

Cropper: Vision-Language Model for Image Cropping through In-Context Learning

Seung Hyun Lee^{2,4*}, Jijun Jiang^{3*}, Yiran Xu^{1,5*}, Zhuofang Li^{3*}, Junjie Ke¹, Yinxiao Li¹, Junfeng He², Steven Hickson¹, Katie Datsenko¹, Sangpil Kim⁶, Ming-Hsuan Yang¹, Irfan Essa¹, Feng Yang^{1†}
Google DeepMind¹, Google Research², Google³,
University of Michigan⁴, University of Maryland⁵, Korea University⁶

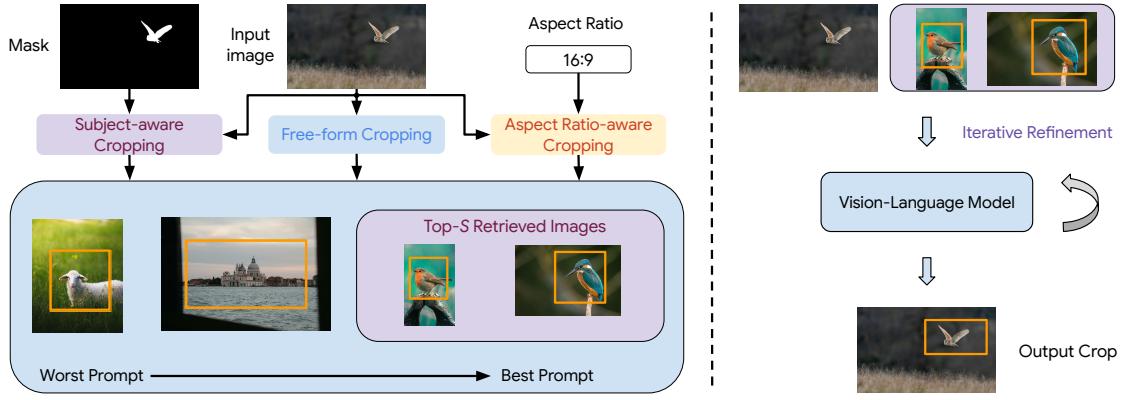


Figure 1. **Cropper** is a unified framework for various cropping tasks, including free-form cropping, subject-aware cropping, and aspect ratio-aware cropping built on top of a pretrained large vision-language model through in-context learning. Given the input image, top-K semantically similar images are retrieved as in-context learning prompt, and fed to pretrained vision-language model to generate crops. The crop candidates are iteratively refined to yield the visually pleasing output crop. All images are from Unsplash [30].

Abstract

The goal of image cropping is to identify visually appealing crops in an image. Conventional methods are trained on specific datasets and fail to adapt to new requirements. Recent breakthroughs in large vision-language models (VLMs) enable visual in-context learning without explicit training. However, downstream tasks with VLMs remain under explored. In this paper, we propose an effective approach to leverage VLMs for image cropping. First, we propose an efficient prompt retrieval mechanism for image cropping to automate the selection of in-context examples. Second, we introduce an iterative refinement strategy to iteratively enhance the predicted crops. The proposed framework, we refer to as **Cropper**, is applicable to a wide range of cropping tasks, including free-form cropping, subject-aware cropping, and aspect ratio-aware cropping. Extensive experiments demonstrate that **Cropper** significantly outperforms state-of-the-art methods across several benchmarks.

1. Introduction

Existing cropping methods [3–7, 16, 19, 24, 28, 29, 33, 34, 41, 42, 45] train neural networks on images and ground-truth crops to localize aesthetic crops. However, these approaches often depend on specially designed networks or features, which struggle to generalize effectively when confronted with new requirements or diverse datasets. Additionally, for specialized cropping tasks such as subject-aware cropping with subject masks or aspect ratio-aware cropping with target aspect ratio, unique networks must be developed and retrained, further complicating the process. This limitation underscores the need for more generalizable and versatile techniques in the field of image cropping.

Recent advancements in large vision-language models (VLM), such as GPT-4o [1] and Gemini [9], have unlocked new potential for various vision tasks. Unfortunately, in a lot of cases, users are not able to fine-tune the VLM for downstream tasks. Effectively adapting large black-box models for downstream tasks is very difficult. Luckily, in-context learning (ICL) ability is observed in large models [25]. Given a test instance and a few in-context example

* Equal contribution. † Work done while working at Google. ‡ Project lead.

demonstrations as input, the model directly infers the output without any parameter update or explicit training for the unseen task. ICL originates from natural language processing (NLP), and it has only recently been explored in the vision realm, mainly in image-to-image tasks [2, 32, 40, 44]. In this paper, we undertake an investigation aimed at harnessing the power of VLMs through ICL for image cropping, which, to our knowledge, has not been explored before.

Despite the remarkable capabilities of VLMs, challenges persist. First, the effectiveness of visual ICL heavily relies on the quality of the in-context examples (i.e. prompts)[40, 44]. Manual selection of these examples would be laborious and difficult to scale. Moreover, how to incorporate aesthetics in VLM for image cropping is not straightforward. Leveraging VLM in-context learning for image cropping presents a novel research area requiring effective strategies.

To address these challenges, we propose an effective framework to adapt VLM for image cropping through in-context learning, referred to as Cropper. It not only addresses the inherent challenges in traditional image cropping methods but also demonstrates versatility across various cropping tasks, including free-form cropping, subject-aware cropping, and aspect ratio-aware cropping. Illustrated in Fig. 1, our approach begins with an efficient prompt retrieval mechanism for image cropping tasks, automating the selection of relevant in-context examples to enhance efficiency. To further improve the performance, we introduce an iterative refinement strategy designed to enhance the quality of the predicted crops produced by VLM. To validate the efficacy of Cropper, we conduct extensive experiments on various benchmark datasets. Cropper significantly outperforms existing state-of-the-art methods across various performance metrics. Notably, with only a few in-context examples, Cropper achieves superior performance without the need for training. It also provides a unified framework for various cropping tasks. Our contributions are:

- We introduce a unified visual in-context-learning framework Cropper for image cropping tasks, including free-form, subject-aware, and aspect ratio-aware cropping.
- Our prompt retrieval strategy automates the effective selection of ICL examples for cropping tasks.
- The proposed iterative refinement strategy enables the model to progressively enhance the output crop.
- With a few in-context examples and no explicit training, Cropper surpasses the existing supervised learning methods across various benchmarks.

2. Related Work

Image Cropping. Image cropping is a critical operation for various photography-related applications. From the perspective of constraints, there are three commonly studied types of cropping problems. The first category is free-

form cropping, where the objective is to directly identify the best crop without imposing additional constraints. Numerous techniques have been explored to tackle this problem, including saliency maps [29], learning-based methods [4, 5, 7, 8, 10, 11, 14, 19, 20, 22, 29, 33, 37, 38], and reinforcement learning [16]. Another cropping task is subject-aware image cropping [12, 36], where an additional subject mask is provided to indicate the subject of interest. The third task is aspect ratio-aware cropping [18], where the crops are expected to adhere to a specified aspect ratio. Most existing image cropping approaches rely on training neural networks on specific datasets, requiring retraining to accommodate different data distributions and requirements. In contrast, our method requires only a few in-context examples and doesn't need explicit training. Moreover, none of the methods are flexible enough to handle all three cropping tasks in unified manner, while our method can do.

In-Context Learning. In-context learning is a recent paradigm originating from NLP, where large-scale models perform inference on unseen tasks by conditioning on a few in-context examples and the test instance. This paradigm is effective because users can directly adapt the model to different downstream tasks without the hassle of fine-tuning or changing the model parameters in any way. Numerous methods based on ICL have been developed for various tasks such as text classification [43] and machine translation [39]. ICL is still relatively new in computer vision, and most visual ICL works focus on using large-scale image-to-image vision models for tasks such as image inpainting [2] and segmentation [32, 40]. Empowered by recent breakthroughs in VLMs such as OpenAI GPT-4o [1] and Google Gemini [9], we investigate effective strategies for in-context learning for cropping for the first time.

Prompt Retrieval. NLP researchers have discovered that the selection and arrangement of in-context examples, also known as prompts, significantly impacts the output performance [21, 44]. These findings have sparked interest in prompt retrieval, where in-context learning examples are retrieved based on similarity metrics given a test instance. Liu et al. [44] have demonstrated success in selecting semantically similar in-context examples based on nearest neighbors measured by embeddings from a pretrained languages model. Rubin et al. [27] propose selecting examples using a supervised prompt retriever to maximize downstream performance. For visual ICL for image-to-image tasks, Zhang et al. [44] use CLIP [25]-based unsupervised embedding similarity measure, and demonstrate further improvement with a supervised prompt retrieval approach.

3. Method

Fig. 2 illustrates the structure of Cropper, which consists of two main steps: visual prompt retrieval on the left and feedback-based iterative refinement on the right. Given the

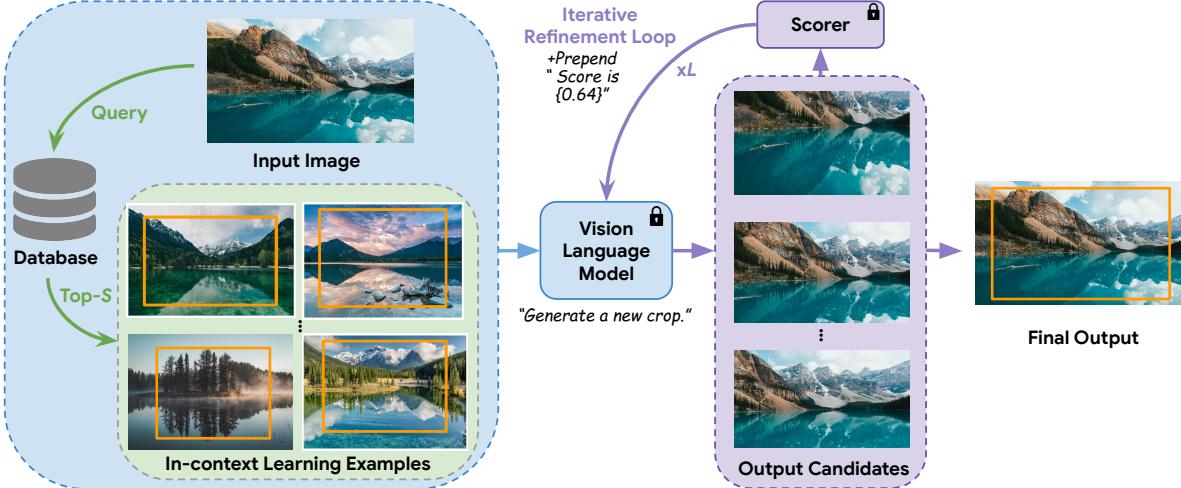


Figure 2. Cropper Overview. Cropper consists of two main steps: visual prompt retrieval and iterative crop refinement. Through visual prompt retrieval, top- S ICL examples are retrieved using an image similarity metric. In the iterative crop refinement stage, the VLM generates candidate crops based on these ICL examples and then these crops are subsequently scored by a scorer which measures aesthetics, content similarity, and area size. The VLM iteratively refines the crop candidates using the feedback from the scorer L times. All images are from Unsplash [30].

input image, Cropper automatically retrieves the top- S suitable in-context learning examples along with their ground-truth crop coordinates. Both the input image and the retrieved in-context learning examples are then fed into the vision-language model. The model is prompted to propose several potential crop candidates represented by their coordinates. In the iterative refinement stage, crops are generated based on the output of the VLM. These candidate crops are then evaluated using an aesthetic scorer, CLIP [25] scorer, and area size, which provides feedback guidance for the VLM. The model then iteratively refines the crop candidates based on this feedback. The iterative refinement process repeats L times to produce the final output.

3.1. Visual Prompt Retrieval for Cropping

The simplest method for retrieving ICL examples is random selection, where one or multiple samples are randomly chosen from the training dataset. However, previous studies have demonstrated that the ICL performance of such random selection is highly sensitive to the chosen samples [21, 44]. In our experiments (Sec. 4), we empirically confirm that random prompt selection often leads to sub-optimal results. Therefore, our objective is to explore an effective strategy for automatically selecting the most suitable examples for various cropping tasks.

Intuitively, similar images are more likely to be cropped similarly. Thus, we aim to retrieve the top- S images and their most relevant ground-truth crops based on some similarity metric. Formally, given an image query z_q and a dataset $\mathcal{D} = (z_i, C_i)_{i=1}^M$ containing M pairs of image z_i and crop ground-truth C_i , where C_i contains multiple crops c_1, \dots, c_s for some datasets, we seek to retrieve the

most relevant in-context examples and crop ground-truth as:

$$\mathcal{Z} = \arg \max_{z_i \in \mathcal{D}} Q(z_q, z_i), \quad |\mathcal{Z}| = S, \quad (1)$$

$$\mathcal{H} = \arg \max_{c_j \in C_j} G(z_q, c_j), \quad z_j \in \mathcal{Z}, \quad |\mathcal{H}| = T, \quad (2)$$

where \mathcal{Z} represents the set of top- S relevant images selected from the dataset \mathcal{D} based on the similarity metric $Q(z_q, z_i)$. $\mathcal{H} = (z_j, c_j)_{j=1}^S$ represents the selected in-context images z_j along with their most relevant T crop ground-truths based on metric $G(z_q, c_j)$. Q and G are designed differently to accommodate different cropping tasks, including free-form cropping, subject-aware cropping, and aspect ratio-aware cropping.

Free-form cropping aims to identify the best crop without additional constraints regarding aspect ratio or target subject. For this cropping task, we use the CLIP [25] image embeddings as an off-the-shelf image feature extractor, and Q corresponds to the cosine similarity between the input image z_q and each training example $z_i \in \mathcal{D}$. In free-form cropping datasets, such as GAICD [38], each image z_i is associated with multiple ground-truth crops c_i , each with its mean opinion score (MOS) aggregated from human evaluation. We use the MOS score as G for selecting the ground-truth crops. Therefore, after obtaining \mathcal{Z} , we select the top-ranked crops based on their MOS. Each crop ground-truth c_i is represented as a 5-tuple, (s, x_1, y_1, x_2, y_2) , indicating the MOS and the leftmost, top, rightmost, and bottom positions, respectively.

Subject-aware cropping intends to identify an aesthetic crop containing the subject of interest, which is represented

Prompt & Output	Instruction
Initial Prompt	<p>“Localize the aesthetic part of the image. (s, x_1, y_1, x_2, y_2) represents the region. x_1 and x_2 are the left and right most positions, normalized into 1 to 1000, where 1 is the left and 1000 is the right. y_1 and y_2 are the top and bottom positions, normalized into 1 to 1000 where 1 is the top and 1000 is the bottom. s is MOS score. We provide several images here.</p> <p>{image 1} $(s_1^1, x_1^{1,1}, y_1^{1,1}, x_2^{1,1}, y_2^{1,1})$, $(s_1^2, x_1^{1,2}, y_1^{1,2}, x_2^{1,2}, y_2^{1,2})$, ..., $(s_1^T, x_1^{1,T}, y_1^{1,T}, x_2^{1,T}, y_2^{1,T})$,</p> <p>{image 2} $(s_2^1, x_1^{2,1}, y_1^{2,1}, x_2^{2,1}, y_2^{2,1})$, $(s_2^2, x_1^{2,2}, y_1^{2,2}, x_2^{2,2}, y_2^{2,2})$, ..., $(s_2^T, x_1^{2,T}, y_1^{2,T}, x_2^{2,T}, y_2^{2,T})$,</p> <p>...</p> <p>{image S} $(s_S^1, x_1^{S,1}, y_1^{S,1}, x_2^{S,1}, y_2^{S,1})$, $(s_S^2, x_1^{S,2}, y_1^{S,2}, x_2^{S,2}, y_2^{S,2})$, ..., $(s_S^T, x_1^{S,T}, y_1^{S,T}, x_2^{S,T}, y_2^{S,T})$,</p> <p>{Query image}.</p>
Output	$(\hat{s}_1, \hat{x}_1^1, \hat{y}_1^1, \hat{x}_2^1, \hat{y}_2^1), (\hat{s}_2, \hat{x}_1^2, \hat{y}_1^2, \hat{x}_2^2, \hat{y}_2^2), \dots, (\hat{s}_R, \hat{x}_1^R, \hat{y}_1^R, \hat{x}_2^R, \hat{y}_2^R)$
Iterative Crop Refinement Prompt	<p>Initial Prompt + {Cropped image 1} $(\hat{s}_1, \hat{x}_1^1, \hat{y}_1^1, \hat{x}_2^1, \hat{y}_2^1)$, Score is {score 1}</p> <p>{Cropped image 2} $(\hat{s}_2, \hat{x}_1^2, \hat{y}_1^2, \hat{x}_2^2, \hat{y}_2^2)$, Score is {score 2}</p> <p>...</p> <p>{Cropped image R} $(\hat{s}_R, \hat{x}_1^R, \hat{y}_1^R, \hat{x}_2^R, \hat{y}_2^R)$, Score is {score R}</p> <p>Propose similar crop that has high score. The region should be represented by (s, x_1, y_1, x_2, y_2).</p>
Output	$(\hat{s}, \hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2)$

Table 1. VLM prompt used for free-form cropping. The goal is to find the most visual pleasing crop $(\hat{s}, \hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2)$. In the initial prompt, we use S in-context (ICL) examples, and T ground-truth crops. The format of image i ’s j -th crop is defined as $(s_i^j, x_1^{i,j}, y_1^{i,j}, x_2^{i,j}, y_2^{i,j})$. Intermediate results of initial prompt are coordinates of R crops. Subsequently, the crop is iteratively refined by accumulating the context into prompts, using refinement prompt. Note that {score} is calculated based on VILA [15] and area size.

by binary masks provided by users. In this task, the query image z_q is accompanied by a binary mask m_q indicating the subject of interest. Similarly, we first use CLIP image embedding similarity as Q for retrieving the top- S relevant images. Since each image in this task is associated with multiple target subject masks and their corresponding ground-truth crops, we further refine it by choosing the most similar mask areas to provide better guidance. G is defined as $-L_2$ distance between the center points of the target mask m_q and mask from image $z \in \mathcal{Z}$ to select the crop with closest masks. As a result the ground-truth crop for the closest mask is provided as the in-context learning example label, and each crop ground-truth c_i is represented (x_1, y_1, x_2, y_2) .

Aspect ratio-aware cropping requires the generated crop to conform to a specified aspect ratio r_q given the query image z_q . Each image in the dataset for this task is associated with ground-truth crops using different aspect ratios, such as 16:9, 3:4, and 1:1. Similarly, CLIP-based image similarity is adopted as Q . G is defined as the similarity between the crop c_i ’s aspect ratio and the target aspect ratio r_q . In other words, for each image $z \in \mathcal{Z}$, only the crop that has the similar target aspect ratio is used as in-context learning ground-truth. Similar to previous tasks, each crop ground-truth c_i is represented (x_1, y_1, x_2, y_2) .

3.2. Iterative Crop Refinement

Without explicit supervision, VLM lacks a deep understanding of the context of the cropping task, such as the provided coordinate system and intended aesthetics. Consequently, it often produces nonsensical outputs even when provided with good in-context learning cropping examples. Empirically, we observe that the initial crop candidates generated by the VLM lack diversity and sometimes fail to make sense (e.g. being too small or too large). Yang et al. [35] have shown that large language models can optimize

the output by iteratively incorporating feedback. Motivated by this, we propose an iterative crop refinement mechanism to further guide the VLM in generating high quality crops. Concretely, the VLM is prompted to generate R crop candidates based on the in-context learning examples retrieved using the method described in Sec. 3.1. Subsequently, we crop the image according to each cropping proposal and feed the cropped images into scorers, such as VILA [15] covering aesthetics, CLIP [25] measuring content preserving, and area size, to obtain corresponding scores. In the refinement phase, we iteratively provide such feedback to the VLM by scoring the crop candidates and prompting it to generate new candidates to improve the score. This iterative process is repeated L times to generate the final output. Tab. 1 shows the prompt design for free-form cropping with the two phases of Cropper. For subject-aware cropping and aspect ratio aware cropping, the only difference depends on whether MOS is predicted together.

4. Experimental Results

We first describe the implementation details and experimental setups before presenting quantitative and qualitative results with comparisons to SOTA as well as ablation studies.

4.1. Implementation Details

Dataset. We evaluate Cropper across three cropping benchmarks: GAICD [38] for free-form cropping, FCDB [4] for free-form cropping and aspect ratio-aware cropping, and SACD [36] for subject-aware cropping.

GAICD [38] dataset has 3,336 images, with 2,636 for training, 200 for validation, and 500 for testing, containing 288,069 densely annotated crops. For evaluation, we retrieve in-context learning examples from the GAICD [38] training set (the validation set is not used). Among the 90 annotations available for each retrieved image, we select the

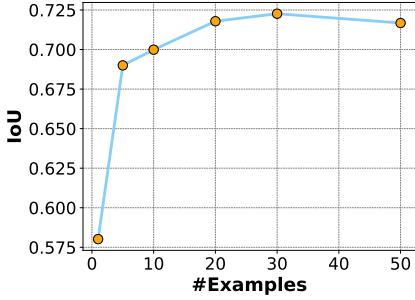


Figure 3. Relationship between number of in-context learning examples S and IoU on the GAICD [38] validation dataset for free-form cropping. We could see when the number of in-context learning examples S is 30, IoU is the best.

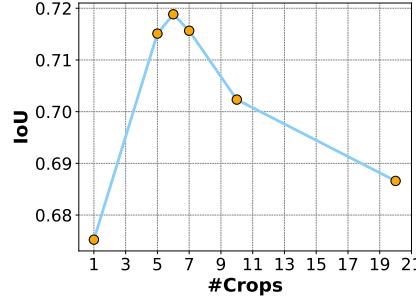


Figure 4. Relationship between number of crops R and IoU on the GAICD [38] validation dataset for free-form cropping. When the number of crops R is 6, IoU is the best.

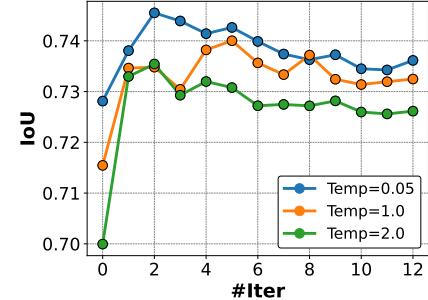


Figure 5. Relationship among the number of refinement iterations L , VLM model temperature and IoU on the GAICD [38] validation dataset. The experiments are based on the optimal number of ICL examples and candidate crops previously determined.

top- T crops ranked by MOS. Additionally, employing the same retrieval strategy on the GAICD [38] training set, we further evaluate the performance on 348 test images from the FCDB [4] dataset, measuring the out-of-domain performance. Following [18], we also use the FCDB [4] dataset to evaluate the performance of aspect ratio-aware cropping, treating the aspect ratio of the user-annotated box as the expected aspect ratio. The subject-aware cropping dataset SACD [36] contains 2,906 images, with 2,326 for training, 290 for validation, and 290 for testing.

Metrics. We use standard metrics in the image cropping community [38], including the Spearman’s rank-order correlation coefficient $SRCC$, the Pearson correlation PCC , and $Acc_{K/N}$, which is also used to evaluate image cropping methods [16] generating arbitrary bounding boxes. These metrics quantify the alignment of the generated crops with aesthetic preferences, using the ground-truth mean opinion score (MOS). Specifically, PCC assesses the linear correlation between the predicted MOS and the ground-truth MOS, whereas $SRCC$ measures the correlation of ranking order. Given that Cropper generates R candidate crops per each iteration step on the GAICD [38] dataset, we compute $SRCC$ and PCC using the best five crops instead of considering all crops. $Acc_{K/N}$ indicates whether top- K from predictions could be involved among top- N crops from the ground-truth based on MOS. $Acc_{1/5}$, $Acc_{2/5}$, $Acc_{3/5}$, $Acc_{4/5}$, $Acc_{1/10}$, $Acc_{2/10}$, $Acc_{3/10}$, $Acc_{4/10}$ are measured with $N \in \{5, 10\}$ and $K \in \{1, 2, 3, 4\}$ to return K of top- N accuracy. Additionally, we use Intersection-over-Union (IoU) and Boundary-displacement-error (Disp) metrics to compare with other approaches on the FCDB [4] and SACD [36] datasets. Disp represents the average L_1 distance between the ground-truth coordinates and the predicted values. To verify the effectiveness of the proposed method, we also conduct user study in Sec. 4.4.

Vision-language model. We adopt publicly available Gemini 1.5 Pro [26] model via the Vertex AI API for our task,

and experiment with GPT-4o [1]. This is because they support visual prompting with many images, which is critical for in-context learning for image cropping.

In-context learning examples. Similarity measurement Q in Eq. 1 is implemented using cosine similarity between image embeddings extracted from the ViT-B/32 variant of CLIP [25]. As shown in Fig. 3, we studied the relationship between the number of in-context examples S and IoU on the GAICD [38] validation dataset for free-form cropping and found the best number of in-context learning examples S is 30. We set S to 30 by default.

Number of crops. To determine the number of crops R , as in Fig. 4, we studied the relationship between the number of crops and IoU on the GAICD [38] validation dataset for free-form cropping and found the best number of crops R is 6. We set the number of crops R to be 6 by default.

Number of refinement iteration and VLM model temperature. Fig. 5 illustrates the effect of the number of refinement iterations L and VLM model temperature on the GAICD [38] validation dataset for free-form cropping. Since performance peaks after approximately two iterations, we set the number of iterations L to 2. The temperature value controls the randomness of VLM reasoning, with higher temperatures leading to more varied reasoning. We could see temperature 0.05 gives better IoU. So we set temperature to 0.05 for our experiments.

Scorer. For evaluating the aesthetics of each crop, we utilize the VILA-R [15] model as our aesthetic scorer. Trained on the AVA [23] dataset, this model specializes in image aesthetic assessment, providing evaluations based on factors such as perspectives, compositions, and color contrast. To measure the content preserving, we calculate the cosine similarity between image CLIP embeddings from cropped image and original image. We also consider the area size of the cropped region as one indicator for content preserving. The area scorer $A = \frac{H_{crop}W_{crop}}{HW} \in [0, 1]$, where H , W , H_{crop} , W_{crop} are the height and width of the input image

Model	$Acc_{1/5}$	$Acc_{2/5}$	$Acc_{3/5}$	$Acc_{4/5}$	\overline{Acc}_5	$Acc_{1/10}$	$Acc_{2/10}$	$Acc_{3/10}$	$Acc_{4/10}$	\overline{Acc}_{10}	$SRCC$	PCC
A2RL [16]	23.2	-	-	-	-	39.5	-	-	-	-	-	-
VPN [33]	36.0	-	-	-	-	48.5	-	-	-	-	-	-
VFN [5]	26.6	26.5	26.7	25.7	26.4	40.6	40.2	40.3	39.3	40.1	0.485	0.503
VEN [33]	37.5	35.0	35.3	34.2	35.5	50.5	49.2	48.4	46.4	48.6	0.616	0.662
GAIC [38]	68.2	64.3	61.3	58.5	63.1	84.4	82.7	80.7	78.7	81.6	0.849	0.874
CGS [19]	63.0	62.3	58.8	54.9	59.7	81.5	79.5	77.0	73.3	77.8	0.795	-
TransView [24]	69.0	66.9	61.9	57.8	63.9	85.4	84.1	81.3	78.6	82.4	0.857	0.880
Chao et al. [31]	70.0	66.9	62.5	59.8	64.8	86.8	84.5	82.9	79.8	83.3	0.872	0.893
Cropper (Ours)	88.9	85.9	83.1	79.4	84.3	98.2	97.2	96.4	94.3	96.5	0.904	0.860

Table 2. Quantitative comparison with existing free-form cropping methods on the GAICD [38] dataset. Cropper demonstrates significant superiority over other baselines despite using only a few in-context learning examples and no explicit training.

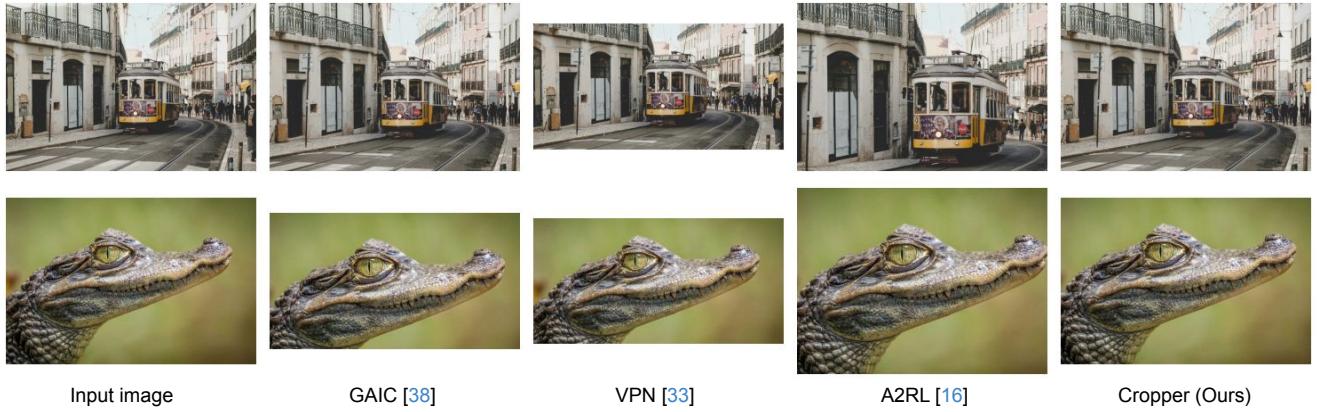


Figure 6. Qualitative comparing Cropper with GAIC [38], VPN [33], A2RL [16] on images from Unsplash [30] for free-form cropping.

and cropped image respectively. These scorers are normalized to the range of 0 to 1, and we use different combinations as final score.

4.2. Comparison with Baselines

Free-form image cropping. Tab. 2 presents a quantitative comparison between Cropper and other training-based baselines on the GAICD [38] dataset. Remarkably, Cropper outperforms training-based methods by a large margin with only a few in-context learning examples and no training.

Fig. 6 shows a visual comparison between Cropper and other free-form image cropping baselines, namely GAIC [38], VPN [33], A2RL [16]. The images are generated using the released codes of these methods. Overall, our approach produces more visually appealing results.

Subject-aware image cropping. Tab. 3 shows the quantitative comparison on the SACD [36] dataset, where Cropper surpasses all other training-based baselines. The reported numbers for other methods are directly taken from the baseline papers. To visually demonstrate the effectiveness of Cropper, we provide visual samples from Cropper in Fig. 7. We show the zero-shot inference results from GPT-4o [1] and Gemini 1.5 Pro [9]. GPT-4o is prompted with chain-of-thoughts to crop out the main subject within the image,

Model	Training-Free	IoU \uparrow	Disp \downarrow
A2RL [16]	✗	0.667	0.0887
VFN [5]	✗	0.669	0.0887
VPN [33]	✗	0.704	0.0699
VEN [33]	✗	0.691	0.0765
LVRN [20]	✗	0.696	0.0765
GAIC [38]	✗	0.712	0.0696
SAC-Net [36]	✗	0.767	0.0491
Cropper (Ours)	✓	0.769	0.0372

Table 3. Quantitative comparison on the SACD [36] dataset in subject-aware cropping task.

such as “*Think step-by-step about finding visually pleasing crops.*”. However, it struggles to generate good crops. For example, the crop from GPT-4o in the first row cuts the main subject, while our cropped image effectively captures the subject of interest.

Aspect ratio-aware image cropping. Tab. 4 shows quantitative comparison results on the FCDB [18] dataset for the aspect ratio-aware cropping task. Cropper outperforms other baselines in both IoU and Disp, indicating that Cropper is more adept at cropping the image according to the desired aspect ratio. Fig. 8 shows example crops from Cropper for different aspect ratio, illustrating qualitatively that Crop-

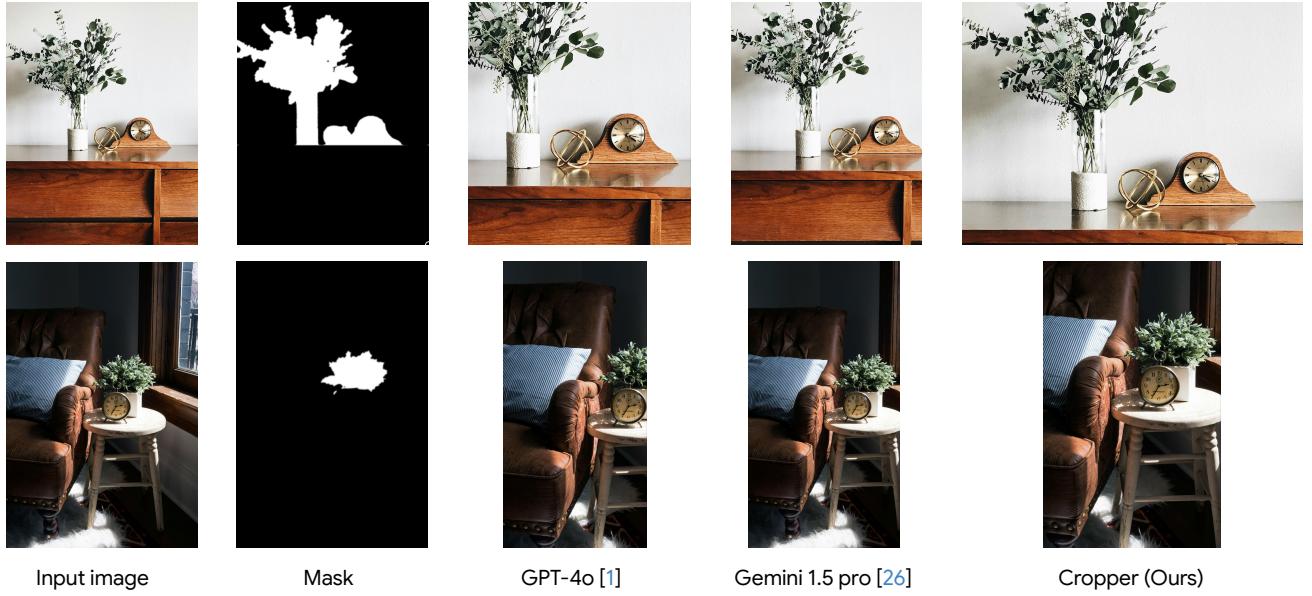


Figure 7. Visual comparisons on the subject-aware image cropping. Cropper preserves the important contents better than directly using VLMs, such as GPT-4o [1] and Gemini 1.5 Pro [26]. All input images are from Unsplash [30].

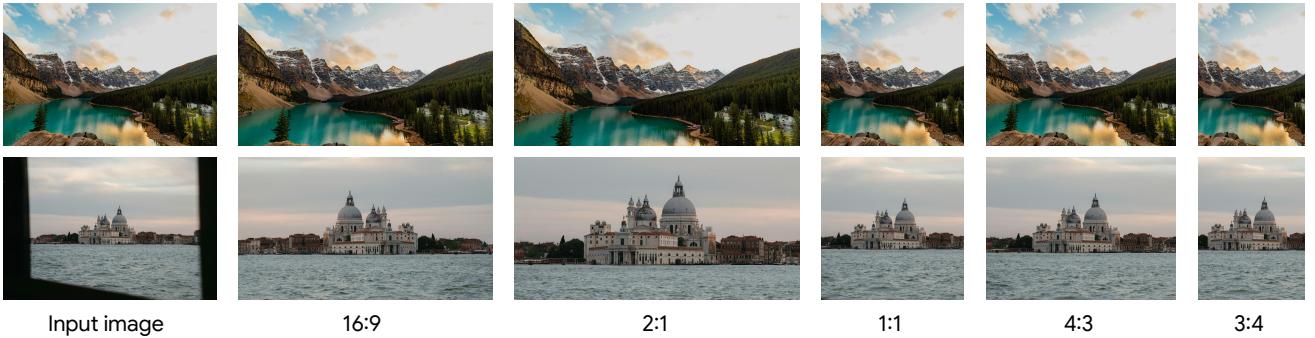


Figure 8. Example crops from Cropper for aspect ratio-aware cropping. This shows that our method can generate crops with the desired aspect ratios. All input images are from Unsplash [30].

per can generate crops that possess both good aesthetics and adhere to the specified aspect ratio.

4.3. Ablation study

Scorer. Tab. 5 presents ablation study for different scorers. We experiment with different combinations of VILA [15], normalized area score, and CLIP [25]. Both “VILA+Area” and “VILA+Area+CLIP” give the best tradeoff for different metrics. We choose “VILA+Area” for simplicity.

In-context learning. We compared the results of our method with another two different approaches: 1) zero-shot Gemini-1.5-pro [26]; 2) Cropper with random retrieval for in-context learning examples on the the GAICD test set [38] for free-form cropping in Tab. 6. Our approach achieves the best performance.

Iterative refinement. We show the ablation study on proposed iterative refinement for free-form cropping on both GAICD test set [38] and FCDB [18] dataset in Tab. 7. The

Model	Training-free	IoU \uparrow	Disp \downarrow
GAIC [38]	\times	0.673	0.064
A2RL [16]	\times	0.695	0.073
VPN [33]	\times	0.716	0.068
Mars [18]	\times	0.735	0.062
Cropper (Ours)	\checkmark	0.756	0.053

Table 4. Quantitative comparison on the FCDB dataset [4] for aspect ratio-aware cropping task. For other methods, we follow [18] to report the modified aspect ratio specified results, which are better than the ones in the original papers.

Scorer	Area	CLIP [25]	Metrics					
			IoU \uparrow	$Acc_5 \uparrow$	$Acc_{10} \uparrow$	$SRCC \uparrow$	$PCC \uparrow$	Avg \uparrow
✓	✗	✗	0.748	83.6	95.7	0.901	0.860	0.852
✗	✓	✗	0.752	83.9	96.0	0.882	0.838	0.854
✗	✗	✓	0.751	81.2	95.1	0.884	0.833	0.846
✓	✓	✗	0.748	84.3	96.5	0.904	0.860	0.864
✓	✗	✓	0.754	82.3	95.8	0.902	0.850	0.858
✗	✓	✓	0.752	83.9	96.1	0.902	0.858	0.862
✓	✓	✓	0.753	83.2	96.0	0.907	0.869	0.864

Table 5. Ablation study for different scorers on the GAICD [38] test dataset for free-form cropping.

Method	IoU \uparrow	Disp \downarrow
Zero-shot Gemini 1.5 Pro [26]	0.509	0.1385
Cropper with random retrieval	0.740	0.0660
Cropper with CLIP top-S	0.748	0.0635

Table 6. Ablation study for in-context learning. Comparing our methods with zero-shot Gemini 1.5 Pro [26], and Cropper with random retrieving in-context learning examples on the GAICD test set [38] for free-form cropping.

performance of our model improves significantly with the iterative refinement.

Dataset	Model	IoU \uparrow	Disp \downarrow
GAICD [38]	Cropper w/o Iter Refine.	0.722	0.0679
	Cropper (ours)	0.748	0.0635
FCDB [4]	Cropper w/o Iter Refine.	0.642	0.0925
	Cropper (ours)	0.667	0.0865

Table 7. Ablation study on iterative refinement.

Final output selection. We investigate the choice between selecting the output from the final iteration and selecting the output with the highest score across all iterations. We report the results in Tab. 8. For free-form cropping and subject-aware cropping, we find that using the prediction from the final iteration yields better performance, while aspect ratio-aware is different.

Method	Free-form		Subject		Aspect-ratio	
	IoU \uparrow	Disp \downarrow	IoU \uparrow	Disp \downarrow	IoU \uparrow	Disp \downarrow
Highest score across all iters.	0.714	0.0843	0.760	0.0381	0.756	0.0529
From final iter.	0.748	0.0635	0.769	0.0372	0.714	0.0632

Table 8. Comparison of selection strategies for free-form cropping on GAICD test set [38], subject-aware cropping on the SACD [36] dataset, and aspect-ratio aware cropping on the FCDB dataset [4].

Vision-language models. We first evaluate the robustness of various VLMs by performing free-form cropping on the GAICD [38] dataset using open-source models, such as Mantis-8B-Idefics2 [13], trained on the Mantis-Instruct [13] dataset to perform multi-image tasks. For Mantis-8B-Idefics2 [13], we use the same prompt as the free-form cropping task with slightly different parameters: 10 training examples, 5 output crops, and 2 iterations. Additionally, we assess Gemini-1.5-flash [26], a lighter variant of

Gemini-1.5-pro [26]. By using Gemini-1.5-flash [26], we reduce latency from 5.83 seconds (Gemini-1.5-pro [26]) to 2.65 seconds. As shown in Tab. 9, while lighter VLMs improve efficiency, we observe that larger model yields better cropping performance. Also better VLMs achieve better performance, i.e., results from Gemini is better than Mantis-8B-Idefics2 [13].

Model	$Acc_5 \uparrow$	$Acc_{10} \uparrow$	$SRCC \uparrow$	$PCC \uparrow$	IoU \uparrow	Avg \uparrow
Cropper with Mantis-8B-Idefics2 [13]	80.2	88.6	0.874	0.797	0.672	0.806
Cropper with Gemini-1.5-flash [26]	87.2	96.7	0.805	0.758	0.781	0.837
Cropper with Gemini-1.5-pro [26]	84.3	96.5	0.904	0.860	0.748	0.864

Table 9. Comparing different vision-language models.

4.4. User study

We conduct a user study on 200 test images from the GAICD dataset [38]. Our methods are compared against A2RL [16], GAIC [38], and CGS [19]. For each test image, we show the input, our result and the result from one of three methods, and ask the user to “*select one image that preserves the most important content from the source image.*”. For each test image, we asked five different users to provide ratings, resulting in a total of $5 \times 200 = 1000$ votes. We show results in Tab. 10. Cropper demonstrates a clear superiority over the baseline methods, outperforming them by a significant margin.

Choice	Baseline (%)	Cropper (%)
A2RL [16] v.s. Cropper	39.2	60.8
GAIC [38] v.s. Cropper	37.8	62.2
CGS [19] v.s. Cropper	36.6	63.4

Table 10. User study on 200 images from the GAICD test set [38]. We compare with A2RL [16], GAIC [38] and CGS [19]. Cropper shows a superiority over other baselines.

5. Conclusion & Limitation

The paper presents Cropper, a novel approach to image cropping that leverages in-context learning and vision-language models to achieve superior performance across various cropping tasks. It presents a novel training-free unified approach for tasks like free-form cropping, subject-aware cropping, and aspect ratio-aware cropping. Through extensive experimentation and comparison with existing baselines, Cropper demonstrates remarkable effectiveness and efficiency, outperforming counterparts with only a few in-context learning examples. Ablation studies show that the proposed visual prompt retrieval strategy and iterative crop refinement approach effectively harness the power of VLMs for effective ICL for cropping. Cropper has shown effectiveness in improving image cropping performance, but its inference speed constrained by finding suitable in-context examples in a dataset D . Advancements of visual retrieval will directly enhance the capabilities of Cropper.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. [1](#), [2](#), [5](#), [6](#), [7](#)
- [2] Amir Bar, Yossi Gandelsman, Trevor Darrell, Amir Globerson, and Alexei Efros. Visual prompting via image inpainting. In *NeurIPS*, 2022. [2](#)
- [3] Jiansheng Chen, Gaocheng Bai, Shaoheng Liang, and Zhengqin Li. Automatic image cropping: A computational complexity study. In *CVPR*, 2016. [1](#)
- [4] Yi-Ling Chen, Tzu-Wei Huang, Kai-Han Chang, Yu-Chen Tsai, Hwann-Tzong Chen, and Bing-Yu Chen. Quantitative analysis of automatic image cropping algorithms: A dataset and comparative study. In *WACV*, 2017. [2](#), [4](#), [5](#), [7](#), [8](#), [11](#)
- [5] Yi-Ling Chen, Jan Klopp, Min Sun, Shao-Yi Chien, and Kwan-Liu Ma. Learning to compose with professional photographs on the web. In *MM*, 2017. [2](#), [6](#)
- [6] Bin Cheng, Bingbing Ni, Shuicheng Yan, and Qi Tian. Learning to photograph. In *MM*, 2010.
- [7] Yang Cheng, Qian Lin, and Jan P Allebach. Re-compose the image by evaluating the crop on more than just a score. In *WACV*, 2022. [1](#), [2](#)
- [8] Yubin Deng, Chen Change Loy, and Xiaoou Tang. Aesthetic-driven image enhancement by adversarial learning. In *MM*, 2018. [2](#)
- [9] Gemini Team Google. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. [1](#), [2](#), [6](#)
- [10] Guanjun Guo, Hanzi Wang, Chunhua Shen, Yan Yan, and Hong-Yuan Mark Liao. Automatic image cropping for visual aesthetic enhancement using deep neural networks and cascaded regression. *IEEE Trans. Multimedia*, 20(8), 2018. [2](#)
- [11] Chaoyi Hong, Shuaiyuan Du, Ke Xian, Hao Lu, Zhiguo Cao, and Weicai Zhong. Composing photos like a photographer. In *CVPR*, 2021. [2](#)
- [12] James Hong, Lu Yuan, Michaël Gharbi, Matthew Fisher, and Kayvon Fatahalian. Learning subject-aware cropping by out-painting professional photos. In *AAAI*, 2024. [2](#)
- [13] Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhui Chen. Mantis: Interleaved multi-image instruction tuning. *TMLR*, 2024. [8](#)
- [14] Yueying Kao, Ran He, and Kaiqi Huang. Automatic image cropping with aesthetic map and gradient energy map. In *ICASSP*, 2017. [2](#)
- [15] Junjie Ke, Keren Ye, Jiahui Yu, Yonghui Wu, Peyman Milanfar, and Feng Yang. Vila: Learning image aesthetics from user comments with vision-language pretraining. In *CVPR*, 2023. [4](#), [5](#), [7](#), [8](#), [14](#)
- [16] Debang Li, Huikai Wu, Junge Zhang, and Kaiqi Huang. A2rl: Aesthetics aware reinforcement learning for image cropping. In *CVPR*, 2018. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#), [11](#), [13](#)
- [17] Debang Li, Huikai Wu, Junge Zhang, and Kaiqi Huang. Fast a3rl: Aesthetics-aware adversarial reinforcement learning for image cropping. *TIP*, 28(10), 2019. [11](#)
- [18] Debang Li, Junge Zhang, and Kaiqi Huang. Learning to learn cropping models for different aspect ratio requirements. In *CVPR*, 2020. [2](#), [5](#), [6](#), [7](#)
- [19] Debang Li, Junge Zhang, Kaiqi Huang, and Ming-Hsuan Yang. Composing good shots by exploiting mutual relations. In *CVPR*, 2020. [1](#), [2](#), [6](#), [8](#), [11](#)
- [20] Weirui Lu, Xiaofen Xing, Bolun Cai, and Xiangmin Xu. Listwise view ranking for image cropping. *IEEE Access*, 7, 2019. [2](#), [6](#)
- [21] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *ACL*, 2022. [2](#), [3](#)
- [22] Long Mai, Hailin Jin, and Feng Liu. Composition-preserving deep photo aesthetics assessment. In *CVPR*, 2016. [2](#)
- [23] Naila Murray, Luca Marchesotti, and Florent Perronnin. Ava: A large-scale database for aesthetic visual analysis. In *CVPR*, 2012. [5](#)
- [24] Zhiyu Pan, Zhiguo Cao, Kewei Wang, Hao Lu, and Weicai Zhong. Transview: Inside, outside, and across the cropping view boundaries. In *ICCV*, 2021. [1](#), [6](#), [11](#)
- [25] 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 *ICML*, 2021. [1](#), [2](#), [3](#), [4](#), [5](#), [7](#), [8](#), [14](#)
- [26] Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricu, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. [5](#), [7](#), [8](#), [11](#)
- [27] Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context learning. In *ACL*, 2022. [2](#)
- [28] Jin Sun and Haibin Ling. Scale and object aware image thumbnailing. *IJCV*, 104, 2013. [1](#)
- [29] Yi Tu, Li Niu, Weijie Zhao, Dawei Cheng, and Liqing Zhang. Image cropping with composition and saliency aware aesthetic score map. In *AAAI*, 2020. [1](#), [2](#), [11](#)
- [30] Unsplash Website. Unsplash. Accessed: March 21, 2025, URL: <https://unsplash.com/>. [1](#), [3](#), [6](#), [7](#), [12](#), [13](#)
- [31] Chao Wang, Li Niu, Bo Zhang, and Liqing Zhang. Image cropping with spatial-aware feature and rank consistency. In *CVPR*, 2023. [6](#), [11](#)
- [32] Xinlong Wang, Wen Wang, Yue Cao, Chunhua Shen, and Tiejun Huang. Images speak in images: A generalist painter for in-context visual learning. In *CVPR*, 2023. [2](#)
- [33] Zijun Wei, Jianming Zhang, Xiaohui Shen, Zhe Lin, Radomir Mech, Minh Hoai, and Dimitris Samaras. Good view hunting: Learning photo composition from dense view pairs. In *CVPR*, 2018. [1](#), [2](#), [6](#), [7](#), [11](#), [13](#)
- [34] Jianzhou Yan, Stephen Lin, Sing Bing Kang, and Xiaoou Tang. Learning the change for automatic image cropping. In *CVPR*, 2013. [1](#)
- [35] Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers. In *ICLR*, 2024. [4](#)

- [36] Guo-Ye Yang, Wen-Yang Zhou, Yun Cai, Song-Hai Zhang, and Fang-Lue Zhang. Focusing on your subject: Deep subject-aware image composition recommendation networks. *CVM*, 9(1), 2023. [2](#), [4](#), [5](#), [6](#), [8](#)
- [37] Hui Zeng, Lida Li, Zisheng Cao, and Lei Zhang. Reliable and efficient image cropping: A grid anchor based approach. In *CVPR*, 2019. [2](#)
- [38] Hui Zeng, Lida Li, Zisheng Cao, and Lei Zhang. Grid anchor based image cropping: A new benchmark and an efficient model. *TPAMI*, 44(3), 2020. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [11](#), [13](#)
- [39] Biao Zhang, Barry Haddow, and Alexandra Birch. Prompting large language model for machine translation: A case study. In *ICML*, 2023. [2](#)
- [40] Jiahao Zhang, Bowen Wang, Liangzhi Li, Yuta Nakashima, and Hajime Nagahara. Instruct me more! random prompting for visual in-context learning. In *WACV*, 2024. [2](#)
- [41] Luming Zhang, Mingli Song, Yi Yang, Qi Zhao, Chen Zhao, and Nicu Sebe. Weakly supervised photo cropping. *TMM*, 16(1), 2013. [1](#)
- [42] Mingju Zhang, Lei Zhang, Yanfeng Sun, Lin Feng, and Weiyi Ma. Auto cropping for digital photographs. In *ICME*, 2005. [1](#)
- [43] Yiming Zhang, Shi Feng, and Chenhao Tan. Active example selection for in-context learning. In *EMNLP*, 2022. [2](#)
- [44] Yuanhan Zhang, Kaiyang Zhou, and Ziwei Liu. What makes good examples for visual in-context learning? In *NeurIPS*, 2023. [2](#), [3](#)
- [45] Lei Zhong, Feng-Heng Li, Hao-Zhi Huang, Yong Zhang, Shao-Ping Lu, and Jue Wang. Aesthetic-guided outward image cropping. *TOG*, 40(6), 2021. [1](#)

Cropper: Vision-Language Model for Image Cropping through In-Context Learning

Supplementary Material

This supplementary material provides:

- Sec. A: We present the implementation details, additional comparison results, additional qualitative results.
- Sec. B: We present the implementation details, additional ablation study, and additional qualitative results for subject-aware cropping task.
- Sec. C: We present the implementation details, additional ablation study, and additional qualitative results for aspect ratio-aware cropping task.
- Sec. D: We present details about user study.

A. Free-form Cropping

A.1. Implementation details

We show the prompt for zero-shot cropping using Gemini 1.5 Pro [26] in Tab. 11.

Prompt & Output	Instruction
Initial Prompt	Localize the aesthetic part of the image. (x_1, y_1, x_2, y_2) represents the region. x_1 and x_2 are the left and right most positions, normalized into 1 to 1000, where 1 is the left and 1000 is the right. y_1 and y_2 are the top and bottom positions, normalized into 1 to 1000 where 1 is the top and 1000 is the bottom. Please propose a new region (x_1, y_1, x_2, y_2)
Output	$(\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2)$

Table 11. Prompt for zero-shot cropping with Gemini 1.5 Pro [26].

Tab. 12 shows further comparison on the FCDB [4] dataset for free-form cropping.

A.2. Additional qualitative results

Iterative update. We showcase some intermediate results of the iterative refinement in Fig. 9. Our method progressively refines the predicted crops, achieving increasing accuracy and better overlap with the ground-truth cropping box in each iteration.

Qualitative comparison. We present additional results in Fig. 10. Our method generates better visual pleasing crops.

Model	Training-Free	Training Set	IoU \uparrow	Disp \downarrow
A2RL [16]	✗	AVA	0.663	0.089
A3RL [17]	✗	AVA	0.696	0.077
VPN [33]	✗	CPC	0.711	0.073
VEN [33]	✗	CPC	0.735	0.072
ASM [29]	✗	CPC	0.749	0.068
GAICD [38]	✗	GAICD	0.672	0.084
CGS [19]	✗	GAICD	0.685	0.079
TransView [24]	✗	GAICD	0.682	0.080
Chao et al. [31]	✗	GAICD	0.695	0.075
Cropper (Ours)	✓	GAICD	0.667	0.087

Table 12. Quantitative comparison among different methods for free-form image cropping on the FCDB [4] dataset. Cropper shows competitive performance as a *training-free* approach.

B. Subject-aware Cropping

B.1. Prompts

We show the details of the prompts for the subject-aware cropping in Tab. 13. The goal is to get accurate coordinates of the crop $(\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2)$. In the initial prompt, we use 30 in-context (ICL) examples for image cropping for 10 iterations. 10 examples are ranked by the scorer and we use top-5 crops for our task, the format of image i 's j -th crop is defined as $(x_1^{i,j}, y_1^{i,j}, x_2^{i,j}, y_2^{i,j})$. Intermediate results of initial prompt are coordinates of 5 crops. Subsequently, the crop is iteratively refined by accumulating the context into prompts, using refinement prompt. Note that scorer is “VILA+Area”.

B.2. Ablation study of scores

We show the ablation study of scorer on the subject-aware cropping in Tab. 14. With “VILA+Area”, our proposed method achieves the best performance.

B.3. Additional qualitative results

We showcase more results in Fig. 11. Our method demonstrates subject awareness, enabling the generation of high-quality cropped images.

C. Aspect-ratio aware Cropping

C.1. Prompts

We show the details of the prompts for the aspect ratio-aware cropping in Tab. 15. For this task, we use the following hyperparameter: number of in-context learning

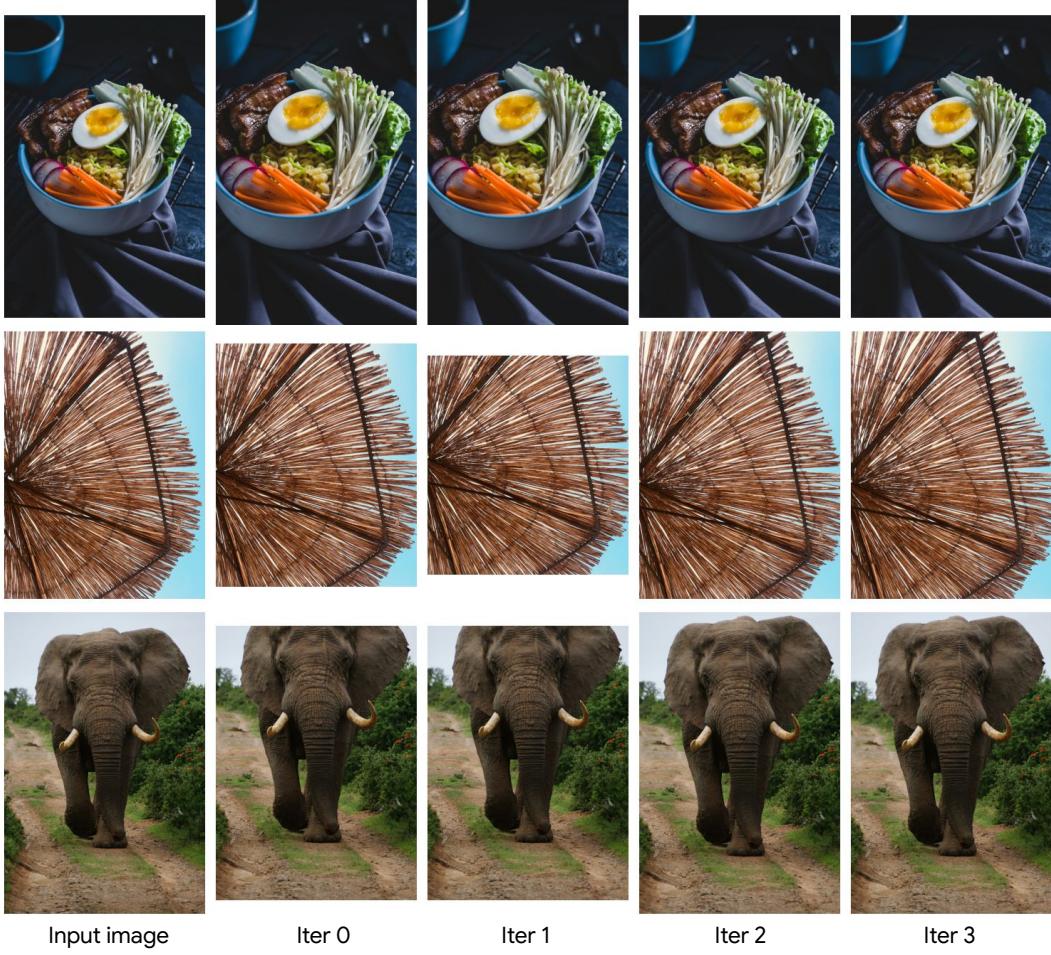


Figure 9. Results from each iteration for free-form cropping using Cropper. The iteration process demonstrated progressive convergence, resulting in improved crop quality. All input images are from Unsplash [30].

examples $S = 10$, number of crops $R = 6$, number of iteration $L = 2$, temperature = 0.05.

C.2. Ablation study of scores

We show the ablation study of scores on the aspect ratio-aware cropping in Tab. 16. With CLIP score only, our proposed method achieves the best performance.

D. User study

We include the instructions for users as follows:

- Your Task: Carefully analyze the source image and the two output images and SELECT one output.
- Content: This refers to the key elements and objects in the image, such as people, buildings, or other recognizable features. The output should keep the important details of these objects as close to the original as possible.
- Aesthetics: The output has a sense of aesthetics. It follows common human natures with proper layout. Select the output that not only preserves the content best and fits

the aesthetics.

- Your Goal: Select the image that, overall, looks the most natural and visually appealing to the source image.

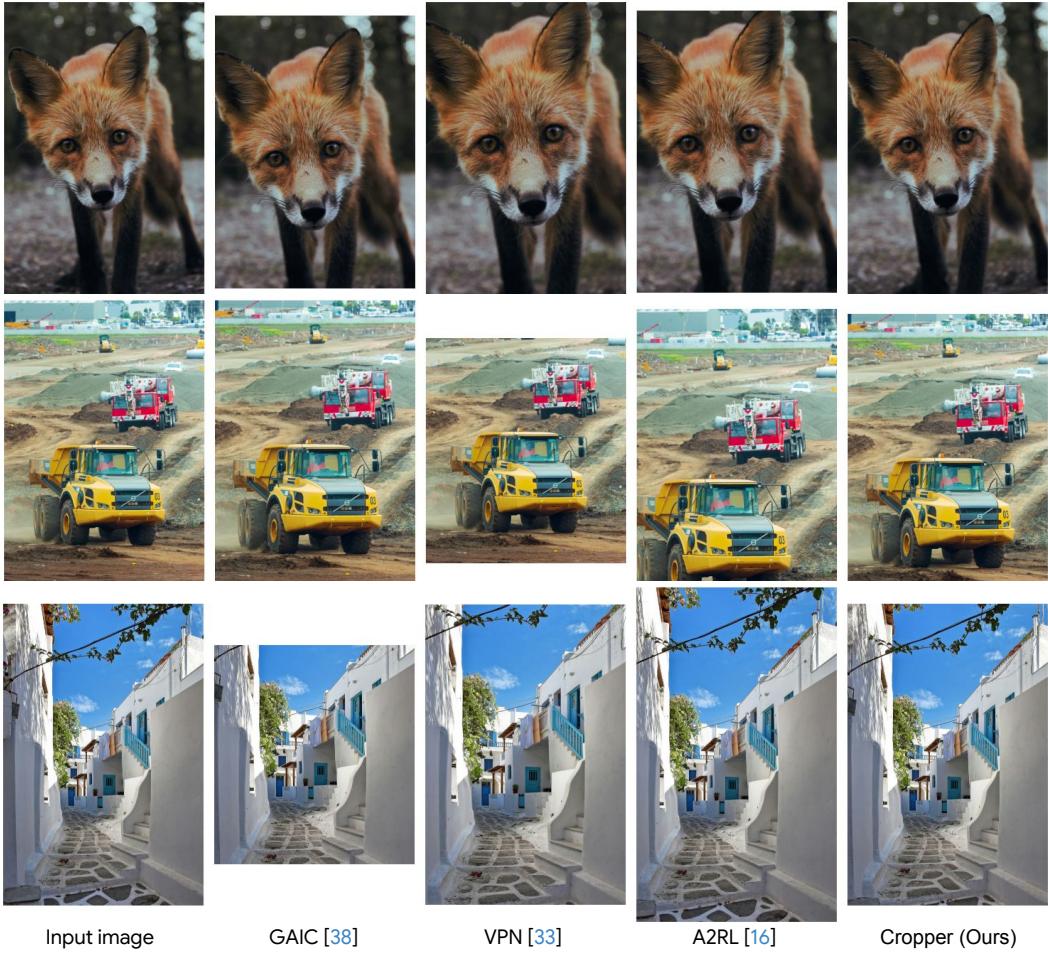


Figure 10. Comparing with GAIC [38], VPN [33] and A2RL [16] for free-form cropping on images from Unsplash [30].

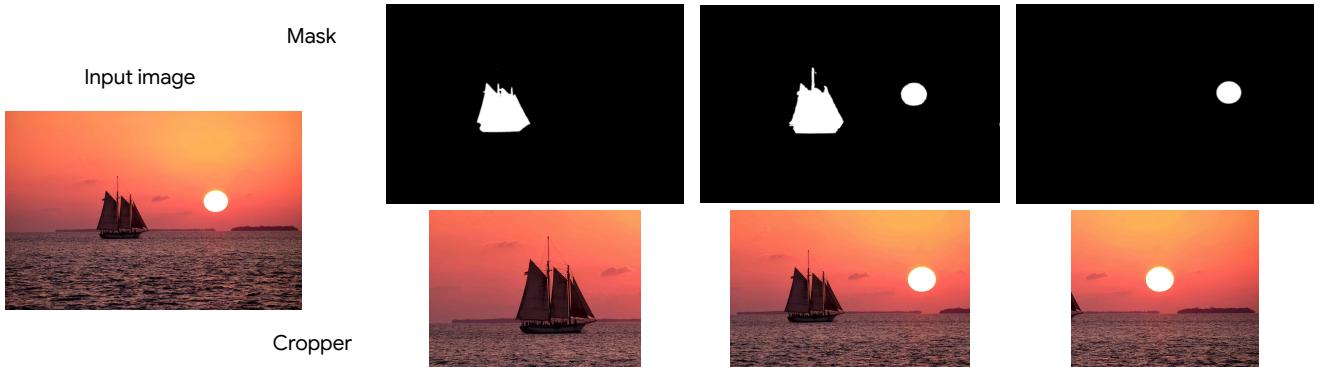


Figure 11. Qualitative results of subject-aware cropping. The result shows that our method can generate crops on different subjects. The input image is from Unsplash [30].

Prompt & Output	Instruction
Initial Prompt	Find visually appealing crop. Each region is represented by (x_1, y_1, x_2, y_2) coordinates. x_1, x_2 are the left and right most positions, normalized into 0 to 1, where 0 is the left and 1 is the right. y_1, y_2 are the top and bottom positions, normalized into 0 to 1 where 0 is the top and 1 is the bottom. $\{\text{image } 1\} ((c_x^1, c_y^1), x_1^1, y_1^1, x_2^1, y_2^1),$ $\{\text{image } 2\}, ((c_x^2, c_y^2), x_1^2, y_1^2, x_2^2, y_2^2),$ \dots $\{\text{image } S\}, ((c_x^S, c_y^S), x_1^S, y_1^S, x_2^S, y_2^S),$ $\{\text{Query image}\}, (c_x, c_y)$ $(\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2)$
Output	
Iterative Crop Refinement Prompt	Localize aesthetic part of image. The region is represented by (x_1, y_1, x_2, y_2) . x_1, x_2 are the left and right most positions, normalized into 0 to 1, where 0 is the left and 1 is the right. y_1, y_2 are the top and bottom positions, normalized into 0 to 1 where 0 is the top and 1 is the bottom. We provide several images here. $\{\text{Cropped image } 1\}$ Output: $(\hat{x}_1^1, \hat{y}_1^1, \hat{x}_2^1, \hat{y}_2^1)$ $\{\text{Cropped image } 2\}$ Output: $(\hat{x}_1^2, \hat{y}_1^2, \hat{x}_2^2, \hat{y}_2^2)$ \dots $\{\text{Cropped image } R\}$ Output: $(\hat{x}_1^R, \hat{y}_1^R, \hat{x}_2^R, \hat{y}_2^R)$ Propose different crop. The region should be represented by (x_1, y_1, x_2, y_2) . Output: $(\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2)$
Output	

Table 13. VLM prompt used for subject-aware cropping.

Prompt & Output	Instruction
Initial Prompt	Find visually appealing crop. Give the best crop in the form of a crop box and make sure the crop has certain width:height. Box is a 4-tuple defining the left, upper, right, and lower pixel coordinate in the form of (x_1, y_1, x_2, y_2) . Here are some example images, its size, and crop w:h triplets and their corresponding crops. $\{\text{image } 1\}$, size (w_1, h_1) , crop ratio (r_1) , output $(x_1^1, y_1^1, x_2^1, y_2^1)$, $\{\text{image } 2\}$, size (w_2, h_2) , crop ratio (r_2) , output $(x_1^2, y_1^2, x_2^2, y_2^2)$, \dots $\{\text{image } S\}$, size (w_S, h_S) , crop ratio (r_S) , output $(x_1^S, y_1^S, x_2^S, y_2^S)$, $\{\text{Now Give the best crop in the form of a crop box for the following image. Give } R \text{ possible best crops.}\}$ $\{\text{Query image}\}$, size (w, h) , crop ratio (r) $(\hat{x}_1^1, \hat{y}_1^1, \hat{x}_2^1, \hat{y}_2^1), (\hat{x}_1^2, \hat{y}_1^2, \hat{x}_2^2, \hat{y}_2^2), \dots, (\hat{x}_1^T, \hat{y}_1^T, \hat{x}_2^T, \hat{y}_2^T)$
Output	
Iterative Crop Refinement Prompt	Initial Prompt + Example Image: $\{\text{Query image}\}$; Crop ratio: r ; Example output: $\{\text{Cropped image } 1\} (\hat{x}_1^1, \hat{y}_1^1, \hat{x}_2^1, \hat{y}_2^1),$ $\{\text{Cropped image } 2\} (\hat{x}_1^2, \hat{y}_1^2, \hat{x}_2^2, \hat{y}_2^2)$ \dots $\{\text{Cropped image } R\} (\hat{x}_1^R, \hat{y}_1^R, \hat{x}_2^R, \hat{y}_2^R)$ Propose a different better crop with the given ratio. Output: $(\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2)$
Output	

Table 15. VLM prompt used for aspect ratio-aware cropping.

VILA [15]	Area	CLIP [25]	IoU \uparrow	Disp \downarrow
✓	✗	✗	0.753	0.0413
✗	✓	✗	0.755	0.0402
✗	✗	✓	0.749	0.0417
✓	✓	✗	0.769	0.0372
✓	✗	✓	0.751	0.0401
✗	✓	✓	0.754	0.0394
✓	✓	✓	0.766	0.0379

Table 14. Ablation study for scores on the subject-aware cropping. Cropper achieves the best performance with VILA [15] + Area score.

VILA [15]	Area	CLIP [25]	IoU \uparrow	Disp \downarrow
✓	✗	✗	0.718	0.0631
✗	✓	✗	0.713	0.0630
✗	✗	✓	0.756	0.0529
✓	✓	✗	0.716	0.0630
✗	✓	✓	0.741	0.0562
✓	✗	✓	0.742	0.0560
✓	✓	✓	0.729	0.0588

Table 16. Ablation study for aspect ratio-aware cropping task. Comparison of combinations of VILA [15], Area, and CLIP [25] components shows that the CLIP-only configuration achieves the best IoU and Disp values.