. (2) **Conflict**: APP-Rec performance degrades with MMU
 321 training steps. (3) **Neutrality**: MM-NER shows no
 322 correlation with MMU training. (4) **Overfitting**:
 323 MT-Plan exhibits initial gains followed by sharp
 324 declines, indicating harmful overfitting beyond op-
 325 timal data volume.
 326

327 Figure 3(b) demonstrates the complex interac-
 328 tion when training on the mixed dataset of APP-
 329 Rec and MMU for the APP-Rec task. The 3D per-
 330 formance surface (X: training steps, Y: APP-Rec
 331 training dataset ratio, Z: APP-Rec bench score) ex-
 332 hibits **non-convex topology with non-monotonic**
 333 **fluctuations** along both axes. This nonlinearity fun-
 334 damentally prevents analytical solutions for Eq. 1
 335 and invalidates conventional exponential and power
 336 functions.

337 Motivated by neural networks’ capacity to model
 338 high-dimensional nonlinearities, we pioneer their
 339 application to DaMo. Our framework implements
 340 f as a multi-layer perceptron (MLP) that directly
 341 maps data mixture and training step to task per-
 342 formance:

$$343 \hat{s} = f_{MLP}(\mathbf{p}, t; \theta_0). \quad (3)$$

344 4.3 Optimal Data Mixture Extrapolation

345 When we define the data mixture space as \mathcal{P}_{fix}
 346 and employ MLP as the fitting function, the opti-
 347 mization objective in Eq. 1 can be reformulated as
 348 follows:

$$349 \mathbf{p}_{fix}^* = \arg \max_{\mathbf{p} \in \mathcal{P}_{fix}, t \leq T} \frac{1}{k} \sum_{j=1}^k f_{MLP}^j(\mathbf{p}, t; \theta_0), \quad (4)$$

350 where j denotes j th downstream task. Given the
 351 negligible inference cost of MLP models, DaMo
 352 can efficiently extrapolate the optimal data mixture.
 353 We first iterate through all possible data mixtures
 354 in the \mathcal{P}_{fix} space to predict downstream task per-
 355 formance scores. Subsequently, we sort these pre-
 356 dicted scores and select the top- k highest-scoring
 357 mixtures to train our MLLM. This approach en-
 358 ables us to systematically identify the optimal data
 359 mixture without exhaustive empirical testing. The
 360 complete algorithm pseudocode is provided in Ap-
 361 pendix D.

362 5 Experiments

363 5.1 Experiments Settings

364 **Training datasets** Please refer to Appendix B.

Implementation Details We used InternVL2.5-
 4B (Chen et al., 2024b) as the base model. Please
 refer to Appendix C.1 for details.

Baseline We selected three heuristic approaches
 commonly used in industry and one representa-
 tive loss-based exponential fitting method. **Uni-
 form Mixture**: All datasets are sampled with equal
 weights. **Natural Mixture**: Sampling weights are
 proportional to the size of each dataset. **Random
 Mixture**: All datasets are sampled with random
 weights. **Data Mixing Laws (DML)** (Ye et al.,
 2024): An exponential function-based loss fitting
 method to predict the optimal mixture.

Downstream Task Evaluation Besides
 PhoneAgentBench, we further evaluated our
 method on four widely used open-source bench-
 marks to verify generalization, including BFCL-
 V3 (Patil et al., 2025a), MME-perception (Fu et al.,
 2023), MME-reasoning (Yuan et al., 2025) and
 OCRBench (Liu et al., 2024d). All metrics of
 these benchmarks are expressed as percentages
 (0-100%), with higher values indicating superior
 performance.

388 5.2 Fitting Score of Neural Network

389 As analyzed theoretically in Section 4.1, for $m =$
 390 12 training datasets and a batch size of $b = 16$,
 391 the discrete space \mathcal{P}_{fix} encompasses approximately
 392 1.3×10^7 potential data mixtures. This combina-
 393 torial explosion poses a significant challenge for
 394 performance prediction. To determine the mini-
 395 mum sample size N_s required to characterize this
 396 vast space, we gradually increased the training sam-
 397 ples for the MLP. As shown in Table 1, even with
 398 only 250 mixture samples, the MLP achieves an
 399 R^2 (Wright, 1921) score of 0.81 in 10-fold cross-
 400 validation. This indicates that the MLP can accu-
 401 rately extrapolate the performance landscape from
 402 a very sparse sampling set.

403 The theoretical justification for this sparse sam-
 404 pling efficiency lies in the approximation prop-
 405 erties of L -Lipschitz functions (Kolmogorov and
 406 Tikhomirov, 1959; Wainwright, 2019). According
 407 to the theory of metric entropy, the sample complex-
 408 ity N_s required to approximate such a function with
 409 a uniform error ε is governed by the covering num-
 410 ber, which scales as $N_s \propto (L/\varepsilon)^m$ (Kolmogorov
 411 and Tikhomirov, 1959; Yarotsky, 2017). Crucially,
 412 this bound is determined by the **intrinsic dimen-
 413 sion** m and the **hypersurface smoothness** (Lips-
 414 chitz constant L), rather than the total cardinality

of the discrete domain $|\mathcal{P}_{\text{fix}}|$ (Anthony and Bartlett, 1999; Wainwright, 2019). Our empirical results (high R^2 with small N_s) suggest that the mapping function between downstream task performance and data mixture is inherently smooth (possessing a small L). This smoothness allows a limited number of representative samples to effectively capture the global topology, decoupling the required sampling size from the search space’s combinatorial explosion.

Mixture number	H20 hours	Score (R^2)
50	872	0.58
100	1817	0.57
150	2581	0.74
200	3521	0.78
250	4225	0.81

Table 1: MLP fitting dynamics.

5.3 Downstream Task Performance of Unseen Data Mixtures

Through selecting mixture with top-1 predicting score on PhoneBenchAgent and open-source benchmarks to train MLLM, we obtain the performance on PhoneBenchAgent and open-source benchmarks, as shown in Table 2. DaMo achieves 23.35% improvement over the native model (without SFT) on PhoneAgentBench. When general capabilities are concurrently considered, DaMo yields an overall average score improvement of 13.73% over the native model. On PhoneAgentBench, DaMo surpasses the previous best-performing Natural mixture by 2.85%. Across overall datasets, it outperforms best-performing method, DML (Ye et al., 2024), by 3.06%. This performance advantage substantiates the effectiveness of fitting the relationship between data mixtures and downstream performance. When predicting the top-1 mixture on a single dataset, DaMo(*) further outperforms DaMo. This is because DaMo(*) is dedicated to predicting optimal data mixtures for single task, thereby mitigating the performance compromises inherent in multiple tasks.

Figure 4(a) shows the score distribution of Random(250). We can observe two critical characteristics: (1) The absence of a right-side long tail indicates that excellent data mixtures are extremely sparse. (2) The performance of random mixture is predominantly mediocre, and baseline methods (vertical dashed line) show no discernible advantage, demonstrating the inefficiency of heuristic

approaches. We used DaMo to predict across \mathcal{P}_{fix} space, selected the top 50 data mixtures with the best predicted performance, and conducted actual training and evaluation on MLLM. The score distribution of their performance is shown in Figure 4(b), which indicates that DaMo successfully identifies data mixtures with significantly higher average scores compared to all the baselines.

To study the generalization of DaMo on general tasks, we employ DaMo to predict MLLM’s performance only on open-source benchmarks, and use the top-1 data mixture to train MLLM, reporting the results in Table 3. It can be observed that our DaMo achieves remarkably superior performance across all open-source benchmarks compared to baselines. Also, it can be observed that focusing on task-specific objectives leads to significantly greater improvements. This is clearly demonstrated by the performance growth from 34.69% obtained by Uniform mixture to 47.43% on the BFCL-V3 benchmark, implemented by DaMo (BFCL-V3) which predicts the performance on BFCL-V3 benchmark only to search optimal data mixture. Crucially, this enhancement is sustained even in the absence of any task-curated training data. We posit that the observed performance benefit is fundamentally driven by DaMo exploring optimal mixtures, which orchestrates a balanced advancement across both specialized and generalizable capabilities.

5.4 Results on Unseen Datasets

To assess DaMo’s performance on unseen data, we evaluated the model trained using its predicted top-1 data mixture on unseen open-source benchmarks. Table 4 lists the performance of DaMo as well as other methods on BFCL-V4 (Patil et al., 2025b), MMBench (Yuan Liu, 2023), DocVQA (Mathew et al., 2021), and MMMU (Yue et al., 2024). Compared to w/o SFT, DaMo achieves substantial improvements on BFCL-V4 and MMMU, comparable performance on MMBench and DocVQA. The result indicates that DaMo improves model performance on unseen datasets. Moreover, DaMo presents significantly superior performance than Uniform, Natural, Random, and DML mixtures on all datasets, which demonstrates the superiority of DaMo predicting optimal data mixture for unseen datasets.

5.5 Extension to Other Models

We are concerned with the effective generalization of the DaMo to other models. Most

Method	MT-Plan	APP-Rec	MM-RR	ACU	MM-NER	Mobile-FC	OS Avg.	PAB Avg.	Overall Avg.
w/o SFT	20.00	6.00	65.38	39.18	84.08	54.31	68.77	44.83	54.40
uniform	54.50	56.00	44.62	86.37	81.71	45.92	55.76	61.52	59.22
natural	47.00	46.00	86.15	83.10	79.83	49.88	59.95	65.33	63.18
random(250)	45.00	43.00	83.08	82.77	80.26	47.32	66.36	63.57	64.89
DML (Ye et al., 2024)	52.00	43.00	85.38	85.72	80.01	42.66	65.48	64.80	65.07
DaMo	55.50	51.00	86.15	85.30	83.34	47.79	68.05	68.18	68.13
DaMo(*)	62.00	67.00	90.00	88.37	85.84	64.80	/	/	/

Table 2: Main results on PhoneAgentBench and open-source benchmarks by using top-1 data mixture to train MLLM, predicted by DaMo on PhoneAgentBench and open-source benchmarks. Random(250) denotes the metric is obtained by getting the best performance among 250 fine-tuned models with different random mixtures.

*: These scores correspond to different checkpoints, which are optimized by DaMo on a single task.

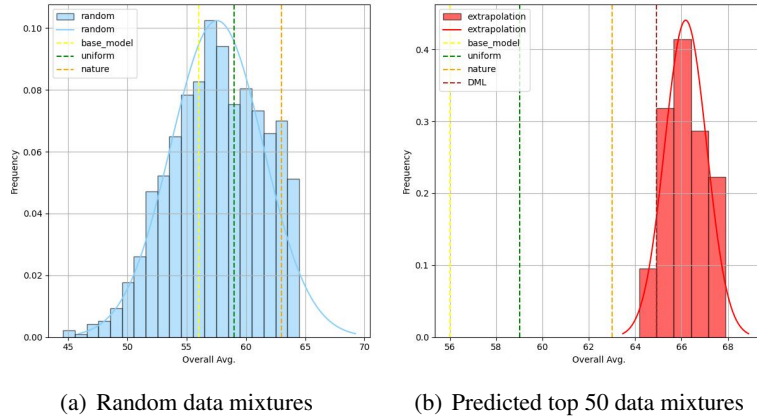


Figure 4: Probability distributions of overall average scores across different checkpoints.

Method	BFCL-V3	MME-perception	MME-reasoning	OCRBench	OS Avg.
w/o SFT	29.32	83.82	79.42	82.50	68.77
Uniform	34.69	58.63	64.91	64.80	55.76
Natural	31.41	75.47	67.01	65.90	59.95
Random(250)	33.59	81.11	74.55	78.20	66.36
DML (Ye et al., 2024)	25.47	83.31	76.34	76.8	65.48
DaMo	43.15	84.53	80.94	83.60	73.06
DaMo (*)	47.43	85.12	82.54	83.90	/

Table 3: Main results on open-source benchmarks of MLLMs trained by predicted optimal data mixture on open-source benchmarks.

Method	BFCL-V4	MMBench	DocVQA	MMMU
w/o SFT	17.86	78.79	91.08	49.33
Uniform	22.12	76.24	88.73	50.22
Natural	22.94	76.70	88.96	51.33
Random(250)	20.07	76.63	88.19	50.67
DML (Ye et al., 2024)	19.01	75.15	86.56	51.33
DaMo	28.10	79.18	91.04	52.44

Table 4: Main results of DaMo on unseen datasets.

Model	w/o SFT	Uniform	Natural	DML	DaMo (orig.)	DaMo (lin.)
Qwen2.5VL-3B-Inst.	56.25	65.15	64.82	65.03	68.02	68.66
Qwen2.5VL-7B-Inst.	59.43	67.48	65.99	66.37	67.79	69.09
InternVL3-14B	67.84	63.56	67.8	66.45	68.86	69.75

Table 5: Main results of transferability testing on PhoneAgentBench and open-source Benchmarks.

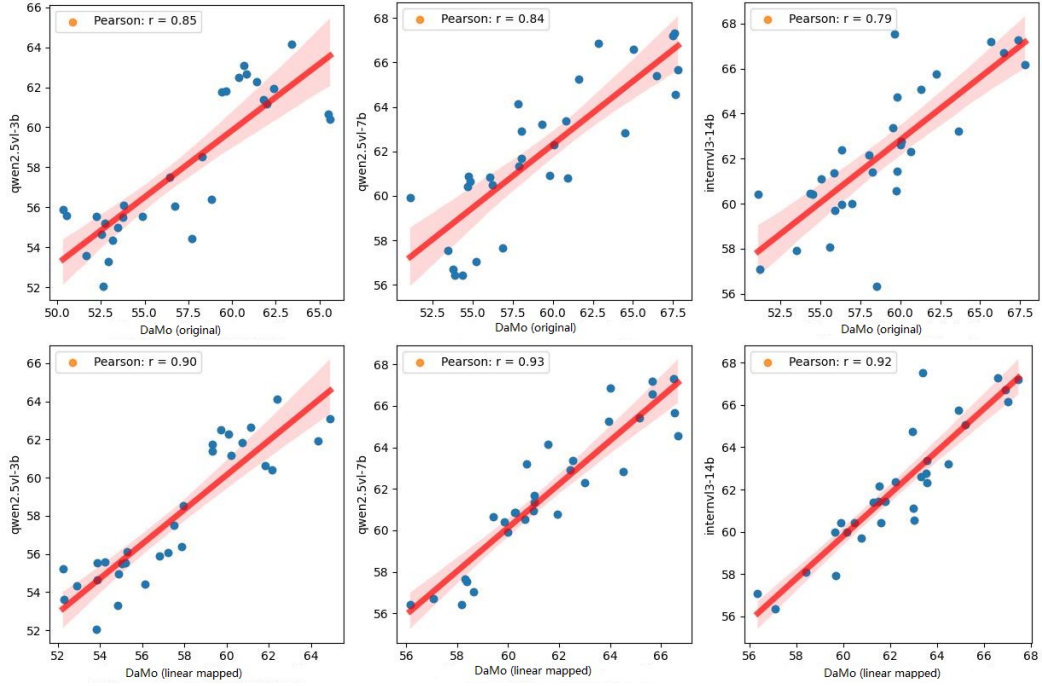


Figure 5: Transferability Analysis. Top: Scatter plot comparing the predicted overall average scores by the original DaMo against the actual scores of target models. Bottom: Apply linear-mapped correction to DaMo.

current works on data mixture during the pre-training phase assume that data mixture strategies can be directly transferred from smaller models to larger ones (Xie et al., 2023a), but their applicability in the supervised fine-tuning phase remains unverified. To this end, we conducted experiments on transferring DaMo obtained from InternVL2.5-4B to Qwen2.5VL-3B-Instruct, Qwen2.5VL-7B-Instruct (Bai et al., 2025b), and InternVL3-14B (Zhu et al., 2025) with zero or minimal extra training cost.

As outlined in Table 5, directly using the predicted best data mixture for InternVL2.5-4B to train Qwen2.5VL-3B, Qwen2.5VL-7B, and InternVL3-14B, DaMo still outperforms Uniform, Natural, and DML by 0.31%~2.87% score for all models, demonstrating the stable transferability of DaMo and it is more efficient than other methods. Note that we omit Random mixture for comparison as it needs large training resource for a new model to get a good mixture. To mitigate biases caused by discrepancies in model capabilities, we fit a compensating linear layer using only 20 calibration samples for each new model. The linear-mapped DaMo is defined as $g = f(\cdot)\mathbf{w} + b$, termed as DaMo (lin). It can be seen in Table 5 that DaMo (lin) further improves Uniform, Natural and DML by 1.61%~3.51% score.

Figure 5 shows the Pearson correlation coeffi-

cients (r) between the predicted overall average scores by DaMo and the actual scores of target models. The coefficients r are generally above 0.75, demonstrating the robust cross-model applicability of DaMo even without extra training cost. This suggests that optimal mixtures identified for the base model likely remain near-optimal for the target models. After applying linear mapping, the discrepancies between models are reduced, leading to a further enhancement in correlation with r increasing to above 0.9.

6 Conclusion

In this paper, we present the Data Mixing Optimizer (DaMo) to optimize data mixtures in multitask fine-tuning of multimodal large language models. By introducing downstream task performance prediction with neural network-based modeling, DaMo can predict model performance for any given data mixture. To support comprehensive evaluation, we introduce PhoneAgentBench for evaluation of multimodal large language models on phone agentic tasks. Moreover, DaMo can be extended to other models and tasks. Experimental results demonstrate the efficacy of DaMo not only on PhoneAgentBench, but also on general benchmarks, outperforming state-of-the-art methods.

562
563
564
565
566

567

568
569
570
571

572
573
574
575
576

577
578
579
580

581
582
583

584
585
586
587
588
589
590

591
592
593
594

595
596
597
598

599
600
601
602

603
604
605
606
607

608
609
610
611
612
613

Limitation

Due to the limited computing resources and time, we can not test the method’s effectiveness on recently released large models, e.g., Qwen3-VL-235B-A22B-Instruct (Bai et al., 2025a).

References

Abhishek Shrivastava. 2023. sentiment-analysis-dataset. <https://www.kaggle.com/datasets/abhi8923shriv/sentiment-analysis-dataset>. Accessed: 2025-03-12.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Alon Albalak, Liangming Pan, Colin Raffel, and William Yang Wang. 2023. Efficient online data mixing for language model pre-training. *arXiv preprint arXiv:2312.02406*.

Martin Anthony and Peter L Bartlett. 1999. *Neural network learning: Theoretical foundations*. Cambridge University Press.

Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, and 45 others. 2025a. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025b. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Kinjal Basu. 2024. Bridging knowledge gaps in llms via function calls. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 5556–5557.

David Castillo-Bolado, Joseph Davidson, Finlay Gray, and Marek Rosa. 2024. Beyond prompts: Dynamic conversational benchmarking of large language models. *Preprint*, arXiv:2409.20222.

Mingyang Chen, Haoze Sun, Tianpeng Li, Fan Yang, Hao Liang, Keer Lu, Bin Cui, Wentao Zhang, Zenan Zhou, and Weipeng Chen. 2024a. Facilitating multi-turn function calling for llms via compositional instruction tuning. *arXiv preprint arXiv:2410.12952*.

Zehui Chen, Weihua Du, Wenwei Zhang, Kuikun Liu, Jiangning Liu, Miao Zheng, Jingming Zhuo, Songyang Zhang, Dahua Lin, Kai Chen, and 1 others. 2023. T-eval: Evaluating the tool utilization capability of large language models step by step. *arXiv preprint arXiv:2312.14033*.

Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, and 1 others. 2024b. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.

Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. 2024. Seeclick: Harnessing gui grounding for advanced visual gui agents. *arXiv preprint arXiv:2401.10935*.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and 1 others. 2023. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*.

Longxi Gao, Li Zhang, Shihe Wang, Shangguang Wang, Yuanchun Li, and Mengwei Xu. 2024. Mobileviews: A large-scale mobile gui dataset. *arXiv preprint arXiv:2409.14337*.

Ce Ge, Zhijian Ma, Daoyuan Chen, Yaliang Li, and Bolin Ding. 2024. Bimix: Bivariate data mixing law for language model pretraining. *arXiv preprint arXiv:2405.14908*.

Longling Geng and Edward Y Chang. 2025. Realm-bench: A real-world planning benchmark for llms and multi-agent systems. *arXiv preprint arXiv:2502.18836*.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Shuhao Gu, Jialing Zhang, Siyuan Zhou, Kevin Yu, Zhaohu Xing, Liangdong Wang, Zhou Cao, Jintao Jia, Zhuoyi Zhang, Yixuan Wang, Zhenchong Hu, Bo-Wen Zhang, Jijie Li, Dong Liang, Yingli Zhao, Yulong Ao, Yaoqi Liu, Fangxiang Feng, and Guang Liu. 2024. Infinity-mm: Scaling multimodal performance with large-scale and high-quality instruction data. *Preprint*, arXiv:2410.18558.

Hamed Rahimi. 2024. Sujetfinancevision10k. <https://huggingface.co/datasets/sujet-ai/Sujet-Finance-Vision-10k>. Accessed: 2025-03-12.

Chen Huang, Shuangfei Zhai, Walter Talbott, Miguel Angel Bautista, Shih-Yu Sun, Carlos Guestrin, and Josh Susskind. 2019. Addressing the loss-metric mismatch with adaptive loss alignment. *Preprint*, arXiv:1905.05895.

Tim Ingold. 2021. *The perception of the environment: essays on livelihood, dwelling and skill*. routledge.

Berivan Isik, Natalia Ponomareva, Hussein Hazimeh, Dimitris Pappas, Sergei Vassilvitskii, and Sanmi

668	Koyejo. 2024. Scaling laws for downstream task performance of large language models. In <i>ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models</i> .	Task planning with clusters across multiple tools. <i>arXiv preprint arXiv:2406.03807</i> .	722
669			723
670			
671			
672	Andrei N Kolmogorov and Vladimir M Tikhomirov. 1959. ε -entropy and ε -capacity of sets in function spaces. <i>Uspekhi Matematicheskikh Nauk</i> , 14(2):3–86.	Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. 2024d. <i>Ocr-bench: on the hidden mystery of ocr in large multimodal models</i> . <i>Science China Information Sciences</i> , 67(12).	724
673			725
674			726
675			727
676	Hao Li, Chenghao Yang, An Zhang, Yang Deng, Xiang Wang, and Tat-Seng Chua. 2024a. Hello again! llm-powered personalized agent for long-term dialogue. <i>arXiv preprint arXiv:2406.05925</i> .		728
677			729
678			
679			
680	Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo, Yuchen Tian, Jing Ma, Zhiyong Huang, and Tat-Seng Chua. 2025a. Screenspot-pro: Gui grounding for professional high-resolution computer use. <i>arXiv preprint arXiv:2504.07981</i> .	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Taffjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. <i>Advances in Neural Information Processing Systems</i> , 35:2507–2521.	730
681			731
682			732
683			733
684			734
685	Lingyu Li, Yixu Wang, Haiquan Zhao, Shuqi Kong, Yan Teng, Chunbo Li, and Yingchun Wang. 2024b. Reflection-bench: probing ai intelligence with reflection. <i>arXiv preprint arXiv:2410.16270</i> .	Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 2200–2209.	735
686			736
687			737
688			738
689	Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023. <i>Api-bank: A comprehensive benchmark for tool-augmented llms</i> . <i>Preprint</i> , arXiv:2304.08244.	Grégoire Mialon, Clémentine Fourier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023. <i>Gaia: a benchmark for general ai assistants</i> . <i>Preprint</i> , arXiv:2311.12983.	739
690			740
691			741
692			742
693			743
694	Yuan Li, Zhengzhong Liu, and Eric Xing. 2025b. Data mixing optimization for supervised fine-tuning of large language models. <i>arXiv preprint arXiv:2508.11953</i> .	Paul Michel, Sebastian Ruder, and Dani Yogatama. 2021. Balancing average and worst-case accuracy in multitask learning. <i>arXiv: Learning, arXiv: Learning</i> .	744
695			745
696			746
697			747
698	Liang Xu. 2024. Superclue-agent. https://github.com/CLUEbenchmark/SuperCLUE-Agent . Accessed: 2025-03-12.	minyang. 2024. invoices-and-receipts_ocr_v1. https://huggingface.co/datasets/mychen76/invoices-and-receipts_ocr_v1 . Accessed: 2025-03-12.	748
699			749
700			750
701	P Liashchynskiy. 2019. Grid search, random search, genetic algorithm: A big comparison for nas. <i>arXiv preprint ArXiv:1912.06059</i> .	Niccolò Zanichelli. 2024. arxiv-ocr-v0.1.2. https://huggingface.co/datasets/nz/arxiv-ocr-v0.1.2 . Accessed: 2025-03-12.	751
702			752
703			753
704	Qian Liu, Xiaosen Zheng, Niklas Muennighoff, Guangtao Zeng, Longxu Dou, Tianyu Pang, Jing Jiang, and Min Lin. 2024a. Regmix: Data mixture as regression for language model pre-training. <i>arXiv preprint arXiv:2407.01492</i> .	Haojie Pan, Zepeng Zhai, Hao Yuan, Yaojia Lv, Ruiji Fu, Ming Liu, Zhongyuan Wang, and Bing Qin. 2023. Kwaiaagents: Generalized information-seeking agent system with large language models. <i>arXiv preprint arXiv:2312.04889</i> .	754
705			755
706			756
707			757
708			758
709	Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, and 1 others. 2023. Agent-bench: Evaluating llms as agents. <i>arXiv preprint arXiv:2308.03688</i> .	Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. 2025a. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In <i>Forty-second International Conference on Machine Learning</i> .	759
710			760
711			761
712			762
713			763
714	Xiao Liu, Tianjie Zhang, Yu Gu, Iat Long Iong, Yifan Xu, Xixuan Song, Shudan Zhang, Hanyu Lai, Xinyi Liu, Hanlin Zhao, and 1 others. 2024b. Visualagent-bench: Towards large multimodal models as visual foundation agents. <i>arXiv preprint arXiv:2408.06327</i> .	Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. 2025b. The berkeley function calling leaderboard (bfcl) v4: From tool use to agentic evaluation of large language models. In <i>Forty-second International Conference on Machine Learning</i> .	764
715			765
716			766
717			767
718			768
719	Yanming Liu, Xinyue Peng, Jiannan Cao, Shi Bo, Yuwei Zhang, Xuhong Zhang, Sheng Cheng, Xun Wang, Jianwei Yin, and Tianyu Du. 2024c. Tool-planner:	qgyd2021. 2024a. chinese_ner_sft. https://huggingface.co/datasets/qgyd2021/chinese_ner_sft . Accessed: 2025-03-12.	769
720			770
721			771
			772
			773
			774
			775

776	qgyd2021. 2024b. few_shot_ner_sft. https://huggingface.co/datasets/qgyd2021/few_shot_ner_sft . Accessed: 2025-03-12.	829
777		830
778		831
779	Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. Toollm: Facilitating large language models to master 16000+ real-world apis . <i>Preprint</i> , arXiv:2307.16789.	832
780		833
781		834
782		835
783		836
784		837
785		838
786	Shiori Sagawa, PangWei Koh, Tatsunori Hashimoto, and Percy Liang. 2020. Distributionally robust neural networks. <i>International Conference on Learning Representations, International Conference on Learning Representations</i> .	839
787		840
788		841
789		842
790		843
791	shibing624. 2023. sharegpt_gpt4. https://huggingface.co/datasets/shibing624/sharegpt_gpt4 . Accessed: 2025-03-12.	844
792		845
793		846
794	Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. 2023. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 2998–3009.	847
795		848
796		849
797		850
798		851
799		852
800	Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2023. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. <i>Advances in Neural Information Processing Systems</i> , 36:38975–38987.	853
801		854
802		855
803		856
804		857
805		858
806	Martin J Wainwright. 2019. <i>High-dimensional statistics: A non-asymptotic viewpoint</i> . Cambridge University Press. See Chapter 5 on Metric Entropy and Covering Numbers.	859
807		860
808		861
809		862
810	Xuehui Wang, Zhenyu Wu, JingJing Xie, Zichen Ding, Bowen Yang, Zehao Li, Zhaoyang Liu, Qingyun Li, Xuan Dong, Zhe Chen, and 1 others. 2025. Mmbench-gui: Hierarchical multi-platform evaluation framework for gui agents. <i>arXiv preprint arXiv:2507.19478</i> .	863
811		864
812		865
813		866
814		867
815		868
816	Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, Haiteng Zhao, Yongfei Liu, Bohan Zhai, Jianbo Yuan, Quanzeng You, and Hongxia Yang. 2024. Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning. <i>arXiv preprint arXiv:2401.06805</i> .	869
817		870
818		871
819		872
820		873
821		874
822		875
823	Sewall Wright. 1921. Correlation and causation. <i>Journal of agricultural research</i> , 20(7):557.	876
824		877
825	Yuxin Xiao, Shujian Zhang, Wenxuan Zhou, Marzyeh Ghassemi, and Sanqiang Zhao. 2024. Sftmix: Elevating language model instruction tuning with mixup recipe. <i>arXiv preprint arXiv:2410.05248</i> .	878
826		879
827		880
828		881
	Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. 2023a. Doremi: Optimizing data mixtures speeds up language model pretraining. <i>Advances in Neural Information Processing Systems</i> , 36:69798–69818.	882
		883
		884
	Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. 2023b. Doremi: Optimizing data mixtures speeds up language model pretraining . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 69798–69818. Curran Associates, Inc.	885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

for open-source multimodal models. *arXiv preprint arXiv:2504.10479*.

Tong Zhu, Daize Dong, Xiaoye Qu, Jiacheng Ruan, Wenliang Chen, and Yu Cheng. 2024. Dynamic data mixing maximizes instruction tuning for mixture-of-experts. *arXiv preprint arXiv:2406.11256*.

A Evaluation datasets

Table 6 summarizes the statistic information of the tasks in PhoneAgentBench. To guarantee the faithfulness of the proposed PhoneAgentBench, we implemented a rigorous workflow encompassing data filtering, synthetic data generation, and manual verification. Details information about our evaluation datasets for PhoneAgentBench are as follows.

Dataset	Evaluation ability	Data size
MT-Plan	Multimodal Task Planning	100
MM-RR	Multimodal Reference Resolution	130
ACU	Agent Context Understand	100
MM-NER	Multimodal Named Entity Recognition	376
APP-Rec	APP Recognition	100
Mobile-FC	Mobile Function Calling	429

Table 6: The statistics of PhoneAgentBench.

A.0.1 MultiModal Task Planning

We introduce the two metrics of complexity and diversity to evaluate the quality of the benchmark for the task planning.

- **Complexity:** The answer of MT-Planning can be viewed as a directed acyclic graph (DAG), where each subtask is a node and the dependency relationship between subtasks are edges. Thus, complexity can be expressed as n_{edge}/n_{node} .
- **Diversity:** The higher the similarity between queries in a dataset, the lower the diversity of that dataset. We use Rough-L to calculate the similarity between every pair of queries, and diversity can be expressed as $1 - \frac{1}{N(N-1)/2} \sum_{i \neq j} \text{Rough-L}(q_i, q_j)$.

Based on this, we compared the data complexity and diversity between MT-Plan and T-Eval planning.

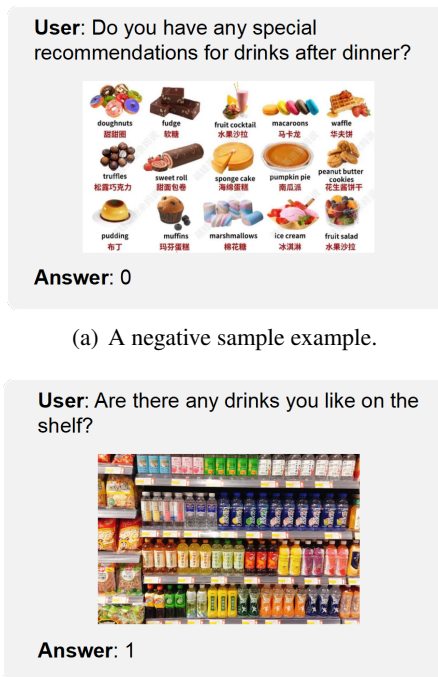
Benchmark	complexity \uparrow	diversity \uparrow
MT-Plan	0.661	0.82
T-Eval (Chen et al., 2023)	0.122	0.73

Table 7: Benchmark metrics.

A.0.2 MultiModal Reference Resolution

The MultiModal Reference Resolution (MM-RR) task requires the model to determine whether the current question refers to information in the image, which is a binary classification task.

As shown in Figure 6, the question in Figure 6(a) does not refer to the content in the image, so the answer is 0; while the question in Figure 6(b) refers to the drinks on the shelf in the image, so the answer is 1.



(a) A negative sample example.

(b) A positive sample example.

Figure 6: Examples of RR dataset.

A.0.3 Multimodal NER

Multimodal NER (MM-NER) benchmark quantitatively measures MLLMs' ability in understanding and extracting key entities. The dataset comprises 376 image-only samples sourced from Baidu's publicly available image repositories, where each image underwent a rigorous curation process: professional annotators manually filtered the raw visual data to retain high-quality, clearly discernible images, which were subsequently annotated with precise labels focusing on seven critical entity categories—temporal references, geographical locations, personal identifiers, contact numbers, tracking Number, flight Number, train Number to establish a structured benchmark for multimodal entity recognition. We adopt the entity F1-score as the evaluation metric. Figure 7 demonstrates time, lo-

945

cation and person extraction from chat logs.

User: Extract the time, location, person's name, telephone number, express number, flight number and train number in the image.

Answer: {
 "time": [
 "11:50",
 "下午5点"
],
 "location": [
 "公司的会议室"
],
 "name": [
 "周总",
 "小榕"
]
}

Figure 7: An example of MM-NER dataset.

A.0.4 Mobile Function Call

The Mobile Function Call (Mobile-FC) task is designed to evaluate the ability of MLLMs to call mobile API functions. The task requires the model to select appropriate functions from a given set of application functions to call according to the user's app instruction questions and output the parameters required for the function calls. We define 50 function call interfaces for different scenarios, such as setting an alarm, checking the weather, and setting navigation. The questions in the data are manually constructed by annotators, simulating real-world scenarios of apps on smartphone operating systems and forming complete multi-round dialogues. The evaluation method mainly compares the predicted function names and parameter names with the annotated results. A perfect match scores 1 point; otherwise, 0 points. As shown in the Figure 8, we define the function name create alarm for setting an alarm, with the time field as the input parameter.

```
User: Set an alarm for 10:00 in the morning.
Assistant: create_alarm(time=10:00)
User: An alarm for 03:00.
Assistant: create_alarm(time=03:00)
```

Figure 8: An example of Mobile-FC dataset.

A.0.5 Agent Context Understanding

The Agent Context Understanding (ACU) task is used to assess the context-aware dialogue comprehension ability of MLLMs. The data is presented in the form of multi-turn conversations (including text and image).

The model is required to resolve the anaphoric information in the user's final question based on multi-turn conversations or image information, and output a question that contains no anaphora. As shown in the Figure 9, the user asks "Do you like his songs?". If no image is provided, the

model needs to determine who "he" refers to based on the historical conversation. Otherwise, the model needs to recognize the person in the image. Model's output is a question that contains no referential information. We use the BLEU of the output answer with the reference answer to evaluate task performance, with scores ranging from 0 to 1.

(a) Pure-text conversation (b) Multimodal conversation sample.

Figure 9: Examples of ACU dataset.

A.0.6 APP Recognition

The APP Recognition (APP-Rec) task, similar to the APP-Rec training set, is used to evaluate the ability of MLLMs to identify mobile applications. The model is required to directly output the APP name based on the content of the input mobile APP interface image, as illustrated in the Figure 11. The performance evaluation is conducted by comparing the overlap between the predicted application name and the annotated result. A correct prediction scores 1 point; otherwise, 0 points.

B Training datasets

The open source data includes: ShareGPT4 (shibing624, 2023), NER (composed of Chinese-NER-SFT (qgyd2021, 2024a), Sentiment-Analysis (Abhishek Shrivastava, 2023), and Few-Shot-NER-SFT (qgyd2021, 2024b)), Infinity-MM (Gu et al., 2024), OCR (consisting of Vision-OCR-Financial-Reports-10K (Hamed Rahimi, 2024), Arxiv-OCR-v0.1-SFT (Niccolò Zanichelli, 2024) and Invoices-and-Receipts-OCR-v1 (minyang, 2024)), and SuperCLUE-Agent (Liang Xu, 2024).

The self-built datasets include MultiModal-Instruction-Evolution (MMIE), APP Recognition (APP-Rec), Reference-Resolution (RR), MultiModal-Understanding (MMU), Function-Calling (FC), Task-Planning (TP), and Image-Text-Relevance (ITR), which are primarily derived from data synthesis and real-world industrial scenarios. The size of samples in all the training data is shown in Table 8.

B.0.1 Multimodal Instruction Evolution

The Multimodal Instruction Evolution (MMIE) task consists of 1.8K pieces of multimodal

13

Dataset	Source	Data size	Dataset	Source	Data size
MMIE	self-built	1.8k	APP-Rec	self-built	22.8k
MMU	self-built	21.1k	RR	self-built	10.5k
TP	self-built	26.8k	FC	self-built	10.4k
ITR	self-built	9.7k	ShareGPT4	open-source	36k
NER	open-source	8k	Infinity-MM	open-source	37.2k
OCR	open-source	33k	SuperCLUE-Agent	open-source	1.5k

Table 8: Training dataset sizes.

question-answering data. As shown in Figure 10, given an initial query and image with several available tools, the methodology requires the model to generate more sophisticated and diversified questions. The generation pipeline comprises six structured phases:

- **Intent analysis:** Analyze the user’s potential needs from multiple perspectives.
- **Scenario expansion:** Expand the scenario to increase the diversity and complexity of the initial question.
- **Task decomposition:** Decompose the scenario into multiple subtasks which can be executed correctly by provided tools.
- **Raise new question:** Propose a new question based on the expanded scenario and subtasks.
- **Iterative Validation:** Evaluate completeness and complexity, where completeness indicates whether the question adequately covers the steps of the subtasks.
- **Naturalization and output:** Refine questions to be more colloquial and output the final result.

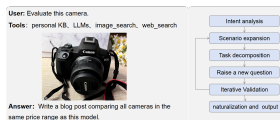


Figure 10: An example of MMIE dataset.

B.0.2 APP Recognition

The APP Recognition (APP-Rec) task consists of 22.8K pieces of multimodal question-answering data, which are composed of images and task instructions. The task requires the model to identify

the interface information of mobile apps in the input images and directly output the app names. To obtain diverse app interface data, we install 100 different applications on a mobile phone, such as WeChat, QQ, Little Red Book, Weibo, Alipay, Pinduoduo, Taobao, and TikTok. Annotators are then required to manually capture screenshots of different functional interfaces of each application, which serve as the image source for the APP-Rec task, as illustrated in Figure 11. The default input task instruction is "Identify which app the screenshot belongs to?", and the answer is the name of the app corresponding to the image.



Figure 11: An example of APP-Rec dataset.

B.0.3 Reference Resolution

The Reference Resolution (RR) task corresponds to the MM-RR task in Section A.0.2, which contains 10.5k pieces of multimodal question-answering data. We collect various images containing text information from the internet, with sources including academic papers, test questions, news, com-

pany official websites, Wikipedia, etc. Annotators design corresponding questions based on the text content in the given images as positive examples, while negative examples are obtained by replacing the images with different types, as shown in the following Figure 6, which provides one positive and one negative example respectively.

B.0.4 Multimodal Understanding

The Multimodal Understanding (MMU) tasks are consistent with ACU tasks in Section A.0.5. It takes the form of multimodal or text-only multi-round dialogues, with 1-4 rounds and a total of 21.1K samples. The images are sourced from publicly available internet data, covering various fields such as people, animals, plants, architecture, and digital products. The dialogue data is manually constructed by annotators based on the given images, focusing on reference problems. The task requires the model to combine the images and historical dialogue content to rewrite the user’s final input text. This is achieved by replacing pronouns or supplementing omitted content to make the text semantically complete.

B.0.5 Function Calling

The Function Calling (FC) task consists of plain text instructions, which requires selecting appropriate tools from given tool set and filling in correct parameters for executing. The tools involve practical mobile applications such as unit conversion, weather inquiry, time calculation, text creation, recipe search, mobile phone bill inquiry, and other 500 types of useful tools. Notably, 90% of the instructions only require the invocation of a single tool.

Here is an example in Figure 12: The user inquires how much 500 US dollars is in Japanese yen, and the answer includes thoughts and actions. The thought process briefly outlines the current step, while the action first provides the name of the selected tool and sets the actual parameters in the action input.

```

User: I am planning a trip to Japan next month and wonder how much Japanese yen I can exchange for 500 USD.
Answer: [
  {
    "Thought": "Check how many Japanese yen can be exchanged for 500 USD.",
    "Action": "exchange_rate",
    "Action-Input": {
      "money": "500",
      "currency": "USD",
      "target_currency": "JPY"
    }
  }
]

```

Figure 12: An output example of FC dataset.

B.0.6 Task Planning

The Task Planning (TP) dataset, also targeting tool calling scenarios, places greater emphasis on multi-stage operations with inter-dependent steps compared to FC. It involves 26.8K pieces of multimodal question-answering data. This dataset requires models to properly decompose complex problems into solvable subtasks while ensuring correct tool selection and execution. In multistep scenarios, managing inter-parameter dependencies becomes critical.

The input is a complex question requiring calling apps on mobile phone. Output contains multistep thinking and actions similar to FC, and symbols start with #E are used to receive parameter for cited in subsequent tasks. (as demonstrated in Figure 2).

B.0.7 Image-Text Relevance



(a) A negative sample example.



(b) A positive sample example.

Figure 13: Examples of ITR dataset.

The Image-Text Relevance (ITR) task involves 9.7K pieces of multimodal question-answering data. The task requires the model to analyze the relevance between the question and the image based on the characteristics of the question. If the image is relevant to the question and can be used to an-

1131 swer the user’s question, the model should answer
 1132 1; otherwise, 0. The images are sourced from pub-
 1133 licly available internet data, the same as those used
 1134 in the MMU task. Annotators manually construct
 1135 questions related to the image content as positive
 1136 examples. For example, for images of people, ques-
 1137 tions about names, works, or family relationships
 1138 can be asked. Negative examples are constructed
 1139 by replacing the images with different types, as
 1140 shown in the following Figure 13, which presents
 1141 one positive and one negative example respectively.

1142 C Details of Experiments Setting

1143 C.1 Implement Details

1144 We applied a series of experiments to verify the
 1145 effectiveness of DaMo. Initially, we conducted
 1146 training and evaluation on InternVL2.5-4B (Chen
 1147 et al., 2024b) to obtain fitting samples for the
 1148 MLP. Specifically, we first sampled 250 random
 1149 data mixtures \mathbf{p} from \mathcal{P}_{fix} . For each mixture,
 1150 training was performed on 8 NVIDIA H20 GPUs,
 1151 and checkpoints were saved at 4 distinct training
 1152 steps—resulting in a total of 1000 checkpoints. All
 1153 1000 checkpoints were then evaluated on down-
 1154 stream tasks, which generated 1000 sample points
 1155 in the format of (\mathbf{p}, t, s) . The hyperparameters for
 1156 training the MLLM are listed in Table 9.

1157 Subsequently, we fitted the MLP on these 1000
 1158 sample points. MLP is structured as a two-
 1159 layer multi-layer perceptron (MLP) built upon
 1160 sklearn.MLPRegressor, where each of the two
 1161 hidden layers contains 100 neurons. To verify the
 1162 model’s fitting score, we assessed the coefficient
 1163 of determination (R^2) (Wright, 1921) of DaMo via
 1164 10-fold cross-validation. More details of MLP are
 1165 provided in Table 9.

1166 Then, we utilized DaMo to predict the down-
 1167 stream task performance of unseen data mixtures.
 1168 Leveraging the low inference cost of the MLP, we
 1169 conducted performance predictions for all mixtures
 1170 $\mathbf{p} \in \mathcal{P}_{\text{fix}}$. Among these, the data mixture with the
 1171 optimal predicted performance was selected for fur-
 1172 ther model training and validation, aiming to obtain
 1173 actual performance metrics.

1174 Finally, to verify the transferability of DaMo
 1175 on other models, we extended DaMo (based
 1176 on InternVL2.5-4B) to Qwen2.5VL-3B-Instruct,
 1177 Qwen2.5VL-7B-Instruct (Bai et al., 2025b), and
 1178 InternVL3-14B (Zhu et al., 2025). For these
 1179 new models, we trained a small number of ran-
 1180 dom mixtures, analyzed the correlation between

1181 DaMo’s predicted performance and the actual train-
 1182 ing performance, and meanwhile used DaMo to
 1183 find the optimal mixtures on the new models to
 1184 verify whether it still maintains competitiveness
 1185 compared with the baselines.

Model	Hyperparameters	setting
MLLMs	AdamW β_1	0.9
	AdamW β_2	0.95
	AdamW ϵ	$1e - 6$
	Max Sequence Length	16384
	Batch Size	16
	Gradient Accumulation Steps	8
	Training Steps	1440
	Warmup Steps	144
	Peak Learning Rate	$1e - 5$
	Weight Decay	0.1
Gradient Clipping	1.0	
MLP	Input Layer Dimension	12
	Hidden Layer 1 Dimension	100
	Hidden Layer 2 Dimension	100
	Output Layer Dimension	10
	Activation Function	ReLU
	Optimizer	Adam
	Learning Rate	$1e - 6$
	Training Steps	1500

Table 9: Hyperparameters of training.

1186 D Algorithm

1187 Algorithm D shows the algorithm of DaMo.

1188 E LLM usage

1189 In this paper, we used LLMs to polish the content
 1190 of the main text and appendices.

Algorithm 1: Algorithm of DaMo

Input: \mathcal{D} : training dataset; \mathcal{D}^{test} : test dataset; θ_0 : initial parameters of MLLM; \mathcal{P} : data mixing space, consisting of data mixture \mathbf{p} ; f_{MLP} : fitted MLP; t : training steps; \mathcal{M} : The data points for fitting MLP, consisting of pairs $\langle (\mathbf{p}, t), \mathbf{s} \rangle$.

Output: θ^* : MLLM trained with the optimal data mixture \mathbf{p}^* .

Initialize $\mathcal{M} \leftarrow \emptyset$

Randomly sample a small subset

$\mathcal{P}_{mlp} \subset \mathcal{P}_{fix}$

foreach $\mathbf{p}^i \in \mathcal{P}_{train}$ **do**

$\theta_t^i \leftarrow \text{Trainer}(\mathcal{D}, \mathbf{p}^i, t, \theta_0)$

$\mathbf{s}^i \leftarrow \text{Evaluator}(\theta_t^i, \mathcal{D}^{test})$

$\mathcal{M} \leftarrow \mathcal{M} \cup \{(\mathbf{p}^i, t, \mathbf{s}^i)\}$

end

$f_{MLP} \leftarrow \text{fit}(\mathcal{M})$

$\mathbf{p}^*, t^* \leftarrow \arg \max_{\mathbf{p} \in \mathcal{P}_{fix}} f_{MLP}(\mathbf{p}, t)$

$\theta^* \leftarrow \text{Trainer}(\mathcal{D}, \mathbf{p}^*, t^*, \theta_0)$

$\hat{\mathbf{s}} \leftarrow \text{Evaluator}(\theta^*, \mathcal{D}^{test})$

return $\theta^*, \hat{\mathbf{s}}$
