

# LOCoRE: Image Re-ranking with Long-Context Sequence Modeling

Zilin Xiao<sup>1</sup>, Pavel Suma<sup>2</sup>, Ayush Sachdeva<sup>1</sup>, Hao-Jen Wang<sup>1</sup>,  
Giorgos Kordopatis-Zilos<sup>2</sup>, Giorgos Tolias<sup>2</sup>, Vicente Ordonez<sup>1</sup>

<sup>1</sup>Rice University    <sup>2</sup>VRG, FEE, Czech Technical University in Prague

## Abstract

We introduce LOCoRE, *Long-Context Re-ranker*, a model that takes as input local descriptors corresponding to an image query and a list of gallery images and outputs similarity scores between the query and each gallery image. This model is used for image retrieval, where typically a first ranking is performed with an efficient similarity measure, and then a shortlist of top-ranked images is re-ranked based on a more fine-grained similarity model. Compared to existing methods that perform pair-wise similarity estimation with local descriptors or list-wise re-ranking with global descriptors, LOCoRE is the first method to perform list-wise re-ranking with local descriptors. To achieve this, we leverage efficient long-context sequence models to effectively capture the dependencies between query and gallery images at the local-descriptor level. During testing, we process long shortlists with a sliding window strategy that is tailored to overcome the context size limitations of sequence models. Our approach achieves superior performance compared with other re-rankers on established image retrieval benchmarks of landmarks ( $\mathcal{ROxf}$  and  $\mathcal{RPar}$ ), products (SOP), fashion items (In-Shop), and bird species (CUB-200) while having comparable latency to the pair-wise local descriptor re-rankers.

## 1. Introduction

Instance-level image retrieval is an important problem in computer vision with many applications. It is usually cast as a problem of metric learning where a model is trained to optimize a distance metric while comparing pairs of examples. Retrieval in large image collections (gallery) is typically performed in two stages. First, images are mapped to a compact global image descriptor that is used for efficient retrieval of a shortlist of candidate images from the gallery. Subsequently, a more powerful but computationally demanding model is used to re-rank the shortlisted images. Local features and their descriptors are typically employed at this stage to enable a more detailed image-to-image comparison and to

provide benefits such as robustness to background clutter and partial visibility due to occlusions. During re-ranking, the common practice is to estimate the improved similarity in a pair-wise manner by comparing the two local descriptor sets of the query and each image in the shortlist. Geometric verification is a common re-ranking method where point correspondences are processed with a RANSAC-like process, and a high number of inliers is expected for images depicting the same object under different viewpoints [9, 45]. Recent methods train a model while optimizing the image-to-image similarity using either dense [28] or sparse [65, 66, 76, 78] local representation as its input. Transformers are becoming a dominant component of such models [65, 66, 76, 78].

Most existing image re-ranking methods perform in a pair-wise way, *i.e.*, the query is separately compared to each of the gallery images. This strategy does not capture interactions between the gallery images, such as objects or object parts that tend to appear more than once. In this work, we introduce LOCoRE, a model that performs image re-ranking by jointly processing the input query and an entire shortlist of gallery images at the local descriptor level. To take into account interactions within the shortlist, we rely on sequence models, in particular transformers, that capture contextual relationships. Processing within a long context comes with computational challenges. We overcome those in two ways. First, we leverage a model design from existing architectures in Natural Language Processing (NLP) that efficiently extends the maximum context length of a standard transformer. Second, we propose a sliding window strategy that reuses LOCoRE multiple times during inference. Figure 1 shows an overview of how the proposed list-wise re-ranker compares to pair-wise re-rankers.

LOCoRE takes inspiration from NLP tasks such as sequence tagging [21, 43, 53] and extractive question-answering [15, 52, 58, 75]. Modern solutions to these problems often involve a sequence model that predicts token-level scores to extract task-specific text spans [6, 20]. In contrast to the current design choice of image re-ranking models, we adopt the common strategy in NLP and optimize all output tokens, *i.e.* many per gallery image, in-

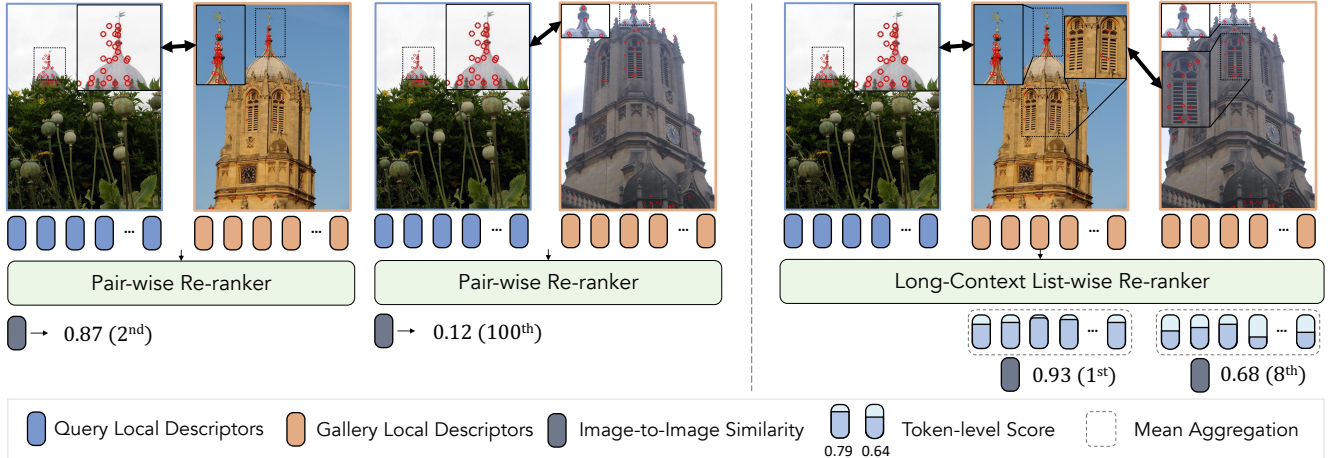


Figure 1. **Overview of pair-wise vs our proposed list-wise re-ranking.** Red circles denote the locations of input local descriptors. *Left:* The pair-wise re-ranker gets a high score for a positive image since it clearly depicts the same structure at the top of the tower as the query, while a different positive image gets a low score because the top of the tower is not as clearly visible. *Right:* Our long-context re-ranker can output a high score for both positive image results since it can exploit the transitive relationship between these images as the two gallery images also share common local descriptors.

stead of a single global token. At inference time, the aggregated score across tokens associated with a gallery image is used as the similarity value. By casting image re-ranking as a span extraction task and modeling long context relations across shortlisted images, we demonstrate superior results on established image retrieval benchmarks. Specifically, LOCORE is the state-of-the-art re-ranker when evaluated with the same global retrieval method across the CUB-200 [71], Stanford Online Products (SOP) [63] and In-shop [32] benchmarks. Moreover, our models trained on Google Landmarks v2 (GLDv2) [73] are achieving leading results with other re-rankers in relative performance gains on the Revisited Oxford and Paris datasets [44, 46, 48]. Code is available at <https://github.com/MrZilinXiao/LongContextReranker>.

## 2. Related Work

**Exhaustive search.** Early image retrieval approaches extract hand-crafted local descriptors that encode visual keypoints in images [7, 34]. The extracted descriptors are compactly quantized into bag-of-words (BoW) representations [12, 37, 62] to enable exhaustive search in large databases. Aggregating local descriptors into a single global image representation [25, 56] allows using simple metrics, such as Euclidean distance, to compute the image-to-image similarity for efficient ranking. Improvements over the vanilla BoW scheme achieve better matching approximation by searching directly with local descriptors [23, 24, 69].

With the rise of deep learning, approaches that leverage neural networks to extract descriptors [2, 9, 28, 38, 77] have dominated over the hand-crafted approaches. The main focus is the optimization of the learning process via dif-

ferent loss functions [10, 13, 41, 57] and network architectures [9, 28, 39]. In early work [55], dense image feature maps are extracted from Convolutional Neural Networks (CNN) as local descriptor sets and compared with chamfer similarity. However, aggregating feature maps into a global descriptor remains the dominant approach [4, 68, 74]. The aggregation can also be learned in an end-to-end fashion via learnable pooling layers [47, 50, 72]. Other works [9, 64, 74] demonstrate advantages in training separate global and local descriptor branches, which are ultimately fused into a global descriptor.

**Re-ranking.** A typical approach for refining the exhaustive search is to apply a re-ranking step that involves more precise similarity estimation. To ensure feasible search times, exhaustive re-ranking is applicable only with global image descriptors. Query expansion derives a new query descriptor by aggregating it with the descriptors of its nearest neighbors from the initial retrieval [11]. Extensions of the base scheme include weighted aggregation based on image ranks [1, 50, 59] or aggregation through a learnable model [17]. An alternative approach is to re-estimate pairwise image similarities by leveraging local descriptors. Due to its high demand in terms of computational complexity, such approaches are only applied to a short list of the top-ranked images provided by the exhaustive search. Geometric verification (GV) of the descriptor correspondences is a long-standing solution, typically involving RANSAC [16]. This approach is applicable to hand-crafted [3, 44] as well as learned features [9, 38, 61]. Recently, the geometric approach is surpassed by deep networks optimized to internally match the local descriptors of an image pair and output a single similarity score. Correlation Verification

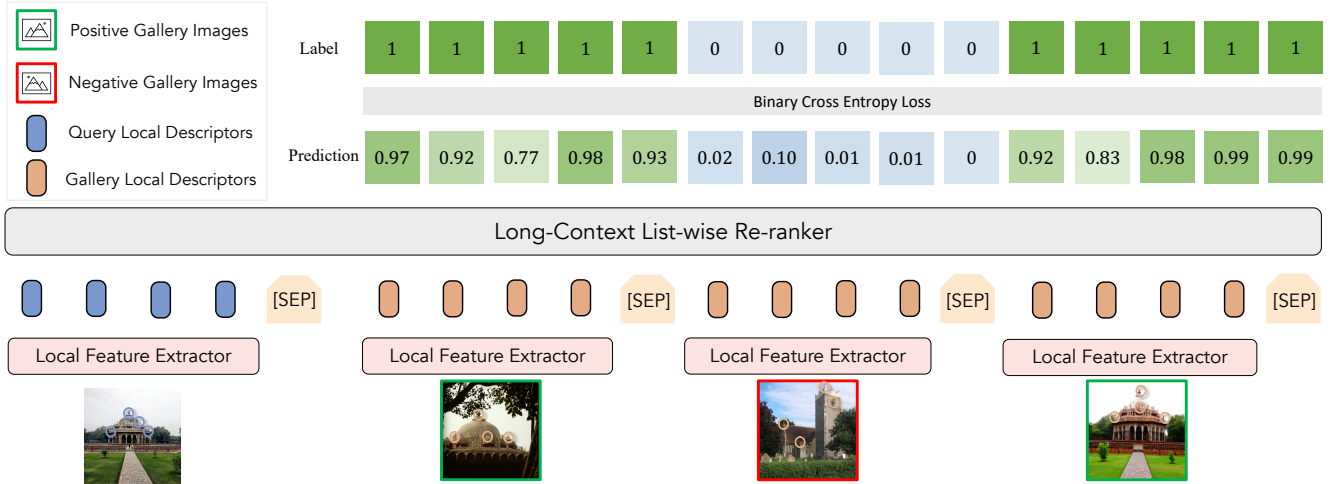


Figure 2. **Overview of training and inference under LOCORE when re-ranking three candidate gallery images.** In practice, we re-rank in one inference step up to 100 gallery images. At training time, the model is trained to optimize a binary cross-entropy loss on each gallery image token. At inference time, the token scores of each gallery image get aggregated to facilitate a re-ranked gallery image list.

Networks (CVNet) [28] incorporates 4D convolutions on top of joint correlations between dense local descriptor sets. Transformer-based architectures, such as Reranking Transformer (RRT) [66] and Asymmetric and Memory-Efficient Similarity (AMES) [65], are designed to be fed with two sets of sparse local descriptors, which are processed as a sequence of tokens to estimate similarity. A closely related work relying on transformers uses global descriptors to re-rank images in a list-wise manner [40] instead of processing each query and gallery image pair separately. It extracts a single affinity image descriptor for each image. Our proposed approach builds on list-wise re-ranking but utilizes local descriptors. This is fundamentally different since each image is represented by multiple descriptors. To this end, we repurpose and fine-tune a transformer network specifically designed for long-context tasks.

### 3. Method

In this section, we describe the model architecture, the training procedure and test-time strategies of our proposed LOCORE image reranker.

#### 3.1. Problem Formulation

Given a query image  $I_q$  and an ordered list of  $K$  gallery images  $I_{g,i}, i \in \{1, \dots, K\}$  obtained by global similarity search, the purpose of image re-ranking is to produce another refined list of gallery images that are reordered based on improved similarity measures to the query image  $I_q$ .

Let  $x_q = \{\mathbf{x}_{q,j} \in \mathbb{R}^d\}_{j=1}^L$  denote  $L$  local descriptors of dimension  $d$  extracted for the query image using a visual backbone and  $x_{g,i} = \{\mathbf{x}_{g,i,j} \in \mathbb{R}^d\}_{j=1}^L$  denote the local descriptors for the  $i$ -th gallery image. These  $L$  local descriptors

are either extracted from a specific layer of a visual backbone model, *e.g.* local descriptors from the last convolutional layer of ResNet [19], or are the top- $L$  multi-scale descriptors weighted by local attention scores, *e.g.* as obtained in DELG [9].

**Pair-wise re-ranker.** A typical neural re-ranker  $f_\theta$  computes a pair-wise confidence score  $S$  for each pair of images  $(I_q, I_{g,i})$  based on their local descriptors:

$$S(I_q, I_{g,i}) = f_\theta(x_q, x_{g,i}),$$

where  $f_\theta$  learns a binary classification objective during training to separate matching (positive) and non-matching (negative) image pairs. Then, the refined gallery list is constructed by sorting the confidence scores of the  $K$  gallery images in descending order. The most common architecture is a transformer [65, 66], and the parameter set  $\theta$  corresponds to the parameters of the transformer layers and of the binary classifier on top of a single output token.

**List-wise re-ranker.** Instead of estimating the similarity to the query separately per gallery image, we propose a model architecture to perform this jointly across all  $K$  images to benefit from interactions across gallery images. The deep network takes the local descriptors of the query and  $K$  gallery images as input to estimate all  $K$  query-to-gallery image similarities via

$$\mathbf{y} = f_\theta(x_q, x_{g,1}, \dots, x_{g,K}),$$

where  $\mathbf{y} \in [0, 1]^K$ . The similarities to perform the sorting and re-ranking are obtained via  $S(I_q, I_{g,i}) = \mathbf{y}(i)$ , with  $\mathbf{y}(i)$  being the  $i$ -th element of vector  $\mathbf{y}$ .

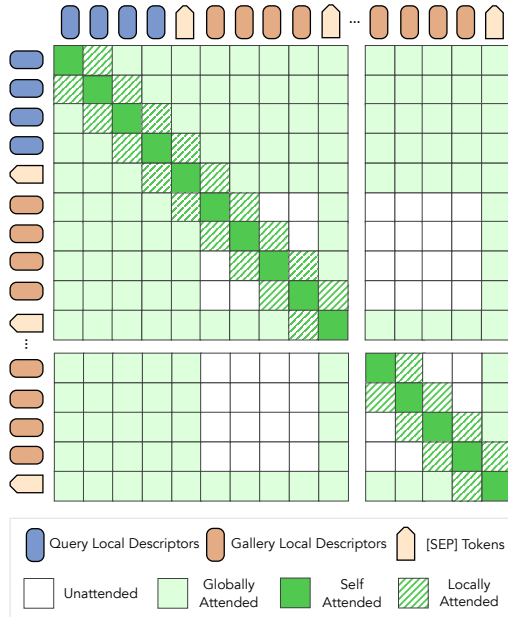


Figure 3. **Attention pattern of LOCORE** when the number of descriptors for each image is  $L = 4$  and the local window size is  $W = 1$ . In our experiments we actually use a larger number of local feature descriptors and window size.

### 3.2. Architecture

**Long context.** We treat the set of input local descriptors as a sequence of vectors by concatenating them. In principle, any sequence model can serve as a backbone for our method, but to rerank up to a hundred gallery images at once we need a sequence model with a very large context window. Similar to pair-wise re-rankers, we rely on transformers, but dense attention matrices restrict the model to very small values of  $K$  and  $L$ . For this reason, we use a pre-trained Longformer [8] as its computational complexity grows linearly with sequence length due to its sliding attention window mechanism. It can capture long-range dependencies, as the positive gallery image can appear in any part of the sequence. Figure 2 presents an overview of our method.

**Input sequence.** The vector sequence is given in a matrix format by  $x = [x_q, \mathbf{x}_s, x_{g,1}, \mathbf{x}_s, \dots, x_{g,K}, \mathbf{x}_s]$ , with  $\mathbf{x}_s \in \mathbb{R}^d$  being a learnable separation token, [SEP], that is meant to act as a global representation of the image that is preceding it. In total there are  $M = (L + 1)(K + 1)$  vectors. We use  $M$  learnable positional encodings to indicate the position in the sequence of each vector. Additionally, we use  $K + 1$  learnable positional encodings to indicate the image from which the vectors come. Both these positional encodings are added to the corresponding columns of  $x$  and then fed to the sequence model, with each column representing a token.

**Query global attention.** Longformer uses a local sliding-window attention layer instead of a full self-attention layer,

which allows it to handle longer sequences. It additionally defines some tokens to perform global attention in a symmetric way, *i.e.* they attend to all other tokens and all other tokens attend to them. To compensate for potential issues handling long-range dependencies, we define all tokens associated with the query image and all [SEP] tokens to perform global attention. In this way, all descriptor tokens attend to other tokens within the same image and of nearby images with a local window size of  $W$ , while the global attention tokens ensure long-range interactions. We illustrate the attention pattern in Figure 3. Note that using global attention for the limited set of query and separation tokens brings only a marginal computation overhead. Additionally, such an attention mechanism is an important inductive bias, without which the model demonstrates trivial re-ranking performance, as indicated in our later experiments.

**Token classifier.** All contextualized output representations are fed to a binary classifier, which aims to predict whether the corresponding input token comes from a positive or a negative image. This resembles the design in a variety of natural language processing tasks that require token-level classification, such as named entity recognition [5, 6], parts-of-speech tagging [36], and extractive question answering [52, 58, 75].

### 3.3. Training

During training, all  $(L + 1) \times K$  classifier outputs are fed to a binary cross-entropy loss according to the ground-truth labels of query-gallery-image pairs. Regarding training batch sampling, we pick a query image from the training set, use global similarity search, and pick the top- $K$  training images as gallery images to form a list-wise training sample. However, directly training LOCORE on list-wise re-ranking supervision comes with inherent issues. Global similarity search tends to have positive gallery images ranked at the top positions of the retrieved gallery list. Feeding the local descriptors for the gallery images in the order given by the global retrieval model introduces a rather counterproductive positional bias. Models can use this bias as a shortcut during training since the model will correlate the positive matching predictions with top positions rather than relying on learning to match local descriptors. Therefore, we shuffle the gallery list in each training step to ensure that each position has an equal probability of being assigned a positive gallery image. We refer to this strategy as *gallery shuffled training*.

### 3.4. Inference

**Image-to-image similarity.** During test time, we compute a single score per gallery image by aggregating all  $L + 1$  classifier outputs that are associated with the specific image. The aggregator function can be defined as one of the following: (i) *the average token score*, (ii) *the first token score* or (iii) *[SEP] token score* of the token span associated with the corresponding gallery images.



Figure 4. **Illustration of sliding window re-ranking** for  $N = 8$  gallery images with a list-wise re-ranker that can re-rank  $K = 4$  images each forward pass. Blue blocks represent the re-ranking window of the current forward pass, and it slides to the next window with a stride of  $S = 2$  images. Brown blocks indicate the re-ranking for the current forward pass is completed.

**Sliding window re-ranking.** Despite using a long-context sequence model, we might still be limited in how many gallery images we can rerank at once. In order to extend the number of gallery images at inference time, we can rerank separate overlapping groups of gallery images to obtain reranking scores for a larger number of gallery images. Figure 4 illustrates our proposed test-time sliding window re-ranking strategy. Let  $N$  be the total number of gallery images pending for re-ranking and  $K$  be the maximum number of images that the reranker can process at once. We first use the model to rank the  $(N - K)$ -th to  $N$ -th gallery images and slide the re-ranking window of size  $K$  forward with a stride step size of  $S$ . This process continues to the start of the shortlist to maximally ensure each gallery image gets to its ranking as accurately as possible. Our experiments demonstrate consistent improvements under this strategy.

## 4. Experiments

We first introduce the experimental setup in Section 4.1 which covers datasets, metrics, and details of training and evaluation. We report the main results and ablation results in Section 4.2 and Section 4.3.

### 4.1. Experimental Setup

**Datasets and Metrics.** To demonstrate the effectiveness of our method, we experiment with large-scale instance-level recognition datasets. We train on the clean Google Landmarks v2 (GLDv2) dataset [73] and evaluate on the Revisited Oxford ( $\mathcal{ROxf}$ ) and Paris ( $\mathcal{RPar}$ ) datasets and their 1M distractor variants ( $\mathcal{ROxf}+1M$ ,  $\mathcal{RPar}+1M$ ) [44, 46, 48]. Following the literature in the relevant field [65, 66], we report mean average precision (mAP) on the Medium and Hard settings. In addition, following established image retrieval benchmarks, we evaluate on the following datasets:

CUB-200 [71], Stanford Online Products (SOP) [63] and In-shop [32]. We choose Recall@ $k$  ( $R@k$ ) and mean average precision at  $R$  ( $mAP@R$ ) as evaluation metrics for these datasets, which are the standard metrics [32, 63, 66].

**Training and Evaluation.** Unless stated otherwise, all experiments follow the configuration of  $L = 50$  and  $K = 100$ , *i.e.* re-ranking is conducted with 1 query image and 100 gallery images, each of them providing 50 descriptors. In inference, we use the [SEP] token score as the default choice to get individual image similarity scores. For each dataset, we follow the official training and validation splits if available. Otherwise, we consider the first half of the training dataset as the training split and the remaining as the validation split. Hyper-parameters are tuned based on the validation split. We utilize a series of models ranging from 19.4M (tiny) to 111.8M (base) parameters. The implementation details employed for each dataset are reported in the supplementary material.

**Competing Methods.** To verify the versatility of our method on the landmarks datasets, we compare LOCoRE with state-of-the-art re-ranker approaches. We evaluate with local descriptors extracted from two different backbones, *i.e.* DELG [9] and DINOv2 [39]. For the DELG experiments, we follow the original process to extract global and local descriptors at 3 scales ( $\{1/\sqrt{2}, 1, \sqrt{2}\}$ ). The top- $L$  local descriptors are selected for each image based on the provided DELG weights. We compare with the RRT [66] and CVNet Reranker [28] which report results on top of the same descriptors. We further present the results of GV [9, 45] as reported in the original DELG paper.

In addition, we experiment with local descriptors from DINOv2-ViT-B/14 [39]. We follow the descriptor extraction of AMES [65] to extract descriptors and to be directly comparable to this re-ranking method. We also adopt their global-local ensemble scheme and combine the local similarities from LOCoRE with the global similarities. In these experiments, we extract global descriptors using the Super-Global [59] ResNet-101 backbone. Our ensemble parameters are tuned on the public test split of the GLDv2 dataset in the same manner as AMES. We compare LOCoRE with the publicly-available AMES model. For a fair comparison, we also train AMES with descriptors of  $d = 768$  dimensionality, *i.e.* same as LOCoRE, using the official implementation.

Regarding the experiments on the other image retrieval benchmarks, we follow previous work [66] and extract global and local descriptors based on a CNN backbone [40, 60]. More precisely, we use ProxyNCA++ [67] to train a ResNet-50 backbone with randomly cropped  $224 \times 224$  images, which is used to extract local descriptors from the dense feature maps of the last convolutional layer. Then, the channel dimension of the feature maps is reduced to the input dimension of re-rankers using a  $1 \times 1$  convolutional layer. On SOP, we directly compare with the results reported in RRT [66].

Method	# desc.	Medium				Hard <sup>1</sup>			
		$\mathcal{ROxf}$	$\mathcal{ROxf}+1M$	$\mathcal{RPar}$	$\mathcal{RPar}+1M$	$\mathcal{ROxf}$	$\mathcal{ROxf}+1M$	$\mathcal{RPar}$	$\mathcal{RPar}+1M$
DELG global & local									
RN50-DELG [9]	-	73.6	60.6	85.7	68.6	51.0	32.7	71.5	44.4
+ GV [9]	1000	78.3	67.2	85.7	69.6	57.9	43.6	71.0	45.7
+ RRT [66]	500	78.1	67.0	86.7	69.8	60.2	44.1	75.1	49.4
+ CVNet Reranker [28]	3,072	78.7	67.7	87.9	72.3	63.0	46.1	76.8	52.5
+ LoCoRE-small	50	80.9	70.9	<b>89.3</b>	77.7	63.3	49.3	<b>77.8</b>	57.9
+ LoCoRE-base	50	<b>81.6</b>	<b>71.7</b>	89.2	<b>77.9</b>	<b>64.2</b>	<b>50.5</b>	77.6	<b>58.2</b>
RN101-DELG [9]	-	76.3	63.7	86.6	70.6	55.6	37.5	72.4	46.9
+ GV [9]	1000	81.2	69.1	87.2	71.5	64.0	47.5	72.8	48.7
+ RRT [66]	500	79.9	-	87.6	-	64.1	-	76.1	-
+ LoCoRE-small	50	81.8	74.1	87.9	75.9	64.2	54.7	75.0	55.2
+ LoCoRE-base	50	<b>83.3</b>	<b>76.1</b>	<b>90.3</b>	<b>80.7</b>	<b>66.9</b>	<b>56.2</b>	<b>81.4</b>	<b>61.8</b>
SuperGlobal global & DINOv2 local									
RN101-SuperGlobal [59]	-	85.3	78.8	92.1	83.9	72.1	61.9	83.5	69.1
+ AMES [65]	50	86.4	80.3	91.7	83.7	73.8	65.4	83.6	70.9
+ AMES* [65]	50	86.4	80.4	91.7	83.8	74.4	65.7	83.7	71.3
+ LoCoRE-small	50	89.0	82.7	91.9	84.0	77.6	68.9	84.9	72.3
+ LoCoRE-base	50	89.7	83.6	92.1	84.3	79.5	71.2	85.2	72.7
+ AMES (top-400) [65]	50	87.4	81.2	92.8	85.7	75.2	67.0	85.6	73.9
+ AMES* (top-400) [65]	50	87.4	81.3	92.8	85.8	75.6	67.3	85.8	74.3
+ SuperGlobal Rerank (top-400) [59]	-	90.9	84.4	93.3	84.9	80.2	71.1	86.7	71.4
+ LoCoRE-small (top-400 <sup>†</sup> )	50	91.0	84.7	93.3	86.4	79.9	70.9	87.3	75.8
+ LoCoRE-base (top-400 <sup>†</sup> )	50	<b>92.0</b>	<b>85.8</b>	<b>93.8</b>	<b>86.8</b>	<b>81.8</b>	<b>73.2</b>	<b>87.7</b>	<b>76.5</b>

Table 1. **Performance (mAP) on  $\mathcal{ROxf}$  and  $\mathcal{RPar}$**  and their 1M distractor variants (+1M) with Medium and Hard evaluation strategy. For a fair comparison, results for re-rankers are reported with their top-100 candidates unless indicated otherwise. Re-ranking for RN50-DELG and RN101-DELG are with DELG local descriptors while re-ranking for RN101-SuperGlobal is based on DINOv2-ViT-B/14 [39] local descriptors. \* indicates AMES trained with the hidden size of 768, serving as a fair comparison with LoCoRE. <sup>†</sup> indicates sliding window re-ranking is enabled with a stride size  $S = 50$ .

For CUB-200 and InShop, we train RRT with the provided implementation and in the identical setup. Furthermore, we compare with re-rankers that operate with global descriptors, *i.e.* SSR Rerank [60], AQE [11], DQE [1], and  $\alpha$ QE [49]. Evaluation for these datasets follows the same evaluation protocols as Jun *et al.* [27].

## 4.2. Results

**Retrieval on landmarks.** We report LoCoRE performance on  $\mathcal{ROxf}$  and  $\mathcal{RPar}$  in Table 1. We observe that LoCoRE-base exhibits clear and consistent improvement over all local-descriptor re-ranking baselines, using different types of descriptors as input and across all settings. The improvement

is most pronounced on  $\mathcal{ROxf}+1M$  and  $\mathcal{RPar}+1M$  in the hard setting, where it achieves significant gains of +17.8 and +13.8 mAP points on RN50-DELG, against the prior state-of-the-art CVNet-Reranker, which achieved +13.4 and +8.1 mAP points, respectively. The base model shows a significant boost compared to the small one. We outperform even the latest state-of-the-art of pair-wise models while re-ranking either 100 or 400 images. The proposed sliding window re-ranking allows us to effectively go beyond the list size used during training.

**Retrieval on metric learning benchmarks.** We report LoCoRE performance on CUB-200 [71], SOP [63] and InShop [32] and compare with previous re-ranking methods in Table 2. Query expansion methods and SSR Rerank lead to marginal or no improvements in retrieval performance, yet we observe inconsistent behaviour across datasets, *e.g.*, the

<sup>1</sup>indicates the hard setting allows *easy* images to be used in the re-ranking and removed before evaluation following [59]. Details in the supplementary clarify how such differences impact the final results.

	CUB-200				SOP				In-Shop			
	R@1	R@2	R@4	mAP	R@1	R@10	R@100	mAP	R@1	R@10	R@20	mAP
Global descriptors	68.9	79.4	87.3	49.8	80.8	92.1	96.9	65.1	88.5	97.7	98.4	74.8
SSR Rerank [60]	69.4	79.0	86.1	54.2	81.2	91.9	95.6	66.7	86.3	97.1	98.2	75.8
AQE [11]	66.9	76.8	82.7	58.4	76.9	89.3	94.5	66.1	81.7	95.9	96.7	73.2
DQE [1]	67.0	75.5	82.0	54.6	67.9	81.5	90.6	47.8	86.4	97.1	98.1	72.4
$\alpha$ QE [49]	70.9	78.8	84.7	56.9	81.1	90.7	96.3	68.1	88.5	97.1	98.1	76.8
RRT [66]	68.7	85.0	95.6	55.6	81.9	92.4	<b>96.9</b>	67.2	88.3	<b>97.9</b>	98.6	77.6
LoCORE-tiny	71.4	86.8	96.4	58.1	82.4	<b>93.1</b>	<b>96.9</b>	68.0	89.1	<b>97.9</b>	98.2	78.4
LoCORE-small	74.6	89.1	97.3	61.0	83.3	92.7	<b>96.9</b>	69.4	<b>89.4</b>	97.7	97.7	<b>78.8</b>
LoCORE-base	<b>78.3</b>	<b>91.9</b>	<b>98.2</b>	<b>64.8</b>	<b>83.8</b>	92.9	<b>96.9</b>	<b>71.0</b>	87.9	<b>97.9</b>	<b>98.7</b>	77.0

Table 2. **Performance (R@k and mAP) on metric learning benchmarks.** Results are reported on the same sets of global descriptors to showcase the effectiveness of re-rankers fairly. mAP refers to the mAP@R. Re-ranking 100 images is used.

N	S	Medium		Hard	
		$\mathcal{R}_{\text{Oxf+1M}}$	$\mathcal{R}_{\text{Par+1M}}$	$\mathcal{R}_{\text{Oxf+1M}}$	$\mathcal{R}_{\text{Par+1M}}$
-	-	78.5	83.6	61.4	68.4
100	N/A	82.7	84.0	68.9	72.3
200	50	83.9	85.1	70.4	74.3
400	25	<b>84.8</b>	<b>86.4</b>	70.8	<b>75.9</b>
400	50	84.7	<b>86.4</b>	<b>70.9</b>	75.8
400	75	84.5	86.1	70.7	75.7
400	100	82.9	84.9	69.1	73.1

Table 3. **Different settings of our sliding window strategy** for re-ranking with LoCORE-small and DINOv2 local descriptors. One forward pass has length of  $K = 100$  for each setting. The sliding windows are used with a stride of  $S$  images, and  $N$  images are re-ranked in total.

DQE method on the SOP dataset severely compromises the original performance of the global retriever. RRT demonstrates improvements across datasets. The re-ranking performance of LoCORE variants is consistently robust, as it shows significant improvements across all datasets and is aligned with increases in both R@k and mAP@R. This is particularly evident from the relative improvements observed on the CUB-200 dataset, where LoCORE-base obtains a +9.4 increase in R@1 and a +15.0 increase in mAP@R, as opposed to the slight decrease of -0.2 in R@1 and smaller gain of +5.8 in mAP@R for RRT. Notably, this is the first work that reports performance improvements with local descriptors re-ranking on CUB and In-Shop. The only previous work reporting such results on SOP is RRT.

### 4.3. Ablation Study and Analysis

**Sliding window re-ranking with LoCORE.** Table 3 presents retrieval performance for different settings of the sliding window approach. This approach effectively re-uses LoCORE to extend the list size defined during training,

Ablation Module	R@1	R@10	mAP
Global descriptors	80.8	92.1	65.1
RRT [66]	81.9	92.4	67.2
LoCORE-tiny	82.4	93.1	68.0
<i>w/o gallery shuffled training</i>	80.7	91.9	65.1
<i>w/o global query attention</i>	60.7	90.0	53.0
<i>w/ shuffled evaluation</i>	81.8 $\pm 0.1$	92.8 $\pm 0.1$	67.5 $\pm 0.0$
LoCORE-small	83.3	92.7	69.4
LoCORE-base	83.8	92.9	71.0
<i>w/ first token score</i>	83.7	92.5	70.9
<i>w/ [SEP] token score</i>	83.7	92.5	70.9

Table 4. **Ablation studies for LoCORE** on the SOP dataset. mAP refers to the mAP@R. Re-ranking 100 images is used.

which is according to the hardware limitations, *i.e.* the GPU memory. We observe that re-ranking with  $N > K$  increases performance compared to  $N = K$  and also that a smaller stride is better up to some extent. We also find that provided  $S < K$ , *i.e.* overlapping the re-ranking window in each sliding process, greatly improves performance, with the performance gains saturating as  $S$  further decreases. This is consistent with our motivation in designing the sliding window re-ranking – to promote interaction across the gallery within the shortlist beyond the context window of the sequence model. As a compromise between performance and runtime,  $S = 50$  is used for the experiment in Table 1.

**Training and evaluation strategies.** Table D shows the results of ablation studies on the SOP benchmark. Without *gallery shuffled training*, LoCORE-tiny exhibits trivial results that are nearly identical to global retrieval. This is because the model learns from positional shortcut as discussed in Section 3.3 and simply repeats the global retrieval results instead of modeling list-wise correspondences. We

Models	LOCORE			RRT [66]	CVNet [28]	AMES* [65]
	tiny	small	base			
<b># of Params</b>	19.4M	58.7M	111.8M	2.2M	7.5M	32.1M
<b>FLOPs</b>	200.0G	517.5G	1035.1G	518.5G	2087.0G	679.2G
<b>FLOP/s (T)</b>	10.8 $\pm$ 0.1	20.9 $\pm$ 0.2	7.4 $\pm$ 0.2	7.0 $\pm$ 0.2	1.8 $\pm$ 0.3	11.1 $\pm$ 0.0
<b>Latency (ms)</b>	18.5 $\pm$ 0.2	24.7 $\pm$ 0.3	140.1 $\pm$ 4.2	74.4 $\pm$ 4.5	1418 $\pm$ 323	61.1 $\pm$ 0.5
<b>Peak Memory</b>	3.2GB	4.5GB	8.9GB	29.0GB	130.0MB	6.0GB

Table 5. **Runtime performance comparison of LOCORE variants with other re-ranking methods.** Floating-point operation (FLOP), latency and memory consumption are measured per 100 re-ranked gallery images with a single query image. \* indicates AMES trained with the hidden size of 768. RRT and CVNet are measured in their original settings.

also observe a significant performance decline when disabling *global query attention*, suggesting that the absence of global query attention makes the model suffer from learning collapse due to challenges in capturing long-distance correspondences.

To investigate whether LOCORE takes advantage of candidate confidence permutations in global retrieval, *i.e.*, the re-ranker preferentially assigns higher scores to top retrieved candidates, we conduct an ablation with *shuffled evaluation*, where we shuffle the candidates in the gallery list even at inference time. We repeat the experiment using 5 different random seeds and report the mean and the standard deviation. Reduced scores demonstrate that the re-ranker has learned to use candidate confidence permutation from global retrieval results, a signal that pair-wise re-rankers cannot exploit.

Finally, we investigate the impact of different strategies to aggregate the individual scores of all tokens of each image. The results with the base model indicate that various aggregators have no substantial impact on re-ranking performance, reflecting the robustness of the proposed method.

**Scalability.** Table 5 shows that the previous largest neural re-ranker CVNet only has 7.5M parameters, in contrast with our smallest model, LOCORE-tiny, which has 19.4M. To explore the potential for scaling up our model, we conduct ablations on model size. The ablation results demonstrate that model performance improves consistently with increases in model size, indicating considerable potential for scalability in our long-context re-ranking method. On the other hand, scaling up pair-wise re-rankers, like RRT, does not yield notable performance improvements. Furthermore, we discuss in the following section a substantial computational cost introduced by large pair-wise re-rankers.

**Runtime Performance.** One major advantage of our approach is that LOCORE can complete the top- $k$  re-ranking with a single forward pass, in contrast to pair-wise re-ranking, which requires  $k$  passes. To quantitatively assess the efficacy of our method, we present runtime performance analysis in Table 5. For models based on the Longformer architecture, even the largest LOCORE-base with 111M parameters ex-

hibits latency close to that of RRT [66]. We can trade off between latency and accuracy with our three model variants. For instance, LOCORE-small offers  $3\times$  lower latency while only marginally sacrificing performance metrics on the  $\mathcal{R}Oxf$  and  $\mathcal{R}Par$  datasets on DELG descriptors. It still offers  $2\times$  lower latency when competing with the efficiency-centric AMES [65] method. All variants of LOCORE exhibit significantly reduced peak memory compared to RRT, with LOCORE-tiny leading at nearly  $10\times$  lower peak memory usage. Long-context models benefit from efficient utilization of the memory hierarchy in modern hardware accelerators, and as a result, all variants of LOCORE achieve higher throughput and increased parallelism as measured in FLOPs per second.

## 5. Conclusion

We present LOCORE, the first image re-ranking framework that leverages list-wise re-ranking supervision at the local descriptor level. With a long-context sequence model, this approach effectively captures dependencies between the query image and each gallery image in addition to the dependencies among the gallery images themselves and implicitly learns to calibrate predictions for a more precise ranking list. As demonstrated through extensive ablation studies, learning from list-wise re-ranking signals enables the model to better leverage the initial rankings obtained from global retrieval. We achieve leading re-ranking results across established image retrieval datasets and obtain state-of-the-art metrics on the  $\mathcal{R}Oxf$  and  $\mathcal{R}Par$  datasets under the same settings.

**Acknowledgements.** This work was partially supported by NSF Award #2201710 and funding from the Ken Kennedy Institute at Rice University. This work was also supported by the Junior Star GACR GM no. 21-28830M, Horizon MSCA-PF no. 101154126, and the Czech Technical University in Prague no. SGS23/173/OHK3/3T/13. This work was supported in part by the NSF Campus Cyberinfrastructure grant ‘‘CC\* Compute: Interactive Data Analysis Platform’’ NSF OAC-2019007, and by Rice University’s Center for Research Computing (CRC).



# LoCORE: Image Re-ranking with Long-Context Sequence Modeling

## Supplementary Material

### A. Implementation details

All training is conducted on 8 NVIDIA A100 PCI-E 40GB GPUs. Training on Google Landmark v2 clean set [73] takes 106 hours on LoCORE-base for 5 epochs. Models are trained with AdamW optimizer [33],  $5e-5$  learning rate and weight decay disabled. Global batch size is set to 128 with 4 gradient accumulation steps. We present the configurations of the different LoCORE variants in Table A. LoCORE-tiny is initialized from `roberta-tiny-cased`<sup>2</sup> by migrating weights and repeatedly copying absolute position embedding along the sequence dimension<sup>3</sup>. LoCORE-small is initialized from the first 6 layers of `longformer-base-4096`<sup>4</sup>, while LoCORE-base is initialized from the full `longformer-base-4096`. To accommodate  $50 \text{ descriptors} \times (1 \text{ query image} + 100 \text{ re-ranking candidates}) = 5,050$  tokens, the position embeddings in the original models are linearly interpolated to extend their length from 4,096 to 5,120.

When experimenting with local descriptors from DINOv2 [39], we use the same training set as AMES [65], which is approximately half the size of the full GLDv2 clean set, *i.e.* 750k images. We adopt the same global-local ensemble scheme as AMES. The ensemble hyper-parameters are selected based on the best-performing configuration on GLDv2 public validation split and applied to  $\mathcal{R}Oxf$  and  $\mathcal{R}Par$  evaluations. For the training of AMES\*, we follow the original training process from AMES, changing only the batch size and learning rate to 150 and  $1e-5$ , respectively.

For baseline results on CUB-200 [71], Stanford Online Products (SOP) [63] and In-shop [32], we reproduce them using their official code releases and identical training configuration, except for ProxyNCA++ [67], we change the training image size from  $256 \times 256$  to  $224 \times 224$  to use the training image size same as the other baselines.

For the performance benchmark in Section 4.3, we use the Deepspeed [54] profiler on a single NVIDIA A100 GPU to measure key performance metrics of the model per 100 re-ranked gallery images as follows: the number of parameters (# of Params), floating-point operations (FLOPs), throughput in FLOPs per second, latency, and peak memory consumption. All dynamic metrics are reported with 10 warmup steps followed by 10 measurements for reporting the mean and standard deviation. Parameters of visual backbones are

<sup>2</sup><https://huggingface.co/haisongzhang/roberta-tiny-cased>

<sup>3</sup>[https://github.com/allenai/longformer/blob/master/scripts/convert\\_model\\_to\\_long.ipynb](https://github.com/allenai/longformer/blob/master/scripts/convert_model_to_long.ipynb)

<sup>4</sup><https://huggingface.co/allenai/longformer-base-4096>

Model Variants	tiny	small	base
Number of Parameters	19.4M	58.7M	111.8M
Number of Layers	4	6	12
Local Attention Window	1024	512	512
Hidden Size	512	768	768
Intermediate Size	2048	3072	3072
Number of Attention Heads	8	12	12
Max Context Length	5120	5120	5120

Table A. Architectural parameters of LoCORE variants.

excluded from # of Params.

We consider the descriptors to be already extracted and exclude I/O from measuring memory, latency, *etc.* For the geometric verification (GV) method, we run RANSAC in OpenCV [22] with 1,000 iterations on AMD EPYC 9354 CPU and measure the wall-clock time as the latency and the maximum resident set size (Max RSS) as the peak memory consumption. All models are benchmarked with batched input except CVNet Reranker [28]. It is worth noting that CVNet Reranker does not support batched inference since it computes pair-wise multi-scale correlation on raw feature maps (without resizing) from query and gallery images of different sizes. Thus, CVNet Reranker heavily underutilizes the GPU and achieves extremely low throughput and high latency. The FLOP, latency, and peak memory are measured assuming query and gallery images of  $512 \times 512$  size in CVNet Reranker.

### B. Additional Experimental Results

#### B.1. Additional comparisons

We present additional experiments with different combinations of global and local features in Table B. We compare with more baseline re-ranking methods, including methods with global, *i.e.* SuperGlobal (SG) Rerank [59], and local, *i.e.* AMES [65], RRT [66], R2Former [78], descriptors. We evaluate the models under different Hard settings, using different global descriptors to generate the shortlist and different backbones for feature extraction. We also test the combination of LoCORE with other re-ranking schemes.

**Variations for Hard setup.** As mentioned in the main paper, there can be two approaches regarding how to handle *easy* images in the hard setup: (i) **Hard**: *easy* images are used to re-rank and removed before the evaluation (typically used in the literature [59]), and (ii) **Hard\***: *easy* images are completely removed from the database. While the two choices (Hard and Hard\*) are equivalent for pair-wise

Global	Local	Re-rank	$\mathcal{R}Oxf+1M$			$\mathcal{R}Par+1M$			
			Medium	Hard	Hard*	Medium	Hard	Hard*	
SG	N/A	N/A <sup>†</sup>	78.8		61.9	83.9		69.1	
		N/A	78.5		61.4	83.6		68.4	
		SG-Rerank <sup>†</sup>	84.4	71.1	N/A	84.9	71.4	N/A	
		SG-Rerank	84.0	69.4	63.9	85.2	72.3	75.7	
	CVNet	R2Former	79.9		63.7	83.8		69.7	
		RRT	79.3		62.7	83.6		69.1	
		AMES	80.7		65.7	84.6		71.8	
		LoCoRE	81.9	68.6	64.9	84.6	71.4	70.7	
		SG + LoCoRE	84.7	71.5	65.6	86.2	74.8	76.1	
		DINOv2	R2Former*	81.0		66.2	84.9		72.1
	RRT*		81.0		66.1	85.5		73.3	
	AMES*		81.3		67.3	85.8		74.3	
	LoCoRE		85.8	75.8	73.2	86.8	75.9	76.5	
	SG + LoCoRE		86.5	76.3	73.7	87.2	76.9	78.2	
	DINOv2		N/A	N/A	59.6		35.2	77.0	
		SG-Rerank		62.2	40.5	31.2	79.8	60.5	65.8
DINOv2		R2Former*	67.8		44.6	78.6		61.3	
		RRT*	68.8		46.0	79.6		64.0	
		AMES*	68.9		46.8	79.9		64.7	
		LoCoRE	73.4	54.9	52.5	80.9	66.4	66.7	
		SG + LoCoRE	71.2	54.4	48.7	81.9	68.7	69.5	

Table B. **Additional results** with re-ranking top-400 candidates. Hard\*: *easy* images are completely removed from the database. Hard: *easy* images are used to re-rank and removed before the evaluation. †: results in the SuperGlobal paper [59]. LoCoRE is reported with the base variant. SG + LoCoRE: re-ranking with SG first and then with LoCoRE. \* indicates models trained with 768 hidden size, serving as a fair comparison with LoCoRE. N/A: not available.

re-ranking methods, this is not the case when interactions between database images are considered (*i.e.* LoCoRE, SG-rerank). In Table B, it is evident that the two setup lead to significantly different results. In most cases, mAP considerably drops in  $\mathcal{R}Oxf$ , comparing results in Hard and Hard\*; whereas, mAP increases in  $\mathcal{R}Par$ .

**Performance with other backbones.** First, we benchmark all models when the shortlist is generated based on DINOv2 global descriptors. It is noteworthy that DINOv2 global descriptors are significantly worse than SG ones. In this setup, LoCoRE outperforms all other re-ranking schemes by a vast margin.

Second, we evaluate LoCoRE using local descriptors extracted from CVNet backbones. CVNet local descriptors have a higher dimension than that of DINOv2, *i.e.* 1024 vs 768; hence, we used a learnable linear projector to match the embedding dimensionality of the transformer. LoCoRE achieves competitive performances compared with the pair-wise re-rankers, with only AMES outperforming it in a few cases. Yet, all local-based re-rankers are outperformed by SG-Rerank. Nevertheless, LoCoRE with DINOv2 outper-

forms SG-Rerank.

**Combination with SG-Rerank.** It is straightforward to combine local and global-based re-ranking. To this end, we combine LoCoRE with SG-Rerank by applying global re-ranking first, followed by local re-ranking. This combination achieves the best performance when SG is used as global descriptor. However, this combination hurts LoCoRE performance on  $\mathcal{R}Oxf$  when DINOv2 is used as global.

**Performance per query.** To highlight the advantages of our proposed list-wise re-ranking over pair-wise re-ranking, we present several scatter plots in Figure A, showing the average precision of each sample in  $\mathcal{R}Oxf+1M$  Hard before and after re-ranking with different re-ranking paradigms. We compare our model with AMES [65], which is considered the state-of-the-art solution for pair-wise re-ranking. In the first two plots, we observe that most data points are concentrated in the upper-left half and above the red reference line, indicating that both re-ranking paradigms improve the ranking accuracy for the majority of query images. However, the list-wise re-ranking method driven by LoCoRE has barely any sample points below the red reference line, meaning

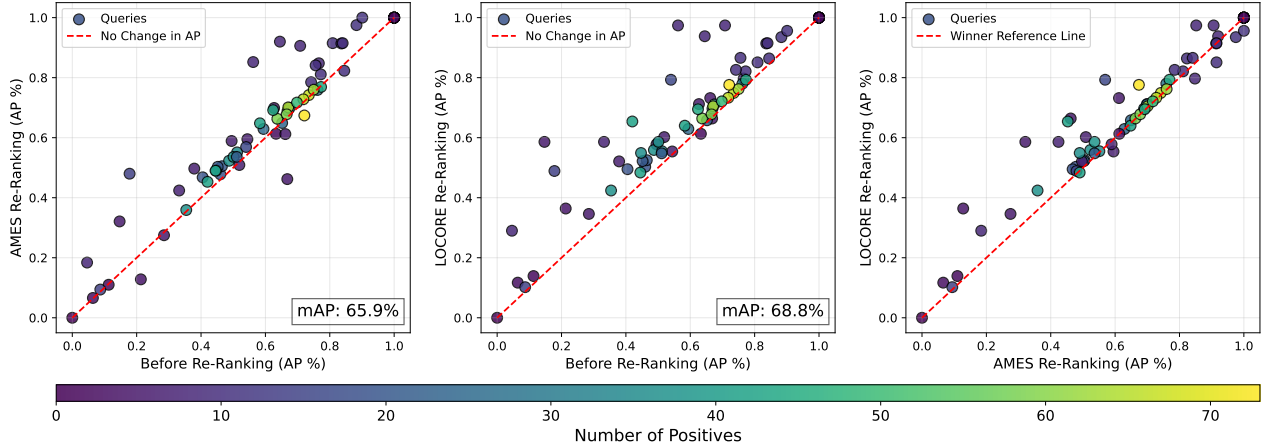


Figure A. **Average precision per query scatter plot** on  $\mathcal{R}Oxf+1M$  Hard for global-only vs. AMES (Left), global-only vs. LoCORE-small (Middle) and AMES vs. LoCORE-small (Right). Global descriptors are from RN101-Superglobal, which by itself achieves mAP=61.4%. Re-ranking is performed for top-100 candidates, and the color bar indicates the number of positive images in the shortlist for each query.

Global	Local	$L$	$K$	$\mathcal{R}Oxf+1M$		$\mathcal{R}Par+1M$	
				Medium	Hard*	Medium	Hard*
		50	100	85.8	73.2	86.8	76.5
SG	DINOv2	100	50	83.9	68.6	85.2	72.5
		25	200	84.5	72.1	85.6	75.1

Table C. Additional results for LoCORE-base with different combinations of the number of local descriptors  $L$  and the number of re-ranking candidates  $K$  on  $N = 400$  candidates.

the re-ranking only improves the retrieval on the individual query level. The distinction between the two models is most prominent in the final plot, where the number of sample points above the winner reference line far exceeds those below, demonstrating that LoCORE outperforms AMES on more query samples. We also observed that the list-wise re-ranking method is relatively robust in terms of the number of positive samples included in the shortlist, as the color distribution of the sample points does not exhibit any discernible pattern. This indicates the general versatility of LoCORE.

## B.2. Additional ablations

**Number of images vs number of descriptors.** We explore the relationship between the number of local descriptors and the number of image candidates within a given context window in Table C. Specifically, we set the context window to 5,120 and examine three configurations of LoCORE: (i) using 100 gallery images with 50 local descriptors per image, *i.e.* the default setup, (ii) using 200 gallery images with 25 local descriptors per image, *i.e.* more candidate images but fewer descriptors per image, and (iii) using 50 database images with 100 local descriptors per image, *i.e.* more descriptors per image but fewer candidate images. The LoCORE in the default settings reports the best results.

**Comparison with other recurrent models.** Other model architectures with no restrictions on context length that could

Ablation Module	$\mathcal{R}@1$	$\mathcal{R}@10$	mAP@ $R$
Global descriptors	80.8	92.1	65.1
LoCORE-tiny	82.4	93.1	68.0
LoCORE-small	83.3	92.7	69.4
LoCORE-base	83.8	92.9	71.0
LoCORE-RWKV	81.4	92.3	66.7
LoCORE-Mamba	80.6	92.1	66.4

Table D. Ablation studies for LoCORE recurrent models on the SOP dataset. Re-ranking is performed with the top 100 candidates.

be employed instead of LongFormer are the recently proposed recurrent models Mamba [18] and RWKV [42]. As the causal nature of the recurrence-based model does not align well with our re-ranking motivation and is strictly less expressive than bi-directional encoders [26, 51], we follow the common practice in recurrent visual encoder community [14, 29, 31] to build a bi-directional variant that serves as an efficient sequence encoder. To ensure recurrent models can still handle long-range interactions and alleviate the inherent information bottleneck in the design of recurrent models, we devise a mechanism that resembles the query global attention in Section 3.2 by interleaving recurrent blocks with uni-directional transformer blocks [70]. These transformer blocks compute attention scores between intermediate hidden states of query image tokens and intermediate hidden states of gallery image tokens and produce fused intermediate representations for the following layers to process. The uni-directional attention guarantees that every gallery image has similar difficulty accessing the query image, irrespective of its position in the sequence relative to the query. Although we find that these recurrence-based models could slightly outperform the base global retrieval model, they do not surpass our transformer-based results, as shown in Table D.

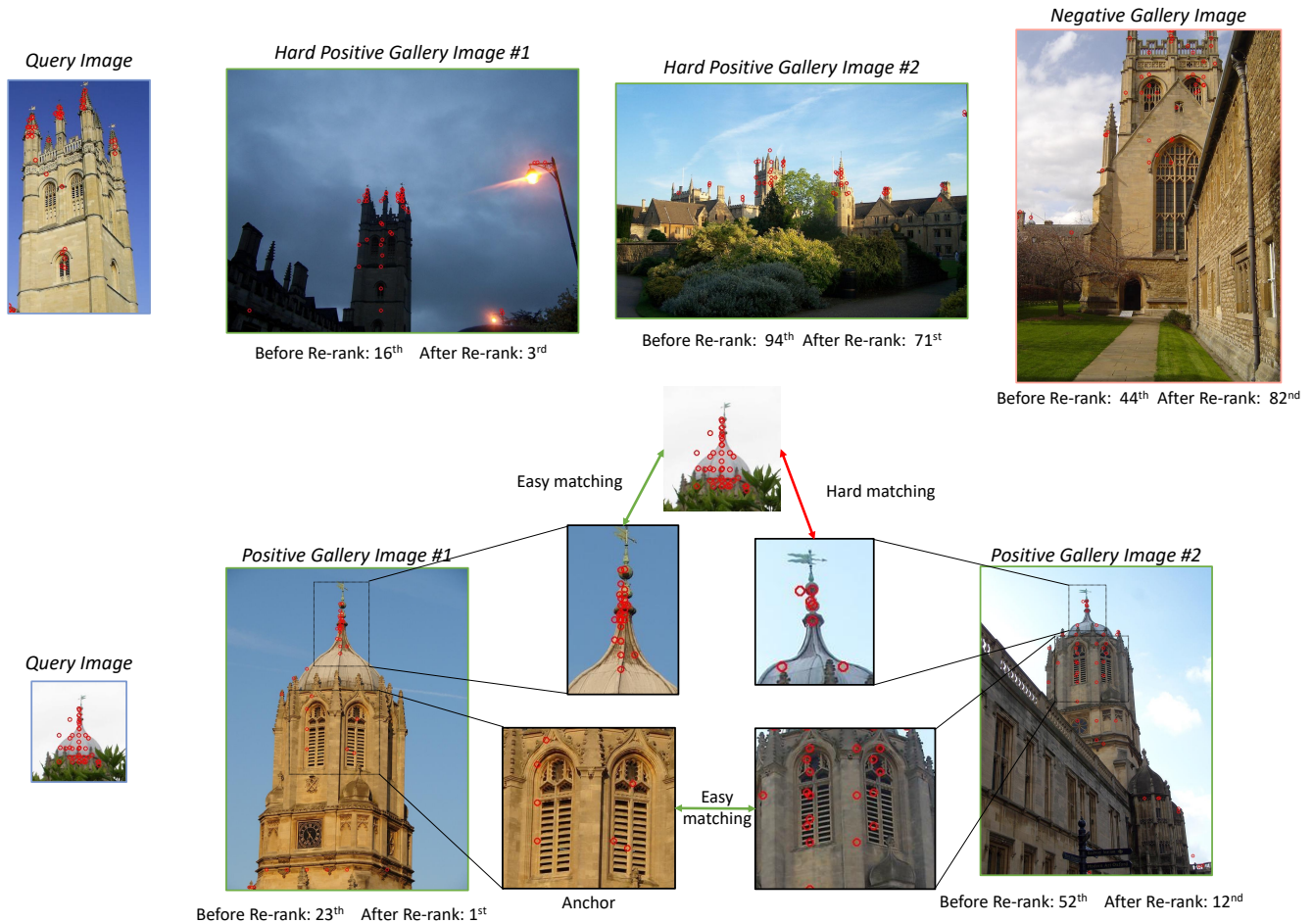


Figure B. **Qualitative analysis** on  $\mathcal{R}$ Oxford dataset of LOCORE-base on RN50-DELG descriptors. **Upper:** two hard positive gallery images get assigned with higher ranks while a negative gallery image is put in lower ranks after re-ranking. **Lower:** the first gallery image can be easily identified as positive due to its dense matching with the query image; it can also serve as a perfect anchor image for refining the ranking of the second gallery image due to their transitive relationship.

### B.3. Qualitative Results

We illustrate the re-ranking performance of LOCORE in Figure B as qualitative results. The upper example underscores the superior performance of our method, demonstrated by its success in elevating the ranking of two hard positive images and lowering that of the negative gallery image. We also show in the lower example that our model is able to capture the transitive relationship between query and gallery images. The transitive relationship is based on the assumption that generally, if two gallery images are similar and one of them is predicted as positive, then the other should be calibrated with higher confidence. In the lower example, the correspondence between the query image and the first gallery image is easy to catch as the common geometric features are evident, resulting in Easy matching in the figure. However, although the global retriever returns the second gallery image as re-

ranking candidates, the sparse local features focused on the top of the tower make it hard for pair-wise re-ranker to assign this gallery image a high confidence score. This misalignment is calibrated by our list-wise re-ranking paradigm since the windows in both gallery images can serve as the anchor to propagate the positive prediction from the easy candidate to the hard one.

Additionally, in Figure C the easy positive gallery has visual overlap with the query (rooftop). The hard positive gallery has little visual overlap with the query, but larger overlap with the first positive (e.g. windows). We wish to answer this question: *Are the local features of the window improving the rank of the hard positive due to a transitive relationship?* We remove local features of the windows (blue crosses), repeat the similarity estimation, and compare the ranks. The decreased similarity score is a sign of LOCORE capturing transitive relationships.



Figure C. **Visualization of LOCORE capturing transitive relationships in gallery images.** We prevent LOCORE from accessing local features of the easy positive corresponding to the windows (blue crosses) and instead randomly sample local features from other negative images. The dropped similarity score indicates LOCORE relies on the transitivity of local features to calibrate predictions for hard positive gallery images.

## C. Limitations and Future Work

Despite the merits in efficiency and re-ranking performance, our model is inherently restricted by the context window of existing encoder-only sequence models. A limited context window limits the number of re-ranking candidates in the gallery and the number of local descriptors that LOCORE can use. While recurrent models offer more flexibility with the context window size, we find that they could not capture list-wise re-ranking dependencies as well as transformer-based models, resulting in sub-optimal performance. Future work could adopt large-scale decoder-only sequence models which typically have longer context windows and greater capacity for list-wise re-ranking. Additionally, context parallelization techniques (e.g., RingAttention [30], Infini-attention [35]) could help expand the context window of current Transformer encoder models. Lastly, extractive re-ranking as proposed in our work could also be seamlessly adopted for other modalities, e.g. document or video re-ranking.

## References

- [1] Relja Arandjelovic and Andrew Zisserman. Three things everyone should know to improve object retrieval. In *2012 IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, USA, June 16-21, 2012*, pages 2911–2918. IEEE Computer Society, 2012. 2, 6, 7
- [2] Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic. NetVLAD: CNN architecture for weakly supervised place recognition. In *CVPR*, 2016. 2
- [3] Yannis Avrithis and Giorgos Tolias. Hough pyramid matching: Speeded-up geometry re-ranking for large scale image retrieval. *Int. J. Comput. Vis.*, 107(1):1–19, 2014. 2
- [4] Artem Babenko and Victor Lempitsky. Aggregating deep convolutional features for image retrieval. In *ICCV*, 2015. 2
- [5] Edoardo Barba, Tommaso Pasini, and Roberto Navigli. ESC: redesigning WSD with extractive sense comprehension. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 4661–4672. Association for Computational Linguistics, 2021. 4
- [6] Edoardo Barba, Luigi Procopio, and Roberto Navigli. Extend: Extractive entity disambiguation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 2478–2488. Association for Computational Linguistics, 2022. 1, 4
- [7] Herbert Bay, Andreas Ess, Tinne Tuytelaars, and Luc Van Gool. Speeded-up robust features (SURF). *Comput. Vis. Image Underst.*, 110(3):346–359, 2008. 2
- [8] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv: Arxiv-2004.05150*, 2020. 4
- [9] Bingyi Cao, Andre Araujo, and Jack Sim. Unifying deep local and global features for image search. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*, pages 726–743. Springer, 2020. 1, 2, 3, 5, 6
- [10] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *CVPR*, 2005. 2
- [11] Ondrej Chum, James Philbin, Josef Sivic, Michael Isard, and Andrew Zisserman. Total recall: Automatic query expansion with a generative feature model for object retrieval. In *IEEE 11th International Conference on Computer Vision, ICCV 2007, Rio de Janeiro, Brazil, October 14-20, 2007*, pages 1–8. IEEE Computer Society, 2007. 2, 6, 7
- [12] Gabriella Csurka, Christopher R. Dance, Lixin Fan, Jutta Katharina Willamowski, and Cédric Bray. Visual categorization with bags of keypoints. In *eccv*, 2004. 2
- [13] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 2
- [14] Yuchen Duan, Weiyun Wang, Zhe Chen, Xizhou Zhu, Lewei Lu, Tong Lu, Yu Qiao, Hongsheng Li, Jifeng Dai, and Wenhai Wang. Vision-rwkv: Efficient and scalable visual perception with rwkv-like architectures, 2024. 3
- [15] Markus Eberts and Adrian Ulges. Span-based joint entity and relation extraction with transformer pre-training. In *ECAI 2020 - 24th European Conference on Artificial Intelligence, 29 August-8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020 - Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020)*, pages 2006–2013. IOS Press, 2020. 1
- [16] Martin A. Fischler and Robert C. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24:381–395, 1981. 2
- [17] Albert Gordo, Filip Radenovic, and Tamara Berg. Attention-based query expansion learning. In *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXVIII*, pages 172–188. Springer, 2020. 2

- [18] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv: 2312.00752*, 2023. 3
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3
- [20] Benjamin Heinzerling and Michael Strube. Sequence tagging with contextual and non-contextual subword representations: A multilingual evaluation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 273–291, Florence, Italy, 2019. Association for Computational Linguistics. 1
- [21] Zhiheng Huang, Wei Xu, and Kai Yu. Bidirectional LSTM-CRF models for sequence tagging. *CoRR*, abs/1508.01991, 2015. 1
- [22] Itseez. Open source computer vision library. <https://github.com/opencv/opencv>, 2015. 1
- [23] Hervé Jégou, Matthijs Douze, and Cordelia Schmid. Hamming embedding and weak geometric consistency for large scale image search. In *ECCV*, 2008. 2
- [24] Hervé Jégou, Matthijs Douze, and Cordelia Schmid. Improving bag-of-features for large scale image search. *IJCV*, 2010. 2
- [25] Hervé Jégou, Matthijs Douze, Cordelia Schmid, and Patrick Pérez. Aggregating local descriptors into a compact image representation. In *The Twenty-Third IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2010, San Francisco, CA, USA, 13-18 June 2010*, pages 3304–3311. IEEE Computer Society, 2010. 2
- [26] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, pages 4904–4916. PMLR, 2021. 3
- [27] HeeJae Jun, ByungSoo Ko, Youngjoon Kim, Insik Kim, and Jongtaek Kim. Combination of multiple global descriptors for image retrieval. *CoRR*, abs/1903.10663, 2019. 6
- [28] Seongwon Lee, Hongje Seong, Suhyeon Lee, and Euntai Kim. Correlation verification for image retrieval. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 5364–5374. IEEE, 2022. 1, 2, 3, 5, 6, 8
- [29] Kunchang Li, Xinhao Li, Yi Wang, Yinan He, Yali Wang, Limin Wang, and Yu Qiao. Videomamba: State space model for efficient video understanding, 2024. 3
- [30] Hao Liu, Matei Zaharia, and Pieter Abbeel. Ring attention with blockwise transformers for near-infinite context. *arXiv preprint arXiv: 2310.01889*, 2023. 5
- [31] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan Liu. Vmamba: Visual state space model, 2024. 3
- [32] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 5, 6, 1
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. 1
- [34] David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vis.*, 60(2):91–110, 2004. 2
- [35] Tsendsuren Munkhdalai, Manaal Faruqui, and Siddharth Gopal. Leave no context behind: Efficient infinite context transformers with infini-attention. *arXiv preprint arXiv: 2404.07143*, 2024. 5
- [36] Dat Quoc Nguyen, Dai Quoc Nguyen, Dang Duc Pham, and Son Bao Pham. A robust transformation-based learning approach using ripple down rules for part-of-speech tagging. *AI communications*, 29(3):409–422, 2016. 4
- [37] D. Nister and H. Stewenius. Scalable recognition with a vocabulary tree. In *CVPR*, 2006. 2
- [38] Hyeonwoo Noh, A. Araújo, Jack Sim, Tobias Weyand, and Bohyung Han. Large-scale image retrieval with attentive deep local features. *IEEE International Conference on Computer Vision*, 2016. 2
- [39] M. Oquab, Timoth’ee Darcet, Théo Moutakanni, Huy Q. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russ Howes, Po-Yao (Bernie) Huang, Shang-Wen Li, Ishan Misra, Michael G. Rabbat, Vasu Sharma, Gabriel Synnaeve, Huijiao Xu, H. Jégou, J. Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. *Trans. Mach. Learn. Res.*, 2023. 2, 5, 6, 1
- [40] Jianbo Ouyang, Hui Wu, Min Wang, Wengang Zhou, and Houqiang Li. Contextual similarity aggregation with self-attention for visual re-ranking. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 3135–3148, 2021. 3, 5
- [41] Yash Patel, Giorgos Tolias, and Jiri Matas. Recall@k surrogate loss with large batches and similarity mixup. In *CVPR*, 2022. 2
- [42] Bo Peng, Eric Alcaide, Quentin Gregory Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Nguyen Chung, Leon Derczynski, Xingjian Du, Matteo Grella, Kranthi Kiran GV, Xuzheng He, Haowen Hou, Przemyslaw Kazienko, Jan Kocon, Jiaming Kong, Bartłomiej Koptyra, Hayden Lau, Jiaju Lin, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Guangyu Song, Xianguo Tang, Johan S. Wind, Stanisław Woźniak, Zhenyuan Zhang, Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. RWKV: Reinventing RNNs for the transformer era. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. 3
- [43] Matthew E. Peters, Waleed Ammar, Chandra Bhagavatula, and Russell Power. Semi-supervised sequence tagging with bidirectional language models. In *Proceedings of the 55th*

- Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1756–1765. Association for Computational Linguistics, 2017. 1
- [44] James Philbin, Ondrej Chum, Michael Isard, Josef Sivic, and Andrew Zisserman. Object retrieval with large vocabularies and fast spatial matching. In *2007 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2007)*, 18-23 June 2007, Minneapolis, Minnesota, USA. IEEE Computer Society, 2007. 2, 5
- [45] James Philbin, Ondrej Chum, Michael Isard, Josef Sivic, and Andrew Zisserman. Object retrieval with large vocabularies and fast spatial matching. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2007. 1, 5
- [46] James Philbin, Ondrej Chum, Michael Isard, Josef Sivic, and Andrew Zisserman. Lost in quantization: Improving particular object retrieval in large scale image databases. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2008)*, 24-26 June 2008, Anchorage, Alaska, USA. IEEE Computer Society, 2008. 2, 5
- [47] Bill Psomas, Ioannis Kakogeorgiou, Konstantinos Karantzas, and Yannis Avrithis. Keep it simple: Who said supervised transformers suffer from attention deficit? In *ICCV*, pages 5350–5360, 2023. 2
- [48] Filip Radenovic, Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Revisiting oxford and paris: Large-scale image retrieval benchmarking. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 5706–5715. Computer Vision Foundation / IEEE Computer Society, 2018. 2, 5
- [49] Filip Radenovic, Giorgos Tolias, and Ondrej Chum. Fine-tuning CNN image retrieval with no human annotation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(7):1655–1668, 2019. 6, 7
- [50] Filip Radenović, Giorgos Tolias, and Ondřej Chum. Fine-tuning cnn image retrieval with no human annotation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(7):1655–1668, 2019. 2
- [51] 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. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, pages 8748–8763. PMLR, 2021. 3
- [52] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 2383–2392. The Association for Computational Linguistics, 2016. 1, 4
- [53] Lance Ramshaw and Mitch Marcus. Text chunking using transformation-based learning. In *Third Workshop on Very Large Corpora*, 1995. 1
- [54] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*, pages 3505–3506. ACM, 2020. 1
- [55] Ali S Razavian, Josephine Sullivan, Stefan Carlsson, and Atsuto Maki. Visual instance retrieval with deep convolutional networks. *ITE Trans. on Media Technology and Applications*, 2016. 2
- [56] Sánchez, Florent Perronnin, Thomas Mensink, and Jakob Verbeek. Image classification with the fisher vector: Theory and practice. *IJCV*, 2013. 2
- [57] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015. 2
- [58] Elad Segal, Avia Efrat, Mor Shoham, Amir Globerson, and Jonathan Berant. A simple and effective model for answering multi-span questions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 3074–3080. Association for Computational Linguistics, 2020. 1, 4
- [59] Shihao Shao, Kaifeng Chen, Arjun Karpur, Qinghua Cui, André Araujo, and Bingyi Cao. Global features are all you need for image retrieval and reranking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11036–11046, 2023. 2, 5, 6, 1
- [60] Xi Shen, Yang Xiao, Shell Xu Hu, Othman Sbai, and Mathieu Aubry. Re-ranking for image retrieval and transductive few-shot classification. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 25932–25943, 2021. 5, 6, 7
- [61] Oriane Simeoni, Yannis Avrithis, and Ondrej Chum. Local features and visual words emerge in activations. In *CVPR*, 2019. 2
- [62] Josef Sivic and Andrew Zisserman. Video Google: A text retrieval approach to object matching in videos. In *ICCV*, 2003. 2
- [63] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 5, 6, 1
- [64] Yuxin Song, Ruolin Zhu, Min Yang, and Dongliang He. Dalg: Deep attentive local and global modeling for image retrieval. In *arxiv*, 2022. 2
- [65] Pavel Suma, Giorgos Kordopatis-Zilos, Ahmet Iscen, and Giorgos Tolias. Ames: Asymmetric and memory-efficient similarity estimation for instance-level retrieval. In *ECCV*, 2024. 1, 3, 5, 6, 8, 2
- [66] Fuwen Tan, Jiangbo Yuan, and Vicente Ordonez. Instance-level image retrieval using reranking transformers. *IEEE International Conference on Computer Vision*, 2021. 1, 3, 5, 6, 7, 8
- [67] Eu Wern Teh, Terrance Devries, and Graham W. Taylor. Proxynca++: Revisiting and revitalizing proxy neighborhood com-

- ponent analysis. *European Conference on Computer Vision*, 2020. 5, 1
- [68] Giorgos Tolias and Hervé Jégou. Visual query expansion with or without geometry: Refining local descriptors by feature aggregation. *Pattern Recognition*, 2014. 2
- [69] Giorgos Tolias, Yannis Avrithis, and Hervé Jégou. Image search with selective match kernels: Aggregation across single and multiple images. *Int. J. Comput. Vis.*, 116(3):247–261, 2016. 2
- [70] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [71] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. Technical report, California Institute of Technology, 2011. 2, 5, 6, 1
- [72] Weinzaepfel, Philippe and Lucas, Thomas and Larlus, Diane and Kalantidis, Yannis. Learning Super-Features for Image Retrieval. In *ICLR*, 2022. 2
- [73] Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2 - a large-scale benchmark for instance-level recognition and retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 5, 1
- [74] Min Yang, Dongliang He, Miao Fan, Baorong Shi, Xuetong Xue, Fu Li, Errui Ding, and Jizhou Huang. DOLG: single-stage image retrieval with deep orthogonal fusion of local and global features. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 11752–11761. IEEE, 2021. 2
- [75] Wonjin Yoon, Richard Jackson, Aron Lagerberg, and Jaewoo Kang. Sequence tagging for biomedical extractive question answering. *Bioinformatics*, 38(15):3794–3801, 2022. 1, 4
- [76] Hao Zhang, Xin Chen, Heming Jing, Yingbin Zheng, Yuan Wu, and Cheng Jin. ETR: an efficient transformer for re-ranking in visual place recognition. In *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2023, Waikoloa, HI, USA, January 2-7, 2023*, pages 5654–5663. IEEE, 2023. 1
- [77] Yidan Zhang, Ting Zhang, Dong Chen, Yujing Wang, Qi Chen, Xing Xie, Hao Sun, Weiwei Deng, Qi Zhang, Fan Yang, Mao Yang, Qingmin Liao, and Baining Guo. Irgen: Generative modeling for image retrieval. *CoRR*, abs/2303.10126, 2023. 2
- [78] Sijie Zhu, Linjie Yang, Chen Chen, Mubarak Shah, Xiaohui Shen, and Heng Wang. R2former: Unified retrieval and reranking transformer for place recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19370–19380, 2023. 1