# Investigating and Enhancing Vision-Audio Capability in Omnimodal Large Language Models

**Rui Hu[1][*], Delai Qiu[2], Shuyu Wei[1], Jiaming Zhang[1],**
**Yining Wang[2], Shengpeng Liu[2], Jitao Sang[1][†]**
[1]Beijing Jiaotong University
[2]Unisound AI Technology Co., Ltd.
*rui.hu@bjtu.edu.cn, jtsang@bjtu.edu.cn*

## Abstract

Omnimodal Large Language Models (OLLMs) have shown significant progress in integrating vision and text, but still struggle with integrating vision and audio, often exhibiting suboptimal performance when processing audio queries compared to text queries. This disparity is primarily due to insufficient alignment between vision and audio modalities during training, leading to inadequate attention to visual information when using audio queries. To mitigate this issue, we propose a Self-Knowledge Distillation (Self-KD) training method where the vision-text component of the OLLM serves as the teacher and the vision-audio component as the student. This enables the model to process audio in a manner analogous to its text processing. Our experimental results demonstrate that Self-KD is an effective method for enhancing the vision-audio capabilities of OLLMs by learning from the vision-text components, which subsequently improves the interaction between audio and images and results in improved performance on multimodal tasks.

## 1 Introduction

Recent years have witnessed significant advancements in large language models (LLMs) (Achiam et al., 2023; Touvron et al., 2023; Yang et al., 2024), which have catalyzed the development of multimodal large language models (MLLMs) (Wang et al., 2024b; Chen et al., 2024a; Liu et al., 2024; Fang et al., 2024; Chu et al., 2024). This progress marks a paradigm shift in how machines understand and interact with the world, with omnimodal large language models (OLLMs) (OpenAI, 2024; Fu et al., 2024; Xie and Wu, 2024; Fu et al., 2025; Li et al., 2024; InfinigenceAI, 2024) emerging as a new frontier. These models, exemplified by GPT-4o, demonstrate advanced capabilities in visual,

---

*This work was done during an internship at Unisound.
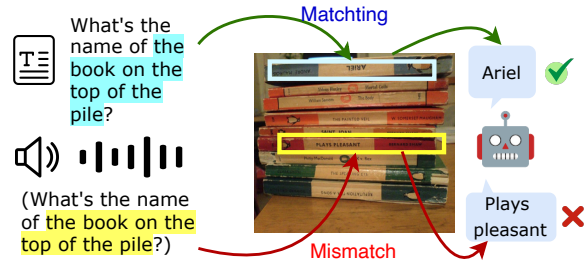†Corresponding author.



Figure 1: An example of the OLLM correctly answering text question but giving incorrect response to the same question in audio form.

linguistic, and auditory functionalities, promising more natural and comprehensive interactions. However, despite these advancements, a critical gap remains in the performance of OLLMs when processing vision-text versus vision-audio inputs. Specifically, OLLMs often exhibit suboptimal performance on vision-audio tasks compared to their vision-language counterparts. For instance, replacing a text question with its audio equivalent can result in contradictory responses from the models. As illustrated in Figure 1, when the text question "What's the name of the book on the top of the pile?" is posed to Megrez (InfinigenceAI, 2024), the model accurately responds with "Ariel". However, when the same question is converted into audio, it erroneously answers "Plays pleasant". This inconsistency is prevalent across various OLLMs, indicating that the models exhibit different behaviors when processing vision-text and vision-audio inputs.

To systematically evaluate this gap, we synthesize text questions from existing vision-language benchmarks into audio using Text-to-Speech (TTS) technology. The results reveal that the vision-audio performance of OLLMs significantly lag behind their vision-text performance. Notably, these incorrect audio responses, as illustrated in Figure 1, share a common thread: they are consistently

image-relevant, despite being factually inaccurate. This observation implies that the models are processing both audio and visual cues but failing to synthesize them into correct answers.

Furthermore, we visualize the attention weights of OLLMs when processing input information and observe that the models show higher attention to query tokens in audio queries than in text queries, while exhibiting lower attention to vision tokens in audio queries compared to text queries. This indicates that OLLMs struggle to effectively integrate visual and audio information. It is hypothesized that this observation arises from a relative deficiency in the alignment between vision and audio compared to vision and text. To evaluate these alignments, we developed a new benchmark MMAlign (See sec 3.2). The results confirm that the alignment between vision and audio is indeed weaker than that between vision and text. This discrepancy stems from the fact that during the alignment phase, OLLM only aligned vision and text as well as audio and text, without directly aligning vision and audio. The model could only learn to process vision-audio inputs during the vision-audio SFT phase. Based on these results, we can conclude that conventional vision-audio SFT alone is insufficient for enabling the model to effectively integrate vision and audio.

To mitigate this issue, we propose a Self-Knowledge Distillation (Self-KD) training framework. In this framework, the vision-text component of the OLLM serves as the teacher model, while the vision-audio component acts as the student model. Unlike conventional vision-audio SFT, Self-KD uses the vision-text outputs of the model as soft labels to guide the training of the vision-audio component. After distillation, the student component learns the behavior of the teacher component, for instance, allocating more attention to vision tokens, thereby enhancing vision-audio performance. In summary, our contributions are as follows:

(1) We identify and analyze the significant gap in performance between vision-language and vision-audio capabilities in OLLMs, attributed to insufficient alignment and between images and audio.

(2) We propose a Self-KD training framework that leverages the vision-text component to guide the training of the vision-audio component, promoting better alignment and integration of visual and audio information.

(3) We conduct extensive experiments on various models and datasets, demonstrating that Self-KD significantly enhances vision-audio performance compared to conventional vision-audio SFT.

## 2 Evaluation of Audio-Vision Capability for OLLM

Currently, the evaluation of OLLMs focuses separately on their vision-language (VL) and audio capabilities, overlooking a holistic assessment of their vision-audio (VA) capability. In this section, we first generate VA benchmarks based on existing VL benchmarks and then conduct a comprehensive evaluation of OLLMs.

### 2.1 Setup

**Datasets Preparation.** We select MME (Fu et al., 2023), HallusionBench (Guan et al., 2024), RealWorldQA (xai, 2024), TextVQA (Singh et al., 2019), ChartQA (Masry et al., 2022), DocVQA (Mathew et al., 2021), and InfographicVQA (Mathew et al., 2022) as the VL evaluation datasets. We then synthesize the text questions in these datasets into audio using TTS (Text-to-Speech) technology. To ensure reproducible evaluation results, we use VLMEvalKit (Duan et al., 2024) uniformly for all evaluations with a zero-shot manner.

**OLLMs.** We select three open-source OLLMs, VITA (Fu et al., 2024), VITA-1.5 (Fu et al., 2025), and Megrez (InfinigenceAI, 2024) for testing, with parameter sizes of 8×7B, 7B, and 3B, respectively.

### 2.2 Performance Gap

**There is a gap between the vision-audio and vision-language capabilities of OLLMs.** We evaluate the OLLMs on both VL and VA datasets, with the results presented in Table 1. All models exhibite relatively strong performance under text-based queries, achieving scores around 70. However, when the same questions are posed in audio form, the performance of all models declined to varying degrees. Specifically, VITA exhibits the most substantial decline, with an average decrease of 62.2, Megrez demonstrates the least decline, but still experiences a reduction of 19.2. These results suggest that current open-source OLLMs generally possess weaker capabilities in integrating images and audio compared to integrating images and text.

**Models exhibit a higher "Yes" bias when using audio to query compared to text:** In Figure 2, we present the "Yes" ratio of OLLMs on the MME and HallusionBench datasets. Both MME and HallusionBench are yes-or-no datasets, and the

Table 1: Vision tasks performance of different OLLMs. In the query, "Text" indicates that the question is posed using text, while "Audio" indicates that the question is posed using audio.

| Model | Query | MME | TextVQA | HalluB | $CQA_H$ | $CQA_A$ | DocVQA | InfoVQA | RWQA | Average |
|---|---|---|---|---|---|---|---|---|---|---|
| VITA-8x7b | Text | 84.81 | 71.52 | 40.98 | 65.60 | 87.60 | 84.49 | 63.85 | 61.44 | 70.04 |
| | Audio | 5.36 | 4.55 | 22.79 | 6.40 | 7.76 | 7.63 | 5.36 | 2.88 | 7.84 |
| | $\Delta Gap$ | 79.45 | 66.97 | 18.19 | 59.20 | 79.84 | 76.86 | 58.49 | 58.56 | 62.20 |
| VITA-1.5-7B | Text | 86.44 | 72.85 | 45.04 | 65.12 | 87.52 | 88.51 | 60.64 | 64.58 | 71.34 |
| | Audio | 32.11 | 44.32 | 14.92 | 27.60 | 67.12 | 47.39 | 23.01 | 33.73 | 36.28 |
| | $\Delta Gap$ | 54.33 | 28.53 | 30.12 | 37.52 | 20.40 | 41.12 | 37.63 | 30.85 | 35.06 |
| Megrez-3B | Text | 80.21 | 90.66 | 52.30 | 48.72 | 82.32 | 78.56 | 47.91 | 70.98 | 68.96 |
| | Audio | 57.52 | 51.25 | 36.48 | 36.88 | 71.60 | 63.38 | 30.57 | 50.07 | 49.72 |
| | $\Delta Gap$ | 22.69 | 39.41 | 15.82 | 11.84 | 10.72 | 15.18 | 17.34 | 20.91 | 19.24 |



Figure 2: The "Yes" ratio of OLLMs in MME and HallusionBench datasets.



Figure 3: Examples of OLLMs provide relevant but inaccurate answers to audio questions. (top) An example from ChartQA. (bottom) An example from TextVQA.
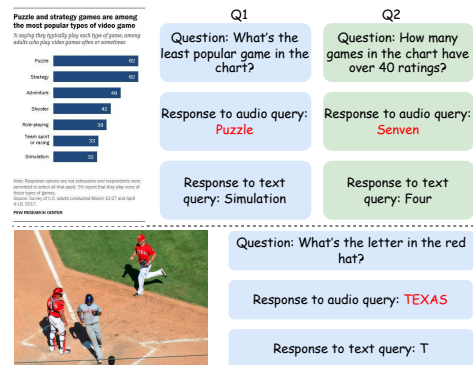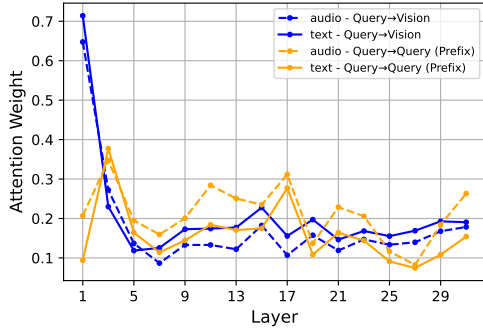
"Yes" ratio reflects the model's output bias. The ground truth "Yes" proportions are 50% for MME and 42% for HallusionBench. For MME, the "Yes" ratio for all models exceeds 50% with audio query, indicating a higher preference for "Yes". In contrast, the "Yes" ratio for text queries is close to 50%, suggesting that the model exhibits no significant bias when using text queries. For HallusionBench, the model demonstrates a moderate of "Yes" bias when using text queries, which is further amplified when using audio queries.

**Models exhibit a tendency to provide relevant but inaccurate answers to audio-based questions:** In the VQA task, models are required to integrate image and question to generate accurate answers. We observe that when using text queries, models can accurately combine the question and the image to produce correct answers. However, when using audio queries, although the answers are relevant to the images and meet the question requirements, they are often inaccurate. For example, as shown in the top part of Figure 3, when querying VITA-1.5 with the audio ques-
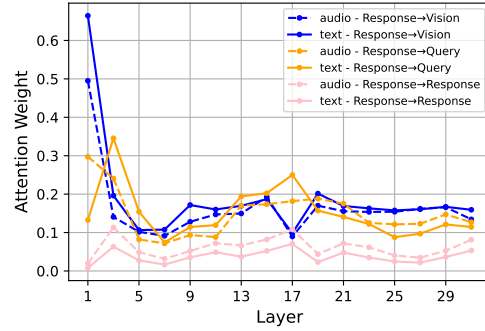
tion "What's the least popular game in the chart?" from ChartQA dataset, the model responded with "Puzzle", which is a game listed in the chart but not the least popular one. Similarly, the response "Seven" represents the total number of games rather than the correct answer to the question "How many games in the chart have over 40 ratings?". The bottom part of Figure 3 and Figure 1 show similar cases of Megrez in the TextVQA dataset, indicating that this phenomenon is widely present in current OLLMs.

## 3  Why is OLLM's Audio-Vision Capability Weaker?

Given that OLLMs exhibit inferior performance on vision-audio tasks compared to vision-text tasks, what factors contribute to this discrepancy? In this section, we first show that when processing vision-audio inputs, the attention weights of query tokens to vision tokens are lower than when processing vision-text inputs. We then build a new benchmark MMAlign to evaluate the alignment be-

(a) Query -> Vision / Query(prefix)

(b) Response -> Vision / Query / Response(prefix)

Figure 4: Layer-wise variation of attention weights assigned to different types of tokens (including query, vision, and response) in OLLMs. "A→B" means the attention weights from A to B.

tween the audio modality and the image modality, as well as between the text modality and the vision modality within OLLMs. The results show that the alignment between audio and vision is significantly weaker than that between text and vision. Finally, we discuss the connection between the model training process and vision-audio capability, suggesting that models need to enhance the integration of vision and audio during training.

## 3.1 Attention Weight Analysis

To find out the behavioral differences of the model in processing vision-text and vision-audio inputs, we measure the attention weights assigned to different token types at each layer. For each sample, we can represent the input and output as "`<system><image><query><response>`", where the `"<query>"` can be in text or audio form. For a causal language model, the model relies solely on the preceding input information when generating sequences. Thus, the assignment of attention weights can reflect the model's behavior in processing sequence information.

**Query tokens pay less attention to vision tokens under audio queries than under text queries:** In Figure 4(a), we present how the attention weights from query tokens to image tokens and to themselves vary across different layers in Megrez (InfinigenceAI, 2024). This reflects how the model processes the input information to prepare for output. Consistent with the findings of Bi et al. (2024); Zhang et al. (2025), we observe that the model's attention to vision tokens is high in the early layers, regardless of whether the query is text or audio. However, in the middle and later layers of the model, when using an audio query,

the attention weights from the query tokens to the vision tokens are consistently lower than those with a text query. In contrast, the model focuses more on the query token itself. This suggests that the model may struggle to effectively integrate audio and visual information in the later layers, leading to the inferior performance on vision-audio tasks compared to vision-text tasks.

**Response tokens show similar attention to input tokens between audio and text queries:** In Figure 4(b), we present how the attention weights from response tokens to image tokens, query tokens, and to themselves vary across different layers. The model's attention to both vision and query tokens shows little difference between audio and text queries. This indicates that the model considers both the image and the query when generating a response to an audio question, consistent with our observations in Section 2.2, where we find that the model's responses to audio queries are relevant to both the image content and the query content. This further suggests that the primary cause of the performance discrepancy lies in the insufficient integration of audio and vision information.

## 3.2 MMAlign: Evaluation on Modality Alignment

According to prior work (Bi et al., 2024), attention distribution to some extent reflects the alignment between different modalities. Therefore, we hypothesize that within OLLMs, the alignment between vision and audio is weaker than that between vision and text. To test this hypothesis, we construct the MMAlign benchmark based on the ARO dataset (Yüksekgönül et al., 2023) to compare the degree of alignment between vision-text and vision-
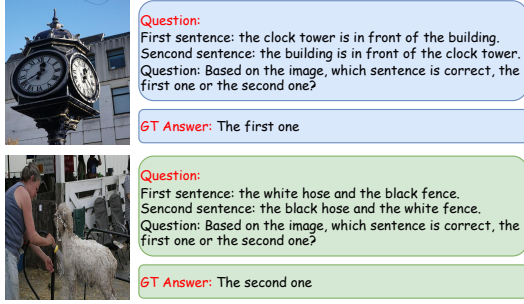
Figure 5: Test samples of MMAlign. The top one is relation type and the bottom one is attribute type.

Table 2: Results on MMAlign.

| Model | Query | Relation | Attribute | Average |
|-------|-------|----------|-----------|---------|
| VITA | Text | 61.33 | 68.00 | 64.67 |
| | Audio | 1.33 | 2.33 | 1.83 |
| VITA-1.5 | Text | 74.00 | 77.33 | 75.67 |
| | Audio | 31.33 | 34.33 | 32.83 |
| Megrez | Text | 54.33 | 59.67 | 57.00 |
| | Audio | 50.00 | 52.00 | 51.00 |

audio within OLLMs. Specifically, ARO (Yük-sekgönül et al., 2023) is a dataset for testing the image understanding capabilities of VLMs, e.g., CLIP (Radford et al., 2021). Each sample contains an image and two short captions, including one correct caption and one perturbed caption. Depending on the type of perturbation, it can be divided into relation perturbation, attribute perturbation, and word order perturbation.

As shown in Figure 5, we build MMAlign by combining the two captions into a single question, asking the model to select the correct one from the two sentences. To ensure the semantic correctness of the sentences, we only consider the relation and attribute types, resulting in a total of 600 samples. Each sample contains a text question, its corresponding audio version, and a correct answer.

Table 2 shows the results of OLLMs on MMA-lign. The results for all models demonstrate better performance with text queries than with audio queries, indicating that the alignment between audio and vision is still not on par with that between text and vision. The models' performance on attributes is slightly better than on relations, indicating that the models' understanding of the relationships between objects in images is weaker than their understanding of attributes.

### 3.3 Limitation of the Training Process of the current OLLMs

The training process for current OLLMs (Fu et al., 2025; Li et al., 2024) can be divided into four steps:

**Vision-Text Alignment:** This step aims to bridge the gap between vision and text, enabling the model to understand visual information and align it with text embeddings.

**Vision-Text SFT:** This step further trains the model to understand image content and answer image-related questions based on instructions,

building on the foundation of visual alignment.

**Audio-Text Alignment:** This step aims to bridge the gap between audio and text, enabling the model to understand audio inputs.

**Vision-Audio SFT:** This step further trains the model to understand audio and answer image-related questions based on audio instructions, building on the foundation of audio alignment.

Unlike vision and text, vision and audio have not been directly aligned at any stage. This is because, due to the characteristics of LLMs, we can only construct the training loss based on text. As a result, we are unable to directly model the alignment task between vision and audio. Therefore, we expect the model to learn to organically integrate vision and audio to complete downstream tasks during the vision-audio SFT stage. However, our experimental results show that the current vision-audio SFT does not achieve the same effect as vision-text SFT.

## 4 A Simple Improvement: Self-Knowledge Distillation from Vision-Text to Vision-Audio

Our analysis shows that vision-text surpasses vision-audio in both modality alignment and downstream task performance. A natural way to bridge this gap is through knowledge distillation (Hinton, 2015), where the vision-text component of the OLLM serves as the teacher and the vision-audio component as the student. Since both originate from the same model, we refer to this method as Self-Knowledge Distillation (Self-KD) of OLLM, which can be used to enhance the effect of vision-audio SFT. Figure 6 illustrates the Self-KD training framework.

**Vision-Audio SFT.** We can represent a vision-text SFT dataset as $[X^T, Y]$, where $X^T$ are inputs and $Y$ are text answers, the current common practice is to convert the text question in $X^T$ into audio to obtain vision-audio inputs $X^A$ and train the model on $[X^A, Y]$. The conventional vision-audio
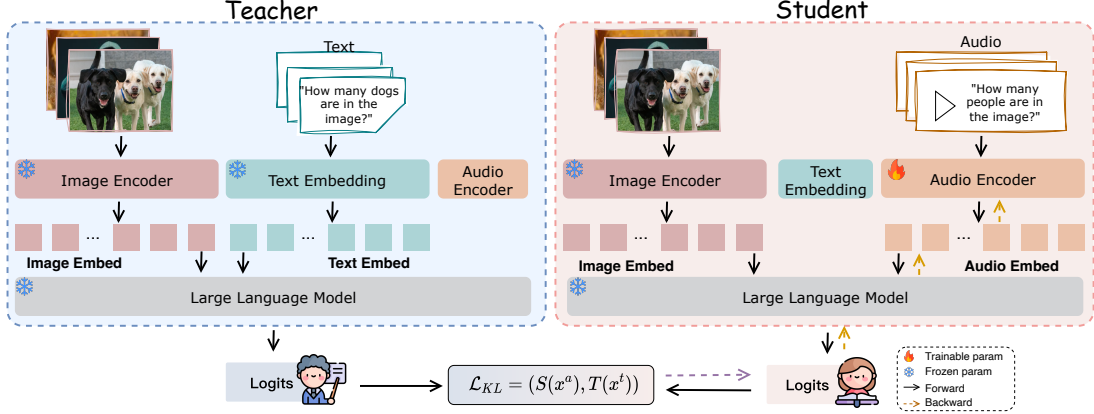
Figure 6: Illustration of our proposed Self-Knowledge Distillation training framework.

SFT loss function can be expressed as:

$$L_{\text{SFT}} = \mathbb{E}_{x^a \sim X^A, y \sim Y} \left[ -\log p_S(y|x^a) \right], \quad (1)$$

where $p_S$ is the vision-audio component of the OLLM, comprising the vision encoder, audio encoder, and the LLM. After vision-audio SFT, the OLLM is expected to learn to process vision-audio inputs. However, our results show that with conventional vision-audio SFT, the model's ability to integrate vision and audio to generate correct responses remains insufficient.

**Self-KD.** We define $p_T$ as the vision-text component of the OLLM, which includes the text embedding layer, the vision encoder and the LLM of the OLLM. Given that $p_T$ outperforms $p_S$, we use $p_T$ as the teacher model and $p_S$ as the student model, employing KL divergence as the loss function for self-knowledge distillation. The formula is as follows:

$$
\begin{aligned}
L_{\text{Self-KD}} &= \text{KL}(p_T \parallel p_S) \\
&= \mathbb{E}_{x^a \sim X^A, x^t \sim X^T, y \sim Y} \left[ \log \frac{p_T(y|x^t)}{p_S(y|x^a)} \right].
\end{aligned}
\tag{2}
$$

As shown in Figure 6, unlike conventional knowledge distillation, where the teacher and student models use the same input, in Self-KD, the teacher model's input $x^t$ is the vision-text sample, while the student model's input $x^a$ is the corresponding vision-audio sample. For the final training, we can combine the SFT loss and the Self-KD loss, and use a hyperparameter to control their proportions:

$$L = \alpha L_{Self-KD} + (1 - \alpha) L_{SFT}. \quad (3)$$

## 5 Experiment

### 5.1 Experimental Setup

To verify the effectiveness of Self-KD, we chose to expand the audio modality on existing LVLMs to obtain OLLMs because they have already completed alignment and SFT on vision-text data.

**Models.** We select the InternVL2 series (Chen et al., 2024b) and Qwen2VL series (Wang et al., 2024a) as our base models due to their excellent performance and the availability of multiple sizes. Following (Li et al., 2024; Chu et al., 2024), we use the Whisper-large-v3 model (Radford et al., 2023) as the audio encoder and a one-layer MLP as the projector to convert audio features to LLM embeddings.

**Training.** For audio-text alignment, we collect ASR datasets such as LibriSpeech (Panayotov et al., 2015), Common Voice (Ardila et al., 2019), GigaSpeech (Chen et al., 2021), and Libriheavy (Kang et al., 2024), totaling 988k samples. For vision-audio SFT and self-KD training, we first sample 50k instruction-following samples from llava-1.5-mix-665k (Liu et al., 2024) and then converte the text questions into audio. See Appendix A for more training details.

### 5.2 Main Results

We conducted extensive experiments on different types and sizes of base models. Based on the results in Table 3, we can draw the following conclusions:

**The gap between VL and VA capabilities is widespread.** After performing audio-text alignment and audio-vision SFT, the gap between VL and VA capabilities persists in various models. This suggests that even with effective audio-text alignment, audio cannot yet fully replace text when in-

Table 3: Performance comparison between conventional vision-audio SFT and Self-KD training (KD ratio=1). The first row for each model shows the performance using text queries.

| Model | Method | MME | TextVQA | HalluB | RWQA | CQA$_H$ | CQA$_A$ | DocVQA | InfoVQA | Average |
|---|---|---|---|---|---|---|---|---|---|---|
| InternVL2-1B | - | 67.2 | 57.62 | 24.64 | 49.93 | 44.48 | 72.4 | 48.32 | 32.87 | 49.68 |
|  | SFT | 27.52 | 16.58 | 16.27 | 25.75 | 17.92 | 29.52 | 17.67 | 18.07 | 21.16 |
|  | Self-KD | **47.63** | **43.09** | **19.23** | **27.19** | **29.84** | **42.96** | **36.14** | **24.63** | **33.84** |
| InternVL2-2B | - | 68.05 | 63.84 | 30.02 | 53.99 | 49.76 | 72.8 | 54.69 | 38.41 | 53.94 |
|  | SFT | **43.47** | 26.16 | 19.45 | 28.76 | 18.32 | 13.2 | 18.21 | 12.61 | 22.52 |
|  | Self-KD | 40.69 | **51.55** | **23.14** | **35.69** | **32.64** | **40.88** | **40.24** | **27.8** | **36.58** |
| InternVL2-4B | - | 76.99 | 64.04 | 35.65 | 57.12 | 59.36 | 80.48 | 56.79 | 44.58 | 59.38 |
|  | SFT | 50.39 | 35.53 | 27.5 | 38.17 | 23.6 | 33.76 | 27.83 | 20.95 | 32.22 |
|  | Self-KD | **54.29** | **53.33** | **28.64** | **38.43** | **39.68** | **48.88** | **43.19** | **31.99** | **42.3** |
| InternVL2-8B | - | 76.74 | 75.31 | 39.73 | 69.52 | 91.44 | 84.99 | 61.66 | 59.87 | 69.91 |
|  | SFT | 44.02 | 45.76 | 27.08 | 28.08 | 43.20 | 48.78 | 36.42 | 36.34 | 38.71 |
|  | Self-KD | **43.78** | **63.37** | **31.49** | **43.60** | **69.76** | **71.36** | **49.52** | **38.69** | **51.45** |
| Qwen2VL-2B | - | 74.98 | 74.82 | 39.69 | 59.22 | 53.44 | 86.08 | 81.08 | 48.89 | 64.77 |
|  | SFT | 54.3 | 54.82 | 28.17 | 40 | 36.32 | 61.6 | 59.45 | 34.98 | 46.21 |
|  | Self-KD | **57.41** | **67.77** | **32.82** | **45.1** | **41.12** | **68.72** | **67.91** | **39.77** | **52.58** |
| Qwen2VL-7B | - | 83.4 | 77.09 | 47.39 | 70.98 | 70.16 | 90.8 | 89.75 | 71.5 | 75.14 |
|  | SFT | **71.14** | 73.27 | 43.01 | **51.37** | 62.88 | 88 | 85.04 | 67.28 | 67.75 |
|  | Self-KD | 70.04 | **73.87** | **43.74** | 50.46 | **64.96** | **89.28** | **85.69** | **68.08** | **68.27** |

teracting with images.

**Model's VL capability is directly proportional to its acquired VA capability after audio-vision SFT.** For example, InternVL2-8B has the best VL performance (69.91) in its series, and after SFT with the same data, its VA performance (38.71) is also the best. This suggests that models with stronger VL capabilities tend to achieve better VA performance after vision-audio SFT. Therefore, when developing OLLMs, it is advisable to prioritize enhancing their VL capabilities.

**Self-KD training can reduce the gap between a model's VL and VA capabilities.** The results in Table 3 show that, with the same training data, using Self-KD compared to conventional SFT can enable the model to achieve better VA performance. Similarly, the effectiveness of Self-KD is also directly proportional to the model's VL capability, which is understandable because Self-KD uses the model's VL component as the teacher. The improvement of Self-KD on the Qwen series is relatively smaller than that on the Intern series. This may be because the Qwen series models have better alignment between vision and text, as indicated by their performance at the same scale. Thus, a standard vision-audio SFT can yield satisfactory results after audio-text alignment.

### 5.3 Further Analysis

**Self-KD aligns the model's behavior when it processes vision-audio and vision-text inputs.** To examine the behavioral differences between models trained with Self-KD and conventional SFT, we visualize the attention weights of the models. We refer to the teacher component as the "base model" and the model with audio-text alignment but without vision-audio SFT as the "ASR model". As shown in Figure 7, the ASR model exhibits higher attention to query tokens and lower attention to vision tokens compared to the base model. After vision-audio SFT, this gap narrows, but only marginally. In contrast, the model trained with Self-KD shows a smaller difference in attention allocation relative to the base model. This indicates that Self-KD effectively brings the model's behavior with vision-audio input closer to its behavior with vision-text input. Figure 9 in the Appendix B further illustrates this behavioral consistency.

**Self-KD enhances the alignment between vision and audio.** As shown in Table 4, compared to conventional vision-audio SFT, models trained with Self-KD achieved better overall results on MMAlign. This indicates that, even though we did not directly align audio and vision during training, learning from the teacher component's behavior can indirectly promote the alignment between audio and vision.
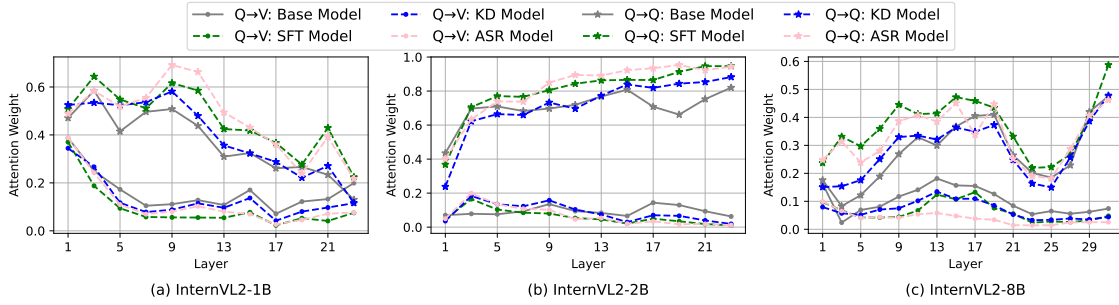
Figure 7: Layer-wise variation of attention weights assigned to different types of token. Q->V means attention from query tokens to vision tokens, Q->Q means query tokens to query tokens.

Table 4: Comparison of conventional vision-audio SFT and Self-KD training on MMAlign.

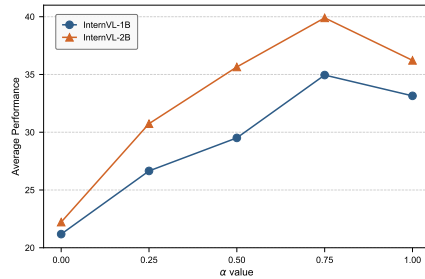| Model | Relation | | Attribute | |
|---|---|---|---|---|
| | SFT | Self-KD | SFT | Self-KD |
| InternVL2-1B | 42.67 | 50.67 | 45.33 | 47.33 |
| InternVL2-2B | 49.00 | 50.00 | 47.67 | 49.67 |
| InternVL2-4B | 47.67 | 52.33 | 50.67 | 50.00 |
| InternVL2-8B | 53.33 | 57.33 | 54.67 | 56.00 |
| Qwen2VL-2B | 50.67 | 51.33 | 51.33 | 55.00 |
| Qwen2VL-7B | 71.00 | 71.33 | 58.33 | 61.67 |
| Average | 52.39 | **55.50** | 51.28 | **53.44** |



Figure 8: Ablation results for KD ratio $\alpha$.

## 5.4 Ablation Study

**KD Loss Ratio.** Figure 8 shows the results for different values of the KD loss ratio $\alpha$ (see Appendix C for detailed results). Performance improves as the KD ratio increases, with the best average results achieved at a KD ratio of 0.75. This indicates that KD and SFT can mutually enhance each other's effectiveness.

## 6 Related Works

**Omnimodal Large Language Models.** Recent advancements in multimodal large models have primarily focused on Vision-Language Models, e.g., CLIP (Radford et al., 2021), followed by models such as Intern-VL (Chen et al., 2024a), which use MLPs to integrate vision encoders and LLMs for enhanced semantic alignment. Audio-Language Models, like Qwen-Audio (Chu et al., 2024), combine audio encoders with LLMs to directly map audio signals to text. Recently, Omnimodal Large Language Models (OLLMs) have emerged, integrating vision, audio, and text by aligning their encoders during training for end-to-end processing. Models such as VITA (Fu et al., 2024, 2025), Mini-Omni2 (Xie and Wu, 2024), MiniCPM-o (MiniCPM-o Team, 2025), and Baichuan-Omni (Li et al., 2024) have demonstrated strong multimodal

performance.

**Knowledge Distillation in MLLMs.** Knowledge distillation (Hinton, 2015) has recently been applied to multimodal large language models (MLLMs). For example, LLaVA-MoD (Shu et al., 2024) and LLaVA-KD (Cai et al., 2024) use knowledge distillation to transfer the performance of large teacher models to smaller student models. This paper proposes a self-knowledge distillation method, dividing the same model into teacher and student components to bring the vision-audio capabilities of OLLMs closer to their vision-text capabilities.

## 7 Conclusions

This paper investigates the issue of integration of audio and vision in OLLMs. We find that, for visual question answering tasks, performance with audio queries is significantly lower than with text queries. Further analysis reveals that this disparity arises from insufficient alignment between images and audio during training, leading to inadequate attention to images when using audio queries. To address this, we propose a Self-Knowledge Distillation training method, where the vision-text component serves as the teacher and the vision-audio component as the student. This approach aims to align the model's vision-audio capability with its

vision-text capability. Experimental results show that our method effectively improves the interaction between audio and images during model inference, outperforming existing baseline models in benchmark performance.

## 8 Limitations

This paper propose a self-knowledge distillation training method for OLLMs, however, there are two limitations in this work. Firstly, under the knowledge distillation paradigm, we need to generate soft labels through teacher model inference, which increases the training cost compared to conventional SFT. Second, the vision-audio (VA) performance of models after Self-KD still falls short of their vision-text (VL) performance, suggesting that there is still room for improvement in OLLM training.

## References

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

Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber. 2019. Common voice: A massively-multilingual speech corpus. *arXiv preprint arXiv:1912.06670*.

Jing Bi, Junjia Guo, Yunlong Tang, Lianggong Bruce Wen, Zhang Liu, and Chenliang Xu. 2024. Unveiling visual perception in language models: An attention head analysis approach. *arXiv preprint arXiv:2412.18108*.

Yuxuan Cai, Jiangning Zhang, Haoyang He, Xinwei He, Ao Tong, Zhenye Gan, Chengjie Wang, and Xiang Bai. 2024. LLaVA-KD: A Framework of Distilling Multimodal Large Language Models. *arXiv e-prints*, arXiv:2410.16236.

Guoguo Chen, Shuzhou Chai, Guanbo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, et al. 2021. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio. *arXiv preprint arXiv:2106.06909*.

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

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

Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, et al. 2024. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*.

Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, et al. 2024. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models. In *Proceedings of the 32nd ACM international conference on multimedia*, pages 11198–11201.

Qingkai Fang, Shoutao Guo, Yan Zhou, Zhengrui Ma, Shaolei Zhang, and Yang Feng. 2024. Llama-omni: Seamless speech interaction with large language models. *arXiv preprint arXiv:2409.06666*.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and Rongrong Ji. 2023. MME: A comprehensive evaluation benchmark for multimodal large language models. *CoRR*, abs/2306.13394.

Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Shaoqi Dong, Xiong Wang, Di Yin, Long Ma, et al. 2024. Vita: Towards open-source interactive omni multimodal llm. *arXiv preprint arXiv:2408.05211*.

Chaoyou Fu, Haojia Lin, Xiong Wang, Yi-Fan Zhang, Yunhang Shen, Xiaoyu Liu, Yangze Li, Zuwei Long, Heting Gao, Ke Li, et al. 2025. Vita-1.5: Towards gpt-4o level real-time vision and speech interaction. *arXiv preprint arXiv:2501.01957*.

Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. 2024. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14375–14385.

Geoffrey Hinton. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.

InfinigenceAI. 2024. Megrez-3b-omni: The first open-source end-side full modality understanding model. https://huggingface.co/Infinigence/Megrez-3B-Omni.

Wei Kang, Xiaoyu Yang, Zengwei Yao, Fangjun Kuang, Yifan Yang, Liyong Guo, Long Lin, and Daniel Povey. 2024. Libriheavy: a 50,000 hours asr corpus with punctuation casing and context. In *ICASSP*

*2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 10991–10995. IEEE.

Yadong Li, Haoze Sun, Mingan Lin, Tianpeng Li, Guosheng Dong, Tao Zhang, Bowen Ding, Wei Song, Zhenglin Cheng, Yuqi Huo, et al. 2024. Baichuan-omni technical report. *arXiv preprint arXiv:2410.08565*, 3(7).

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.

Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*.

Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. 2022. Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1697–1706.

Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209.

OpenBMB MiniCPM-o Team. 2025. Minicpm-o 2.6: A gpt-4o level mllm for vision, speech, and multimodal live streaming on your phone. https://huggingface.co/openbmb/MiniCPM-o-2_6.

OpenAI. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.

Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5206–5210. IEEE.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR.

Fangxun Shu, Yue Liao, Le Zhuo, Chenning Xu, Lei Zhang, Guanghao Zhang, Haonan Shi, Long Chen, Tao Zhong, Wanggui He, Siming Fu, Haoyuan Li, Bolin Li, Zhelun Yu, Si Liu, Hongsheng Li, and Hao Jiang. 2024. LLaVA-MoD: Making LLaVA Tiny via MoE Knowledge Distillation. *arXiv e-prints*, arXiv:2408.15881.

Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024b. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.

xai. 2024. Grok-1.5 vision preview. https://x.ai/blog/grok-1.5v.

Zhifei Xie and Changqiao Wu. 2024. Mini-omni2: Towards open-source gpt-4o with vision, speech and duplex capabilities. *arXiv preprint arXiv:2410.11190*.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.

Mert Yüksekgönül, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. 2023. When and why vision-language models behave like bags-of-words, and what to do about it?

Shaolei Zhang, Qingkai Fang, Zhe Yang, and Yang Feng. 2025. Llava-mini: Efficient image and video large multimodal models with one vision token. *arXiv preprint arXiv:2501.03895*.
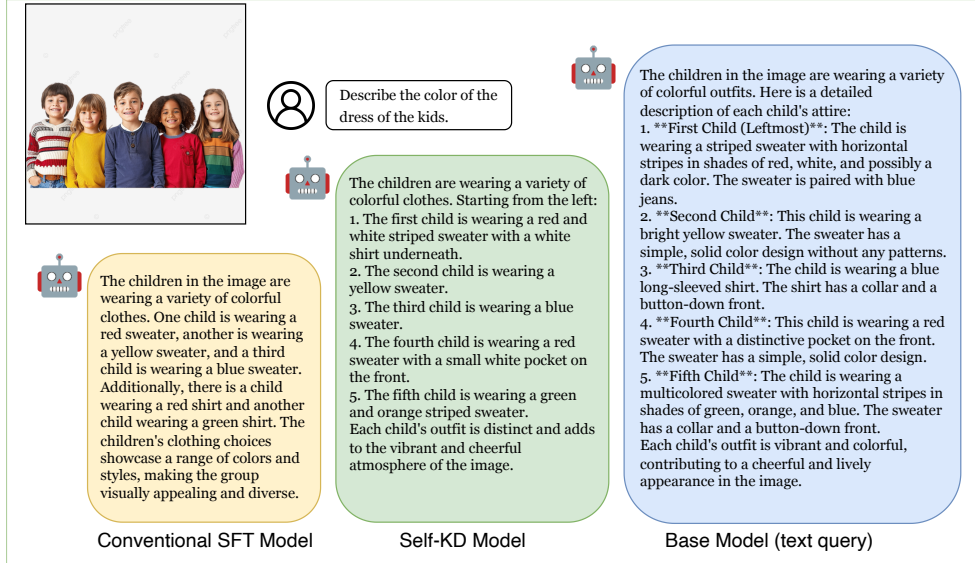
Figure 9: An example shows the output differences between three models: conventional SFT model, Self-KD model, and base model. The Self-KD model has very similar output to the base model.

Table 5: Ablation results of different KD loss ratio.

| Model | KD ratio | MME | TextVQA | HalluB | RWQA | CQA_H | CQA_A | DocVQA | InfoVQA | Average |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 27.52 | 16.58 | 16.27 | 17.92 | 29.52 | 17.67 | 18.07 | 25.75 | 21.16 |
| | 0.25 | 26.52 | 22.70 | 16.78 | 21.12 | 35.12 | 37.48 | 24.45 | 28.76 | 26.62 |
| InternVL2-1B | 0.5 | 27.31 | 29.23 | 17.80 | 24.08 | 41.04 | 42.82 | 25.50 | 27.97 | 29.47 |
| | 0.75 | 46.04 | 40.68 | 18.22 | 27.20 | **45.52** | **49.62** | **28.77** | **28.24** | **35.53** |
| | 1.0 | **47.63** | **43.09** | **19.23** | **29.84** | 42.96 | 36.14 | 24.63 | 27.19 | 33.84 |
| | 0 | 43.47 | 26.16 | 19.45 | 18.32 | 13.20 | 18.21 | 12.61 | 28.76 | 22.52 |
| | 0.25 | 47.35 | 37.56 | 19.62 | 23.52 | 26.48 | 43.67 | 21.60 | 31.63 | 31.43 |
| InternVL2-2B | 0.5 | 38.71 | 47.36 | 22.61 | 26.00 | 37.52 | 52.73 | 27.17 | 33.99 | 35.76 |
| | 0.75 | **43.85** | **54.04** | **25.07** | 28.80 | **45.04** | **58.45** | **31.66** | 34.25 | **40.15** |
| | 1.0 | 40.69 | 51.55 | 23.14 | **32.64** | 40.88 | 40.24 | 27.80 | **35.69** | 36.58 |

## A   Training Details

The entire training process is completed on eight A100 GPUs. For audio-text alignment, we set the batch size to 128, which takes about 4 hours. For vision-audio SFT and Self-KD, we set the batch size to 64, and each training session takes approximately half an hour and one hour, respectively. The learning rate is set to 4e-5 throughout the training process, and we employ a cosine-type learning rate decay strategy. Both training stages are conducted for only one epoch. To avoid degrading the model's vision-language performance, we freeze the LLM and vision encoder, and only train the audio encoder and its corresponding MLP layer.

## B   Case Study

In Figure 9, we present an example comparing the output differences between models trained with

conventional vision-audio SFT and trained with Self-KD. We use the output of the base model as a reference. Faced with the request "Describe the color of the dress of the kids", the base model can accurately describe the dress of each kid. The Self-KD model also describes each child but with less detail compared to the base model, while the SFT model can only provide a general description of the overall image.

## C   Ablation Study

Table 5 shows the detailed results of different KD loss ratios. When the KD ratio is relatively high, the models achieve better results. Specifically, when the KD ratio is set to 0.75, the models achieve the best average performance.