
Algorithm 1 Chain-of-Merging (CoM) Framework

Require: Vision-Language Model M_{vlm} with parameters Θ_{vlm} , Large Language Model M_{llm} with parameters Θ_{llm} .

Require: Input Image I , Question Q .

Require: Hyperparameters: Penalty coefficient λ , Matching threshold τ_{match} , Weight range $[w_{min}, w_{max}]$.

Ensure: Merged Model Parameters $\Theta_{\mathcal{M}}$, Final Answer A .

```
1: // Stage 1: Initialization
2: Generate caption  $C \leftarrow M_{vlm}(I)$ 
3: Extract hidden states  $H_{vlm} \leftarrow f_{vlm}(I, Q; \Theta_{vlm})$ 
    $\{H_{vlm} = \{h_{vlm}^i\}_{i=1}^L\}$ 
4: Extract hidden states  $H_{llm} \leftarrow f_{llm}(C, Q; \Theta_{llm})$ 
    $\{H_{llm} = \{h_{llm}^j\}_{j=1}^L\}$ 
5: // Stage 2: Adaptive Layer Matching
6: Initialize Merging Plan  $\mathcal{P} \leftarrow \emptyset$ 
7: Target Sequence Length  $L_{seq} \leftarrow$ 
    $\lfloor \min(\text{len}(H_{vlm}), \text{len}(H_{llm}))/2 \rfloor$ 
8: for  $i = 1$  to  $L$  do
9:    $j^* \leftarrow -1, S_{max} \leftarrow -\infty$ 
10:  for  $j = 1$  to  $L$  do
11:    {Dimensionality Reduction for CKA}
12:     $\tilde{h}_{vlm}^i, \tilde{h}_{llm}^j \leftarrow \text{AdaptivePool}(h_{vlm}^i, h_{llm}^j, L_{seq})$ 
13:    {Calculate Scores (Eq. 3, 6, 7)}
14:     $S_{struct} \leftarrow \text{CKA}(\tilde{h}_{vlm}^i, \tilde{h}_{llm}^j)$ 
15:     $S_{sem} \leftarrow \text{CosineSim}(\text{MeanPool}(h_{vlm}^i), \text{MeanPool}(h_{llm}^j))$ 
16:     $S_{base} \leftarrow (S_{struct} + S_{sem})/2$ 
17:     $S_{match}^{(i,j)} \leftarrow S_{base} - \lambda \cdot \frac{|i-j|}{L}$ 
18:    if  $S_{match}^{(i,j)} > S_{max}$  then
19:       $S_{max} \leftarrow S_{match}^{(i,j)}, j^* \leftarrow j$ 
20:    end if
21:  end for
22:  if  $S_{max} > \tau_{match}$  then
23:     $\mathcal{P} \leftarrow \mathcal{P} \cup \{(i, j^*, S_{max})\}$ 
24:  end if
25: end for
26: // Stage 3: Dynamic Weight Merging
27: Initialize  $\Theta_{\mathcal{M}} \leftarrow \Theta_{vlm}$  {Copy VLM parameters as base}
28: for each layer  $i \in \{1, \dots, L\}$  do
29:  if  $\exists (i, j^*, S_{match}) \in \mathcal{P}$  then
30:    {Calculate Dynamic Weight (Eq. 9)}
31:     $S_{norm} \leftarrow (S_{match} - \tau_{match}) / (1.0 - \tau_{match})$ 
32:     $w \leftarrow w_{min} + (w_{max} - w_{min}) \cdot S_{norm}$ 
33:    {Perform SLERP (Eq. 10-12)}
34:     $\theta_{merged}^i \leftarrow \text{SLERP}(\theta_{vlm}^i, \theta_{llm}^{j^*}, w)$ 
35:     $\Theta_{\mathcal{M}}^i \leftarrow \theta_{merged}^i$ 
36:  else
37:     $\Theta_{\mathcal{M}}^i \leftarrow \theta_{vlm}^i$  {Keep original VLM layer}
38:  end if
39: end for
40: // Inference
41:  $A \leftarrow \text{Generate}(I, Q; \Theta_{\mathcal{M}})$ 
42: return  $A$ 
```

A. Chain-of-Merging (CoM) Algorithm Pseudocode

In this section, we present the comprehensive algorithmic details of the proposed Chain-of-Merging (CoM) framework. As illustrated in Algorithm 1, the framework is designed to operate in a strictly training-free manner, dynamically integrating the reasoning capabilities of a Large Language Model (LLM) into a Vision-Language Model (VLM) during inference.

The execution flow is structured into three distinct phases to ensure optimal parameter integration:

- 1. Initialization:** The process begins by extracting hidden representations from both models. Crucially, this stage generates a textual caption from the image to bridge the modality gap, allowing the LLM to process visual context alongside the textual question.
- 2. Adaptive Layer Matching:** Unlike conventional methods that enforce rigid one-to-one mapping, our algorithm dynamically identifies the optimal layer correspondence. It computes a composite score derived from structural consistency (via CKA) and semantic similarity, penalized by layer distance to maintain architectural coherence.
- 3. Dynamic Weight Merging:** Finally, the algorithm constructs the merged model parameters. Instead of simple linear averaging, we employ Spherical Linear Interpolation (SLERP) weighted by the normalized matching scores. This ensures that the geometric properties of the parameter space are preserved, minimizing the risk of feature collapse.

This algorithmic design allows CoM to adaptively adjust its merging strategy for every specific image-question pair, maximizing reasoning performance while maintaining computational feasibility.

B. Dimensionality Reduction for Efficient CKA Computation

In the Adaptive Layer Matching stage, we employ Centered Kernel Alignment (CKA) to evaluate the structural consistency between the hidden representations of the Vision-Language Model (VLM), denoted as $H_{vlm} \in \mathbb{R}^{n_{vlm} \times d}$, and the Large Language Model (LLM), denoted as $H_{llm} \in \mathbb{R}^{n_{llm} \times d}$. A fundamental challenge in this process is the discrepancy in sequence lengths: n_{vlm} includes both image and text tokens, whereas n_{llm} consists solely of caption and text tokens.

Standard CKA necessitates the computation of Gram matrices $K = H_{vlm}H_{vlm}^T$ and $L = H_{llm}H_{llm}^T$. When $n_{vlm} \neq n_{llm}$, the resulting matrices have mismatched dimensions, rendering the Hilbert-Schmidt Independence Criterion (HSIC) mathematically undefined for direct comparison. Furthermore, computing Gram matrices for high-

resolution VLMs incurs a computational complexity of $O(n^2)$, which is prohibitively expensive. To resolve these issues, we implement an **Adaptive Mean Pooling** strategy that aligns sequence lengths while preserving global topological structure.

B.1. Target Length Determination

To enable valid kernel operations, we first establish a unified target sequence length, L_{target} . This length is determined dynamically to balance the preservation of information granularity with computational efficiency. We define L_{target} based on the shorter of the two sequences, reduced by a factor of 2:

$$L_{target} = \max\left(1, \left\lfloor \frac{\min(n_{vlm}, n_{llm})}{2} \right\rfloor\right) \quad (14)$$

The inclusion of the $\max(1, \cdot)$ operation ensures numerical stability by preventing a zero target length in edge cases, while the reduction factor serves to accelerate the subsequent matrix operations.

B.2. Adaptive Pooling Logic

We apply 1D Adaptive Average Pooling along the sequence dimension to map an arbitrary input hidden state tensor $H \in \mathbb{R}^{n \times d}$ to a reduced representation $\tilde{H} \in \mathbb{R}^{L_{target} \times d}$.

For the i -th output vector in the reduced sequence (where $0 \leq i < L_{target}$), the pooling window is calculated dynamically. This approach handles non-integer stride lengths uniformly, ensuring that every token in the original sequence contributes to the pooled representation without overlap or omission gaps. The window indices are defined as:

$$\text{start_idx} = \lfloor i \times \frac{n}{L_{target}} \rfloor \quad (15)$$

$$\text{end_idx} = \min\left(n, \lceil (i + 1) \times \frac{n}{L_{target}} \rceil\right) \quad (16)$$

The pooled feature vector $\tilde{H}[i]$ is computed as the mean of the hidden states within this window:

$$\tilde{H}[i] = \frac{1}{\text{end_idx} - \text{start_idx}} \sum_{k=\text{start_idx}}^{\text{end_idx}-1} H[k] \quad (17)$$

B.3. HSIC Computation

Following the adaptive pooling, both the reduced VLM tensor \tilde{H}_{vlm} and LLM tensor \tilde{H}_{llm} share the dimensions $L_{target} \times d$. We subsequently compute the aligned Gram matrices $K = \tilde{H}_{vlm} \tilde{H}_{vlm}^T$ and $L = \tilde{H}_{llm} \tilde{H}_{llm}^T$, both of size $L_{target} \times L_{target}$. This alignment allows for the efficient computation of the HSIC value:

$$\text{HSIC}(K, L) = \frac{1}{(L_{target} - 1)^2} \text{tr}(KHLH) \quad (18)$$

By focusing on the global topological correspondence of the feature spaces rather than strict token-to-token alignment, this method effectively bridges the modality gap between visual and textual representations.