

SEEING THROUGH WORDS: CONTROLLING VISUAL RETRIEVAL QUALITY WITH LANGUAGE

Anonymous authors

Paper under double-blind review

ABSTRACT

Text-to-image retrieval is a fundamental task in vision-language learning, yet in real-world scenarios it is often challenged by short and underspecified user queries. Such queries are typically only one or two words long, making them semantically ambiguous, prone to collisions across diverse visual interpretations, and lacking explicit control over the quality of retrieved images. To address these issues, we propose a new paradigm of *quality-controllable retrieval*, which enriches short queries with contextual details while incorporating explicit notions of image quality. Our key idea is to leverage a generative large language model as a query completion function, extending underspecified queries into descriptive forms that capture fine-grained visual attributes such as pose, scene, and aesthetics. We introduce a training framework that conditions query completion on discretized quality levels, derived from relevance and aesthetic scoring models, so that query enrichment is not only semantically meaningful but also quality-aware. The resulting system provides three key advantages: ① *flexibility*, as it is compatible with any pretrained vision-language model without modification; ② *transparency*, since enriched queries are explicitly interpretable by users; and ③ *controllability*, enabling retrieval results to be steered toward user-preferred quality levels. Extensive experiments demonstrate that our proposed approach significantly improves retrieval results and provides effective quality control, bridging the gap between the expressive capacity of modern vision-language models and the underspecified nature of short user queries.

1 INTRODUCTION

Text-to-image retrieval (T2IR) aims to return the most relevant images from a gallery given a textual query. Recent progress in this task has been largely driven by vision-language models (VLMs) (Radford et al., 2021; Jia et al., 2021; Yu et al., 2022; Li et al., 2022; Yang et al., 2022; Li et al., 2023; Yang et al., 2024; Lu et al., 2024), which learn joint representations of text and images through large-scale pretraining on web-scale image-text pairs (Schuhmann et al., 2021; 2024; Liu et al., 2023a). Such models significantly narrow the semantic gap between modalities and achieve strong alignment across diverse benchmarks (Ilharco et al., 2021; Singh et al., 2022; Gao et al., 2022; Khan & Fu, 2023; Wang et al., 2024; Li et al., 2024).

Despite these advances, retrieval performance often degrades in realistic scenarios where user queries are very short (typically just one or two words, e.g., “a dog”). Short queries encode only limited semantics, which results in large and ambiguous search subspaces and less discriminative results. This issue becomes more pronounced in large-scale galleries, where underspecified queries yield many candidate matches and cause semantic collisions among visually diverse results.

Another limitation of existing retrieval systems is their singular focus on semantic alignment. Naïve retrieval approaches simply return the top- k images with the highest similarity scores, overlooking other critical aspects of user satisfaction such as aesthetics, interestingness, or popularity. In practice, *retrieval quality* is context-dependent: art students may prefer visually inspiring images, architects may seek unique and creative references, and shoppers may favor popular or visually appealing products. However, conventional systems lack mechanisms for steering retrieval toward these quality dimensions.

To address these limitations, we introduce the task of *quality-controllable retrieval* (QCR). Formally, given a frozen VLM and a short textual query, the objective is to retrieve images that not only align semantically but also satisfy user-specified quality requirements. This setting is feasible because short queries naturally span a broad subspace that contains images of varying quality levels. With appropriate conditioning, this subspace can be partitioned into perceptually distinct regions, enabling fine-grained quality-aware retrieval.

In this work, we define retrieval quality along two widely applicable dimensions: *relevance* (semantic consistency) (Cherti et al., 2023) and *aesthetics* (visual appeal) (Yi et al., 2023). For each image in the gallery, we construct auxiliary annotations consisting of a textual description, a relevance score, and an aesthetic score. We discretize these continuous scores into categorical quality levels (e.g., Low, Medium, High) and associate each description with its corresponding quality condition.

The central challenge is how to steer retrieval results toward specific quality levels given only a short query. We propose a simple yet effective solution: *quality-conditioned query completion* (QC²). QC² enriches short queries with quality-aware details by leveraging a generative large language model (LLM). Trained on the quality-augmented dataset, the LLM learns to append appropriate descriptive phrases that capture both semantic and quality-related attributes. Conditioning on different quality levels guides retrieval toward the desired regions of the search space. This is particularly valuable because, in practice, users often do not know how to formulate queries that precisely reflect their preference or may not be aware of what constitutes “high” or “low” quality within the dataset. By learning from how textual descriptions vary across quality scores, our approach bridges this gap and enables more controllable retrieval through query completion. To summarize, our key contributions are summarized as follows:

- *Problem*: we introduce quality-controllable retrieval, a new setting where retrieval can be explicitly conditioned on user-defined quality requirements.
- *Methodology*: we propose QC², a query completion framework that leverages LLMs to enrich short queries with quality-aware descriptive details.
- *Validation*: we conduct extensive experiments to show that QC² effectively steers retrieval outcomes according to quality preferences and is readily adaptable to multiple VLMs.

2 PRELIMINARIES

2.1 MOTIVATION

We study the problem of text-to-image retrieval, where the goal is to return the desired images from a large gallery given a set of natural language queries. Specifically, let $\mathcal{Q} := \{Q_1, \dots, Q_m\}$ denote a collection of m text queries and $\mathcal{I} := \{I_1, \dots, I_n\}$ an image gallery of size n . We consider a state-of-the-art VLM as the retrieval backbone, equipped with a text encoder $g : \mathcal{Q} \rightarrow \mathbb{R}^d$ and an image encoder $f : \mathcal{I} \rightarrow \mathbb{R}^d$, both producing d -dimensional normalized embeddings. Given a query set \mathcal{Q} , the system returns the top- η relevant images according to

$$\mathcal{X} := \text{sort}(f(\mathcal{I}), g(\mathcal{Q}), \eta), \quad (1)$$

where $\mathcal{X} \subseteq \mathcal{I}$ denotes the top- η matches of queries \mathcal{Q} . The sort function typically operates on the similarity scores $\mathbf{S} \in \mathbb{R}^{m \times n}$ with $S_{ij} := g(Q_i)^\top f(I_j)$.

Although modern VLMs achieve strong cross-modal alignment, retrieval performance deteriorates in realistic scenarios where user queries are usually very short (typically just one or two words, e.g., “a dog”). Given such short queries, naive retrieval system faces several challenges: ① *Semantic ambiguity*: a few words can refer to a wide range of possible images, leading to a large and diffuse search subspace with less discriminative retrieval results. ② *Semantic collisions*: short queries tend to yield similar similarity scores for visually diverse images (e.g., realistic vs. cartoon dogs). These collisions confuse ranking and are particularly problematic in large-scale galleries where many candidate images satisfy the vague query. ③ *Lack of quality control*: the quality of retrieved images is not explicitly enforced during retrieval. At best, one can apply post-retrieval filtering, but the system itself provides no mechanism to ensure that high-quality results consistently appear among the top matches. These issues highlight a fundamental gap between the expressive capacity of modern VLMs and the underspecified nature of user queries, motivating the need for enriched query representations and controllable retrieval mechanisms.

2.2 PROBLEM SETTING

To address the above limitations, we propose to enrich short queries with additional descriptive details that potentially capture more distinguishable attributes of images. Formally, let h denote a query completion function that maps \mathcal{Q} to enriched queries $h(\mathcal{Q})$. Retrieval is then performed as

$$\tilde{\mathcal{X}} := \text{sort}(f(\mathcal{I}), g(h(\mathcal{Q})), \eta), \quad (2)$$

where $h(\mathcal{Q})$ augments the short queries with contextual details. The enriched queries are expected to capture not only object categories but also additional information such as pose, scene, action, and fine-grained attributes. To be effective, the completion function should *be aware of* the retrieval gallery, so that it generates meaningful context rather than irrelevant or hallucinated content.

To achieve this, we implement h using a generative large language model (LLM). However, simply training the LLM on image descriptions is insufficient, since it cannot guarantee that retrieval results satisfy user expectations of quality. Instead, we partition the textual descriptions into non-overlapped quality levels \mathcal{C} that reflect different image quality categories. We then finetune the LLM with these quality levels, enabling it to generate query completions conditioned on quality preferences. This yields the following *quality-controllable retrieval* (QCR) formulation:

$$\tilde{\mathcal{X}} := \text{sort}(f(\mathcal{I}), g(\text{LLM}(\mathcal{Q}; \mathcal{C})), \eta), \quad (3)$$

where $\text{LLM}(\mathcal{Q}; \mathcal{C})$ expands the short queries based on the specified quality constraint \mathcal{C} . The extended queries thus steer retrieval toward images that align with the desired quality criteria.

This approach offers several practical benefits: ① *Flexibility*: it requires no modification to pre-trained VLMs and remains compatible with any VLMs; ② *Transparency*: the generated query completions are human-readable, allowing users to review and select preferred options. ③ *Controllability*: the LLM can produce different query completions with different quality conditions \mathcal{C} , enabling explicit quality control during retrieval.

2.3 THEORETICAL ANALYSIS

Before implementing the completion function, we justify why enriching short queries may help to improve retrieval. Let $\mathcal{Q}^+ = \{Q_1^+, \dots, Q_m^+\} := h(\mathcal{Q})$ denote the extended queries by h , where $Q_i^+ := Q_i + \text{suffix}_i$, $\forall i \in \{1, \dots, m\}$, and suffix_i denotes additional descriptive details appended to query Q_i . Let $\mathbf{C} \in \mathbb{R}^{n \times d}$ be the image embedding matrix with rows $\mathbf{c}_j := f(I_j) \in \mathbb{R}^d$, $\forall j \in \{1, \dots, n\}$, and $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times d}$ be two sets of text embeddings with a strict one-to-one pairing of rows, with rows $\mathbf{a}_i := g(Q_i) \in \mathbb{R}^d$ and $\mathbf{b}_i := g(Q_i^+) \in \mathbb{R}^d$, $\forall i \in \{1, \dots, m\}$. Let $r := \text{rank}(\mathbf{A})$ be the rank of \mathbf{A} , $\sigma_r(\mathbf{A})$ be the smallest nonzero singular value of \mathbf{A} , and $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ denote its singular value decomposition (SVD). We then partition the right singular vectors as $\mathbf{V} = [\mathbf{V}_S \ \mathbf{V}_\perp]$, where $\mathbf{V}_S \in \mathbb{R}^{d \times r}$ and $\mathbf{V}_\perp \in \mathbb{R}^{d \times (d-r)}$ satisfy $\text{span}(\mathbf{V}_S) = \mathcal{R}(\mathbf{A})$ and $\text{span}(\mathbf{V}_\perp) = \mathcal{R}(\mathbf{A})^\perp$, with $\mathcal{R}(\mathbf{A}) := \text{span}\{\mathbf{a}_1^\top, \dots, \mathbf{a}_m^\top\} \subseteq \mathbb{R}^d$ the row space of \mathbf{A} .

Definition 1. We define a perturbation matrix $\mathbf{\Delta} := \mathbf{B} - \mathbf{A} \in \mathbb{R}^{m \times d}$, score matrices $\mathbf{S}_A := \mathbf{A}\mathbf{C}^\top \in \mathbb{R}^{m \times n}$ and $\mathbf{S}_B := \mathbf{B}\mathbf{C}^\top \in \mathbb{R}^{m \times n}$ for the queries \mathcal{Q} and \mathcal{Q}^+ , $\mathbf{A}_S := \mathbf{A}\mathbf{V}_S$, $\mathbf{\Delta}_S := \mathbf{\Delta}\mathbf{V}_S$, $\mathbf{\Delta}_\perp := \mathbf{\Delta}\mathbf{V}_\perp$, $\mathbf{C}_S := \mathbf{C}\mathbf{V}_S$, $\mathbf{C}_\perp := \mathbf{C}\mathbf{V}_\perp$, $\mathbf{X} := (\mathbf{A}_S + \mathbf{\Delta}_S)\mathbf{C}_S^\top$, $\mathbf{Y} := \mathbf{\Delta}_\perp\mathbf{C}_\perp^\top$, $\mathcal{U} := \text{col}(\mathbf{X})$, and $\mathbf{P} := \mathbf{P}_X$ as the orthogonal projector onto \mathcal{U} .

Proposition 1. Assume that: i) $\|\mathbf{\Delta}\|_2 < \sigma_r(\mathbf{A})$; ii) there exists $I \subseteq \{1, \dots, n\}$ with $|I| = r$ such that the columns \mathbf{X}_I form a basis of \mathcal{U} ; iii) $\|\mathbf{X}_I^\dagger \mathbf{P} \mathbf{Y}_I\|_2 < 1$; and iv) there exists disjoint index set $K \subseteq \{1, \dots, n\} \setminus I$ such that $k := \text{rank}\left((\mathbf{I} - \mathbf{P}_{Z_I})\mathbf{Z}_K\right) \geq 1$, where $\mathbf{Z} := (\mathbf{I} - \mathbf{P})\mathbf{Y}$, $\mathbf{Z}_I := \mathbf{Z}_{:,I}$, and $\mathbf{Z}_K := \mathbf{Z}_{:,K}$. Then, $\text{rank}(\mathbf{S}_B) \geq r + k > r = \text{rank}(\mathbf{S}_A)$.

Remark 1. We decompose $\mathbf{\Delta}$ into two parts: one ($\mathbf{\Delta}_S$) that lies in the original row space of \mathbf{A} , and another ($\mathbf{\Delta}_\perp$) that introduces directions outside this space. Assumption (i) ensures the in-span perturbation $\mathbf{\Delta}_S$ is not too large (controlled by $\sigma_r(\mathbf{A})$) so the original r query directions in \mathbf{A} are not destroyed by completion. Assumption (ii) asserts that we can select r columns from \mathbf{X} that span \mathcal{U} . This fixes a stable r -dimensional basis for the existing subspace. Assumption (iii) claims that adding the projected perturbation $\mathbf{P}\mathbf{Y}_I$ does not reduce the independence of these r columns.

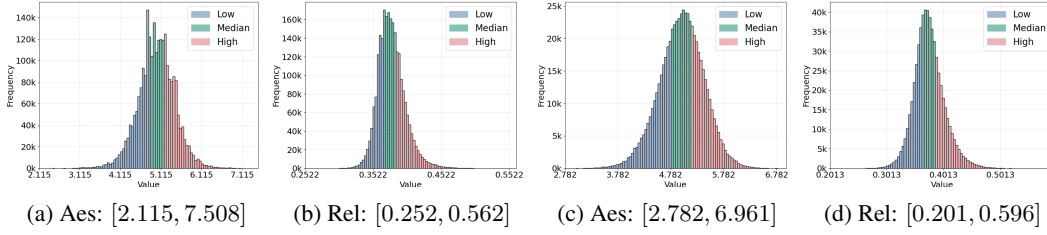


Figure 1: Aesthetic and relevance score distributions of Flickr2.4M in (a) and (b), and of MS-COCO in (c) and (d). It is worth noting that the numbers of high-quality and low-quality images are limited, which leads to the average scores of any two random sets being very close.

Thus the original r -dimensional structure is preserved. Assumption (iv) requires that there exist $k \geq 1$ columns outside I whose orthogonal components (after removing projections onto both \mathcal{U} and $\text{col}(\mathbf{Z}_I)$) are linearly independent. These contribute k genuinely new directions in \mathcal{U}^\perp . Together, these assumptions ensure that the rank of \mathbf{S}_B contains at least the r preserved dimensions from \mathcal{U} plus the k fresh orthogonal ones. Consequently, \mathbf{S}_B can express more independent scoring patterns and has the ability to potentially make finer-grained distinctions.

3 QUALITY-CONDITIONED QUERY COMPLETION

This section first provides the definition of quality, and then shows how to construct the training data, and finally illustrates how to implement and train the query completion function LLM.

3.1 QUALITY DEFINITION

For the proposed QCR task, we require a clear notion of *quality*. In this work, we characterize quality along two primary dimensions: ① *Relevance*, which measures the semantic consistency between textual queries and their corresponding images; and ② *Aesthetics*, which reflects the visual appeal or attractiveness of retrieved images. Note that our framework is inherently flexible, permitting the incorporation of arbitrary quality metric, provided that corresponding and reliable scoring models are available and applicable to general image datasets. Other notions of quality, such as *interestingness* (Gygli et al., 2013; Abdullahu & Grabner, 2024) can also be adopted in a similar manner and are left for future exploration. To facilitate user control over retrieved results, we discretize the quality dimensions into non-overlapping conditions. Specifically, we define \mathcal{C}^R for relevance and \mathcal{C}^A for aesthetics, each partitioned into perceptually distinct and user-friendly levels. For example, both can be represented as $\mathcal{C}^R, \mathcal{C}^A := \{\text{Low}, \text{Medium}, \text{High}\}$.

3.2 DATA GENERATION

To ensure the completion function LLM can perceive the retrieval gallery, we construct an augmented training dataset for each gallery \mathcal{I} . The dataset integrates three complementary components: textual descriptions $\mathcal{T} = \{T_i\}_{i=1}^n$ of images, relevance scores $\mathbf{s}^r \in \mathbb{R}^n$, and aesthetic scores $\mathbf{s}^a \in \mathbb{R}^n$.

Textual Descriptions. For each image I_i , we generate a textual description T_i using an image caption model $\text{CAP}(\cdot)$, i.e., $T_i = \text{CAP}(I_i), \forall i \in \{1, \dots, n\}$. In our experiments, we utilize strong pretrained captioning models without additional fine-tuning for description generation. Each T_i is a concise sentence summarizing the main content of the image.

Aesthetic Scores. We assign an aesthetic score s_i^a to each image I_i using an aesthetic evaluation model $\text{EV}_A(\cdot)$, i.e., $s_i^a = \text{EV}_A(I_i), \forall i \in \{1, \dots, n\}$. The aesthetic scores represent the visual quality of the images, with higher scores indicating greater visual appeal.

Relevance Scores. For each image-description pair $\{I_i, T_i\}$, we compute a relevance score using a pretrained VLM. Specifically, we extract the image feature $f(I_i)$ and text feature $g(T_i)$, then calculate their cosine similarity as their relevance score, i.e., $s_i^r = \cos(f(I_i), g(T_i)), \forall i \in \{1, \dots, n\}$.

Table 1: Query completions with their retrieved images and quality scores on MS-COCO

Rel: Low , Aes: Low	Rel: Median , Aes: Median	Rel: High , Aes: High
		
<i>a train that is sitting on the tracks in gravel</i> Aes 4.715, Rel 0.347	<i>a train sitting on the tracks with black smoke coming out of it</i> Aes 4.818, Rel 0.382	<i>a train is traveling near some water and houses</i> Aes 5.935, Rel 0.394
		
<i>a bird standing on the ground near some leaves</i> Aes 4.616, Rel 0.346	<i>a bird flying above some brown water</i> Aes 5.079, Rel 0.374	<i>a bird flying across some water at the beach</i> Aes 5.120, Rel 0.386
		
<i>a teddy bear wearing eye glasses and laying on a bed</i> Aes 4.788, Rel 0.359	<i>a teddy bear that is sitting on a tree</i> Aes 5.649, Rel 0.388	<i>a teddy bear sitting on a wall next to an old stone house</i> Aes 5.818, Rel 0.437

3.3 TRAINING FRAMEWORK

Score Discretization. To simulate quality-controlled retrieval, we discretize the continuous quality scores of images into categorical levels that are more intuitive for users. Given a score vector \mathbf{r} (either aesthetics \mathbf{s}^a or relevance \mathbf{s}^r), each score r_i is mapped into one of three descriptive levels by partitioning the score distribution into three percentiles:¹

$$l(r_i) = \begin{cases} \text{Low}, & r_i \leq \text{perc}(\mathbf{r}, p_1), \\ \text{High}, & r_i > \text{perc}(\mathbf{r}, p_2), \\ \text{Median}, & \text{otherwise.} \end{cases} \quad (4)$$

Here, r_i is the score of the i -th sample and $\text{perc}(\mathbf{r}, p)$ calculates the $p\%$ percentile of \mathbf{r} as




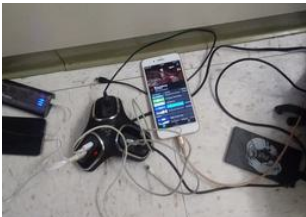


$$\text{perc}(\mathbf{r}, p) = \tilde{\mathbf{r}}[\lfloor \xi \rfloor] + (\xi - \lfloor \xi \rfloor) \cdot (\tilde{\mathbf{r}}[\lfloor \xi \rfloor + 1] - \tilde{\mathbf{r}}[\lfloor \xi \rfloor]), \quad (5)$$

where $\xi = \frac{p}{100} \cdot (n - 1)$, $\tilde{\mathbf{r}}$ is the sorted version of \mathbf{r} , and $\lfloor \cdot \rfloor$ is the floor function. Figure 1 illustrates the distributions of aesthetics and relevance scores and their discretized partitions.

Instruction Design. We train the completion function LLM on the augmented training set $\mathcal{D} = \{\mathcal{T}, \mathbf{s}^a, \mathbf{s}^r\}$. The discretized quality levels serve as explicit conditions within instructions, enabling

¹Our framework is general and supports arbitrary numbers of levels depending on the desired granularity. In Table 5, Sec. 4.5, we provide an example with five quality levels.

Table 2: Query completions with their retrieved images and quality scores on Flickr2.4M

Aes: Low , Rel: Low	Aes: Median , Rel: Median	Aes: High , Rel: High
		
<i>a chair with wires on it</i> Aes 4.019, Rel 0.362	<i>a chair with red and black ropes on it</i> Aes 4.847, Rel 0.379	<i>a chair on a stage in a field</i> Aes 5.257, Rel 0.387
		
<i>a cell phone with wires attached to it</i> Aes 4.531, Rel 0.362	<i>a cell phone with an acoustic guitar on it</i> Aes 5.035, Rel 0.390	<i>a cell phone on a tripod in front of a waterfall in yellowstone national park</i> Aes 5.441, Rel 0.429

LLM to generate quality-aware query completions. For each image I_i , we design a concise instruction P_i of the form:

$$\text{"Relevance: } l(s_i^r), \text{ Aesthetic: } l(s_i^a), \text{ Query: "}$$

where $l(s_i^r)$ and $l(s_i^a)$ represent the categorical quality levels defined in Eq. (4). During training, this instruction provides a lightweight yet effective mechanism to condition query generation on specified quality preferences.

Model Training. To stimulate the quality control process during model training, we use the descriptive levels $l(s_i^r)$ and $l(s_i^a)$ of image I_i as the quality conditions, which are incorporated into the instruction P_i . We then concatenate instructions P_i with the textual description T_i for each image I_i , and then train the completion model LLM with the standard autoregressive next-token prediction loss. In this way, LLM learns to generate query completions that are not only semantically relevant but also controllable according to the given quality constraints.

Inference Strategy. During the inference stage, we concatenate a similar instruction with each testing query. To simulate user preferences, we evaluate various relevance-aesthetic combinations, such as “low relevance, low aesthetic” and “high relevance, high aesthetic”. Then the model generates completed queries based on the instructions, testing queries, and specified quality conditions. For efficient similarity search on large-scale galleries, we utilize the FAISS library (Johnson et al., 2019) to identify the nearest images for the queries.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETTINGS

Datasets. We evaluate our method on two image datasets: one with real textual descriptions and one without. For the image-only one, we construct a large dataset sourced from the Openverse website (Openverse, 2025). We refer to this dataset as Flickr2.4M, which contains over 2.4 million CC0-licensed images randomly selected from the Flickr subset of Openverse. For image

Table 3: Retrieval quality of various methods on `Flickr2 . 4M`. CoCa and Blip2 are used to generate textual descriptions; **L** (Low), **M** (Median), and **H** (High) indicate the quality conditions; and Ctrl specifies whether the method enables controllable retrieval over quality. For both average relevance (Ave Rel) and average aesthetics (Ave Aes), higher values indicate better retrieval quality.

Quality	VLM	Aes Cond Rel Cond	L L	L M	L H	M L	M M	M H	H L	H M	H H	Ctrl ?
Prefix	--	Ave Aes Ave Rel	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	×
LLaMA3	--	Ave Aes Ave Rel	4.730 0.351	4.822 0.351	4.831 0.351	4.823 0.353	4.837 0.351	4.784 0.350	4.798 0.354	4.722 0.354	4.842 0.352	×
GPT-4o	--	Ave Aes Ave Rel	4.359 0.378	4.651 0.361	4.728 0.357	4.712 0.358	4.668 0.360	4.791 0.356	4.791 0.361	4.816 0.357	5.056 0.361	×
PT	--	Ave Aes Ave Rel	4.776 0.345	4.556 0.346	4.722 0.349	4.811 0.349	4.781 0.346	4.757 0.348	4.693 0.345	4.751 0.350	4.746 0.350	×
FT	CoCa	Ave Aes Ave Rel	4.756 0.365	4.834 0.368	4.777 0.364	4.838 0.363	4.863 0.369	4.882 0.368	4.821 0.365	4.905 0.364	4.770 0.365	×
Ours	CoCa	Ave Aes Ave Rel	4.458 0.355	4.615 0.366	4.530 0.391	4.934 0.354	4.852 0.360	4.841 0.386	5.222 0.353	5.170 0.368	5.270 0.390	✓
FT	Blip2	Ave Aes Ave Rel	4.795 0.370	4.871 0.368	4.890 0.367	4.894 0.367	4.844 0.371	4.856 0.366	4.901 0.367	4.847 0.371	4.888 0.369	×
Ours	Blip2	Ave Aes Ave Rel	4.541 0.353	4.523 0.370	4.455 0.397	4.940 0.354	4.906 0.366	4.922 0.396	5.309 0.355	5.222 0.372	5.191 0.390	✓

datasets with real textual descriptions, we adopt the widely-used MS-COCO dataset for experiments, which includes both images and human-annotated descriptions. Specifically, we utilize the training subset of MS-COCO, which consists of 118, 287 samples, each sample containing one image and five corresponding descriptions. In total, approximately 0.6 million descriptions are used for training.

Model Selection. For the backbone of our method, we evaluate two different LLMs: GPT2-1.5B (Radford et al., 2019) and Qwen2.5-0.5B (Yang et al., 2024). Other LLMs can be validated similarly and we leave them for future study. We implement the caption models CAP(\cdot) using a pretrained CoCa (Yu et al., 2022) and a pretrained Blip2 (Li et al., 2023) model, respectively. For feature extraction, we adopt a pretrained VLM OpenCLIP (ViT-H-14-quickgelu) (Cherti et al., 2023; Ilharco et al., 2021). The relevance score is computed as the cosine similarity between the features of each image-description pair. The aesthetic evaluation model $EV_A(\cdot)$ is realized using a pretrained aesthetic predictor (Schuhmann, 2022).

Implementation Details. All experiments are conducted on a node with 8 NVIDIA A100 GPUs. For GPT2-1.5B (Radford et al., 2019), we set the learning rate, warmup steps, number of epochs, and batch size to $2e-3$, 100, 50, 150, respectively. For Qwen2.5-0.5B (Yang et al., 2024), these hyperparameters are set to $2e-5$, 100, 30, and 80, respectively. For score discretization, we set $p_1 = 33$ and $p_2 = 66$ to divide the score distribution into three evenly spaced percentiles (examples of five-level cases are also considered). Note that we only train LLM for query completion, while the quality evaluation model $EV_A(\cdot)$, the caption models CAP(\cdot), and the retrieval model VLM are all pretrained without additional fine-tuning. Since the pretrained caption models may occasionally generate non-English characters, we clean these characters directly before training to prevent potential issues for query completion. Before training, we prepend a start token to the instructions and append an end token to the descriptions. The training loss is computed only on the tokens of the descriptions and the end tokens, while excluding those of the instructions.

Evaluation Strategy. For performance evaluation, we use the 80 class names from MS-COCO dataset as query objectives. These include common objects such as trains, cars, and animals, as well as more specific categories like teddy bear, fire hydrant, and toothbrush. Based on the capitalization of each class name, we prepend either "a" or "an" to form the input queries. Since we focus on controlling the quality of retrieved images, we use aesthetic and relevance scores as the evaluation metrics. We calculate and report the average aesthetic and relevance scores of the retrieved images across all

Table 4: Retrieval quality of various methods on MS-COCO, where **L** (Low), **M** (Median), and **H** (High) indicate the quality conditions for retrieval, and Ctrl specifies whether the method enables controllable retrieval over image quality. For both average relevance (Ave Rel) and average aesthetics (Ave Aes), higher values indicate better retrieval quality.

Quality	Aes Cond Rel Cond	L L	L M	L H	M L	M M	M H	H L	H M	H H	Ctrl ?
Prefix	Ave Aes Ave Rel	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	×
LLaMA3	Ave Aes Ave Rel	4.903 0.348	4.891 0.349	4.855 0.347	4.916 0.348	4.875 0.349	4.880 0.347	4.871 0.348	4.858 0.350	4.911 0.344	×
GPT-4o	Ave Aes Ave Rel	4.673 0.371	4.754 0.357	4.686 0.354	4.782 0.360	4.808 0.358	4.880 0.350	4.838 0.359	5.075 0.352	5.048 0.351	×
PT	Ave Aes Ave Rel	4.819 0.343	4.793 0.340	4.789 0.344	4.829 0.348	4.810 0.339	4.828 0.344	4.794 0.346	4.826 0.343	4.820 0.340	×
FT	Ave Aes Ave Rel	4.925 0.370	4.845 0.367	4.848 0.366	4.882 0.368	4.934 0.368	4.990 0.365	4.849 0.371	4.941 0.371	4.929 0.367	×
FT-CoCa	Ave Aes Ave Rel	4.878 0.346	4.852 0.351	4.859 0.356	4.902 0.349	4.858 0.350	4.941 0.354	4.952 0.345	4.961 0.352	4.944 0.352	×
FT-Blip2	Ave Aes Ave Rel	4.828 0.350	4.815 0.352	4.785 0.356	4.932 0.344	4.894 0.351	4.893 0.353	5.034 0.345	4.948 0.351	4.933 0.347	×
Ours	Ave Aes Ave Rel	4.811 0.356	4.790 0.370	4.773 0.382	4.911 0.354	4.873 0.370	4.862 0.387	5.016 0.352	4.983 0.365	5.024 0.387	✓

input queries as the final evaluation performance. We also test the results using the recall metric for further validation, which can be seen in the appendix.

4.2 QUALITATIVE VALIDATION

We first perform qualitative analysis to validate whether our approach effectively achieves quality control in retrieval. In Tables 1 and 2, we present three retrieved images per query, along with their corresponding completed queries and quality scores under three different quality conditions. As shown, our method generates distinct query completions for different quality conditions. From left to right, as the quality level improves, both aesthetic and relevance scores increase accordingly. This demonstrates that our proposed method effectively controls the quality of the retrieved images. Due to space limitations, we provide more qualitative results in the Appendix A.5.

4.3 QUANTITATIVE VALIDATION

Since no existing retrieval methods are directly applicable to the proposed QCR task, we design the following baselines for quantitative comparison: a) *Prefix*: using the input query prefix directly without query completion; b) *PT (Pretrained)*: using a pretrained LLM for query completion without finetuning; c) *FT (Finetuned)*: finetuning a pretrained LLM on textual descriptions while conditioning on randomly generated quality scores; and d) *general-purpose LLMs*, including pretrained LLaMA-3 (LLaMA-3-8B-Instruct) and GPT-4o (via API). Tables 3 and 4 report the retrieval performance of the baseline models and our proposed method with Qwen2.5 on the two datasets, respectively. The key observations are summarized as follows: ① Prefix-only retrieval yields unsatisfactory quality performance, highlighting the necessity of query completion. ② Pretrained models for query completion degrade retrieval quality, performing worse than using only the query prefix in most cases. This is because these pretrained models tend to generate irrelevant words, negatively impacting retrieval performance. ③ Finetuning on textual descriptions improves both relevance and aesthetics compared to prefix-only and pretrained models. However, models finetuned on randomly assigned scores fail to effectively control the quality of retrieved images. ④ Our method not only enhances retrieval under high-quality conditions but also excels in quality control, demonstrating strong adaptability regardless of whether it is trained on real or generated captions.

Table 5: Results with five quality levels.

\mathcal{M}		Relevance (Red \rightarrow Red)				
		VL	L	M	H	VH
Aesthetics (Green)	VL	4.597	4.507	4.610	4.529	4.445
		0.355	0.364	0.375	0.382	0.397
	L	4.805	4.765	4.825	4.729	4.761
		0.353	0.366	0.369	0.380	0.392
Aesthetics (Green)	M	4.909	4.961	4.878	4.889	4.901
		0.355	0.3642	0.370	0.3754	0.390
	H	5.028	4.967	5.045	4.952	5.009
		0.355	0.365	0.370	0.374	0.387
Aesthetics (Green)	VH	5.282	5.153	5.263	5.245	5.121
		0.355	0.363	0.371	0.378	0.389

Table 6: Comparison with post-retrieval filtering, where the *rerank* method first retrieves the top- k images based on relevance and then reorders the candidates by aesthetic scores to identify the best result.

		k	1	2	3	5	10
Rerank	Aes		4.735	4.947	5.014	5.198	5.313
	Rel		0.350	0.348	0.347	0.345	0.341
LLaMA3	Aes		4.842	5.071	5.154	5.298	5.377
	Rel		0.352	0.349	0.347	0.342	0.337
GPT-4o	Aes		5.056	5.205	5.293	5.393	5.518
	Rel		0.361	0.356	0.353	0.349	0.343
Ours	Aes		5.236	5.320	5.364	5.432	5.533
	Rel		0.387	0.385	0.381	0.376	0.366

4.4 DATASET DEPENDENCE

To achieve quality control in retrieval, the model should be tailored to the specific dataset, as different datasets exhibit varying quality characteristics. To illustrate this, we conduct cross-dataset retrieval experiments. Specifically, we evaluate retrieval quality on MS-COCO using queries completed by the model finetuned on Flickr2.4M. In Table 4, we assess FT-CoCa and FT-Blip2, which are finetuned on descriptions generated by CoCa and Blip2, respectively. The results indicate that both models achieve higher aesthetic scores as quality conditions improve, suggesting that aesthetically relevant semantic cues may be universal across natural images. Nevertheless, they consistently exhibit low relevance across all quality conditions. This limitation stems from the dataset mismatch between the query completion and image retrieval stages, since the two datasets encode different semantic information. See Appendix A.2 for additional analysis and results.

4.5 FURTHER VALIDATION

Table 5 presents the results on the Flickr2.4M dataset across five quality levels: VL (Very-Low), L (Low), M (Median), H (High), and VH (Very-High). As shown, our method effectively enables fine-grained control over the quality of retrieved images, adhering to more nuanced descriptive constraints. We also compare against a post-retrieval filtering baseline that first retrieves images based on relevance and then re-ranks the results by aesthetic scores. The comparison results are listed in Table 6. As shown, this two-stage strategy is unreliable for short queries, which typically offer vague representations and limited descriptive cues. As a result, the initial retrieval set tends to be semantically broad and aesthetically subpar, leaving little room for the re-ranking step to improve. While increasing k can surface images with higher aesthetic quality, it typically comes at the cost of reduced semantic relevance, illustrating a trade-off between these two dimensions. In contrast, our method performs quality control during the query stage, which inherently guides retrieval toward the desired quality level. This quality-aware conditioning cannot be achieved by the two-step baseline, which lacks knowledge of the dataset’s quality distribution and operates in a detached, post-hoc manner. See Appendix A.5 for more experimental results.

5 CONCLUSION

We presented a quality-controllable retrieval framework to address the limitations of short and underspecified text queries in text-to-image retrieval. Our approach enriches queries using a generative language model conditioned on discretized quality levels, enabling retrieval that is both semantically expressive and aligned with user preferences. Extensive experiments demonstrate that our method effectively improves and controls retrieval quality, serving as a flexible augmentation to existing VLMs while improving quality control in retrieval. Future work will extend our method to other dimensions of quality beyond relevance and aesthetics, such as interestingness, diversity, or user personalization. We hope this work inspires further research on integrating controllable language-based query enrichment with large-scale multimodal retrieval systems.

STATEMENTS

ETHICS STATEMENT

This work investigates large language models for query completion in text-to-image retrieval, where image quality information is integrated into the training process. The study relies on publicly available datasets and does not involve human subjects, private information, or sensitive content. We acknowledge that retrieval models may inherit biases present in the underlying vision-language datasets; however, our approach does not introduce new data collection and instead focuses on methodological contributions. The models and results are intended solely for academic research, and no harmful or deceptive applications are pursued. We adhere to the ICLR Code of Ethics and confirm that this research complies with principles of fairness, transparency, and responsible use.

REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our work. All code and pretrained checkpoints used in our experiments will be released upon acceptance. Theoretical results are stated with all necessary assumptions in the main text, and their complete proofs are provided in the appendix. Experimental settings are included in the main body and supplementary materials. Together, these resources are intended to enable full replication and verification of our results.

REFERENCES

- Fitim Abdullahu and Helmut Grabner. Commonly interesting images, 2024. URL <https://arxiv.org/abs/2409.16736>.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- Ziv Bar-Yossef and Naama Kraus. Context-sensitive query auto-completion. In *Proceedings of the 20th International Conference on World Wide Web, WWW ’11*, pp. 107–116. Association for Computing Machinery, 2011.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, pp. 1877–1901, 2020.
- Fei Cai, Maarten De Rijke, et al. A survey of query auto completion in information retrieval. *Foundations and Trends® in Information Retrieval*, 10(4):273–363, 2016.
- Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2818–2829, 2023.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint*, 2018.
- Yuting Gao, Jinfeng Liu, Zihan Xu, Jun Zhang, Ke Li, and Chunhua Shen. Pyramidclip: hierarchical feature alignment for vision-language model pretraining. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, 2022.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, and Ahmad Al-Dahle et al. The llama 3 herd of models, 2024.

- Michael Gygli, Helmut Grabner, Hayko Riemenschneider, Fabian Nater, and Luc Van Gool. The interestingness of images. In *Proceedings of the IEEE international conference on computer vision*, pp. 1633–1640, 2013.
- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pp. 4904–4916. PMLR, 2021.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.
- Zaid Khan and Yun Fu. Contrastive alignment of vision to language through parameter-efficient transfer learning. *arXiv preprint*, 2023.
- Dong-Ho Lee, Zhiqiang Hu, and Roy Ka-Wei Lee. Improving text auto-completion with next phrase prediction. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 4434–4438. Association for Computational Linguistics, 2021.
- Yibin Lei, Yu Cao, Tianyi Zhou, Tao Shen, and Andrew Yates. Corpus-steered query expansion with large language models. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 393–401. Association for Computational Linguistics, 2024.
- Jinhao Li, Haopeng Li, Sarah Monazam Erfani, Lei Feng, James Bailey, and Feng Liu. Visual-text cross alignment: Refining the similarity score in vision-language models. In *Forty-first International Conference on Machine Learning*, 2024.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pp. 19730–19742. PMLR, 2023.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 34892–34916. Curran Associates, Inc., 2023a.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023b.
- Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng Xi, Shiyi Cao, Yuxian Gu, Dacheng Li, Xiuyu Li, Yunhao Fang, Yukang Chen, Cheng-Yu Hsieh, De-An Huang, An-Chieh Cheng, Vishwesh Nath, Jinyi Hu, Sifei Liu, Ranjay Krishna, Daguang Xu, Xiaolong Wang, Pavlo Molchanov, Jan Kautz, Hongxu Yin, Song Han, and Yao Lu. Nvlla: Efficient frontier visual language models, 2024.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. Deepseek-vl: Towards real-world vision-language understanding, 2024.
- Mayug Maniparambil, Raiymbek Akshulakov, Yasser Abdelaziz Dahou Djilali, Mohamed El Amine Seddik, Sanath Narayan, Karttikeya Mangalam, and Noel E. O’Connor. Do vision and language encoders represent the world similarly? In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14334–14343, 2024.

- Bhaskar Mitra and Nick Craswell. Query auto-completion for rare prefixes. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, CIKM '15, pp. 1755–1758. Association for Computing Machinery, 2015. ISBN 9781450337946. doi: 10.1145/2806416.2806599.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, and Lama Ahmad et al. Gpt-4 technical report, 2024.
- Openverse. Openverse: Openly licensed images, audio and more. <https://openverse.org/>, 2025.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Christoph Schuhmann. Improved aesthetic predictor: Clip + mlp aesthetic score predictor. <https://github.com/christophschuhmann/improved-aesthetic-predictor>, 2022.
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv*, 2021.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: an open large-scale dataset for training next generation image-text models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2024.
- Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15638–15650, 2022.
- Andy Sun, Tianqi Zheng, Aakash Kolekar, Rohit Patki, Hossein Khazaei, Xuan Guo, George Cai, David Liu, Ruirui Li, Yupin Huang, Dante Everaert, Hanqing Lu, Garima Patel, and Monica Cheng. A product-aware query auto-completion framework for e-commerce search via retrieval-augmented generation method. *SIGIR 2024 Workshop on Information Retrieval's Role in RAG Systems (IR-RAG)*, 2024.
- Liang Wang, Nan Yang, and Furu Wei. Query2doc: Query expansion with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 9414–9423. Association for Computational Linguistics, December 2023.
- Xiyao Wang, Jiahai Chen, Zhaoyang Wang, Yuhang Zhou, Yiyang Zhou, Huaxiu Yao, Tianyi Zhou, Tom Goldstein, Parminder Bhatia, Furong Huang, et al. Enhancing visual-language modality alignment in large vision language models via self-improvement. *arXiv preprint*, 2024.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint*, 2024.
- Jinyu Yang, Jiali Duan, Son Tran, Yi Xu, Sampath Chanda, Liqun Chen, Belinda Zeng, Trishul Chilimbi, and Junzhou Huang. Vision-language pre-training with triple contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15671–15680, 2022.

- Ran Yi, Haoyuan Tian, Zhihao Gu, Yu-Kun Lai, and Paul L Rosin. Towards artistic image aesthetics assessment: a large-scale dataset and a new method. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22388–22397, 2023.
- Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint*, 2022.
- Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models, 2024.
- Hongyi Zhu, Jia-Hong Huang, Stevan Rudinac, and Evangelos Kanoulas. Enhancing interactive image retrieval with query rewriting using large language models and vision language models. In *Proceedings of the 2024 International Conference on Multimedia Retrieval*, pp. 978–987, 2024.

A APPENDIX

A.1 RELATED WORK

A.1.1 VISION-LANGUAGE MODELS

Vision-language models (VLMs) have become the de facto foundation for image-text tasks, demonstrating exceptional potential across a variety of applications (Alayrac et al., 2022; Liu et al., 2024; Zhang et al., 2024; Maniparambil et al., 2024). Pioneering work such as CLIP (Radford et al., 2021) and ALIGN (Jia et al., 2021) learn directly from raw texts about images by aligning them in a shared embedding space. CoCa (Yu et al., 2022) combines contrastive loss with captioning loss to train an image-text encoder-decoder model, effectively integrating capabilities from both contrastive and generative approaches. Blip2 (Li et al., 2023) bridges the modality gap with a lightweight Q-Former to improve pretraining efficiency. In this paper, we adopt joint-embedding VLMs like CLIP as foundation models for text-to-image retrieval. Instead of fine-tuning the VLMs on target datasets, we keep them frozen and focus on refining textual queries to achieve both quality improvement and control over the retrieved images. Improving existing VLMs for retrieval quality control is orthogonal to our approach and represents a promising direction for our future research.

A.1.2 LARGE LANGUAGE MODELS

Large language models (LLMs) are a class of foundation models designed to process, understand, and generate natural language at scale (Devlin, 2018; Radford et al., 2019; Brown et al., 2020). With fine-tuning and prompting, these models excel across a variety of tasks, including text generation, summarization, reasoning, translation, and coding (Liu et al., 2023b; Zhao et al., 2024). Notable examples, such as LLaMA3 (Grattafiori et al., 2024), GPT-4o (OpenAI et al., 2024), and Qwen2.5 (Yang et al., 2024), contain billions of parameters and are trained on extensive textual datasets. The large-scale pretraining enables them to capture complex contextual, semantic, and syntactic relationships in natural language. To tackle the proposed QCRR task, we utilize pretrained LLMs for query modification. By integrating quality information as conditions, the LLMs autonomously learn to generate quality-aware details for query extension. This provides users with multiple visible query suggestions, allowing them to explore diverse retrieval results.

A.1.3 TEXT-TO-IMAGE RETRIEVAL

Text-to-image retrieval aims to identify the most relevant images from a database given a natural language query. It plays a critical role in applications such as visual search, e-commerce, and content-based recommendation. Recent advances in VLMs (Radford et al., 2021; Gao et al., 2022; Li et al., 2022; Yu et al., 2022; Jia et al., 2021) have significantly improved performance on this task by learning powerful cross-modal representations. These models map images and texts into a shared embedding space, typically through contrastive learning on web-scaled image-text pairs. However, existing retrieval systems are primarily optimized to return the top-k images that are semantically aligned with the input query. They overlook other crucial dimensions—such as aesthetic appeal, interestingness, or popularity—that strongly affect user satisfaction in practical scenarios. In this work, we advocate for incorporating quality control into retrieval. By allowing users to explicitly influence the quality attributes of the returned results, we enable a more personalized and controllable search experience, moving beyond simple semantic matching toward a more adaptive, user-centric paradigm.

A.1.4 QUERY COMPLETION

Query completion (QC) aims to extend user short inputs, referred to as *query prefixes*, by generating longer and more informative *query completions*. It is a widely used technique that helps users better articulate their intent and resolve potential query ambiguity. Traditional QC methods rely on factors such as user profiles, query libraries, and prior search history to extend prefixes into query completions, which limits their applicability to unforeseen prefixes (Bar-Yossef & Kraus, 2011; Mitra & Craswell, 2015; Cai et al., 2016). Recently, several generative approaches have been proposed for query completion with arbitrary prefixes, primarily for text generation and document retrieval tasks (Lee et al., 2021; Wang et al., 2023; Lei et al., 2024). QC-based methods for text-to-image retrieval remain scarce, with only a few related works. Zhu *et al.* (Zhu et al., 2024) enhance

interactive image retrieval through query rewriting based on user relevance feedback, while Sun *et al.* (Sun et al., 2024) leverage LLMs to generate product-aware query completions. However, these approaches primarily focus on query suggestion and refinement rather than achieving control over the quality of retrieved images. In contrast, we tailor query completion to enhance retrieval quality, making the first attempt to adapt it to a given search corpus for quality-controllable retrieval.

A.2 ANALYSIS OF DATASET DEPENDENCY

Our work focuses on text-to-image retrieval, where the goal is to retrieve relevant images from a fixed dataset based on a textual query. This task is inherently dataset-dependent, as the retrieval process relies entirely on the available images within the dataset. Therefore, the query is crucial in this task: the more specific and detailed the query, the easier the retrieval system can match it to the corresponding image. Conversely, short or vague queries make it significantly more difficult for the system to identify the intended image. That’s why our proposed query completion method aims to enrich the original short queries with more specific, quality-aware details. We’d like to emphasize that these details are not randomly generated. Instead, they are learned directly from the dataset itself (by fine-tuning the LLM to fit the captions). As a result, the completed queries remain dataset-dependent and contextually relevant. The additional details are not unnecessary, as they provide essential guidance to the retrieval system, helping it to more accurately identify images with desired quality (as demonstrated by our experimental results).

A.3 ANALYSIS OF SCORE DIFFERENCES

In Tables 1 and 2, the quality scores across low, medium, and high conditions may appear close for a single query. This is expected due to dataset limitations. As shown in Figure 2, the similarity scores across both Flickr2.4M and MS-COCO are not uniformly distributed, and images with extremely low or high scores are rare. For instance, on Flickr2.4M, the similarity scores range from 0.252 to 0.562, and the entire span across the whole dataset is only about 0.3 (where the extreme values correspond to two images from different classes). When retrieving images for a single query, the available results often fall within a narrower score range (much smaller than 0.3) because the dataset lacks images at both ends of the quality spectrum (sparsely distributed). For example, if all images retrieved from the query “a dog” have aesthetic scores between 3.8 and 4.7 (due to dataset limitations), even under the “High” condition, the best available image might score 4.7—which lies in a low range (given the whole range : [2.782, 6.961]). But it is still higher than the score of 3.8 under the “Low” condition. Thus, the method is still effectively ranking and retrieving better images within the constraints of the dataset.

Despite this dataset-level constraint that limits the score differences, our method demonstrates effective ranking ability and a consistent, meaningful trend. As shown from left to right in Tables 1 and 2, both the retrieved image scores and their visual appeal improve progressively as the quality condition increases. This pattern is further supported by quantitative results in Tables 3 and 4, where the average quality scores clearly increase across the low, median, and high conditions. This behavior cannot be reproduced by baseline methods that lack quality consideration in retrieval.

A.4 PROOF OF PROPOSITION 1

Lemma 1. *If $\text{rank}(\mathbf{X}_I) = r$ and $\mathbf{I}_r + \mathbf{X}_I^\dagger \mathbf{P} \mathbf{Y}_I$ is invertible, then*

$$\text{rank}(\mathbf{X}_I + \mathbf{P} \mathbf{Y}_I) = r.$$

Proof of Lemma 1. Since $\text{col}(\mathbf{X}_I) = \mathcal{U}$, we have $\mathbf{P} = \mathbf{X}_I \mathbf{X}_I^\dagger$. Hence

$$\mathbf{X}_I + \mathbf{P} \mathbf{Y}_I = \mathbf{X}_I + \mathbf{X}_I \mathbf{X}_I^\dagger \mathbf{P} \mathbf{Y}_I = \mathbf{X}_I (\mathbf{I}_r + \mathbf{X}_I^\dagger \mathbf{P} \mathbf{Y}_I).$$

As \mathbf{X}_I has rank r and the factor in parentheses is invertible, the product has rank r . \square

Proof of Proposition 1. Since right-multiplication by the orthogonal matrix $V = [V_S \ V_\perp]$ is rank-preserving, we analyze the following matrices:

$$\begin{aligned} A' &:= AV = A[V_S \ V_\perp] = [AV_S \ AV_\perp] = [A_S \ 0], \\ \Delta' &:= \Delta V = \Delta[V_S \ V_\perp] = [\Delta V_S \ \Delta V_\perp] = [\Delta_S \ \Delta_\perp], \\ B' &:= BV = (A + \Delta)V = AV + \Delta V = A' + \Delta' = [A_S + \Delta_S \ \Delta_\perp], \\ C' &:= CV = C[V_S \ V_\perp] = [CV_S \ CV_\perp] = [C_S \ C_\perp]. \end{aligned} \quad (6)$$

Then, for the score matrices, we have:

$$S_A = AC^\top = (AV)(CV)^\top = A'C'^\top = [A_S \ 0] \begin{bmatrix} C_S^\top \\ C_\perp^\top \end{bmatrix} = A_S C_S^\top, \quad (7)$$

$$S_B = BC^\top = B'C'^\top = [A_S + \Delta_S \ \Delta_\perp] \begin{bmatrix} C_S^\top \\ C_\perp^\top \end{bmatrix} = (A_S + \Delta_S)C_S^\top + \Delta_\perp C_\perp^\top =: X + Y.$$

By the SVD construction, A_S has full column rank r and $\sigma_{\min}(A_S) = \sigma_r(A) > 0$. Since $\Delta_S = \Delta V_S$ and V_S is orthogonal (i.e., $\|V_S\|_2 = 1$), it follows that

$$\|\Delta_S\|_2 \leq \|\Delta\|_2. \quad (8)$$

Given that $\|\Delta_S\|_2 \leq \|\Delta\|_2 < \sigma_r(A) = \sigma_{\min}(A_S)$, the standard minimum-singular-value perturbation argument (or Weyl's inequality in spectral norm form) yields that $A_S + \Delta_S$ remains full column rank r . Since left multiplication by a full-column-rank matrix does not change rank, it follows that:

$$\begin{aligned} \text{rank}(X) &= \text{rank}((A_S + \Delta_S)C_S^\top) = \text{rank}(C_S^\top) = \text{rank}(C_S), \\ \text{rank}(S_A) &= \text{rank}(A_S C_S^\top) = \text{rank}(C_S) = \text{rank}(X). \end{aligned} \quad (9)$$

Consider the linear operator

$$T = \begin{bmatrix} P \\ (I - P_{Z_I})(I - P) \end{bmatrix}. \quad (10)$$

Since left multiplication cannot increase rank,

$$\text{rank}(S_B) \geq \text{rank}(TS_B) \geq \text{rank}((TS_B)_{:,I \cup K}). \quad (11)$$

Now

$$TS_B = \begin{bmatrix} X + PY \\ (I - P_{Z_I})Z \end{bmatrix}. \quad (12)$$

Restricting to $I \cup K$ gives the block form

$$(TS_B)_{:,I \cup K} = \begin{bmatrix} X_I + PY_I & X_K + PY_K \\ \mathbf{0} & (I - P_{Z_I})Z_K \end{bmatrix}. \quad (13)$$

By the lemma, the top-left block has rank r . By assumption (4), the bottom-right block has rank $k \geq 1$. Thus block-triangular rank additivity yields

$$\text{rank}((TS_B)_{:,I \cup K}) \geq r + k. \quad (14)$$

Therefore

$$\text{rank}(S_B) \geq r + k > r = \text{rank}(S_A). \quad (15)$$

□

A.5 ADDITIONAL EXPERIMENTAL RESULTS

Recall metrics such as R@1, R@5, and R@10 are standard in retrieval evaluation. However, it's important to note that recall is also derived from similarity—that is, images are ranked by similarity, and recall is computed based on their rank positions. Thus, recall metrics and similarity scores are inherently connected, especially when comparing methods built on the same retrieval backbone. To provide a complementary view of effectiveness, we conduct additional experiments on MS-COCO using R@1, R@5, and R@10 for evaluation. The results are shown in Table 10.

In addition, Table 8 presents the quantitative results on MS-COCO datasets using GPT2 as the backbone. Tables 7-9 provide more qualitative results on the two datasets.

Table 7: Retrieval quality of various methods on Flickr2.4M. CoCa and Blip2 are used to generate textual descriptions; **L** (Low), **M** (Median), and **H** (High) indicate the quality conditions; and Ctrl specifies whether the method enables controllable retrieval over quality. For both average relevance (Ave Rel) and average aesthetics (Ave Aes), higher values indicate better retrieval quality.

Quality	VLM	Aes Cond Rel Cond	L L	L M	L H	M L	M M	M H	H L	H M	H H	Ctrl ?
Prefix	--	Ave Aes Ave Rel	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	4.735 0.350	×
LLaMA3	--	Ave Aes Ave Rel	4.730 0.351	4.822 0.351	4.831 0.351	4.823 0.353	4.837 0.351	4.784 0.350	4.798 0.354	4.722 0.354	4.842 0.352	×
GPT-4o	--	Ave Aes Ave Rel	4.359 0.378	4.651 0.361	4.728 0.357	4.712 0.358	4.668 0.360	4.791 0.356	4.791 0.361	4.816 0.357	5.056 0.361	×
PT	--	Ave Aes Ave Rel	4.681 0.351	4.639 0.344	4.673 0.350	4.688 0.350	4.504 0.346	4.654 0.347	4.610 0.349	4.556 0.352	4.692 0.352	×
FT	CoCa	Ave Aes Ave Rel	4.848 0.366	4.818 0.365	4.864 0.367	4.847 0.365	4.827 0.363	4.876 0.366	4.829 0.367	4.896 0.366	4.853 0.368	×
Ours	CoCa	Ave Aes Ave Rel	4.646 0.354	4.674 0.372	4.632 0.382	4.878 0.355	4.921 0.369	4.894 0.386	5.182 0.357	5.095 0.366	5.124 0.385	✓
FT	Blip2	Ave Aes Ave Rel	4.838 0.369	4.674 0.360	4.744 0.369	4.592 0.365	4.599 0.362	4.772 0.365	4.727 0.373	4.749 0.359	4.818 0.368	×
Ours	Blip2	Ave Aes Ave Rel	4.528 0.355	4.560 0.374	4.470 0.393	4.948 0.354	4.946 0.374	4.885 0.391	5.266 0.354	5.160 0.367	5.236 0.387	✓

Table 8: Retrieval quality of various methods (GPT2) on MS-COCO, where **L** (Low), **M** (Median), and **H** (High) indicate the quality conditions for retrieval, and Ctrl specifies whether the method enables controllable retrieval over image quality. For both average relevance (Ave Rel) and average aesthetics (Ave Aes), higher values indicate better retrieval quality.

Quality	Aes Cond Rel Cond	L L	L M	L H	M L	M M	M H	H L	H M	H H	Ctrl ?
Prefix	Ave Aes Ave Rel	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	4.817 0.349	×
LLaMA3	Ave Aes Ave Rel	4.903 0.348	4.891 0.349	4.855 0.347	4.916 0.348	4.875 0.349	4.880 0.347	4.871 0.348	4.858 0.350	4.911 0.344	×
GPT-4o	Ave Aes Ave Rel	4.673 0.371	4.754 0.357	4.686 0.354	4.782 0.360	4.808 0.358	4.880 0.350	4.838 0.359	5.075 0.352	5.048 0.351	×
PT	Ave Aes Ave Rel	4.742 0.347	4.731 0.345	4.855 0.350	4.821 0.349	4.775 0.344	4.854 0.345	4.830 0.351	4.726 0.347	4.847 0.344	×
FT	Ave Aes Ave Rel	4.785 0.369	4.820 0.369	4.866 0.373	4.813 0.373	4.852 0.369	4.888 0.373	4.833 0.367	4.919 0.376	4.960 0.372	×
FT-CoCa	Ave Aes Ave Rel	4.890 0.347	4.889 0.348	4.793 0.356	4.885 0.346	4.939 0.349	4.903 0.352	4.950 0.347	5.004 0.349	4.898 0.351	×
FT-Blip2	Ave Aes Ave Rel	4.776 0.349	4.883 0.351	4.824 0.352	4.914 0.344	4.968 0.349	4.873 0.350	5.039 0.343	4.967 0.349	5.053 0.349	×
Ours	Ave Aes Ave Rel	4.896 0.354	4.809 0.365	4.719 0.385	4.973 0.356	4.879 0.368	4.916 0.387	5.017 0.353	5.020 0.368	5.109 0.391	✓

A.6 LIMITATION

In rare cases, the completed queries may not align with the semantics of the query prefixes. This occurs when the query completion model generates a sentence referencing different objects. Additionally, the relevance and aesthetic quality of the retrieved images depend on the reliability of the

Table 9: Retrieval quality with five quality levels on CoCa.

\mathcal{M}		Relevance (Red \rightarrow Red)				
		VL	L	M	H	VH
Aesthetics (Green \leftarrow Green)	VL	4.581	4.551	4.559	4.579	4.507
		0.355	0.364	0.372	0.376	0.382
	L	4.870	4.792	4.784	4.748	4.718
		0.357	0.363	0.370	0.376	0.383
	M	4.882	4.954	4.863	4.849	4.820
		0.356	0.366	0.371	0.377	0.381
	H	5.054	5.048	5.005	5.019	4.998
		0.355	0.362	0.371	0.370	0.381
	VH	5.159	5.166	5.161	5.135	5.084
		0.352	0.366	0.369	0.373	0.386

Table 10: Comparison with post-retrieval filtering

	R@1	R@5	R@10
Finetuned	0.8500	0.8875	0.9125
F-CoCa	0.8375	0.9375	0.9750
F-Blip2	0.7375	0.8750	0.9250
ours	0.8750	0.9625	0.9750

VLMs and aesthetic evaluation models. If these models are not sufficiently reliable, retrieval performance can be significantly affected. Refer to Table 14 for examples of such cases. As mentioned before, the model needs to perceive image quality within the datasets to achieve quality control in retrieval. However, the retrieval datasets may sometimes lack the granularity needed to differentiate between high-quality and low-quality images. In some instances, the retrieval database may not contain high-quality or low-quality images that match specific queries.

A.7 THE USE OF LARGE LANGUAGE MODELS (LLMs)

In preparing this paper, large language models were used only as writing assistants for grammar checking and minor sentence rephrasing. All technical aspects of the work, including the design, implementation, and verification of experiments and analyses, were carried out by the authors.

Table 11: Query completions with their retrieved images and quality scores on Flickr2.4M

Rel: Low , Aes: Low	Rel: Median , Aes: Median	Rel: High , Aes: High
		
<i>a bowl of soup with meat and vegetables in it</i> Aes 4.648, Rel 0.343	<i>a bowl on display</i> Aes 4.980, Rel 0.379	<i>a bowl with flowers on it</i> Aes 5.386, Rel 0.387
		
<i>a train on a track next to a grassy field</i> Aes 4.585, Rel 0.369	<i>a train station with people waiting to board a bus</i> Aes 4.910, Rel 0.380	<i>a train in the desert</i> Aes 5.488, Rel 0.391
		
<i>a horse drawn carriage on a dirt road</i> Aes 4.718, Rel 0.357	<i>a horse drawn carriage with people on it</i> Aes 5.023, Rel 0.3773	<i>a horse is grazing in a field under a cloudy sky</i> Aes 5.207, Rel 0.390
		
<i>a truck is parked in front of a building</i> Aes 4.800, Rel 0.331	<i>a truck is parked under a bridge</i> Aes 5.070, Rel 0.370	<i>a truck is parked in front of the washington monument</i> Aes 5.335, Rel 0.389
		
<i>a boat docked in the water next to other boats</i> Aes 3.814, Rel 0.358	<i>a boat is in the water near a castle</i> Aes 4.839, Rel 0.371	<i>a boat on the water with buildings in the background</i> Aes 5.113, Rel 0.396

Table 12: Query completions with their retrieved images and quality scores on MS-COCO

Rel: Low , Aes: Low	Rel: Median , Aes: Median	Rel: High , Aes: High
		
<i>an aeroplane flying in the air with a big blue sky behind it</i> Aes 4.536, Rel 0.354	<i>an aeroplane flying high on a clear sky</i> Aes 4.739, Rel 0.360	<i>Query: an aeroplane flying over the beach and two guys standing on it</i> Aes 5.516, Rel 0.425
		
<i>a fire hydrant stands in front of a bald eagle wall mural</i> Aes 4.991, Rel 0.355	<i>a fire hydrant sitting in front of a sign for a cafe</i> Aes 5.511, Rel 0.368	<i>a fire hydrant is painted to look like a dalmatian</i> Aes 5.799, Rel 0.447
		
<i>a toilet with a raised lid in some lavatory</i> Aes 4.457, Rel 0.361	<i>a toilet and sink in a small bathroom with a seat up</i> Aes 4.502, Rel 0.372	<i>a toilet is sitting outside with a sign on it</i> Aes 5.264, Rel 0.394
		
<i>a sheep is standing on a white fence</i> Aes 4.557, Rel 0.358	<i>a sheep and baby sheep standing in a field</i> Aes 5.014, Rel 0.378	<i>a sheep dog herding sheep through a grass field</i> Aes 5.213, Rel 0.388

Table 13: Query completions with their retrieved images and quality scores on MS-COCO









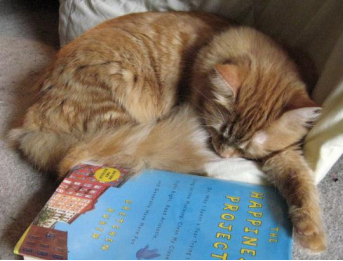


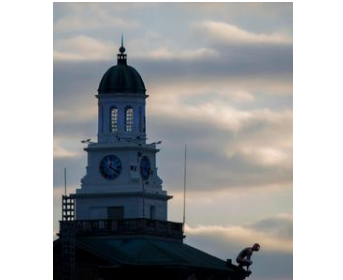

Rel: Low , Aes: Low	Rel: Median , Aes: Median	Rel: High , Aes: High
 <p><i>a wine glass next to a plate with some fish and veggies on it</i> Aes 4.790, Rel 0.350</p>	 <p><i>a wine glass next to a plate with some meat and vegetables on it</i> Aes 5.209, Rel 0.376</p>	 <p><i>a wine glass and three clocks all set at different times</i> Aes 5.555, Rel 0.417</p>
 <p><i>a toothbrush on a table with a bunch of scissors</i> Aes 4.149, Rel 0.345</p>	 <p><i>a toothbrush that is on down on the counter</i> Aes 4.837, Rel 0.370</p>	 <p><i>a toothbrush with a smiley face sitting on a sink</i> Aes 5.184, Rel 0.412</p>
 <p><i>a backpack and a line of supplies laying out</i> Aes 3.937, Rel 0.355</p>	 <p><i>a backpack some water rocks and plants</i> Aes 5.130, Rel 0.374</p>	 <p><i>a backpack with rollers is sitting unattended in the middle of this forested dirt road</i> Aes 5.231, Rel 0.470</p>

Table 14: Some bad retrieval cases on the two datasets

Rel: Low , Aes: Low	Rel: Median , Aes: Median	Rel: High , Aes: High
		
<i>a toilet sits next to a shower an sink</i> Aes 4.192, Rel 0.361	<i>a toilet with a wooden seat on top of it</i> Aes 5.226, Rel 0.371	<i>a toilet in between two trash cans</i> Aes 5.551, Rel 0.434
		
<i>an apple phone and some other type of machine</i> Aes 3.586, Rel 0.361	<i>an apple and other fruit are sitting together</i> Aes 5.007, Rel 0.373	<i>an apple with a knife stuck into it dripping blood</i> Aes 5.484, Rel 0.395
		
<i>an orange and blue bath-room with a tub sink and toilet</i> Aes 4.187, Rel 0.357	<i>an orange cat with its eyes closed sitting next to books</i> Aes 4.434, Rel 0.359	<i>an orange and black fire hy-drant sitting close to a curb</i> Aes 5.603, Rel 0.393
		
<i>a clock on the wall of a room</i> Aes 4.578, Rel 0.355	<i>a clock tower with a statue in front of it</i> Aes 5.187, Rel 0.376	<i>a clock tower with a cross on top</i> Aes 5.425, Rel 0.388