

# Disentangling Identity from Clothing: A Semantically-Supervised Decoupling Framework for Robust Person Re-Identification

Anonymous ICME submission

**Abstract**—Text-to-Image Person Re-Identification is critically hampered by the difficulty of fine-grained semantic alignment, as retrieval accuracy is degraded by clothing-induced interference and a persistent modality gap. In this paper, we propose a novel framework to resolve this issue, centered on feature decoupling guided by a Multimodal Large Language Model (MLLM). Our framework introduces two core components: a Bidirectional Decoupled Alignment Module and a Mamba State Space Model (SSM) for efficient fusion. To obtain high-quality, fine-grained supervision, we first employ MLLM to automatically generate separate identity and clothing descriptions. These descriptions then guide our decoupling module, which utilizes bidirectional attention and a gated weighting strategy to meticulously disentangle visual features into identity and clothing subspaces. To enforce this separation and ensure identity purity, we design a multi-task loss strategy comprising an alignment loss that actively suppresses the influence of clothing-related features, and a kernel-based orthogonal constraint that ensures statistical independence. Furthermore, we pioneer the integration of the Mamba SSM into cross-modal Re-ID as an efficient fusion module. By leveraging its linear-time complexity and proficiency in modeling long-range dependencies, it facilitates deep contextual interactions across modalities while avoiding the quadratic complexity of Transformers. Comprehensive experiments on multiple benchmark datasets reveal that our proposed method achieves superior performance compared to leading contemporary methods, proving its effectiveness and robustness.

**Index Terms**—Multimodal Learning, Text-to-Image Re-Identification, Feature Decoupling, Semantic Supervision.

## I. INTRODUCTION

Text-to-Image Person Re-Identification (T2I-ReID) retrieves a target pedestrian from a large-scale image gallery given a natural-language description [1], [2]. It is valuable for video surveillance [3], intelligent security [4], public safety, and social media. Despite recent progress, practical deployment remains challenging due to image factors (pose, viewpoint, illumination) obscuring identity cues and a persistent modality gap hampering fusion. These issues are exacerbated at the fine-grained level, where semantic alignment is particularly difficult.

A core challenge in T2I-ReID is the semantic gap between images and text. Early work attempted to reduce this discrepancy by projecting global visual and textual features into a shared space [5], [6], but struggled with high intra-class and low inter-class variance. To overcome this, subsequent studies introduced feature disentanglement [7], [8]. These are broadly explicit, using auxiliary modules for part-alignment [2], [9], or implicit, using regularizers to associate noun phrases with

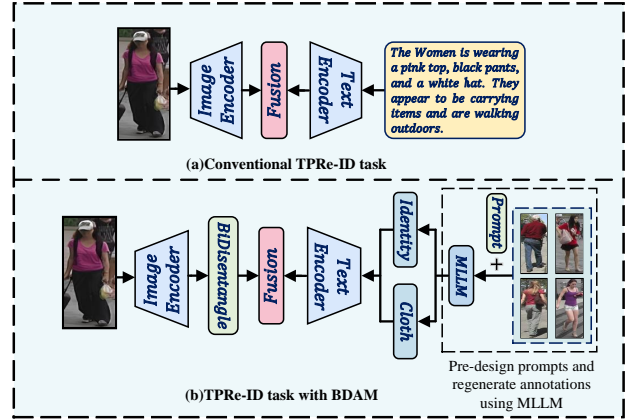


Fig. 1. Comparison of ReID methods. (a) Traditional: Direct fusion of image and text features without distinguishing identity from non-identity information limits alignment. (b) Proposed (BDAM): Introduces a decoupling module that aligns separated identity/clothing features with MLLM-generated descriptions for fine-grained matching.

regions [1], [10]. This progression highlights that distinguishing identity-relevant from irrelevant semantics is essential for advancing T2I-ReID.

This pursuit of disentanglement has been propelled by powerful backbones. Models employing ViT [11]–[13] capture fine-grained details, while methods leveraging CLIP [14]–[16] learn a well-aligned joint space. Despite these advances, a critical limitation persists: both lines of work commonly treat the textual description holistically. This overlooks the semantic distinction between content relevant to identity (e.g., gender, body shape) and content irrelevant to identity (e.g., clothing, hairstyle). This coarse-grained treatment forces the model to entangle these factors, often prioritizing salient clothing details over stable identity cues, which blurs identity features and degrades matching robustness.

To address this challenge, we propose a novel framework with explicit, fine-grained semantic supervision. Inspired by the style-clustering paradigm [17], our approach employs MLLM to guide feature decoupling by automatically generating distinct descriptions for identity and clothing. These decoupled annotations provide precise supervision for our BDAM to meticulously separate and align identity and clothing information. This separation is enforced by a multi-task loss strategy,

including an alignment loss and an HSIC-based orthogonality constraint. Furthermore, we pioneer the integration of the Mamba SSM as an efficient fusion module, adept at capturing long-range cross-modal dependencies with linear complexity. Our main contributions are threefold: (1) An automatic prompt construction pipeline that combines style clustering with an MLLM to produce fine-grained, decoupled identity and clothing descriptions; (2) The BDAM, which achieves precise decoupling and alignment reinforced by a multi-task loss strategy combining an alignment loss and an orthogonality constraint based on HSIC; and (3) The novel integration of a Mamba SSM fusion module that models long-range cross-modal dependencies with linear complexity.

## II. RELATED WORK

### A. Feature Disentanglement

Feature disentanglement aims to separate semantically distinct factors [18]. Strategies include adversarial training [19], metric learning [20], and orthogonal projections [21]. In the ReID domain, this is used to separate identity signals from nuisances [22], [23], such as countering clothing changes [22]. While these efforts improve semantic purity, critical limitations persist. Without effective interaction, isolated factors can lead to brittle alignment. Furthermore, reliance on manual annotations or external detectors constrains scalability, and implicit regularizers can yield spurious separations on unseen data. Our framework addresses this by pairing model-level disentanglement with explicit supervision from MLLM-generated descriptions, ensuring controllable separation while preserving cross-modal synergy.

### B. Feature Fusion

Feature fusion is central to T2I-ReID, often relying on Transformers or CLIP [24]. However, Transformers built on multi-head attention suffer from quadratic complexity in sequence length [25]. Pipelines built upon CLIP [26] benefit from pretraining but their global pooling and holistic processing often blur identity versus clothing cues. Other approaches like dynamic fusion [27] or graph-based models [28] introduce trade-offs in computational cost or adaptability.

Contemporary evidence suggests that accuracy improves only when semantic disentanglement and efficient fusion advance in tandem. An ideal fusion module should respect factorized semantics (i.e., identity and clothing), capture long-range dependencies, and scale with linear or near-linear complexity. This motivates our design: BDAM supplies factor-aware representations, while a Mamba SSM fusion module models long-horizon interactions with linear complexity, enabling precise alignment without the memory and efficiency bottlenecks of standard Transformers.

## III. METHOD

### A. Overview

To learn pedestrian representations robust to variations in clothing, pose, and environment, this paper proposes the BDAM. This module disentangles and extracts robust features

via contrastive and supervised learning, guided by encoded textual features. As illustrated in fig. 2, our framework comprises two primary feature extraction modules for vision and text, our core BDAM, and an efficient Mamba SSM Fusion Module.

Specifically, given a pedestrian image  $I \in \mathbb{R}^{B \times C \times H \times W}$ , a visual encoder first extracts image features  $f_i$ . To obtain semantic guidance, we use an MLLM with pre-designed prompts to generate corresponding descriptions for identity and clothing. A text encoder subsequently encodes these into  $f_{id}^t$  and  $f_{clo}^t$ . During disentanglement, BDAM leverages these textual features to guide the image feature learning process. To ensure the quality of this separation, we introduce a loss based on HSIC to constrain the two resulting feature types towards orthogonality. We also employ an alignment clothing loss, denoted as  $\mathcal{L}_{aln}$ , to supervise the learning of visual clothing features using clothing descriptions. Finally, to achieve a deep fusion of visual identity and textual semantics, we introduce the Mamba SSM as a fusion module. It dynamically models and facilitates interaction among the disentangled multimodal features, enhancing the model’s overall representation capability.

### B. Bidirectional Decoupled Alignment Module

Some studies directly adopt CLIP [14] as a feature extractor for both modalities, aligning global embeddings for retrieval or discrimination. However, this approach presents two key limitations. First, its limited capacity for fine-grained semantics hinders the separation of identity from clothing. Second, its holistic encoding of images and text lacks the modeling of cross-modal structure at the token level, which reduces robustness in complex scenes.

In this paper, we use a pre-trained ViT (ViT-B/16) as the visual encoder  $E_v$  [11]. Given an image  $I_i$ ,  $E_v$  outputs token features  $f_i \in \mathbb{R}^{B \times L \times D}$  that entangle cues relevant to identity and cues irrelevant to identity. A linear projection with two branches then yields preliminary identity features  $f'_{id} \in \mathbb{R}^{B \times L \times D}$  and clothing features  $f'_{clo} \in \mathbb{R}^{B \times L \times D}$ . These are followed by multi-layer self-attention in each branch to enhance local consistency and contextual awareness.

Instead of using the ViT [CLS] token as a global descriptor, we exploit the full patch sequence and introduce cross-attention between the branches to exchange information. In the identity stream, the clothing stream provides auxiliary context, and vice versa, reinforcing semantic distinctions. Each stream then applies global average pooling to produce  $\hat{f}_{id}$  and  $\hat{f}_{clo}$ . To enable soft disentanglement that is adaptive to the input, we design a gating mechanism. The two global vectors are concatenated and fed to a lightweight linear network with a Sigmoid output, producing a gate  $g \in \mathbb{R}^{B \times D}$ . We obtain the final gated features  $f_{id}^i = g \odot \hat{f}_{id}$  and  $f_{clo}^i = (1 - g) \odot \hat{f}_{clo}$ , where  $\odot$  denotes element-wise multiplication. This weighting at the dimension level provides a fine degree of control;  $f_{id}^i$  is further sent to the fusion module.

To train BDAM and enforce separation, we introduce two specialized loss functions. The first is a clothing alignment

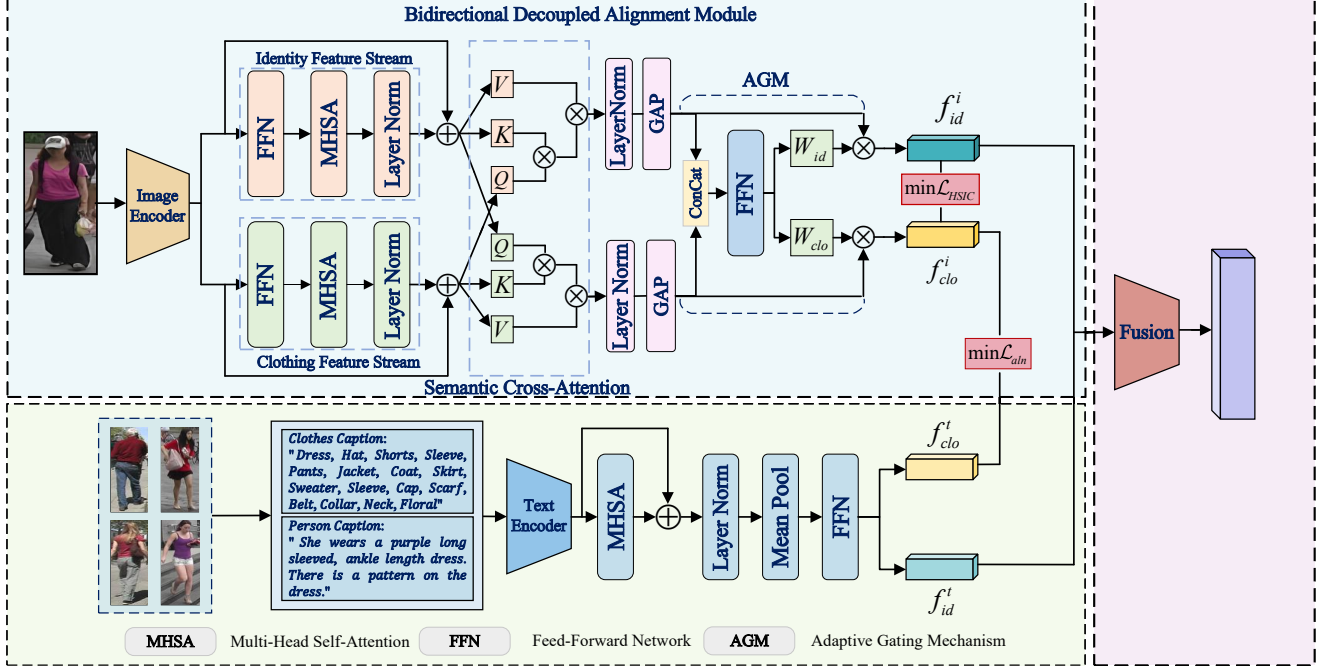


Fig. 2. First, we use an MLLM with predefined prompts to generate identity-related and clothing-related descriptions from the pedestrian image, which are encoded as  $f_{clo}^t$  and  $f_{id}^t$ , respectively. Then, the input pedestrian image is encoded into a visual feature  $f_i$ , which still resides in an entangled feature space. Subsequently, the BDAM module, composed of a dual-branch attention mechanism, decouples  $f_i$  into identity feature  $f_{id}^i$  and non-identity feature  $f_{clo}^i$ . The decoupling process is supervised and optimized through disentanglement loss and corresponding textual descriptions via contrastive learning. Finally, the fusion module integrates  $f_{id}^i$  and  $f_{id}^t$  to generate the final fused feature representation.

loss, and the second is a decoupling loss based on HSIC to enforce independence between identity and clothing. The clothing alignment loss supervises the visual clothing features with the MLLM-generated clothing descriptions, ensuring that the model accurately captures clothing semantics:

$$\mathcal{L}_{aln} = -\mathbb{E}_i \left[ \log \frac{\exp(s_{ii}/\tau)}{\sum_j \exp(s_{ij}/\tau)} \right] \quad (1)$$

where  $s_{ij} = \hat{f}_{clo}^i \cdot (f_{clo}^t)^j$  is the dot-product similarity between the visual clothing feature of sample  $i$  and the textual clothing feature of sample  $j$ , and  $\tau$  is a temperature parameter. This formulation encourages high similarity ( $s_{ii}$ ) for positive pairs (same sample) and low similarity ( $s_{ij}, i \neq j$ ) for negative pairs (different samples). In practice, clothing features are linearly projected to the text dimension and normalized using L2 for stable similarity estimation. This alignment objective explicitly ensures the clothing stream learns accurate representations under semantic supervision, which indirectly enhances the purity of the identity features by providing clear guidance on what constitutes clothing information. This works in conjunction with the cross-attention mechanism described earlier, which implicitly sharpens the separation via interaction.

To further encourage statistical independence, we minimize a

decoupling loss based on HSIC:

$$\begin{aligned} \mathcal{L}_{Decouple} &= \text{HSIC}(f_{id}^i, f_{clo}^i) \\ &= \frac{1}{(N-1)^2} \text{tr}(K_{id} H K_{clo} H) \end{aligned} \quad (2)$$

Here,  $f_{id}^i \in \mathbb{R}^{B \times D}$  and  $f_{clo}^i \in \mathbb{R}^{B \times D}$  are the gated identity and clothing features, respectively.  $K_{id} = f_{id}^i (f_{id}^i)^T$  and  $K_{clo} = f_{clo}^i (f_{clo}^i)^T$  are their respective kernel matrices.  $H = I_N - (1/N) \mathbf{1}_N \mathbf{1}_N^T$  is the centering matrix, where  $I_N$  is the  $N$ -dimensional identity matrix and  $\mathbf{1}_N$  is a column vector of all ones. HSIC measures the statistical dependence between features by calculating the mean trace of the product of their kernel matrices and the centering matrix. By minimizing this value, the loss encourages the features to be statistically independent.

### C. Semantic enhancement

We employ an MLLM, as detailed in Section IV, to automatically generate fine-grained identity and clothing descriptions for pedestrians. This approach reduces the burden of manual annotation and enriches the available supervision. Prior work, notably HAM [17], shows that modeling annotator styles can steer an MLLM to produce diverse texts. Adapting this core insight, we extend the pipeline to meet our model's design goals.

We first use the CLIP text encoder to embed the original descriptions into vectors of a fixed dimension. Using prompts,

an MLLM generalizes and substitutes entity attributes to emphasize expression style rather than content. We then cluster these style embeddings with DBSCAN [29], which adaptively discovers dense regions without predefining the cluster count. To stabilize the clusters, we reassign noise points and merge small clusters. This setup aligns the learned style categories with the identity and clothing disentanglement expected by BDAM.

A dual prompt generator guides the MLLM to output two distinct texts per image: one description for identity, covering biological traits, and another for clothing, detailing apparel, colors, and patterns. We control the generation process with length and temperature constraints. We also apply syntax checks and validation for attribute coverage to ensure the outputs remain grammatical, structured, and parsable.

In summary, our adaptations deliver flexible style modeling via DBSCAN and a decoupled dual description mechanism. This process strengthens the distinctiveness and diversity of identity and clothing semantics. The resulting supervision improves data expressiveness and provides richer training signals for BDAM.

#### D. Feature Fusion

For efficient feature fusion that is sensitive to semantics, we introduce the Mamba SSM. The core objective is to preserve the semantic integrity of the purified identity features from both image and text, enabling a fusion that is robust to clothing variations previously isolated by the BDAM. The process begins with an FFN performing dimensional alignment to mitigate distributional discrepancies between modalities. It processes the decoupled visual features  $f_{id}^t$  and the textual features  $f_{id}^t$  to generate aligned features,  $f_{img}$  and  $f_{txt}$ . Notably, the decoupled visual clothing feature  $f_{clo}^t$  is *intentionally discarded* during fusion. This design is central to our goal: the BDAM, supervised by  $\mathcal{L}_{aln}$  and  $\mathcal{L}_{Decouple}$ , is tasked with purging information irrelevant to identity into  $f_{clo}^t$ . By excluding this feature from the final fusion, the model is forced to learn a representation based purely on stable identity semantics.

Following this alignment, a gating mechanism achieves dynamic weighted fusion. It outputs a weight vector  $g_{fus} \in \mathbb{R}^{B \times 2}$ , which is normalized via a SoftMax layer to produce image  $W_{img}$  and text  $W_{txt}$  weights, satisfying  $W_{img} + W_{txt} = 1$ . The resulting fusion is computed as:  $f_{fusion} = W_{img} \cdot f_{img} + W_{txt} \cdot f_{txt}$ . This mechanism allows the model to adaptively balance modal contributions based on context. This fusion gate is distinct from the one in the disentanglement module; it outputs a global, two-dimensional weight vector  $g_{fus} \in \mathbb{R}^{B \times 2}$  for the modalities, whereas the disentanglement gate provides a vector  $g \in \mathbb{R}^{B \times D}$  for feature control at the dimension level. The resulting  $f_{fusion}$  features are then processed by the Mamba SSM to enhance interaction between modalities. Leveraging its capability to model dependencies over long ranges, Mamba effectively captures complex sequential relationships. We employ a stack of Mamba layers, where each layer updates its input  $f_{fusion}^{(l)}$  using a residual connection:  $f_{fusion}^{(l+1)} = \text{Mamba}(f_{fusion}^{(l)}) + f_{fusion}^{(l)}$ . This structure mitigates

the vanishing gradient problem and improves information flow. Finally, the output from the last Mamba layer is projected to produce the final representation,  $f_{final} \in \mathbb{R}^{B \times D_{out}}$ . The resulting feature is highly adaptive in its modal weighting and benefits from Mamba’s semantic modeling, providing robust support for downstream tasks like person re-identification.

#### E. Loss Function

To achieve alignment between modalities at a fine-grained level, we adopt the InfoNCE loss. This loss maximizes similarity for positive image and text pairs (representing the same identity) while separating negatives. It is defined as:

$$\mathcal{L}_{info} = -\log \frac{\exp(v_i^\top t_i / \tau)}{\sum_j \exp(v_i^\top t_j / \tau)} \quad (3)$$

Here,  $v_i$  is the final fused representation, normalized using the L2 norm;  $t_i$  is the text feature for the  $i$ -th identity; and  $\tau$  controls distribution sharpness. Negatives within the batch help reduce the semantic gap between modalities and promote alignment in a shared space.

To enhance identity discrimination within a single modality, we include a triplet loss. This loss enforces compactness within classes and separation between classes:

$$\mathcal{L}_{triplet} = \mathbb{E}_{(a,p,n)} \left[ \max \left( \|f_a - f_p\|_2^2 - \|f_a - f_n\|_2^2 + m, 0 \right) \right] \quad (4)$$

Here,  $f_a, f_p, f_n$  are decoupled visual identity features of the anchor, positive, and negative samples, respectively;  $\|\cdot\|_2$  denotes the L2 norm; and  $m$  is the margin parameter used to enforce a minimum distance gap between positive and negative pairs.

In training with multiple tasks, differing loss scales can cause one task to dominate. We adopt GradNorm to balance the training process by dynamically adjusting the gradient norm of each task:

$$\mathcal{L}_{GradNorm} = \sum_k |\nabla_\theta(w_k \mathcal{L}_k) - \tilde{r}_k G_{ref}| \quad (5)$$

Here,  $\nabla_\theta(w_k \mathcal{L}_k)$  represents the gradient of the weighted loss of task  $k$ ,  $w_k \mathcal{L}_k$ , with respect to the shared parameters  $\theta$ . The term  $|\cdot|_1$  denotes the L1 norm, which emphasizes a linear penalty on the deviation.  $\tilde{r}_k = (\mathcal{L}_k / \mathcal{L}_k^0) / \bar{r}$  is the normalized loss ratio, where  $\mathcal{L}_k^0$  is the initial loss of task  $k$  at the start of training, serving as a baseline, and  $\bar{r}$  is the average of the loss ratios over all tasks.  $G_{ref}$  is a reference gradient norm, typically set to the gradient norm of the first task,  $\|\nabla_\theta(w_1 \mathcal{L}_1)\|$ .

This mechanism enforces a gradient balance among tasks during training by minimizing the L1 deviation between the actual gradient norm and a target value,  $\tilde{r}_k G_{ref}$ . Furthermore, to prevent instability caused by the abnormal scaling of task weights  $w_k$ , we add a regularization term,  $\lambda \sum_k (\log w_k)^2$ , to the total loss. This term effectively suppresses large discrepancies among the task weights by penalizing the square of

their logarithms, thereby improving training stability. Finally, the overall loss function is defined as follows:

$$\mathcal{L}_{\text{Total}} = \sum_k w_k \mathcal{L}_k + \alpha \mathcal{L}_{\text{GradNorm}} + \lambda \sum_k (\log w_k)^2 \quad (6)$$

Here,  $\mathcal{L}_k$  represents the loss term for task  $k$  (which includes  $\mathcal{L}_{\text{info}}$ ,  $\mathcal{L}_{\text{triplet}}$ ,  $\mathcal{L}_{\text{aln}}$ , and  $\mathcal{L}_{\text{Decouple}}$ ), and  $w_k = \exp(s_k)$  is the task weight, where  $s_k$  is a learnable parameter initialized to zero and optimized during training to capture task-specific uncertainty. The hyperparameter  $\alpha$  is used to control the strength of the GradNorm loss, while  $\lambda$  serves as the regularization coefficient.

#### IV. EXPERIMENTS

##### A. Implementation Details

**Datasets and Metrics.** We evaluate our method on three standard benchmarks: CUHK-PEDES [30], ICFG-PEDES [31], and RSTPReid [2]. We follow their official identity-based splits and report mean Average Precision (mAP) and Rank-k accuracy (R-1, R-5, R-10).

**Model and Training.** Our model employs a pre-trained ‘bert-base-uncased’ as the text encoder and ‘vit-base-patch16-224’ as the visual encoder. All images are resized to  $224 \times 224$  pixels. The BDAM module disentangles visual tokens into 768-dim identity and clothing features. A 2-layer Mamba fusion module (256-dim I/O, 16-dim state, 4-conv kernel) integrates the representations. Dropout is 0.1 throughout. We use the Adam optimizer (LR  $1 \times 10^{-4}$ , WD  $1 \times 10^{-3}$ ) with cosine annealing. The total loss combines InfoNCE, triplet, clothing alignment, and HSIC decoupling terms, which are dynamically balanced using GradNorm ( $\alpha = 1.5$ ).

**Data Augmentation and Reporting.** We use ‘clip-vit-base-patch32’ and DBSCAN to find style clusters, then use ‘ChatGPT-4’ to generate decoupled identity and clothing descriptions. All experiments are repeated with three random seeds (0, 1, 2), and we report the mean results.

##### B. Parameter and Efficiency Analysis

**Gate Weight Analysis.** We analyze the learned gate weights to validate our design. The BDAM identity weights ( $W_{\text{id}}$ ) exhibit a mean of 0.61 ( $\sigma = 0.17$ ), significantly exceeding clothing weights ( $W_{\text{clothing}} = 0.38$ ) [cite: 299]. This asymmetry confirms that identity-relevant features dominate the representation [cite: 300]. Conversely, the fusion module maintains balanced weights (means of 0.52 and 0.48) [cite: 302], indicating stable cross-modal alignment without modality collapse [cite: 303]. Detailed distributions are provided in the supplementary material.

**Efficiency of Mamba Fusion.** Our Mamba-based fusion achieves an optimal trade-off between accuracy and efficiency. Compared to a standard 4-layer Transformer, our module delivers comparable retrieval accuracy while significantly reducing parameters and FLOPs [cite: 382, cite: 385]. Against a simple baseline, it demonstrates marked mAP improvements with minimal computational overhead [cite: 386].

##### C. Ablation Study

We conduct systematic ablation studies on CUHK-PEDES to validate each component [cite: 307].

TABLE I  
ABLATION STUDY ON THE BDAM MODULE.

Method	mAP(%) $\uparrow$	R-1(%) $\uparrow$	R-5(%) $\uparrow$	R-10(%) $\uparrow$
Baseline (w/o BDAM)	59.81	70.54	85.49	91.26
+ BDAM	66.74	76.27	89.30	94.02
w/o Cross-Attn	62.56	71.39	87.05	92.98
w/o Gate	65.11	74.63	88.77	93.56
Shallow (3-layer)	64.27	73.74	88.09	93.32

**Disentanglement Module.** Table I shows that BDAM yields substantial gains [cite: 308]. Removing cross-attention, ablating the gating mechanism, or reducing depth all degrade performance, confirming that semantic interaction and adaptive control are essential for robust disentanglement [cite: 313-314].

TABLE II  
ABLATION STUDY ON THE FUSION MODULE.

Method	mAP(%) $\uparrow$	R-1(%) $\uparrow$	R-5(%) $\uparrow$	R-10(%) $\uparrow$
Baseline (w/o Fusion)	59.81	70.54	85.49	91.26
Full Fusion	69.58	78.42	90.74	95.11
w/o Mamba	66.89	75.73	89.06	93.92
w/o Gate	68.64	77.58	90.11	94.87
w/o Alignment	68.15	77.09	89.84	94.53

**Fusion Module.** Table II confirms the synergy of our fusion components. The Mamba SSM is critical; its removal causes the largest drop [cite: 344-345], highlighting the value of long-range dependency modeling. The gating and alignment layers provide further essential gains [cite: 346-347].

TABLE III  
ABLATION STUDY ON INDIVIDUAL LOSS COMPONENTS.

Method	mAP(%) $\uparrow$	R-1(%) $\uparrow$	R-5(%) $\uparrow$	R-10(%) $\uparrow$
Full Model	72.61	79.93	92.95	96.47
w/o InfoNCE	28.14	36.55	55.21	65.83
w/o Triplet	67.22	74.89	88.15	93.12
w/o Alignment	69.15	76.92	89.53	94.22
w/o Decoupling	70.03	77.81	90.11	94.98
w/o Gate Reg.	71.98	79.23	91.35	95.71

**Loss and Prompting.** Table III validates our multi-task loss. InfoNCE [cite: 362] and triplet losses [cite: 363] are foundational. Crucially, the degradation without alignment [cite: 364] or HSIC decoupling losses [cite: 365] proves that explicit constraints are necessary to enforce identity-clothing separation. For prompt generation, we found density-based clustering (DBSCAN) significantly outperforms K-Means and random sampling, as it adaptively handles irregular style distributions [cite: 353-354]. Full clustering results and t-SNE visualizations are in the supplementary material.

##### D. Comparisons with State-of-the-Art Methods

Table IV compares our method against state-of-the-art approaches. On CUHK-PEDES and RSTPReid, our method establishes new benchmarks, substantially surpassing both ViT-based methods and the strong CLIP-based HAM baseline. On

TABLE IV  
PERFORMANCE COMPARISON WITH STATE-OF-THE-ART METHODS ON THREE BENCHMARK DATASETS.

Method	Backbone	CUHK-PEDES				ICFG-PEDES				RSTPReid			
		R-1	R-5	R-10	mAP	R-1	R-5	R-10	mAP	R-1	R-5	R-10	mAP
Methods with CLIP backbone:													
IRRA [10]	CLIP-ViT	73.38	89.93	93.71	66.10	63.36	80.82	85.82	38.06	60.20	81.30	88.20	47.17
IRLT [32]	CLIP-ViT	73.67	89.71	93.57	65.94	63.57	80.57	86.32	38.34	60.51	82.85	89.71	47.64
CFAM [33]	CLIP-ViT	74.46	90.19	94.01	-	64.72	81.35	86.31	-	61.49	82.26	89.23	-
Propot [34]	CLIP-ViT	74.89	89.90	94.17	67.12	65.12	81.57	86.97	42.93	61.87	83.63	89.70	47.82
RDE [35]	CLIP-ViT	75.94	90.14	94.12	67.56	67.68	82.47	87.36	40.06	65.35	83.95	89.90	50.88
HAM [17]	CLIP-ViT	77.99	91.34	95.03	69.72	69.95	83.88	88.39	42.72	72.50	87.70	91.95	55.47
Methods with ViT backbone:													
CPCL [36]	ViT	70.03	87.28	91.78	63.19	62.60	79.07	84.46	36.16	58.35	81.05	87.65	45.81
PDReid [37]	ViT	71.59	87.95	92.45	65.03	60.93	77.96	84.11	36.44	56.65	77.40	84.70	45.27
SSAN [2]	ViT	61.37	80.15	86.73	-	54.23	72.63	79.53	-	43.50	67.80	77.15	-
CFine [38]	ViT	69.57	85.93	91.15	-	60.83	76.55	82.42	-	50.55	72.50	81.60	-
IVT [39]	ViT	65.59	83.11	89.21	-	56.04	73.60	80.22	-	46.70	70.00	78.80	-
Ours	ViT	79.93	92.95	96.47	72.61	68.68	84.29	89.74	41.78	74.33	88.85	92.95	57.68

ICFG-PEDES, we remain highly competitive. Notably, this high performance is achieved without leveraging large-scale pre-trained vision-language models, validating that our explicit identity-clothing disentanglement and efficient Mamba fusion provide superior cross-modal alignment through task-specific design.

## V. CONCLUSION AND LIMITATIONS

In this work, we proposed a novel framework to address interference caused by clothing in text-to-image person re-identification. Our approach is centered on feature decoupling guided by a Multimodal Large Language Model. We successfully demonstrated that an MLLM can provide supervision at a fine-grained level to guide our Bidirectional Decoupling Alignment Module. By design, this module explicitly isolates features relevant to identity while actively suppressing interference related to clothing through a combined alignment and orthogonal loss strategy based on a kernel function. Furthermore, our integration of a Mamba State Space Model proved to be an effective and efficient fusion strategy, adept at capturing dependencies between modalities. Our method's effectiveness was validated by achieving new state-of-the-art results on the CUHK-PEDES and RSTPReid benchmarks. Nonetheless, limitations remain. These primarily concern the potential for noise in descriptions generated by the MLLM and the need for further optimization for deployment at a large scale. Future work will focus on improving description reliability to mitigate noise and on enhancing inference efficiency. We also plan to explore the framework's applicability to more challenging scenarios, such as person Re-ID involving clothing changes, to further test the robustness of our decoupling mechanism.

## REFERENCES

- [1] Zhiyin Shao, Xinyu Zhang, Meng Fang, Zhifeng Lin, Jian Wang, and Changxing Ding, "Learning Granularity-Unified Representations for Text-to-Image Person Re-identification," 2022.
- [2] Zefeng Ding, Changxing Ding, Zhiyin Shao, and Dacheng Tao, "Semantically Self-Aligned Network for Text-to-Image Part-aware Person Re-identification," 2021.
- [3] Maryam Bukhari, Sadaf Yasmin, Sheneela Naz, Muazzam Maqsood, Jehyeok Rew, and Seungmin Rho, "Language and vision based person re-identification for surveillance systems using deep learning with LIP layers," *Image and Vision Computing*, vol. 132, pp. 104658, 2023.
- [4] Hiren Galiyawala and Mehul S. Raval, "Person Retrieval in Surveillance Using Textual Query: A Review," 2021.
- [5] Liwei Wang, Yin Li, and Svetlana Lazebnik, "Learning Deep Structure-Preserving Image-Text Embeddings," 2016.
- [6] Yuyu Wang, Chunjuan Bo, Dong Wang, Shuang Wang, Yunwei Qi, and Huchuan Lu, "Language Person Search with Mutually Connected Classification Loss," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Brighton, United Kingdom, 2019, pp. 2057-2061.
- [7] Chenyang Gao, Guanyu Cai, Xinyang Jiang, Feng Zheng, Jun Zhang, Yifei Gong, Pai Peng, Xiaowei Guo, and Xing Sun, "Contextual Non-Local Alignment over Full-Scale Representation for Text-Based Person Search," 2021.
- [8] Zijie Wang, Aichun Zhu, Jingyi Xue, Xili Wan, Chao Liu, Tian Wang, and Yifeng Li, "CAIBC: Capturing All-round Information Beyond Color for Text-based Person Retrieval," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 5314-5322.
- [9] Zhe Wang, Zhiyuan Fang, Jun Wang, and Yezhou Yang, "ViTAA: Visual-Textual Attributes Alignment in Person Search by Natural Language," 2020.
- [10] Ding Jiang and Mang Ye, "Cross-Modal Implicit Relation Reasoning and Aligning for Text-to-Image Person Retrieval," 2023.
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," 2021.
- [12] Wentao Tan, Changxing Ding, Jiayu Jiang, Fei Wang, Yibing Zhan, and Dapeng Tao, "Harnessing the Power of MLLMs for Transferable Text-to-Image Person ReID," 2024.
- [13] Shuyu Yang, Yinan Zhou, Yaxiong Wang, Yujiao Wu, Li Zhu, and Zhedong Zheng, "Towards Unified Text-based Person Retrieval: A Large-scale Multi-Attribute and Language Search Benchmark," 2023.
- [14] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever, "Learning Transferable Visual Models From Natural Language Supervision," 2021.
- [15] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu, "FILIP: Fine-grained Interactive Language-Image Pre-Training," 2021.

- [16] Qun Niu, Tao Chen, Xing Zhang, Yifan Wang, and Ning Liu, "LLM-Loc: Bootstrap single-image indoor localization with large language model," *Expert Systems with Applications*, vol. 291, pp. 128368, 2025.
- [17] Jiayu Jiang, Changxing Ding, Wentao Tan, Junhong Wang, Jin Tao, and Xiangmin Xu, "Modeling Thousands of Human Annotators for Generalizable Text-to-Image Person Re-identification," 2025.
- [18] Xin Wang, Hong Chen, Si'ao Tang, Zihao Wu, and Wenwu Zhu, "Disentangled Representation Learning," 2024.
- [19] Yang Liu, Zhaowen Wang, Hailin Jin, and Ian Wassell, "Multi-task Adversarial Network for Disentangled Feature Learning," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Salt Lake City, UT, USA, 2018, pp. 3743–3751.
- [20] Hao Cheng, Yufei Wang, Haojiang Li, Alex C. Kot, and Bihan Wen, "Disentangled Feature Representation for Few-shot Image Classification," 2021.
- [21] Joanna Materzynska, Antonio Torralba, and David Bau, "Disentangling visual and written concepts in CLIP," 2022.
- [22] Yubo Li, De Cheng, Chaowei Fang, Changzhe Jiao, Nannan Wang, and Xinbo Gao, "Disentangling Identity Features from Interference Factors for Cloth-Changing Person Re-identification," in *Proceedings of the 32nd ACM International Conference on Multimedia*, Melbourne VIC Australia, 2024, pp. 2252–2261.
- [23] Shehreen Azad and Yogesh Singh Rawat, "Activity-Biometrics: Person Identification from Daily Activities," 2024.
- [24] Johann Schmidt and Sebastian Stober, "Robust Canonicalization through Bootstrapped Data Re-Alignment," 2025.
- [25] Haoli Yin, Jiayao Li, Eva Schiller, Luke McDermott, and Daniel Cummings, "GraFT: Gradual Fusion Transformer for Multimodal Re-Identification," 2023.
- [26] Seoyeon Kim, Minguk Kang, Dongwon Kim, Jaesik Park, and Suha Kwak, "Extending CLIP's Image-Text Alignment to Referring Image Segmentation," 2024.
- [27] Xinyu Feng, Yuming Lin, Lihua He, You Li, Liang Chang, and Ya Zhou, "Knowledge-Guided Dynamic Modality Attention Fusion Framework for Multimodal Sentiment Analysis," 2024.
- [28] Jiawei Li, Jiansheng Chen, Jinyuan Liu, and Huimin Ma, "Learning a Graph Neural Network with Cross Modality Interaction for Image Fusion," 2023.
- [29] Nooshin Hanafi and Hamid Saadatfar, "A fast DBSCAN algorithm for big data based on efficient density calculation," *Expert Systems with Applications*, vol. 203, pp. 117501, 2022.
- [30] Shuang Li, Tong Xiao, Hongsheng Li, Bolei Zhou, Dayu Yue, and Xiaogang Wang, "Person Search with Natural Language Description," 2017.
- [31] Aichun Zhu, Zijie Wang, Yifeng Li, Xili Wan, Jing Jin, Tian Wang, Fangqiang Hu, and Gang Hua, "DSSL: Deep Surroundings-person Separation Learning for Text-based Person Retrieval," 2021.
- [32] Yu Liu, Guihe Qin, Haipeng Chen, Zhiyong Cheng, and Xun Yang, "Causality-Inspired Invariant Representation Learning for Text-Based Person Retrieval," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, pp. 14052–14060, 2024.
- [33] Jialong Zuo, Hanyu Zhou, Ying Nie, Feng Zhang, Tianyu Guo, Nong Sang, Yunhe Wang, and Changxin Gao, "UFineBench: Towards Text-based Person Retrieval with Ultra-fine Granularity," 2024.
- [34] Shuanglin Yan, Jun Liu, Neng Dong, Liyan Zhang, and Jinhui Tang, "Prototypical Prompting for Text-to-image Person Re-identification," 2024.
- [35] Yang Qin, Yingke Chen, Dezhong Peng, Xi Peng, Joey Tianyi Zhou, and Peng Hu, "Noisy-Correspondence Learning for Text-to-Image Person Re-identification," 2024.
- [36] Yanwei Zheng, Xinpeng Zhao, Chuanlin Lan, Xiaowei Zhang, Bowen Huang, Jibin Yang, and Dongxiao Yu, "CPCL: Cross-Modal Prototypical Contrastive Learning for Weakly Supervised Text-based Person Re-Identification," 2024.
- [37] Weihao Li, Lei Tan, Pingyang Dai, and Yan Zhang, "Prompt Decoupling for Text-to-Image Person Re-identification," 2024.
- [38] Shuanglin Yan, Neng Dong, Liyan Zhang, and Jinhui Tang, "CLIP-Driven Fine-grained Text-Image Person Re-identification," 2022.
- [39] Xiujun Shu, Wei Wen, Haoqian Wu, Keyu Chen, Yiran Song, Ruizhi Qiao, Bo Ren, and Xiao Wang, "See Finer, See More: Implicit Modality Alignment for Text-based Person Retrieval," 2022.