# MULTIWAY-ADAPTER: ADAPTING MULTIMODAL LARGE LANGUAGE MODELS FOR SCALABLE IMAGE-TEXT RETRIEVAL

*Zijun Long, George Killick, Richard McCreadie, Gerardo Aragon Camarasa*

The University of Glasgow, Scotland, UK

## ABSTRACT

As Multimodal Large Language Models (MLLMs) grow in size, adapting them to specialized tasks becomes increasingly challenging due to high computational and memory demands. Indeed, traditional fine-tuning methods are costly, due to the need for extensive, task-specific training. While efficient adaptation methods exist that aim to reduce these costs, in practice they suffer from shallow inter-modal alignment, which severely hurts model effectiveness. To tackle these computational challenges and improve inter-modal alignment, we introduce the MultiWay-Adapter (MWA), a novel framework featuring an 'Alignment Enhancer'. This enhancer deepens inter-modal alignment, enabling high transferability with minimal tuning effort. Our experiments show that unlike prior efficient tuning approaches, MWA maintains model effectiveness, while reducing training time by up-to 57%. MWA is also lightweight, increasing model size by only 2-3% (in terms of parameters) for state-of-the-art foundation models like BEiT-3 Large. These results demonstrate that MWA provides an efficient and effective adaptation method for MLLMs, significantly broadening their applicability.

***Index Terms***— Multimodal Large Language Models, Image-Text Retrieval, Adapter, Transformers, Transfer Learning

## 1. INTRODUCTION

Recent advancements in Multimodal Large Language Models (MLLMs), such as BLIP2 [1] and BEiT-3 [2], have demonstrated state-of-the-art performance in multimodal tasks, exemplified by their capabilities in Visual Question Answering. However, the adaptation of these MLLMs to specialized downstream tasks remains a substantial challenge, particularly for image-text retrieval, a common use-case in multimodal learning. Traditional full fine-tuning requires isolated, exhaustive retraining for each new task, demanding intensive computational resources and thus limiting practical applications. For instance, training BLIP2-Giant on an Nvidia A100 GPU takes 144 days [1].

Given the challenge of fine-tuning MLLMs, there is a growing need to develop efficient adaptation methods for MLLMs [3, 4]. While progress has been made in unimodal domains using adapter modules, these methods remain largely underexplored in multimodal contexts, particularly for image-text retrieval. Furthermore, existing adaptation methods for MLLMs [4, 5, 6] focus on information extraction from downstream datasets but neglect the critical need for inter-modal alignment. The goal of inter-modal alignment is to bring different modalities into a common feature space where they can be effectively compared, combined, or related. With shallow alignment, the model would fail to capture the complex inter-relations between different modalities, thereby impacting its effectiveness in multi-modal tasks [7, 8, 9].

To address the issue of shallow inter-modal alignment while preserving the efficiency advantages of adapter approaches, we introduce the MultiWay-Adapter (MWA), a lightweight yet effective framework designed explicitly for MLLMs adaptation. Additional components of MWA are small in size but bring a significant performance boost in transfer learning with minimal fine-tuning cost. Our key contributions include:

- We propose MWA that incorporates a dual-component approach, namely the New Knowledge Extractor and the Modality Enhancer. MWA not only extracts new knowledge from downstream datasets but also ensures deep inter-modal alignment, which is crucial for superior performance in vision-language tasks. To the best of our knowledge, this paper is the first work that mitigates the issue of shallow inter-modal alignment in adapter approaches for MLLMs.

- Through comprehensive experiments, we demonstrate that MWA achieves superior zero-shot performance on the Flickr30k dataset by tuning merely an additional 2.58% of parameters to the BEiT-3 Large model, saving up to 57% in fine-tuning time compared to full-model fine-tuning. MWA also demonstrates no statistically significant decreases in performance in other settings, compared to full fine-tuning, requiring significantly fewer resources.

- Experimental results demonstrate the robustness of MWA when parameters scale up, making it ready for MLLMs that are continually increasing in size.

- Our ablation study confirms the effectiveness of both MWA components, substantiating our design choices.
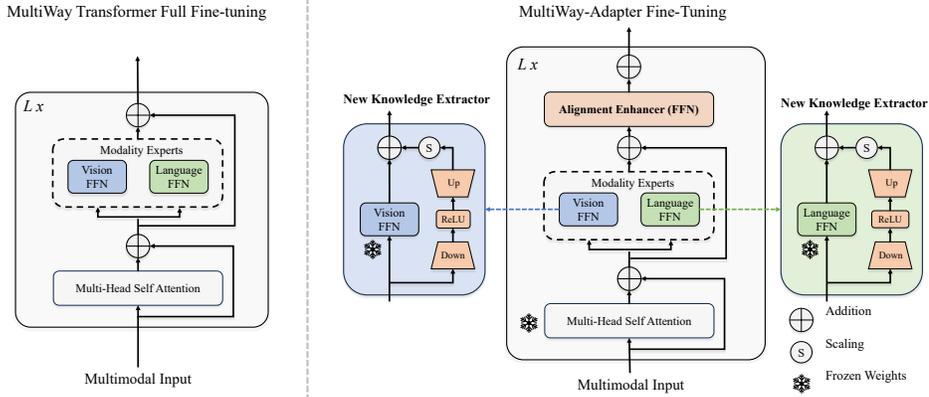
**Fig. 1**. **Comparison of a MultiWay Transformer and our MultiWay-Adapter fine-tuning.** MultiWay-Adapter uses a dual-component design, including New Knowledge Extractor and Alignment Enhancer. We replace the original FFN with New Knowledge Extractor: frozen branch (left) and the trainable bottleneck module (right). Moreover, we add a Alignment Enhancer upon the original FFN to enhance the inter-modal alignment.

## 2. RELATED WORK

**Challenges in Adapting Large Multimodal Models**. Recently, increasing model size has been shown to be an effective strategy for improving performance. Models such as BEiT-3 [2] and BLIP-2 [1], with up to 1.9 billion and 12.1 billion parameters, respectively, have set new state-of-the-art results in multimodal tasks such as Visual Question Answering. However, their application to specialized downstream tasks is often limited by computational constraints [10, 11, 12, 13, 14, 15]. For instance, the requirement for large GPU memory in full fine-tuning limits their adaptations for specialized tasks on commodity hardware, e.g., 45GB for full fine-tuning of the BEiT-3 Large model.

**Efficient Transfer Learning Methods**. The challenge of computational efficiency in fine-tuning MLLMs has given rise to Parameter-Efficient Transfer Learning (PETL) methods. These are broadly categorized into partial parameter updates [16] and modular additions [17, 4]. The former is resource-intensive and model-specific, while the latter adds new modules to architectures, updating only these components. However, most studies only focus on unimodal tasks in domains such as vision [18], text [19] or audio [20, 21, 22], neglecting multimodal tasks. A few works [4, 5, 16, 23] target multimodal tasks but suffer from shallow inter-modal alignment. Our work introduces the MultiWay-Adapter, designed for efficient MLLM transfer learning and enhanced inter-modal alignment.

## 3. APPROACH

We introduce MultiWay-Adapter (MWA), designed for the efficient transfer of Multimodal Large Language Models (MLLM) to downstream tasks. Although the primary focus

of this paper is on image-text retrieval tasks, the potential applicability of the MWA is broader, such as video text retrieval and image captioning.

**Preliminaries**. The overall framework is constructed on the basis of a popular architecture of MLLM, which utilizes a MultiWay Transformer design [2]. As depicted on the left of Figure 1, each MultiWay Transformer block comprises a shared self-attention module and a pool of feed-forward networks (i.e., modality experts) tailored for different modalities. This design is similar to the dual-backbones architecture of multimodal models, e.g., one encoder for vision input and another encoder for language input, yet differs by sharing the weights within each self-attention module. This design choice reduces the parameter count and enhancing inter-modal alignment—an essential quality for high-performance multimodal tasks [2].

### 3.1. MultiWay-Adapter

**Overall Architecture.** Our proposed MWA uses a dual-component approach: the New Knowledge Extractor and the Alignment Enhancer, as illustrated on the right of Figure 1.

**New Knowledge Extractor.** The New Knowledge Extractor is designed for extracting new knowledge from the target downstream tasks. In contrast to the conventional full fine-tuning of MultiWay Transformers, we replace both feed-forward networks (FFNs) in the transformer block with a New Knowledge Extractor. This extractor comprises two branches: the left branch, identical to the original network, and an additional right branch introduced for task-specific fine-tuning. The latter utilizes a bottleneck structure to limit the number of parameters and includes a down-projection layer and an up-projection layer. Formally, for a specific input feature $x_i'$, the right branch of the New Knowledge Extractor produces the

adapted features, $\tilde{x}_i$, as:

$$\tilde{x}_i = \text{ReLU}(\text{LN}(x_i{}') \cdot \mathbf{W}_{\text{down}}) \cdot \mathbf{W}_{\text{up}} \tag{1}$$

Here, $\mathbf{W}_{\text{down}} \in \mathbb{R}^{d \times \check{d}}$ and $\mathbf{W}_{\text{up}} \in \mathbb{R}^{\check{d} \times d}$ denote the down-projection and up-projection layers, respectively. $\check{d}$ is the bottleneck middle dimension and satisfies $\check{d} \ll d$. LN denotes LayerNorm. This bottleneck module is connected to the original FFN (left branch) through a residual connection via a scale factor $\alpha$. Then, these features, $x_i{}'$ and $\tilde{x}_i$, are fused with the original one, $x_i$, through a residual connection:

$$x_i = FFN(LN(x_i{}')) + \alpha \cdot \tilde{x}_i + x_i{}' \tag{2}$$

**Alignment Enhancer.** After extracting new knowledge from the target downstream task, to maintain and improve inter-modal alignment, an Alignment Enhancer module is added atop the pool of feed-forward networks. This module mimics the architecture of the New Knowledge Extractor but uses a larger middle dimension to facilitate better feature fusion and alignment.

During the fine-tuning phase, only the parameters of these newly added modules are optimized, while the rest of the model is frozen (as indicated by the frozen sign in Figure 1). This strategy makes MWA a plug-and-play module, applicable to other MLLM, such as CLIP [26], VLMo [27], and ALIGN [25].

## 4. EXPERIMENTS

**Setup**. We conducted experiments on two state-of-the-art MLLMs, BEiT-3 Base and BEiT-3 Large, across two widely-used image-text retrieval datasets: MSCOCO [28] and Flickr30K [29]. We use the 5k test set of MSCOCO and 1k test set of Flickr30k to report metrics, in accordance with previous studies [28, 29]. We initialized the backbone, excluding our additional modules, with pre-trained weights, which were frozen during the fine-tuning process when employing MultiWay-Adapter. For fine-tuning, the batch size is 512 for the Large model and 1024 for the Base model, over 20 epochs with an initial learning rate of 0.001. Middle dimensions for the New Knowledge Extractor and the Alignment Enhancer were set to 64 and 128, respectively. All the code used in our experiments can be found in `https://github.com/longkukuhi/MultiWay-Adapter`.

**Experimental Results**. The objective of this experiment is to assess the efficiency and efficacy of our MWA framework in comparison to traditional full fine-tuning methods. We compared our MWA approach with full fine-tuning in two distinct settings: fine-tuning performance and zero-shot performance.

**Fine-Tuning Performance**: As shown in Table 1, our MWA method demonstrates superior computational efficiency. Specifically, it utilizes a mere $3.21\%$ and $2.58\%$ of the trainable parameters for the Base and Large variants of BEiT-3, respectively, in contrast to conventional full fine-tuning. This

leads to a substantial reduction in GPU memory consumption—by 7GB and 9GB for the Base and Large variants, respectively. Furthermore, MWA significantly reduces the time required for fine-tuning. For instance, fine-tuning MWA with the BEiT-3 Base model is reduced by $57\%$ compared to full fine-tuning.

Regarding effectiveness, the performance decrement when utilizing MWA is statistically insignificant for both the Base and Large BEiT-3 variants, with deviations falling within a margin of less than $1\%$. Synthesizing these efficiency and effectiveness attributes demonstrates that MWA, when applied to the BEiT-3 Large model, consumes merely $86\%$ of the time required for full fine-tuning of the BEiT-3 Base model, yet surpasses its performance. This suggests that MWA enables enhanced performance with reduced computational time, particularly for larger models. Additionally, as the model size increases, the performance disparity between MWA and full fine-tuning diminishes, indicating a positive correlation between MWA's effectiveness and model size.

**Zero-Shot Performance**: To evaluate the transfer capabilities of MWA and full fine-tuned methods, we conducted experiments in a zero-shot setting. In this setting, the model is evaluated on Flickr30k (1k test set), with which it has no prior knowledge of, thereby necessitating reliance on intrinsically learned knowledge to simulate the handling of previously unseen samples. These models were initially fine-tuned on the MSCOCO dataset. As shown in Table 2, MWA surpasses the performance of full fine-tuning when employed with the BEiT-3 Large model. We hypothesize that this enhancement is attributable to the preservation of generalizable knowledge in the frozen weights, knowledge potentially lost during the full fine-tuning process. This retained knowledge augments the model's ability to adeptly manage unseen instances. Thus, MWA not only match the performance of full fine-tuning method but also distinguishes itself in terms of resource efficiency and transferability.

In summary, the experimental results demonstrate that MWA serves as an effective and resource-efficient fine-tuning method for MLLMs, especially when computational resources are constrained.

## 5. ANALYSIS

**Scaling Tunable Parameters Up**: The primary aim of this section is to investigate the impact of varying the number of tunable parameters on performance and to identify the optimal value for additional parameters. The "mid-dimension" of the New Knowledge Extractor largely controls the number of tunable parameters. We conducted an empirical evaluation across a range of mid dimensions $\{0, 1, 16, 32, 64, 128\}$ on the MSCOCO dataset using the BEiT-3 Base model. The results are summarized in Figure 2. The data reveals a noticeable increase in performance as the dimension grows, plateauing at 64. Specifically, we observed a peak perfor-

| | | | | | MSCOCO (5k test set) | | Flickr30k (1k test set) | |
|---|---|---|---|---|---|---|---|---|
| Model | FT-Way | Tunable params (M) | GPU Mem (GB) | Time (Min) | IR@1 | TR@1 | IR@1 | TR@1 |
| ALBEF [24] | Full Fine-tune | 196 | N/A | N/A | 60.7 | 77.6 | 85.6 | 95.9 |
| ALIGN [25] | Full Fine-tune | 825 | N/A | N/A | 59.9 | 77.0 | 84.9 | 95.3 |
| BEiT-3-Base | Full Fine-tune | 222 (100%) | 37GB | 225 | 61.4 | 79.0 | 86.2 | 96.3 |
| BEiT-3-Large | Full Fine-tune | 675 (304%) | 45GB | 353 | 63.4 (+2.0) | 82.1 (+3.1) | 88.1 (+1.9) | 97.2 (+0.9) |
| BEiT-3-Base | MultiWay-Adapter | 7.13 (**3.21%**) | **30GB** | **130** | 60.7 (-0.7) | 78.3 (-0.7) | 85.4 (-0.8) | 95.4 (-0.9) |
| BEiT-3-Large | MultiWay-Adapter | 17.40 (**2.58%**) | **36GB** | **194** | 63.3 (+1.9) | 82.1 (+3.1) | 88.0 (+1.8) | 97.1 (+0.8) |

**Table 1**. **Comparative Analysis of Full Fine-Tuning and the MultiWay-Adapter**: The table shows Top-1 recall metrics on COCO and Flickr30k datasets, presented as both absolute values and relative gaps to the BEiT-3 Base full fine-tuning Model. Metrics for Text-to-Image Retrieval (IR) and Image-to-Text Retrieval (TR) are provided. GPU memory usage and training time are also included. Training time is measured using a single NVIDIA A6000 GPU with 48GB memory for one epoch.

| | | Flickr30k | |
|---|---|---|---|
| Model | FT-Way | IR@1 | TR@1 |
| BEiT-3-Large | Full fine-tune | 85.99 | 95.48 |
| BEiT-3-Large | MultiWay-Adapter | 86.26 | 95.51 |

**Table 2**. **Zero-shot performance on Flickr30k**.

| | | | MSCOCO | |
|---|---|---|---|---|
| Model | KE | AE | IR@1 | TR@1 |
| BEiT-Base | | | 61.40 | 79.00 |
| BEiT-Base | ✓ | | 57.32 | 73.92 |
| BEiT-Base | | ✓ | 57.88 | 74.61 |
| BEiT-Base | ✓ | ✓ | 60.72 | 78.26 |

**Table 3**. **Ablation study of two modules of MultiWay-Adapter**. KE refers to the New Knowledge Extractor and AE refers to the Alignment Enhancer.

was performed on the MSCOCO dataset using the BEiT-3 Base model. The performance metrics for each component, both in isolation or in combination, are detailed in Table 3. Our findings demonstrate that omitting either component leads to a significant decline in performance, approximately 3%, for image to text retrieval and around 4%, for text to image retrieval. Importantly, the Alignment Enhancer, a novel element distinct from previous Adapter methods, validates its critical role in maintaining deep alignment between modalities through observed performance gains. In summary, both components not only significantly contribute to the overall performance but also complement each other effectively.
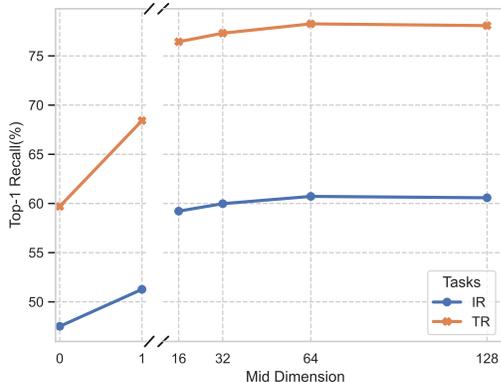


**Fig. 2**. **Evaluation of different sizes of mid-dimension New Knowledge Extractor on MSCOCO**.

mance gain of 9.45%, in text to image retrieval when increasing the dimension from 1 to 64. This indicates that increasing the number of parameters in the adapter does not guarantee performance improvement. When the dimension is set to zero, it represents the zero-shot performance of the BEiT-3 Base model without MWA. Notably, MWA delivers superior performance compared to the zero-shot performance of the BEiT-3 Base model, even when the mid-dimension is as low as one. Furthermore, performance variability is relatively small when increasing the dimension from 16 to 64, indicating that MWA is stable in tuning and not sensitive to changes in size.

**Ablation on MultiWay Adapter's Components**: In this section, our focus is to quantify the individual contributions of our two newly introduced components: the New Knowledge Extractor and the Alignment Enhancer. An ablation study

## 6. CONCLUSION

We introduce the MultiWay-Adapter (MWA), an effective framework designed for the efficient adaptation of Multimodal Large Language Models (MLLM) to downstream tasks. Addressing the issue of shallow inter-modal alignment in existing methods, MWA employs a dual-component approach, utilizing both the New Knowledge Extractor and the Alignment Enhancer. This strategy enables MWA to not only extract novel information from downstream datasets but also to secure deep inter-modal alignment. Our empirical findings reveal that with the addition of a mere 2.58% in extra parameters, there is no statistically significant decline in performance across all tested settings while reducing the fine-tuning time by up to 57%. Our research paves the way for future studies on efficient multimodal fine-tuning methods and holds potential for extension into other vision-language tasks.

# 7. REFERENCES

[1] Junnan Li, Dongxu Li, Silvio Savarese, et al., "BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models," in *Proc. ICML, 23-29 July 2023, Honolulu, Hawaii, USA*.

[2] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, and et al., "Image as a foreign language: Beit pretraining for all vision and vision-language tasks," in *Proc. CVPR*, 2023.

[3] Edward J. Hu, Yelong Shen, et al., "Lora: Low-rank adaptation of large language models," in *Proc. ICLR, 2022*.

[4] Yi-Lin Sung, Jaemin Cho, and Mohit Bansal, "VL-ADAPTER: parameter-efficient transfer learning for vision-and-language tasks," in *Proc. CVPR, 2022*.

[5] Xiaohua Zhai, Xiao Wang, Basil Mustafa, et al., "Lit: Zero-shot transfer with locked-image text tuning," in *Proc. CVPR*.

[6] Shoufa Chen, Chongjian Ge, Zhan Tong, et al., "Adaptformer: Adapting vision transformers for scalable visual recognition," in *Proc. NeurIPS*, 2022.

[7] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency, "Multimodal machine learning: A survey and taxonomy," *TPAMI*, vol. 41, no. 2, pp. 423–443, 2018.

[8] Weijie Su, Xizhou Zhu, Yue Cao, et al., "VL-BERT: pre-training of generic visual-linguistic representations," in *Proc. ICLR, 2020*.

[9] Zijun Long and Richard Mccreadie, "Is multi-modal data key for crisis content categorization on social media?," in *Proc. ISCRAM*. May 2022, Tarbes, France.

[10] Zijun Long, Richard Mccreadie, and Imran Muhammad, "Crisisvit: A robust vision transformer for crisis image classification," in *Proc. ISCRAM*, May 2023.

[11] Zijun Long, George Killick, et al., "Robollm: Robotic vision tasks grounded on multimodal large language models," *arXiv preprint arXiv:2310.10221*, 2023.

[12] Zixuan Yi, Zijun Long, Iadh Ounis, et al., "Large multimodal encoders for recommendation," *arXiv preprint arXiv:2310.20343*, 2023.

[13] Zijun Long, George Killick, Lipeng Zhuang, et al., "Elucidating and overcoming the challenges of label noise in supervised contrastive learning," *arXiv preprint arXiv:2311.16481*, 2023.

[14] Zijun Long, Zaiqiao Meng, et al., "Lacvit: A label-aware contrastive training framework for vision transformers," in *Proc. ICASSP, 2024, Seoul*.

[15] Zijun Long and Richard Mccreadie, "Automated crisis content categorization for covid-19 tweet streams," in *Proc. ISCRAM*. May 2021, Blacksburg, VA, USA.

[16] Zhao Song, Ke Yang, Naiyang Guan, et al., "Vppt: Visual pre-trained prompt tuning framework for few-shot image classification," in *Proc. ICASSP*, 2023, pp. 1–5.

[17] Junyi Peng, Themos Stafylakis, Rongzhi Gu, et al., "Parameter-efficient transfer learning of pre-trained transformer models for speaker verification using adapters," in *Proc. ICASSP*, 2023, pp. 1–5.

[18] Sylvestre-Alvise Rebuffi, Hakan Bilen, et al., "Learning multiple visual domains with residual adapters," in *Proc. NeurIPS, December 4-9, 2017, USA*.

[19] Neil Houlsby, Andrei Giurgiu, et al., "Parameter-efficient transfer learning for NLP," in *Proc. ICML, 9-15 June 2019, Long Beach, California, USA*.

[20] Bethan Thomas, Samuel Kessler, et al., "Efficient adapter transfer of self-supervised speech models for automatic speech recognitio," in *Proc. ICASSP*, 2022.

[21] Steven Vander Eeckt and Hugo Van hamme, "Using adapters to overcome catastrophic forgetting in end-to-end automatic speech recognition," in *Proc. ICASSP*. 2023, pp. 1–5, IEEE.

[22] Samuel Kessler, Bethan Thomas, et al., "An adapter based pre-training for efficient and scalable self-supervised speech representation learning," in *Proc. ICASSP*. 2022, pp. 3179–3183, IEEE.

[23] Odysseas S. Chlapanis, Georgios Paraskevopoulos, et al., "Adapted multimodal bert with layer-wise fusion for sentiment analysis," in *Proc. ICASSP*, 2023, pp. 1–5.

[24] Junnan Li, Ramprasaath Selvaraju, et al., "Align before fuse: Vision and language representation learning with momentum distillation," in *Proc. NeurIPS*, 2021.

[25] Chao Jia, Yinfei Yang, et al., "Scaling up visual and vision-language representation learning with noisy text supervision," in *Proc. ICML*. PMLR, 2021.

[26] Alec Radford, Jong Wook Kim, et al., "Learning transferable visual models from natural language supervision," in *Proc. ICML*. PMLR, 2021.

[27] Hangbo Bao, Wenhui Wang, et al., "Vlmo: Unified vision-language pre-training with mixture-of-modality-experts," in *Proc. NeurIPS,*, 2022.

[28] Tsung-Yi Lin, Michael Maire, et al., "Microsoft coco: Common objects in context," in *Proc. ECCV, 2014*.

[29] Bryan A Plummer, Wang, et al., "Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models," in *Proc. ICCV*, 2015.