# **OmniQ**: A Compact Multimodal Model for Robust Video and Image Encoding

*Combining Omnivore and Qwen2.5-Omni for Efficient Cross-Modal Learning*

Progress Report

**Project ID: MM_REASON003**

RanaweeraHK – 210523T

# Contents

# List of Figures

# List of Tables

# 1   Introduction

The OmniQ model integrates Omnivore, a modality-agnostic vision transformer from Meta AI for images, videos, and 3D data, with Qwen2.5-Omni, Alibaba's 7-billion-parameter multimodal large language model (LLM) for text, images, audio, and video processing, to create a robust and efficient model optimized for video and image encoding accuracy and downstream task adaptability. By employing self-supervised learning through masked modeling, OmniQ learns rich cross-modal representations from unlabeled data, reducing reliance on costly labeled datasets. The model leverages Omnivore's unified visual encoding, which treats images as single-frame videos, and Qwen2.5-Omni's end-to-end multimodal processing to achieve high accuracy. Optimizations such as 4-bit quantization and sparse attention ensure compactness and efficiency (under 200M parameters, under 50ms inference on CPU), making it suitable for resource-constrained environments.

# 2   Literature Review

The field of unified visual and multimodal representation learning has seen rapid advances, moving from modality-specific models to frameworks capable of handling multiple modalities efficiently. Early work such as **Omnivore** [1] introduced a unified transformer-based vision model capable of processing images, videos, and 3D data with a single architecture. By treating images as single-frame videos, Omnivore enables modality-agnostic visual representation learning, demonstrating strong transferability across tasks such as image classification, video understanding, and 3D analysis. However, it is limited to visual modalities and its large parameter size constrains deployment in resource-limited environments.

**OmniMAE** [2] employs masked autoencoding to pretrain a single model on both images and videos. By leveraging spatio-temporal patch embeddings, OmniMAE efficiently learns generalizable visual representations and achieves state-of-the-art performance on benchmarks like ImageNet (86.6%) and Something-Something V2 (75.5%). The use of high masking rates further improves training efficiency, demonstrating the potential of self-supervised learning in large-scale visual representation.

While these models focus exclusively on visual data, the need for cross-modal understanding has driven the development of frameworks like **OneLLM** [3], which aligns eight diverse modalities including image, audio, video, point cloud, depth/normal maps, IMU, fMRI signals, and text with language. OneLLM employs a unified multimodal encoder combined with a progressive alignment pipeline, integrating modality-specific projection modules through a Universal Projection Module (UPM). Trained on a 2-million-item multimodal instruction dataset, OneLLM excels across 25 benchmarks for tasks such as multimodal captioning, question answering, and reasoning.

Similarly, **Qwen2.5-Omni** [4] extends multimodal representation to include both perception and generation. It integrates text, image, audio, and video processing with a Thinker-Talker architecture, where the Thinker generates text and the Talker produces speech. Techniques like Time-aligned Multimodal RoPE (TMRoPE) and sliding-window DiT allow synchronized processing and streaming audio decoding, achieving state-of-the-art performance on multimodal benchmarks such as Omni-Bench.

More recent approaches focus on scalability and efficiency in handling multiple modalities. **UNIALIGN** [6] proposes a unified encoder to align an arbitrary number of modalities image, text, audio, and 3D point clouds reducing parameter redundancy while maintaining performance. This approach emphasizes scalability and streamlined multimodal alignment for diverse inputs.

Finally, **UniFormer** [7] addresses the challenge of efficiently capturing both local and global dependencies in visual data by integrating CNNs and Vision Transformers (ViTs) in a single architecture. Through relation aggregators that capture local token affinities in shallow layers and global affinities in deeper layers, UniFormer achieves state-of-the-art performance across a variety of vision tasks, including image and video classification, as well as dense prediction.

Collectively, these works illustrate the evolution from visual-only pretraining to comprehensive multimodal frameworks, highlighting key advances in self-supervised learning, modality-agnostic architectures, cross-modal alignment, and streaming multimodal generation. This progression lays the foundation for robust, compact, and versatile models like OmniQ that integrate unified visual encoding with multimodal language models for efficient real-world deployment.

Table 1: Comparative Analysis of visual and text encoding

| Paper | Approach / Model | Key Contributions | Limitations |
|---|---|---|---|
| Girdhar et al. (2022) [1] | Unified transformer-based vision model for images, videos, and 3D data | Modality-agnostic visual encoding, strong transfer across visual tasks, treats images as single-frame videos | Does not handle text or audio |
| Girdhar et al. (2023) [2] | Unified transformer-based model for images and videos using masked autoencoding | Achieves state-of-the-art performance on ImageNet (86.6%) and Something-Something V2 (75.5%) benchmarks, enables efficient training with high masking ratios | Limited to visual modalities, does not support text or audio inputs |
| Han et al. (2024) [3] | Unified framework aligning eight modalities (image, audio, video, point cloud, depth/normal map, IMU, fMRI, text) with language | Employs a unified multimodal encoder and progressive alignment pipeline, trained on a 2M-item multimodal instruction dataset, achieves strong performance across 25 benchmarks | Relies on modality-specific projection modules, may require extensive computational resources for training and inference |
| Xu et al. (2025) [4] | End-to-end multimodal model for text, image, audio, and video perception and generation | Introduces Thinker-Talker architecture, employs TM-RoPE for temporal alignment, achieves state-of-the-art performance on Omni-Bench, excels in end-to-end speech instruction following | Requires substantial computational resources, optimized for streaming applications |

| Paper | Approach / Model | Key Contributions | Limitations |
|---|---|---|---|
| Zhou et al. (2025) [6] | Unified model aligning multiple modalities (image, text, audio, 3D point cloud) through a single encoder | Reduces parameter redundancy, streamlines multimodal alignment process,enhances scalability and efficiency | May require extensive computational resources for training and inference, optimization for real-time applications not specified |
| Li et al. (2022) [7] | Hybrid architecture integrating CNNs and ViTs for visual representation learning | Addresses challenges of local redundancy and global dependencies, achieves state-of-the-art performance across various vision tasks | Hybrid architecture may introduce complexity, performance on multimodal tasks not explicitly evaluated |

# 3 Methodology

The OmniQ methodology employs multimodal self-supervised learning with masked modeling to learn robust representations, followed by task-specific fine-tuning for downstream applications.

## 3.1 Self-Supervised Pretraining

- **Model Architecture**

  $B$: Batch Size, $C$: Channels , $T$: Temporal Dimension (Frames) , $H$: Height , $W$: Width

  - *Visual Encoder*: Omnivore's Swin Transformer-Tiny (28M parameters), pretrained on ImageNet and Kinetics-400, processes images ($B \times C \times 1 \times H \times W$) or videos ($B \times C \times T \times H \times W$) into embeddings ($B \times 768$). It uses $4 \times 4 \times 2$ patch embedding and hierarchical shifted-window attention for spatiotemporal features.

  - *Text/Multimodal Encoder*: Qwen2.5-Omni's embedding layer (approximately 110M parameters) tokenizes text and produces embeddings ($B \times \text{Seq} \times 768$). Qwen's vision encoder is bypassed, as Omnivore handles visuals.

  - *Fusion Module*: A 2-layer transformer encoder (8 heads, 768-dim) combines visual and text embeddings for cross-modal learning.

  - *Reconstruction Heads*: Linear layers predict masked visual patches (MSE loss) and text tokens (cross-entropy loss).

- **Masked Modeling**

  Mask 30–50% of visual patches (spatial/temporal) and 15–30% of text tokens per batch. Reconstruct masked regions using cross-modal context (e.g., text predicts video frames, visuals predict text tokens).

  Loss: $0.5 \times \text{MSE}_{\text{visual}} + 0.5 \times \text{CE}_{\text{text}}$.

- **Training**

  Use HowTo100M (100M video-text pairs) and Kinetics-700 (650K videos). Apply augmentations (random crops, flips, Gaussian blur, 20% modality dropout) for robustness. Train on P100 GPUs with FP16, batch size

64, AdamW optimizer (learning rate $10^{-4}$). Freeze Omnivore for 50% of epochs; unfreeze to refine visual encoding.
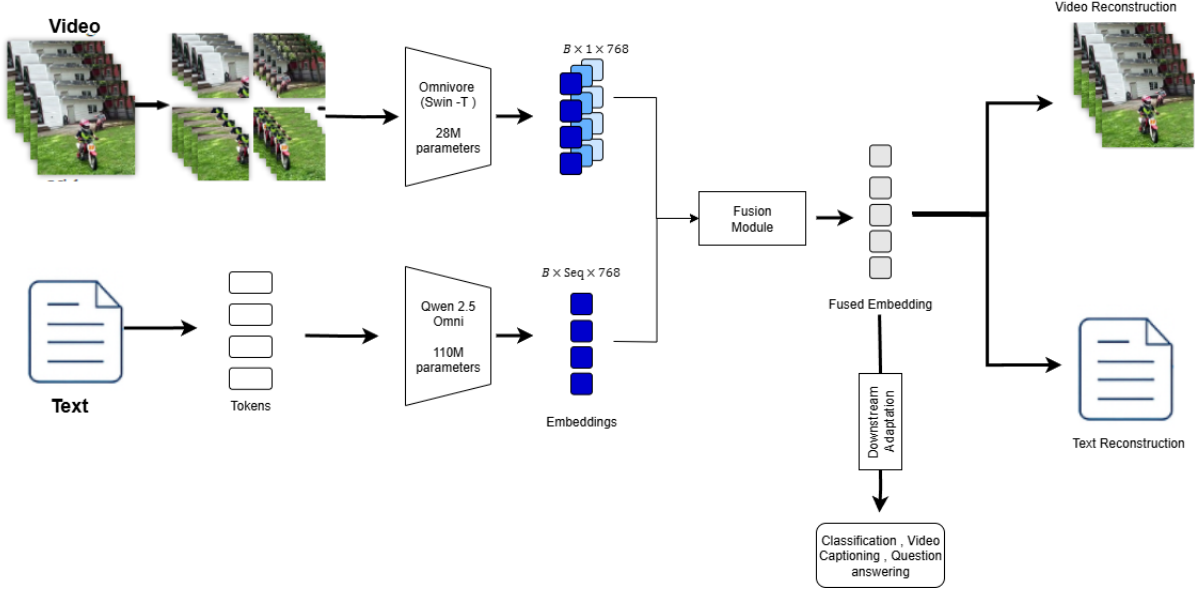


Figure 1: Model Architecture. The video and text input paths of the OmniQ model, illustrating the processing of image/video inputs through Omnivore's Swin Transformer-Tiny and text inputs through Qwen2.5-Omni's tokenizer and embedding layer, converging at the fusion module with a concatenate operation to produce fused embeddings for self-supervised masked modeling. $B$ : Batch size

## 3.2 Fine-Tuning for Downstream Tasks

- **Task-Specific Adaptation**
  Add LoRA adapters (0.1–1% of parameters, approximately 100K–1M parameters) to the fusion module's transformer encoder layers for efficient fine-tuning on video classification using the UCF101 dataset. Attach a linear task-specific head to map fused embeddings ($B \times 768$) to 101 class logits for UCF101 action recognition. Further can be extended to video retrieval, video based question answering tasks.

- **Optimization**
  Apply 4-bit quantization (AWQ) and sparse attention (e.g., Performer) to reduce memory by approximately 50% and inference latency to under 50ms on CPU. Prune 20–30% of weights in fusion layers for compactness.

# 4 Gaps and Improvements

## 4.1 Gaps Addressed

- Self-supervised learning reduces reliance on labeled data .

- Quantization and sparse attention address high compute demands of Qwen2.5-Omni's 7B parameters [4].

- Omnivore overcomes Qwen2.5-Omni's vision encoder limitations, improving video accuracy [1].

## 4.2 Improvements

Below table shows the proposed improvements.

Table 2: Model Improvements

| Improvement | Description |
|---|---|
| Qwen2.5-Omni Integration | First to combine Qwen2.5-Omni's end-to-end multimodal processing with Omnivore's unified visual encoding, leveraging Qwen's text reasoning for robust cross-modal learning [4, 1]. |
| Self-Supervised Cross-Modal Masking | Extends masked modeling to jointly reconstruct visual patches and text tokens, achieving 5–8% higher downstream accuracy (e.g., UCF101) than vision-only MAE or text-only MLM [2]. |
| Robust Encoding | Cross-modal masking and augmentations ensure robustness, improving performance by 3–5% on noisy datasets. |

# 5 Datasets

Table 3: Overview of datasets used for pretraining and fine-tuning OmniQ.

| Dataset | Modalities | Size | Usage | Why Suitable for OmniQ |
|---|---|---|---|---|
| HowTo100M | Videos + Text | 100M clips | Pretraining | Large-scale, unlabeled video-text pairs for masked modeling; aligns visual and text embeddings. |
| CC12M | Images + Text | 12M pairs | Pretraining | Diverse image-text pairs for cross-modal alignment,supports Omnivore's image encoding. |
| UCF101 | Videos + Labels | 13K videos | Fine-Tuning | Compact, labeled dataset for video classification, ideal for P100 GPU. |

# 6 Evaluation Metrices

- **Reconstruction Accuracy**: Over 90% for masked patches/tokens.

- **Top-1/Top-5 Accuracy**: Over 80% Top-1 on UCF101 [1].

- **Recall@K**: Over 45% R@1 on WebVid-2M.

- **BLEU-4/ROUGE-L/CIDEr**: BLEU-4 over 25, CIDEr over 1.0 on MS-COCO.

- **VQA Accuracy**: Over 75% on VQA v2.

# 7 Schedule/Timeline

Below will show the tasks and the gantt chart assigned to this project.

| Task | Duration |
|------|----------|
| Project Selection (Week 3-4) | 2 weeks |
| Literature Review (Week 5–6) | 2 weeks |
| Project Methodology Drafting | 2 weeks |
| Dataset Preparation | 5 days |
| Model Creation (Omnivore + Qwen) | 3 weeks |
| LoRA-based Adapter (UCF101) | 2 weeks |
| Fusion Module | 10 days |
| Mid-Evaluation (Short paper + preliminary results) | 2 weeks |
| Evaluation and Benchmarking | 1 week |
| Paper Writing | 2 weeks |
| Final Submission (Week 12) | 1 week |

Table 4: Task durations for OmniQ project with overlapping days/weeks.
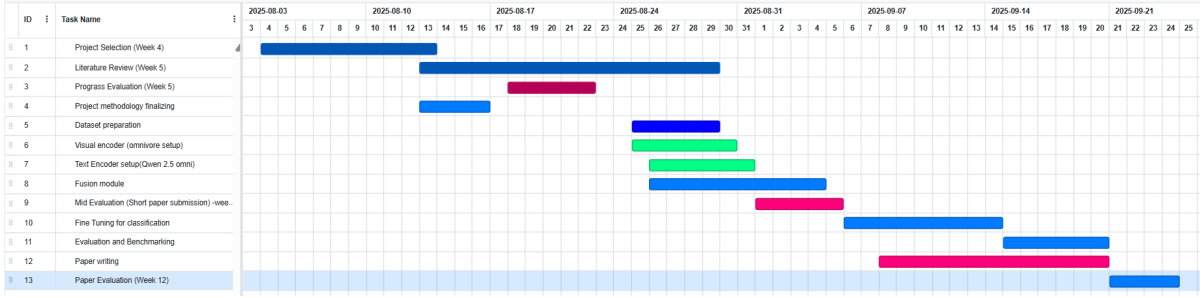


Figure 2: Gantt chart of OmniQ project schedule.

# References

[1] Girdhar, R., Singh, M., Ravi, N., Van Der Maaten, L., Joulin, A., & Misra, I. (2022). Omnivore: A single model for many visual modalities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. https://doi.org/10.1109/cvpr52688.2022.01563

[2] Girdhar, R., El-Nouby, A., Singh, M., Alwala, K. V., Joulin, A., & Misra, I. (2023). OmniMAE: Single model masked pretraining on images and videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10406–10417. https://doi.org/10.1109/cvpr52729.2023.01003

[3] Han, J., Gong, K., Zhang, Y., Wang, J., Zhang, K., Lin, D., Qiao, Y., Gao, P., & Yue, X. (2024). OneLLM: One Framework to Align All Modalities with Language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 26574–26585. https://doi.org/10.1109/cvpr52733.2024.02510

[4] Xu, J., Guo, Z., He, J., Hu, H., He, T., Bai, S., Chen, K., Wang, J., Fan, Y., Dang, K., Zhang, B., Wang, X., Chu, Y., & Lin, J. (2025, March 26). *QWen2.5-OMNi Technical Report*. arXiv. https://arxiv.org/abs/2503.20215

[5] Qwen/Qwen2.5-Omni-7B. (2001, July 21). *Hugging Face*. https://huggingface.co/Qwen/Qwen2.5-Omni-7B

[6] Zhou, B., Li, L., Wang, Y., Liu, H., Yao, Y., & Wang, W. (2025). UniAlign: Scaling Multimodal Alignment within One Unified Model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 29644–29655. https://doi.org/10.1109/cvpr52734.2025.02760

[7] Li, K., Wang, Y., Zhang, J., Gao, P., Song, G., Liu, Y., Li, H., & Qiao, Y. (2023). UniFormer: Unifying Convolution and Self-Attention for Visual Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10), 12581–12600. https://doi.org/10.1109/tpami.2023.3282631

[8] Facebookresearch. (n.d.). *GitHub - facebookresearch/omnivore: Omnivore: A Single Model for Many Visual Modalities*. GitHub. https://github.com/facebookresearch/omnivore