

DRAGE: Dual-Gaps Robust Alignment for Text-Based Person Search

Anonymous ACL submission

Abstract

Text-based Person Search (TBPS) faces two critical and intertwined challenges: the semantic gap caused by noisy image-text correspondences, and the modality gap stemming from the structural heterogeneity between dense visual features and sparse textual attributes. Crucially, we identify that this heterogeneity leads to intra-modal disorganization, where embeddings—particularly sparse text representations—lack internal structure, hindering robust alignment. To address these limitations, we propose Dual-Gaps Robust Aligned General Embedding (DRAGE), a unified framework comprising two synergistic mechanisms. First, Semantic Embedding Reliability Division (SERD) dynamically partitions training data into reliable and noisy subsets using a Beta Mixture Model, providing clean supervision to mitigate the semantic gap. Second, Semantic Intra-Modality Regularization (SIMR) explicitly addresses intra-modal disorganization by enforcing semantic-aware distance constraints within each modality *before* cross-modal alignment. This transforms disorganized embeddings into coherent semantic clusters, establishing a stable geometric foundation for matching. Extensive experiments on CUHK-PEDES, ICFG-PEDES, and RSTPReid demonstrate that DRAGE significantly outperforms state-of-the-art methods, achieving **78.58%** Rank-1 accuracy on CUHK-PEDES and exhibiting superior robustness in cross-domain settings.

1 Introduction

Text-based Person Search (TBPS) aims to retrieve a specific person from a large image gallery using natural language descriptions (Lei et al., 2022; Miech et al., 2021; Sun et al., 2021). This task has critical applications ranging from video surveillance to social media analysis (Li et al., 2025; Zhao et al., 2024). Despite significant progress, deploying TBPS models in real-world scenarios remains constrained by two fundamental and intertwined

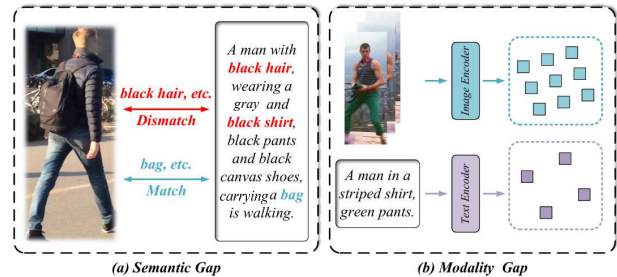


Figure 1: Illustration of the semantic gap and modality gap in Text-based Person Search. (a) **Semantic gap**: training image-text pairs often exhibit content inconsistencies due to annotation biases or visual ambiguities, resulting in noisy supervision. (b) **Modality gap**: the fundamental heterogeneity between dense visual features and sparse textual features, exacerbated by disorganized internal embedding spaces that hinder robust cross-modal alignment.

challenges: the semantic gap and the modality gap, as illustrated in Figure 1.

Recent research makes significant progress in TBPS. Many methods explore both global (Jiang and Ye, 2023b; Xie et al., 2025a) and local matching (Zuo et al., 2024; Park et al., 2024) strategies. Studies leverage large pre-trained models (Li et al., 2025; Liu et al., 2025b) by adapting vision-language models such as CLIP (Radford et al., 2021), or employ generative models for data augmentation through diffusion-based image synthesis (Song et al., 2024) and language-model-driven text enrichment (Tan et al., 2024; Xie et al., 2025b). Despite these advancements, existing TBPS methods remain fundamentally constrained by two critical yet rarely co-addressed challenges.

The first is the **semantic gap**, which arises from noisy and inconsistent image-text correspondences caused by real-world data variability such as pose variations, camera angles, and illumination changes. The second is the **modality gap**, stemming from a deeper structural disparity operating on two levels: (1) *cross-modal heterogeneity*—the

fundamental difference between **dense visual features** (characterized by redundant spatial information) and **sparse textual attributes** (composed of discrete, abstract tokens); and (2) *intra-modal disorganization*—the lack of structural organization within each modality’s embedding space. Existing methods focus exclusively on cross-modal alignment while neglecting this internal organization. As a result, embeddings within a single modality (especially the sparse text modality) tend to be distributed chaotically, where semantically similar samples scatter and dissimilar ones cluster together. This internal disorder sets a low upper bound for the subsequent cross-modal alignment.

To address these limitations, we propose **Dual-Gaps Robust Aligned General Embedding (DRAGE)**, a unified framework that systematically tackles both challenges through two synergistic mechanisms.

First, to address the *semantic gap*, we propose **Semantic Embedding Reliability Division (SERD)**. Leveraging a Beta Mixture Model, SERD dynamically separates reliable image-text pairs from noisy ones. By prioritizing high-confidence data and suppressing noise, it ensures the model focuses on valid semantic correspondences.

Second, to bridge the *modality gap*, DRAGE employs **Semantic Intra-Modality Regularization (SIMR)**. This module addresses the overlooked issue of internal disorganization. Rather than immediately forcing cross-modal alignment, SIMR first organizes each modality’s internal embedding structure by enforcing semantic-aware distance constraints: pulling similar samples closer while pushing dissimilar ones apart. This intra-modal organization is particularly crucial for sparse text representations, transforming disorganized embeddings into coherent semantic clusters and compensating for their inherent lack of structure. By establishing stable, well-organized foundations within each modality *before* cross-modal alignment, SIMR enables more effective and robust mapping into a unified space.

By integrating SERD and SIMR, DRAGE systematically addresses both gaps to achieve state-of-the-art performance. The main contributions of this work are:

- We propose DRAGE, a unified framework that addresses the semantic and modality gaps in TBPS. Unlike prior works that treat these challenges separately, DRAGE integrates SERD

and SIMR to jointly mitigate noisy correspondences and structural heterogeneity.

- We introduce SERD, a dynamic reliability-aware mechanism. By leveraging a Beta Mixture Model to model per-sample loss distributions, SERD effectively filters noisy supervision signals caused by real-world data inconsistencies, ensuring the model learns from semantically reliable pairs.
- We propose SIMR to address the overlooked issue of intra-modal disorganization. By enforcing semantic-aware distance constraints, SIMR structures the embedding spaces of dense visual and sparse textual features *before* cross-modal alignment, providing a stable geometric foundation for matching.
- Extensive experiments on three benchmarks demonstrate that DRAGE achieves state-of-the-art performance. **78.58%** Rank-1 accuracy on CUHK-PEDES. Furthermore, the framework exhibits superior robustness in cross-domain settings, validating the effectiveness of our dual-gap mitigation strategy.

2 Related Work

Text-based Person Search. TBPS aims to match textual descriptions with images to identify specific individuals. Early TBPS approaches focus on extracting global features by mapping both image and text data into a shared embedding space (Ge et al., 2019). However, these methods struggle to capture the fine-grained details necessary to differentiate similar-looking individuals. To improve accuracy, researchers develop local-matching methods (Chen et al., 2018; Guo et al., 2019; Ding et al., 2020) that explicitly associate specific regions in images with corresponding words in text descriptions. Although these methods yield more precise alignments, they incur significantly higher computational costs due to their complexity. Recent advances using pre-trained models such as CLIP (Radford et al., 2021) and ViT (Ge et al., 2019) further advance TBPS by leveraging large-scale multimodal learning to generate enhanced representations. Nevertheless, these models still encounter two key challenges: the semantic gap and the modality gap. Current methods have not fully resolved these challenges—the semantic gap results in misalignments even when text and images refer to the same person, while the

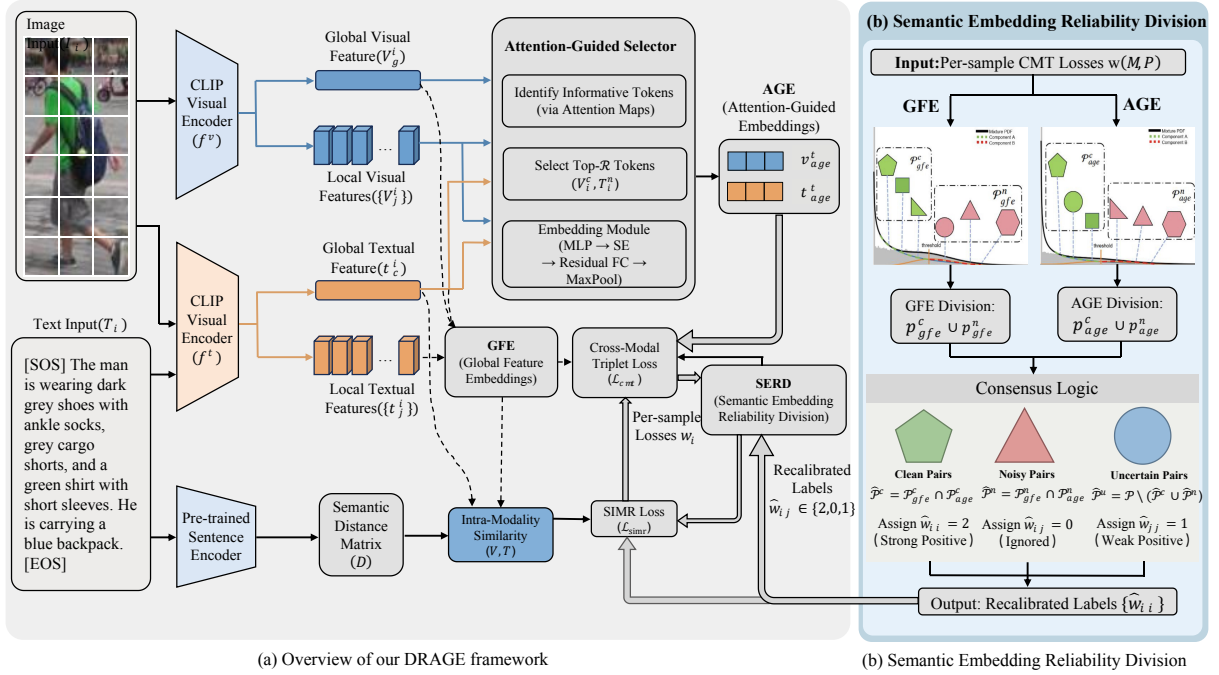


Figure 2: Overview of the DRAGE framework. DRAGE addresses the semantic gap through SERD, which identifies reliable image-text pairs for robust supervision, and bridges the modality gap through SIMR, which organizes each modality’s internal embedding structure before cross-modal alignment. The framework integrates Attention-Guided Selector for discriminative features, Cross-Modal Triplet Loss for robust matching for alignment optimization.

165 modality gap persists due to fundamental differ-
 166 ences between the two modalities. Consequently,
 167 more sophisticated techniques are required to sys-
 168 tematically address these gaps and enhance TBPS
 169 performance in real-world applications.

170 **Modality Gap.** Early work such as
 171 CPM (Zhang and Lu, 2018a) tackles this
 172 challenge by aligning global features across modal-
 173 ities. More recent efforts enhance this alignment
 174 through local attribute modeling (Yan et al., 2022a;
 175 Zuo et al., 2024) and refined loss functions (Jiang
 176 and Ye, 2023a; Qin et al., 2024), yielding better
 177 semantic correspondence across views. However,
 178 these methods typically treat each modality’s
 179 embedding space as a black box—focusing only
 180 on cross-modal pairing while largely ignoring the
 181 internal structure of each modality. As a result,
 182 embeddings within the same modality can remain
 183 disorganized, leading to unstable representations
 184 that hinder alignment consistency, especially in
 185 challenging or noisy scenarios. In contrast, our
 186 study explicitly addresses this gap by proposing a
 187 dedicated intra-modal regularization strategy that
 188 enforces meaningful distance constraints within
 189 each modality, stabilizing the feature space for
 190 robust cross-modal alignment.

3 Method

3.1 Preliminary

192 The goal of TBPS is to retrieve pedestrian im-
 193 ages matching a textual description. Let the im-
 194 age set be $\mathcal{V} = \{I_i\}_{i=1}^N$ and the text set be
 195 $\mathcal{T} = \{T_i\}_{i=1}^N$. We aim to learn a joint em-
 196 bedding space where matched pairs (I_i, T_i) are close.
 197 We employ CLIP (Radford et al., 2021) as the
 198 backbone, consisting of a visual encoder f^v and
 199 a textual encoder f^t . For an image I_i , f^v out-
 200 puts features $V_i = \{v_g^i, v_1^i, \dots\}^\top$, where v_g^i
 201 is the global [CLS] token. For a text T_i , f^t out-
 202 puts $T_i = \{t_s^i, \dots, t_e^i\}^\top$, where t_e^i represents the global
 203 [EOS] feature.
 204

3.2 Attention-Guided Selector

205 Global features often miss fine-grained details. To
 206 address this, we use Attention-Guided Selector to
 207 extract discriminative local tokens. First, we se-
 208 lect the top- \mathcal{R} most informative local tokens from
 209 last attention layer, denoted as \hat{V}_i^a and \hat{T}_i^a . Then,
 210 we project them into compact vectors, termed
 211 Attention-Guided Embedding (AGE) vectors:
 212

$$v_{age}^i = \text{MaxPool}(\text{SE}(\text{MLP}(\hat{V}_i^a)) + \text{FC}(\hat{V}_i^a)), \quad (1)$$

$$t_{age}^i = \text{MaxPool}(\text{SE}(\text{MLP}(\hat{T}_i^a)) + \text{FC}(\hat{T}_i^a)). \quad (2)$$

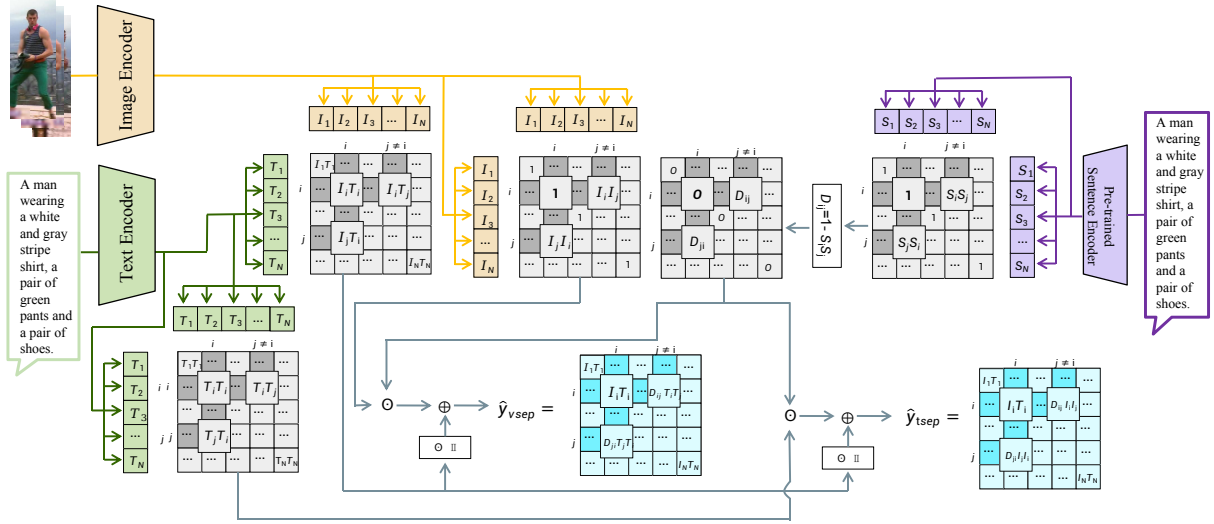


Figure 3: Overview of Semantic Intra-Modality Regularization (SIMR) method. This approach is implemented through vision-to-vision (v2v) and text-to-text (t2t) contrastive losses, adjusting embedding distances within each modality based on semantic distances between samples. Note: The input embeddings for SIMR are purely intra-modal (image-to-image or text-to-text).

The AGE similarity S_{ij}^a is the cosine similarity between these vectors, complementing the global similarity S_{ij}^g .

3.3 Semantic Embedding Reliability Division

Real-world data often contains noisy image-text pairs (e.g., mismatched descriptions), causing the *semantic gap*. SERD identifies and filters these pairs to ensure reliable supervision.

We leverage the "small-loss trick" (Arpit et al., 2017), observing that clean data exhibits lower loss values during early training. We first compute the per-sample loss $\ell_i = \mathcal{L}_{\text{cmt}}(I_i, T_i)$. To satisfy the input requirement of the Beta Mixture Model (BMM), we apply min-max normalization to map losses to $[0, 1]$:

$$\bar{\ell}_i = \frac{\ell_i - \min(\ell)}{\max(\ell) - \min(\ell)}. \quad (3)$$

We then fit a two-component BMM to $\bar{\ell}_i$. Based on the posterior probability $p(\text{clean}|\bar{\ell}_i)$ and a threshold $\delta = 0.6$, we first divide pairs into Reliable and Unreliable sets for both global and AGE views. By computing the consensus between these views, we assign a **recalibrated** reliability weight w_{ii} to each pair (I_i, T_i) :

$$w_{ii} = \begin{cases} 2, & \text{if deemed Reliable (Strong Positive),} \\ 1, & \text{if deemed Uncertain (Weak Positive),} \\ 0, & \text{if deemed Unreliable (Ignored).} \end{cases} \quad (4)$$

For negative pairs ($i \neq j$), we set $w_{ij} = 1$. This ensures the model learns only from semantically valid correspondences.

3.4 Cross-Modal Triplet Loss

To learn robust features against noisy correspondence, we introduce a reliability-aware Cross-Modal Triplet (CMT) loss. Unlike standard contrastive losses, we incorporate the reliability weights w_{ij} to prevent noisy positives from misguiding the gradient optimization.

Formally, the matching probability p_{ij} between image I_i and text T_j is defined as:

$$p_{ij} = \frac{w_{ij} \exp(S(I_i, T_j)/\tau)}{\sum_{k=1}^N w_{ik} \exp(S(I_i, T_k)/\tau)}, \quad (5)$$

where τ denotes the temperature parameter. Crucially, the inclusion of w_{ik} in the denominator ensures that p_{ij} remains a valid probability distribution over the batch.

Based on this, the reliability-aware CMT loss from image to text is formulated as:

$$\mathcal{L}_{\text{i2t}} = \frac{1}{N} \sum_{i=1}^N \left[\alpha - \hat{S}_{i,i} + \tau \log \left(\sum_{j=1}^N w_{ij} \exp \left(\frac{S(I_i, T_j)}{\tau} \right) \right) \right]_+, \quad (6)$$

where $[\cdot]_+ = \max(\cdot, 0)$ denotes the hinge function, and α is the margin hyperparameter. The

term $\hat{S}_{i,i} = \sum_j p_{ij} S(I_i, T_j)$ represents the expected similarity score of reliable positive pairs. This formulation encourages the model to maximize the similarity of reliable positive pairs while suppressing the scores of negative ones relative to the margin. The total CMT loss integrates both image-to-text and text-to-image directions across both Global and AGE feature views:

$$\mathcal{L}_{\text{cmt}} = \mathcal{L}_{\text{cmt}}^{\text{gfe}} + \mathcal{L}_{\text{cmt}}^{\text{age}}. \quad (7)$$

3.5 Semantic Intra-Modality Regularization

As illustrated in Figure 3, the *modality gap* is exacerbated by *intra-modal disorganization*. SIMR addresses this by enforcing semantic constraints *within* each modality to organize the embedding space before cross-modal alignment.

3.5.1 Target Topology Construction

We employ a pre-trained **Sentence Encoder** (e.g., MPNet) to define the ground-truth semantic topology, as shown in the purple blocks of Figure 3. Let $\mathbf{S} = [s_1, \dots, s_N]^\top \in \mathbb{R}^{N \times d_s}$ denote the batch of sentence embeddings derived from the raw descriptions. We compute the target semantic distance matrix \mathcal{D} by performing strictly row-wise normalization to ensure valid cosine distances:

$$\mathcal{S}_{ij}^s = \frac{s_i \cdot s_j}{\|s_i\| \|s_j\|}, \quad \mathcal{D}_{ij} = 1 - \mathcal{S}_{ij}^s. \quad (8)$$

3.5.2 Intra-Modal Alignment

Next, we align the internal structure of our trainable features to this target \mathcal{D} . To maintain consistency with the global features defined in Sec. 3.1, let $\mathbf{I} = [v_g^1, \dots, v_g^N]^\top$ and $\mathbf{T} = [t_e^1, \dots, t_e^N]^\top$ denote the batch matrices of global visual and textual embeddings, corresponding to the yellow and orange blocks in Figure 3, respectively.

We define calibrated logits \hat{y} by combining the cross-modal matching confidence (Bias) with intra-modal topological constraints (Topology):

$$\hat{y}_v = s \cdot \left[\underbrace{\text{diag}(\mathcal{S}^{i2t}) \cdot \mathbf{1}^\top}_{\text{Alignment Bias}} + \underbrace{\mathcal{S}^{i2i} \odot (1 - \mathcal{D})}_{\text{Topology Constraint}} \right], \quad (9)$$

$$\hat{y}_t = s \cdot \left[\text{diag}(\mathcal{S}^{t2i}) \cdot \mathbf{1}^\top + \mathcal{S}^{t2t} \odot (1 - \mathcal{D}) \right], \quad (10)$$

where s is a scaling factor. Here, $\mathcal{S}^{i2i} = \bar{\mathbf{I}}\bar{\mathbf{I}}^\top$ and $\mathcal{S}^{t2t} = \bar{\mathbf{T}}\bar{\mathbf{T}}^\top$ represent the intra-modal similarity matrices (the $I \cdot I$ and $T \cdot T$ grids in Figure 3), while $\mathcal{S}^{i2t} = \bar{\mathbf{I}}\bar{\mathbf{T}}^\top$ is the cross-modal similarity matrix.

The term $\text{diag}(\mathcal{S}^{i2t}) \cdot \mathbf{1}^\top$ involves **broadcasting** the diagonal elements (scores of positive pairs) across the rows. This acts as an alignment-aware bias, ensuring that the intra-modal clustering remains consistent with the cross-modal matching objective. The second term, $\mathcal{S}^{i2i} \odot (1 - \mathcal{D})$, enforces the target topology: features with small semantic distance ($\mathcal{D}_{ij} \approx 0$) are encouraged to have high similarity scores.

Finally, the SIMR loss is formulated as a classification task using the standard Cross-Entropy loss $H(\cdot, \cdot)$:

$$\mathcal{L}_{\text{simr}} = H(\hat{y}_v, Y_{\text{id}}) + H(\hat{y}_t, Y_{\text{id}}), \quad (11)$$

where Y_{id} represents the ground-truth identity labels.

3.6 Overall Optimization Objective

Our framework integrates SERD, CMT, and SIMR to address the dual gaps systematically. The total optimization objective is:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{simr}} + \mathcal{L}_{\text{cmt}}. \quad (12)$$

By jointly optimizing these losses, DRAGE mitigates the semantic gap via reliable supervision (SERD) and bridges the modality gap through structured intra-modal regularization (SIMR).

4 Experiments

4.1 Datasets and Performance Measurements

Datasets. We evaluate our method on three widely-used TBPS datasets: CUHK-PEDES (Li et al., 2017), ICFG-PEDES (Ding et al., 2021a), and RSTPReid (Zhu et al., 2021a). CUHK-PEDES is a dataset containing over 40,000 images of 13,003 identities. ICFG-PEDES includes over 54,000 images of 4,102 identities, and RSTPReid consists of over 20,000 images of 4,101 identities. A detailed statistical breakdown of each dataset is provided in the **supplementary material**.

Evaluation Metrics. In accordance with standard practices, search performance is assessed using the metrics Rank-k (with k set to 1, 5, and 10), mean Average Precision (mAP), and mean Inverse Negative Penalty (mINP) (Ye et al., 2022). Higher values of these metrics signify enhanced search performance.

4.2 Implementation Details.

Consistent with previous studies (Jiang and Ye, 2023a; Li et al., 2023; Liu et al., 2025a), DRAGE

Table 1: Performance comparisons on the CUHK-PEDES dataset. The best results are in **bold**.

Methods	Ref.	Image Enc.	Text Enc.	R-1	R-5	R-10	mAP	mINP
Lapscore (Wu et al., 2021)	ICCV'21	RN50	BERT	63.40	-	87.80	-	-
AXM-Net (Farooq et al., 2022)	ACMMM'22	RN50	BERT	64.44	80.52	86.77	58.73	-
LGUR (Shao et al., 2022)	ACMMM'22	DeiT-Small	BERT	65.25	83.12	89.00	-	-
DCEL (Qin et al., 2022)	ACMMM'22	CLIP-ViT	CLIP-Xformer	71.36	88.11	92.48	64.25	48.26
CFine (Yan et al., 2022a)	TIP'23	CLIP-ViT	BERT	69.57	85.93	91.15	-	-
IRRA (Jiang and Ye, 2023b)	CVPR'23	CLIP-ViT	CLIP-Xformer	73.38	89.93	93.71	66.13	50.24
PBSL (Shen et al., 2023)	ACMMM'23	RN50	BERT	65.32	83.81	89.26	-	-
BEAT (Ma et al., 2023)	ACMMM'23	RN101	BERT	65.61	83.45	89.54	-	-
LCR ² S (Yan et al., 2023)	ACMMM'23	RN50	TextCNN	67.36	84.19	89.62	59.24	-
DCEL (Li et al., 2023)	ACMMM'23	CLIP-ViT	CLIP-Xformer	75.02	90.89	94.52	-	-
APTM (Yang et al., 2023)	ACMMM'23	Swin-Transformer	Bert	76.53	90.04	94.15	66.91	-
UniPT (Shao et al., 2023)	ICCV'23	CLIP-ViT	CLIP-Xformer	68.50	84.67	-	-	-
RaSa (Bai et al., 2023)	IJCAI'23	ALBEF	ALBEF	76.51	90.29	94.25	69.38	-
SEN-XL (Liu et al., 2025a)	Neural Networks'25	CLIP-ViT	CLIP-Xformer	76.64	91.33	94.66	69.19	53.88
TBPS (Cao et al., 2024)	AAAI'24	CLIP-ViT	CLIP-Xformer	73.54	88.19	92.35	65.38	-
DP (Song et al., 2024)	AAAI'24	CLIP-ViT	CLIP-Xformer	75.66	90.59	94.07	66.58	-
IRRA+IRLT (Liu et al., 2024)	AAAI'24	CLIP-ViT	CLIP-Xformer	74.46	90.19	94.01	-	-
UASA (Zhao et al., 2024)	AAAI'24	CLIP-ViT	CLIP-Xformer	74.25	89.83	93.58	66.15	-
RDE (Qin et al., 2024)	CVPR'24	CLIP-ViT	CLIP-Xformer	75.94	90.14	94.12	67.56	51.44
CFAM (Zuo et al., 2024)	CVPR'24	CLIP-ViT	CLIP-Xformer	75.60	90.53	-	67.27	-
MGRL (Lv et al., 2024)	ICASSP'24	CLIP-ViT	CLIP-Xformer	73.91	90.68	-	67.28	-
DM-Adapter (Liu et al., 2025b)	AAAI'25	CLIP-ViT	CLIP-Xformer	72.17	88.74	92.85	64.33	-
OCDL (Li et al., 2025)	ICASSP'25	CLIP-ViT	CLIP-Xformer	75.10	89.43	-	68.18	-
HAM(1.0 M) (Jiang et al., 2025)	CVPR'25	CLIP-ViT	CLIP-Xformer	77.71	91.42	94.57	69.68	-
DRAGE (Ours)	-	CLIP-ViT	CLIP-Xformer	78.58	91.68	95.19	70.68	53.89

utilizes CLIP-ViT as the image encoder, the CLIP text transformer as the text encoder, and in addition, all-mpnet-base-v2 as the sentence encoder. Standard data augmentation techniques, including random horizontal flipping, random cropping, and random erasing, are applied. An Attention-Guided Selector with a sampling ratio of 0.4 is employed. The images are resized to 384×128 pixels, and text sequences are truncated to 77 tokens. Training is conducted for 60 epochs using the Adam optimizer, with an initial learning rate of 1×10^{-6} , which decays according to a cosine schedule after a 2-epoch warm-up. Furthermore, the α parameter in CMT is set to 0.1, and the temperature τ is set to 0.015. All experiments are carried out on a single NVIDIA RTX 4090 GPU.

4.3 Comparison with State-of-the-Art

In this section, we compare our method with current sota approaches on widely-used datasets.

Performance Comparisons on CUHK-PEDES.

Table 1 presents a comprehensive comparison of

different methods on the CUHK-PEDES dataset, a well-established benchmark for Text-Based Person Search (TBPS). Our method, DRAGE, demonstrates superior performance across all key metrics, achieving **78.58%** in Rank-1 accuracy, 91.68% in Rank-5, and 95.19% in Rank-10. While several state-of-the-art methods, including DCEL, APTM, and RaSa, also yield competitive results, DRAGE's performance significantly surpasses theirs, underscoring its enhanced capability on this benchmark. Crucially, these results are achieved without employing post-processing techniques like re-ranking (e.g., ICL (Qin et al., 2025) and MTI (Xie et al., 2025b)) or leveraging external data sources, unlike competing methods such as APTM and HAM.

Performance Comparisons on ICFG-PEDES.

We present the results on the ICFG-PEDES dataset in Table 4. Our method achieves a Rank-1 accuracy of 69.63%, a Rank-5 accuracy of 83.40%, and a mAP of 42.08%. Compared to the existing methods, our approach outperforms several recent models, including RaSa and TBPS, which show perfor-

mance gains of up to 4.56% on Rank-1 and 4.67% on mAP. Despite the relatively low mINP score of 8.97%, which reflects the challenge in identifying the hardest samples, our method demonstrates solid retrieval performance.

Performance Comparisons on RSTPReid. We report results on the RSTPReid dataset in Table 5. Our method achieves a Rank-1 accuracy of 67.55%, Rank-5 accuracy of 85.00%, and a mAP of 52.88%, significantly outperforming other methods, including TBPS and CFAM, by a considerable margin. In comparison to the global-matching method IVT (Shu et al., 2022), our approach shows an improvement of 13.5% on Rank-1, 11.3% on Rank-5, and 9.4% on Rank-10. Additionally, when compared to the local-matching method Cfine (Yan et al., 2022a), we achieve a notable improvement in performance with gains of up to 9.65 on Rank-1.

Table 2: Performance comparisons on the ICFG-PEDES dataset. The best results are in **bold**.

Methods	R-1	R-5	R-10	mAP	mINP
Dual Path (Zheng et al., 2020)	38.99	59.44	68.41	-	-
CMPM/C (Zhang and Lu, 2018b)	43.51	65.44	74.26	-	-
ViTAA (Wang et al., 2020)	50.98	68.79	75.78	-	-
SSAN (Ding et al., 2021b)	54.23	72.63	79.53	-	-
IVT (Shu et al., 2022)	56.04	73.60	80.22	-	-
ISANet (Yan et al., 2022b)	57.73	75.42	81.72	-	-
CFine (Yan et al., 2022a)	60.83	76.55	82.42	-	-
IRRA (Jiang and Ye, 2023b)	63.46	80.25	85.82	38.06	7.93
BiLMA (Fujii and Tarashima, 2023)	63.83	80.15	85.74	38.26	-
PBSL (Shen et al., 2023)	57.84	75.46	82.15	-	-
BEAT(Ma et al., 2023)	58.25	75.92	81.96	-	-
LCR ² S (Yan et al., 2023)	57.93	76.08	82.40	38.21	-
DCEL (Li et al., 2023)	64.88	81.34	86.72	-	-
UniPT (Shao et al., 2023)	60.09	76.19	-	-	-
RaSa (Bai et al., 2023)	65.28	80.40	85.12	41.29	-
TBPS (Cao et al., 2024)	65.05	80.34	85.47	39.83	-
CFAM (Zuo et al., 2024)	65.38	81.17	-	39.42	-
MGRL (Lv et al., 2024)	67.28	63.87	-	82.34	-
OC DL (Li et al., 2025)	64.53	80.23	-	40.76	-
DRAGE (Ours)	69.63	83.40	88.04	42.08	8.97

Table 3: Performance comparisons on the RSTPReid dataset. The best results are in **bold**.

Methods	R-1	R-5	R-10	mAP	mINP
DSSL (Zhu et al., 2021b)	39.05	62.60	73.95	-	-
SSAN (Ding et al., 2021b)	43.50	67.80	77.15	-	-
LBUL (Wang et al., 2022)	45.55	68.20	77.85	-	-
IVT (Shu et al., 2022)	46.70	70.00	78.80	-	-
CFine (Yan et al., 2022a)	50.55	72.50	81.60	-	-
IRRA (Jiang and Ye, 2023b)	60.20	81.30	88.20	47.17	25.28
BiLMA (Fujii and Tarashima, 2023)	61.20	81.50	88.80	48.51	-
PBSL (Shen et al., 2023)	47.80	71.40	79.90	-	-
BEAT(Ma et al., 2023)	48.10	73.10	81.30	-	-
LCR ² S (Yan et al., 2023)	54.95	76.65	84.70	40.92	-
DCEL (Li et al., 2023)	61.35	83.95	90.45	-	-
RaSa (Bai et al., 2023)	66.90	86.50	91.35	52.31	-
TBPS (Cao et al., 2024)	61.95	83.55	88.75	48.26	-
CFAM (Zuo et al., 2024)	62.45	83.55	-	49.50	-
OC DL (Li et al., 2025)	61.60	82.35	-	49.77	-
DRAGE (Ours)	67.55	85.00	90.90	52.88	30.05

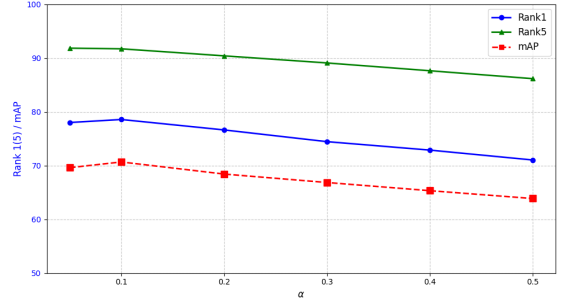


Figure 4: Performance variation with different values of the margin parameter α .

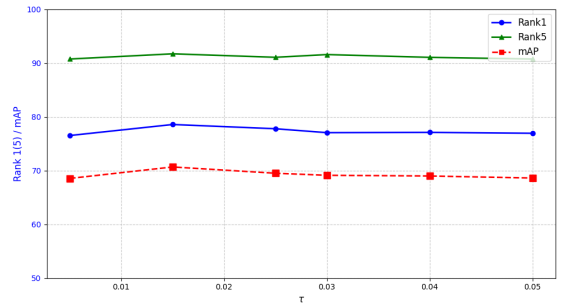


Figure 5: Performance variation with different values of the temperature parameter τ .

4.4 Ablation Study

We conduct ablation experiments on the CUHK-PEDES dataset to evaluate the contribution of each component in DRAGE. As shown in Table 5, removing any module leads to significant performance degradation, demonstrating that each component is essential. Specifically, removing SERD causes the largest drop (-6.40%), confirming its critical role in mitigating the semantic gap by identifying noisy pairs and providing reliable supervision. Without SIMR, performance decreases by 3.16%, validating its effectiveness in bridging the modality gap by organizing internal structures before cross-modal alignment. Eliminating AGE results in a 2.48% decline, showing that local features effectively complement global representations. The loss \mathcal{L}_{cmt} also contributes to robust alignment. When removing both SERD and SIMR, performance drops sharply to 68.45%, exceeding the sum of their individual effects. This demonstrates strong complementarity: SERD provides reliable supervision that helps SIMR build stable internal structures, while SIMR’s organized spaces help SERD better identify noisy pairs.

Table 4: Performance of methods and components in cross-domain tests. CUHK-PEDES is abbreviated as CUHK, ICFG-PEDES as ICFG, and RSTPReid as RSTP. The arrow \rightarrow indicates that the method was trained on the dataset on the left and tested on the dataset on the right. The best results are in bold.

Method	Domain	R1	R5	mAP	Domain	R1	R5	mAP
IRRA (Jiang and Ye, 2023a)	CUHK \rightarrow RSTP	53.30	77.15	39.63	CUHK \rightarrow ICFG	42.42	62.11	21.77
	ICFG \rightarrow RSTP	45.30	69.35	36.83	ICFG \rightarrow CUHK	33.46	56.30	31.56
	RSTP \rightarrow ICFG	32.30	49.69	20.54	RSTP \rightarrow CUHK	32.80	55.25	30.29
SEN (Liu et al., 2025a)	CUHK \rightarrow RSTP	55.50	77.85	45.29	CUHK \rightarrow ICFG	45.34	63.45	23.26
	ICFG \rightarrow RSTP	47.45	71.95	39.86	ICFG \rightarrow CUHK	37.88	60.48	35.07
	RSTP \rightarrow ICFG	36.23	53.31	22.32	RSTP \rightarrow CUHK	35.40	57.71	33.41
DRAGE (Ours)	CUHK \rightarrow RSTP	62.10	81.40	47.31	CUHK \rightarrow ICFG	52.45	69.61	28.46
	ICFG \rightarrow RSTP	54.40	76.00	42.55	ICFG \rightarrow CUHK	45.27	66.72	41.26
	RSTP \rightarrow ICFG	43.31	58.94	27.18	RSTP \rightarrow CUHK	42.58	63.35	38.84
w/o <i>SIMR</i>	CUHK \rightarrow RSTP	58.65	81.05	45.55	CUHK \rightarrow ICFG	49.03	67.04	25.90
w/o <i>SERD</i>	CUHK \rightarrow RSTP	50.25	72.40	38.44	CUHK \rightarrow ICFG	39.15	56.92	21.57

Table 5: Ablation studies on the CUHK-PEDES dataset.

No.	Method	R-1	R-5	R-10	mAP	mINP
0	DRAGE (Full)	78.58	91.68	95.19	70.68	53.89
<i>Single component ablation:</i>						
1	w/o \mathcal{L}_{cmt}	77.24	91.11	94.96	68.98	52.67
2	w/o SIMR	75.42	89.86	93.51	67.35	51.22
3	w/o SERD	72.18	87.94	92.36	64.82	48.91
4	w/o AGE	76.10	90.81	94.40	67.91	51.60
<i>Joint component ablation:</i>						
5	w/o AGE + \mathcal{L}_{cmt}	75.94	90.69	94.54	67.87	51.51
6	w/o SIMR + SERD	68.45	84.72	90.18	61.29	46.08
7	w/o AGE + SERD	70.82	86.53	91.44	63.15	47.36

4.5 Parametric Analysis

To investigate the impact of hyperparameter selection on model performance, we conduct a systematic sensitivity analysis on the CUHK-PEDES dataset, focusing on two critical hyperparameters: the margin parameter α and the temperature parameter τ . The results are illustrated in Figure 4 and Figure 5. First, regarding the margin parameter α , Figure 4 reveals a clear optimal range. Setting α to excessively large values leads to significant performance degradation, particularly across the Rank-1, Rank-5, and mAP metrics. Based on these empirical observations, we set $\alpha = 0.1$ for all subsequent experiments. Second, for the temperature parameter τ (see Figure 5), we observe that extremely low values result in training instability. Conversely, while increasing τ initially improves Rank-1 and Rank-5 accuracy, excessively high values cause a gradual decline in performance, especially in mAP. Consequently, we determine that $\tau = 0.015$ yields the optimal trade-off between training stability and discriminative capability, achieving the best results.

4.6 Robustness Analysis

To evaluate our method’s robustness, we conduct cross-domain tests by training the model on one dataset and testing it on another. As presented in Table 4, DRAGE demonstrates superior generalization capabilities, consistently outperforming prior state-of-the-art methods across all six settings. Notably, the improvements in Rank-1 accuracy reach up to 7.39%. This robust performance suggests that our strategy of addressing the fundamental semantic and modality gaps enables the model to learn more transferable representations. This conclusion is further validated by cross-domain ablation studies, where removing either the SERD or the SIMR module leads to a significant performance drop, confirming their crucial role in enhancing the model’s robustness.

5 Conclusion

In this paper, we propose DRAGE, a unified framework that systematically addresses the semantic and modality gaps in TBPS. DRAGE introduces two synergistic mechanisms: SERD mitigates the semantic gap by identifying noisy pairs to provide robust supervision, while SIMR bridges the modality gap by organizing internal embedding structures. Notably, DRAGE consistently achieves superior performance, demonstrating a state-of-the-art 78.58% Rank-1 accuracy on CUHK-PEDES and significantly enhancing cross-domain generalization. By improving both internal organization and supervision reliability, DRAGE offers a robust solution for real-world person search applications.

6 Limitations

While DRAGE achieves state-of-the-art performance in TBPS, several limitations merit acknowledgment. First, the effectiveness of SIMR depends on the quality of the pre-trained sentence encoder used to construct semantic distance matrices. We adopt all-mpnet-base-v2, whose capability to capture fine-grained nuances in pedestrian descriptions inherently constrains SIMR’s performance. Future work could explore domain-specific fine-tuning or alternative encoders for improved semantic structure modeling. Although DRAGE demonstrates strong cross-domain generalization, our evaluation is limited to standard TBPS benchmarks. These datasets may not fully capture real-world complexities such as extreme weather, heavy occlusion, or out-of-distribution language descriptions. Additionally, all experiments use English-language data; DRAGE’s performance in multilingual or low-resource settings remains unexplored.

References

Devansh Arpit, Stanisław Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S. Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, and Simon Lacoste-Julien. 2017. [A closer look at memorization in deep networks](#). In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 233–242. PMLR.

Yang Bai, Min Cao, Daming Gao, Ziqiang Cao, Chen Chen, Zhenfeng Fan, Liqiang Nie, and Min Zhang. 2023. Rasa: relation and sensitivity aware representation learning for text-based person search. In *IJCAI*, pages 555–563.

Min Cao, Yang Bai, Ziyin Zeng, Mang Ye, and Min Zhang. 2024. An empirical study of clip for text-based person search.

Tianlang Chen, Chenliang Xu, and Jiebo Luo. 2018. [Improving text-based person search by spatial matching and adaptive threshold](#). In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1879–1887. IEEE.

Changxing Ding, Kan Wang, Pengfei Wang, and Dacheng Tao. 2020. [Multi-task learning with coarse priors for robust part-aware person re-identification](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(3):1474–1488.

Zefeng Ding, Changxing Ding, Zhiyin Shao, and Dacheng Tao. 2021a. [Semantically self-aligned network for text-to-image part-aware person re-identification](#). *Preprint*, arXiv:2107.12666.

Zefeng Ding, Changxing Ding, Zhiyin Shao, and Dacheng Tao. 2021b. [Semantically self-aligned network for text-to-image part-aware person re-identification](#). *arXiv preprint arXiv:2107.12666*.

Ammarah Farooq, Muhammad Awais, Josef Kittler, and Syed Safwan Khalid. 2022. Axm-net: Implicit cross-modal feature alignment for person re-identification. In *AAAI*, pages 4477–4485.

Takuro Fujii and Shuhei Tarashima. 2023. Bilma: Bidirectional local-matching for text-based person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2786–2790.

Jing Ge, Guangyu Gao, and Zhen Liu. 2019. [Visual-textual association with hardest and semi-hard negative pairs mining for person search](#). In *arXiv preprint arXiv:1912.03083*.

Jianyuan Guo, Yuhui Yuan, Lang Huang, Chao Zhang, Jin-Ge Yao, and Kai Han. 2019. [Beyond human parts: Dual part-aligned representations for person re-identification](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3642–3651. IEEE.

Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor W. Tsang, and Masashi Sugiyama. 2018. [Co-teaching: Robust training of deep neural networks with extremely noisy labels](#). *Neural Information Processing Systems, Neural Information Processing Systems*.

D. Jiang and M. Ye. 2023a. Cross-modal implicit relation reasoning and aligning for text-to-image person retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2787–2797.

Ding Jiang and Mang Ye. 2023b. Cross-modal implicit relation reasoning and aligning for text-to-image person retrieval. In *CVPR*, pages 2787–2797.

Jiayu Jiang, Changxing Ding, Wentao Tan, Junhong Wang, Jin Tao, and Xiangmin Xu. 2025. Modeling thousands of human annotators for generalizable text-to-image person re-identification. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 9220–9230.

Jie Lei, Xinlei Chen, Ning Zhang, Mengjiao Wang, Mohit Bansal, Tamara L. Berg, and Licheng Yu. 2022. [Loopitr: Combining dual and cross encoder architectures for image-text retrieval](#). *Preprint*, arXiv:2203.05465.

Haiwen Li, Delong Liu, Fei Su, and Zhicheng Zhao. 2025. Object-centric discriminative learning for text-based person retrieval. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

594	Shenshen Li, Xing Xu, Yang Yang, Fumin Shen, Yijun Mo, Yujie Li, and Heng Tao Shen. 2023. Dcel: Deep cross-modal evidential learning for text-based person retrieval. In <i>Proceedings of the 31st ACM International Conference on Multimedia</i> , pages 6292–6300.	647	Yang Qin, Yingke Chen, Dezhong Peng, Xi Peng, Joey Tianyi Zhou, and Peng Hu. 2024. Noisy-correspondence learning for text-to-image person re-identification. In <i>IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	648
595		649		650
596		651		
597				
598				
599	Shuang Li, Tong Xiao, Hongsheng Li, Bolei Zhou, Dayu Yue, and Xiaogang Wang. 2017. Person search with natural language description . In <i>2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 5187–5196.	652	Yang Qin, Dezhong Peng, Xi Peng, Xu Wang, and Peng Hu. 2022. Deep evidential learning with noisy correspondence for cross-modal retrieval. In <i>Proceedings of the 30th ACM International Conference on Multimedia</i> , pages 4948–4956.	653
600		654		655
601		656		
602				
603				
604	Delong Liu, Haiwen Li, Zhicheng Zhao, and Yuan Dong. 2025a. Text-guided image restoration and semantic enhancement for text-to-image person retrieval. <i>Neural Networks</i> , 184:107028.	657	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 <i>ICML</i> , pages 8748–8763. PMLR.	658
605		659		660
606		661		662
607		663		
608	Yating Liu, Zimo Liu, Xiangyuan Lan, Wenming Yang, Yaowei Li, and Qingmin Liao. 2025b. Dm-adapter: Domain-aware mixture-of-adapters for text-based person retrieval. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 5703–5711.	664	Zhiyin Shao, Xinyu Zhang, Changxing Ding, Jian Wang, and Jingdong Wang. 2023. Unified pre-training with pseudo texts for text-to-image person re-identification. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 11174–11184.	665
609		666		667
610		668		
611				
612				
613				
614	Yu Liu, Guihe Qin, Haipeng Chen, Zhiyong Cheng, and Xun Yang. 2024. Causality-inspired invariant representation learning for text-based person retrieval. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 14052–14060.	669	Zhiyin Shao, Xinyu Zhang, Meng Fang, Zhifeng Lin, Jian Wang, and Changxing Ding. 2022. Learning granularity-unified representations for text-to-image person re-identification. In <i>ACM MM</i> , pages 5566–5574.	670
615		671		672
616		673		
617				
618				
619	Tianle Lv, Shuang Li, Jiaxu Leng, and Xinbo Gao. 2024. Mgrl: Mutual-guidance representation learning for text-to-image person retrieval. In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 2895–2899. IEEE.	674	Fei Shen, Xiangbo Shu, Xiaoyu Du, and Jinhui Tang. 2023. Pedestrian-specific bipartite-aware similarity learning for text-based person retrieval. In <i>Proceedings of the 31st ACM International Conference on Multimedia</i> , pages 8922–8931.	675
620		676		677
621		678		
622				
623				
624				
625	Yiwei Ma, Xiaoshuai Sun, Jiayi Ji, Guannan Jiang, Weilin Zhuang, and Rongrong Ji. 2023. Beat: Bi-directional one-to-many embedding alignment for text-based person retrieval. In <i>Proceedings of the 31st ACM International Conference on Multimedia</i> , pages 4157–4168.	679	Xiujun Shu, Wei Wen, Haoqian Wu, Keyu Chen, Yiran Song, Ruizhi Qiao, Bo Ren, and Xiao Wang. 2022. See finer, see more: Implicit modality alignment for text-based person retrieval. In <i>European Conference on Computer Vision</i> , pages 624–641. Springer.	680
626		681		682
627		683		
628				
629				
630				
631	Antoine Miech, Jean-Baptiste Alayrac, Ivan Laptev, Josef Sivic, and Andrew Zisserman. 2021. Thinking fast and slow: Efficient text-to-visual retrieval with transformers . In <i>2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 9821–9831.	684	Zifan Song, Guosheng Hu, and Cairong Zhao. 2024. Diverse person: Customize your own dataset for text-based person search. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 4943–4951.	685
632		686		687
633		688		
634				
635				
636				
637	Jicheol Park, Dongwon Kim, Boseung Jeong, and Suha Kwak. 2024. Plot: Text-based person search with part slot attention for corresponding part discovery. In <i>ECCV</i> , pages 474–490. Springer.	689	Siqi Sun, Yen-Chun Chen, Linjie Li, Shuohang Wang, Yuwei Fang, and Jingjing Liu. 2021. Lightningdot: Pre-training visual-semantic embeddings for real-time image-text retrieval . <i>Preprint</i> , arXiv:2103.08784.	690
638		691		692
639		693		
640				
641	Yang Qin, Chao Chen, Zhihang Fu, Dezhong Peng, Xi Peng, and Peng Hu. 2025. Human-centered interactive learning via mllms for text-to-image person re-identification. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 14390–14399.	694	Wentan Tan, Changxing Ding, Jiayu Jiang, Fei Wang, Yibing Zhan, and Dapeng Tao. 2024. Harnessing the power of mllms for transferable text-to-image person reid. In <i>CVPR</i> , pages 17127–17137.	695
642		696		697
643				
644				
645				
646				
			Zhe Wang, Zhiyuan Fang, Jun Wang, and Yezhou Yang. 2020. Vitaa: Visual-textual attributes alignment in person search by natural language. In <i>Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XII 16</i> , pages 402–420. Springer.	698
				699
				700
				701
				702
				703

704	Zijie Wang, Aichun Zhu, Jingyi Xue, Xili Wan, Chao	Zhiwei Zhao, Bin Liu, Yan Lu, Qi Chu, and Nenghai	759
705	Liu, Tian Wang, and Yifeng Li. 2022. Look before	Yu. 2024. Unifying multi-modal uncertainty model-	760
706	you leap: Improving text-based person retrieval by	ing and semantic alignment for text-to-image person	761
707	learning a consistent cross-modal common manifold.	re-identification. In <i>Proceedings of the AAAI Con-</i>	762
708	In <i>ACM MM</i> , pages 1984–1992.	<i>ference on Artificial Intelligence</i> , volume 38, pages	763
		7534–7542.	764
709	Yushuang Wu, Zizheng Yan, Xiaoguang Han, Guan-	Zhedong Zheng, Liang Zheng, Michael Garrett,	765
710	bin Li, Changqing Zou, and Shuguang Cui. 2021.	Yi Yang, Mingliang Xu, and Yi-Dong Shen. 2020.	766
711	Lapscore: language-guided person search via color	Dual-path convolutional image-text embeddings with	767
712	reasoning. In <i>Proceedings of the IEEE/CVF Interna-</i>	instance loss. <i>ACM Transactions on Multime-</i>	768
713	<i>tional Conference on Computer Vision</i> , pages 1624–	<i>dia Computing, Communications, and Applications</i> ,	769
714	1633.	16(2):1–23.	770
715	Zequn Xie, Haoming Ji, and Lingwei Meng. 2025a.	Aichun Zhu, Zijie Wang, Yifeng Li, Xili Wan, Jing Jin,	771
716	Dynamic uncertainty learning with noisy correspon-	Tian Wang, Fangqiang Hu, and Gang Hua. 2021a.	772
717	dence for text-based person search. <i>arXiv preprint</i>	<i>Dssl: Deep surroundings-person separation learning</i>	773
718	<i>arXiv:2505.06566</i> .	for text-based person retrieval. In <i>Proceedings of the</i>	774
		<i>29th ACM International Conference on Multimedia</i> .	775
719	Zequn Xie, Chuxin Wang, Sihang Cai, Yeqiang Wang,	Aichun Zhu, Zijie Wang, Yifeng Li, Xili Wan, Jing Jin,	776
720	Shulei Wang, and Tao Jin. 2025b. Chat-driven	Tian Wang, Fangqiang Hu, and Gang Hua. 2021b.	777
721	text generation and interaction for person retrieval .	<i>Dssl: Deep surroundings-person separation learning</i>	778
722	<i>Preprint</i> , arXiv:2509.12662.	for text-based person retrieval. In <i>Proceedings of the</i>	779
		<i>29th ACM International Conference on Multimedia</i> ,	780
723	Shuanglin Yan, Neng Dong, Jun Liu, Liyan Zhang, and	pages 209–217.	781
724	Jinhui Tang. 2023. Learning comprehensive repre-	Jialong Zuo, Hanyu Zhou, Ying Nie, Feng Zhang,	782
725	sentations with richer self for text-to-image person	Tianyu Guo, Nong Sang, Yunhe Wang, and Changxin	783
726	re-identification. In <i>Proceedings of the 31st ACM</i>	Gao. 2024. Ufinebench: Towards text-based person	784
727	<i>international conference on multimedia</i> , pages 6202–	retrieval with ultra-fine granularity. In <i>Proceedings</i>	785
728	6211.	<i>of the IEEE/CVF Conference on Computer Vision</i>	786
		<i>and Pattern Recognition</i> , pages 22010–22019.	787
729	Shuanglin Yan, Neng Dong, Liyan Zhang, and Jin-		
730	hui Tang. 2022a. Clip-driven fine-grained text-		
731	image person re-identification. <i>arXiv preprint</i>		
732	<i>arXiv:2210.10276</i> .		
733	Shuanglin Yan, Hao Tang, Liyan Zhang, and Jinhui		
734	Tang. 2022b. Image-specific information suppression		
735	and implicit local alignment for text-based person		
736	search. <i>arXiv preprint arXiv:2208.14365</i> .		
737	Shuyu Yang, Yinan Zhou, Yaxiong Wang, Yujiao Wu,		
738	Li Zhu, and Zhedong Zheng. 2023. Towards uni-		
739	fied text-based person retrieval: A large-scale multi-		
740	attribute and language search benchmark. In <i>Pro-</i>		
741	<i>ceedings of the 2023 ACM on Multimedia Confer-</i>		
742	<i>ence</i> .		
743	Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling		
744	Shao, and Steven C. H. Hoi. 2022. Deep learning		
745	for person re-identification: A survey and outlook .		
746	<i>IEEE Transactions on Pattern Analysis and Machine</i>		
747	<i>Intelligence</i> , page 2872–2893.		
748	Ji Zhang, Jingkuan Song, Lianli Gao, Ye Liu, and		
749	Heng Tao Shen. 2022. Progressive meta-learning		
750	with curriculum . <i>IEEE Transactions on Circuits and</i>		
751	<i>Systems for Video Technology</i> , 32(9):5916–5930.		
752	Ying Zhang and Huchuan Lu. 2018a. Deep cross-		
753	modal projection learning for image-text matching.		
754	In <i>ECCV</i> .		
755	Ying Zhang and Huchuan Lu. 2018b. Deep cross-modal		
756	projection learning for image-text matching. In <i>Pro-</i>		
757	<i>ceedings of the European conference on computer</i>		
758	<i>vision (ECCV)</i> , pages 686–701.		

A Supplementary Material

In this supplementary material, we provide additional information for DRAGE.

A.1 Datasets

To comprehensively evaluate the performance of our proposed method under both clean and noisy supervision, we conduct experiments on three widely adopted benchmarks for text-to-image person retrieval: CUHK-PEDES (Zhang et al., 2022), ICFG-PEDES (Han et al., 2018), and RSTPReid (Zhu et al., 2021b). Each dataset presents unique characteristics in terms of scale, annotation style, and visual diversity.

CUHK-PEDES. CUHK-PEDES (Zhang et al., 2022) is the first and most well-established benchmark for the text-based person re-identification (ReID) task. It comprises a total of 40,206 pedestrian images corresponding to 13,003 unique identities. Each image is annotated with two independent natural language descriptions, collected via crowdsourcing, providing rich semantic diversity. Following the standard evaluation protocol, the dataset is divided into three subsets: a training set with 34,054 images from 11,003 identities, a validation set containing 3,078 images from 1,000 identities, and a testing set with 3,074 images from another 1,000 identities. The average caption length is approximately 23 words, encompassing a variety of appearance-related cues such as clothing, accessories, and actions.

ICFG-PEDES. ICFG-PEDES (Han et al., 2018) is a relatively newer but larger-scale dataset designed to reflect more complex and fine-grained textual descriptions. It contains 54,522 images corresponding to 4,102 person identities. Unlike CUHK-PEDES, each image in ICFG-PEDES is annotated with a single sentence, which on average contains around 37 words—making the descriptions more elaborate and detailed. The official data split consists of 34,674 image-text pairs from 3,102 identities for training, and 19,848 image-text pairs from the remaining 1,000 identities for testing. This dataset presents greater challenges due to longer and potentially more ambiguous captions.

RSTPReid. RSTPReid (Zhu et al., 2021b) is a large-scale and challenging benchmark constructed under real-world surveillance scenarios. It contains 20,505 images from 4,104 identities captured by 15 different cameras, introducing substantial variations in pose, lighting, and background. Each

identity is associated with five images taken from distinct camera views. Importantly, each image is paired with two distinct textual descriptions, which significantly increases the textual diversity. According to the official splitting protocol, the training set includes image-text pairs from 3,701 identities. Both the validation and testing sets are composed of 200 identities each, ensuring consistent and reliable benchmarking. Compared to CUHK-PEDES and ICFG-PEDES, RSTPReid reflects more realistic retrieval conditions with strong cross-view challenges.

B Training Process

Algorithm 1 The training process of our model

- 1: **Input:** The training data \mathcal{P} with N image-text pairs, maximal epoch N_e , the model $\mathcal{M}(\Theta)$, and the hyper-parameters $\mathcal{R}, \alpha, \tau$
 - 2: Initialize the backbones with the weights of the pretrained model CLIP;
 - 3: Initialize the Pre-trained Sentence Encoder: `sbert = SBERT(sentence-transformers/all-mpnet-base-v2)`;
 - 4: **for** $e = 1, 2, \dots, N_e$ **do**
 - 5: Calculate the per-sample loss $\ell(\mathcal{M}, \mathcal{P})$ with Equation (3);
 - 6: Divide the training data with the predictions of GFE and AGE using BMM;
 - 7: SERD: Deriving and Recalibrating the Labels $\{\ell_i\}_{i=1}^N$;
 - 8: **for** each x in mini-batches $\{x_m\}_{m=1}^M$ of x **do**
 - 9: Extract the GFE & AGE features of x ;
 - 10: Compute the similarities between K image-text pairs in x with above features;
 - 11: Calculate the \mathcal{L}_{v2v} for vision-to-vision and \mathcal{L}_{t2t} for text-to-text using SIMR;
 - 12: Calculate \mathcal{L}_{total} with Equation (17);
 - 13: $\Theta = \text{Optimizer}(\Theta, \mathcal{L}_{total})$;
 - 14: **end for**
 - 15: **end for**
 - 16: **Output:** The optimized parameters $\hat{\Theta}$;
-

C Visualization of Modality Gap Reduction

To visualize the modality gap reduction, we project the learned image and text embeddings from CUHK-PEDES into a 2D space using Principal Component Analysis. Figure 6 illustrates the pro-

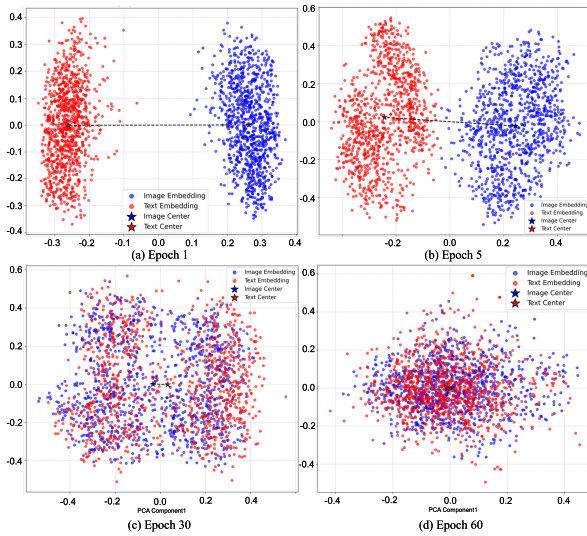


Figure 6: Visualization of modality gap reduction between image and text embeddings in CUHK-PEDES dataset .

857 gressive alignment throughout the training process.
 858 At epoch 1, the embeddings form distinct clusters.
 859 These clusters merge by epoch 30 and become in-
 860 distinguishable by epoch 60, demonstrating the
 861 effective bridging of the modality gap.

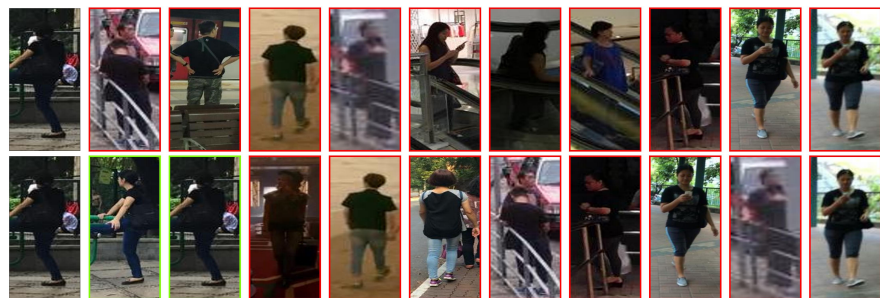
862 D Qualitative Results

863 To illustrate the advantages of our DRAGE, some
 864 retrieved examples for TBPS are presented in Fig-
 865 ure 7.

A woman with long hair is holding a small wood table. She wears a white dress.



The woman has short dark hair, a black short sleeved t-shirt and denim pants. She is leaning on a green rail with one leg raised.



The lady wears a black and white shirt a pink towel wrapped around her neck blue jean pants with pink flipflops carries she carries a black shoulder bag holding a black book in her hand.



The woman is wearing a tan long coat, a scarf, gray jeans, and brown boots.



Figure 7: Qualitative comparison of the top-10 retrieved results on the CUHK-PEDES dataset. For each text query, we compare the baseline IRRA (first row) with our DRAGE (second row). The **matched** and **mismatched** person images are highlighted with **green** and **red** rectangles, respectively. The first column of each row represents the ground truth.