# MIXTURE-OF-QUERIES TRANSFORMER: CAMOU-FLAGED INSTANCE SEGMENTATION VIA QUERIES CO-OPERATION AND FREQUENCY ENHANCEMENT

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Due to the high similarity between camouflaged instances and the surroundings and the widespread camouflage-like scenarios, the recently proposed camouflaged instance segmentation (CIS) is a challenging and relevant task. Previous approaches achieve some progress on CIS, while many overlook camouflaged objects' color and contour nature and then decide on each candidate instinctively. In this paper, we contribute a Mixture-of-Queries Transformer (MoQT) in an end-toend manner for CIS which is based on two key designs (a Frequency Enhancement Feature Extractor and a Mixture-of-Queries Decoder). First, the Frequency Enhancement Feature Extractor is responsible for capturing the camouflaged clues in the frequency domain. To expose camouflaged instances, the extractor enhances the effectiveness of contour, eliminates the interference color, and obtains suitable features simultaneously. Second, a Mixture-of-Queries Decoder utilizes multiple experts of queries (several queries comprise an expert) for spotting camouflaged characteristics with cooperation. These experts collaborate to generate outputs, refined hierarchically to a fine-grained level for more accurate instance masks. Coupling these two components enables MoQT to use multiple experts to integrate effective clues of camouflaged objects in both spatial and frequency domains. Extensive experimental results demonstrate our MoQT outperforms 18 state-of-the-art CIS approaches by 2.69% on COD10K and 1.93% on NC4K in average precision.

031 032 033

034 035

006

008 009 010

011 012 013

014

015

016

017

018

019

021

023

025

026

027

028

029

### 1 INTRODUCTION

Camouflage is a naturally evolved strategy for animals to hide themselves via adapting their body's coloring to match the surroundings, which is used for hunting prey or avoiding detection by natural enemies, as shown in Figure 1(a). Since there is a lot of demand for understanding the widespread camouflage-like scenarios, (*e.g.*, polyp segmentation Fan et al. (2020b), lung infection segmentation Fan et al. (2020c), search-and-rescue work Turić et al. (2010), manipulated image/video detection and segmentation Zhou et al. (2020)), the task of predicting the location and instance-level masks of camouflaged objects (*i.e.*, Camouflaged Instance Segmentation, CIS) has been proposed. Therefore, CIS is worth studying and has gradually received more attention in recent years. However, it also has challenges due to high intrinsic similarities between the target objects and the background.

044 Compared to the tremendous development in generic instance segmentation Bolya et al. (2019); Wang et al. (2020a;b); Ren et al. (2015); He et al. (2017a); Cai & Vasconcelos (2018); Chen et al. 046 (2019), camouflaged instance segmentation remains an under-explored issue, and only a few efforts 047 have been made to study it in the past three years Pei et al. (2022); Luo et al. (2023); Dong et al. 048 (2023); Li et al. (2024); Le et al. (2023). CFL Le et al. (2021) is a first attempt. It is a two-stage method that fuses general instance segmentation methods for camouflaged instance segmentation but has limited performance. Subsequently, OSFormer Pei et al. (2022) is proposed as the first 051 one-stage method for CIS. It takes advantage of a transformer network, which achieves a flexible framework that can be trained end-to-end for camouflaged instance segmentation. Recently, DCNet 052 Luo et al. (2023) has been proposed to segment camouflaged instances via explicit de-camouflaging and achieves CIS by jointly modeling pixel-level camouflage decoupling and instance-level camou-

056

059

060

061

062

063

064

065

067

068 069



Figure 1: The Nature of Camouflaged Objects. Careful contrast of the camouflaged inputs (a) and the corresponding ground truth (b), color (c) (reconstructed with only amplitude component of Fourier transformation) and contour (d) (reconstructed with only phrase component of Fourier transformation) information shows the important priori principle of camouflaged objects: Low-level statistics like color contain more information from the surroundings while high-level semantics like contour tend to preserve more camouflaged characteristics.

flage suppression. In the same period, UQFormer Dong et al. (2023) adopts a typical DETR-like
 architecture Carion et al. (2020) and exploits the information on object edges.

073 Although the aforementioned works have made some progress, two fundamental inspirations have 074 not been taken into account: the priori of camouflage principles and the human habit of segmenting 075 camouflaged instances. (1) The priori of camouflage principles: Only when you know how to cam-076 ouflage can you see through camouflage. Many years ago, zoologists discovered that animals can 077 camouflage themselves by matching their colors or patterns with the background. To look deeper into camouflage, we sample some images of camouflaged animals and analyze them thoroughly. Since it is hard to spot camouflaged objects in the surroundings, we perform the Fourier transform 079 on these images to discover some clues in the frequency domain. We first decompose these camouflaged images into phrase and amplitude components and reconstruct images from only phrase 081 component and amplitude component, respectively, presented in Figure 1. It is easy to find that the phase component of the Fourier spectrum preserves high-level semantics (contours and semantics) 083 of original images, while the amplitude component contains low-level statistics (colors and styles). 084 Therefore, enhancing the influence of contours and eliminating the interference of colors would certainly benefit the performance on CIS. (2) The human habit of segmenting camouflaged instances: When humans segment a camouflaged image, their visual system instinctively sweeps across the 087 scene and determines some candidates. Then, the visual system gradually searches for valuable clues throughout the scene to obtain accurate segmentation masks. For some heavily camouflaged scenes with highly accurate segmentation like some medical image datasets Fu et al. (2019), it may even combine the masks labeled by multiple experts. Gradually refining and integrating the deci-090 sions of multiple experts are also potentially effective for CIS. Therefore, it makes sense to take full 091 advantage of both inspirations of (1) and (2) for improving the performance of the CIS task. 092

093 Motivated by the above discussions, we proposed a Mixture-of-Queries Transformer (MoOT) trained in an end-to-end manner for CIS, which includes a Frequency Enhancement Feature Extrac-094 tor (FEFE) based on modeling the colors and contours of camouflaged instances and a Mixture-of-095 Queries Decoder (MoQ Decoder) in transformer architecture referring to the segmentation process 096 of multi-experts collaboration. First, inspired by the camouflage principles discussed above, we design a Frequency Enhancement Feature Extractor to capture more clues of camouflaged instances in 098 the frequency domain. Specifically, we propose to adopt Fourier spectrum amplitude and phase to model image color information and contour information, respectively, as shown in Figure 1. With 100 the help of color and contour information, we design a contour enhancement module and a color 101 removal module, which can increase the contour effect while eliminating color interference. This 102 mechanism in the frequency domain is suitable for debunking the principle of animal camouflag-103 ing, which is reasonable for providing gains on CIS. Second, for the Mixture-of-Queries Decoder, 104 which is different from the standard DETR framework, we design multiple expert groups of queries 105 according to the success of the Mixture-of-Experts (MoE). In transformer-based architecture, object queries are wonderful designs in transformer decoders, which have two roles: a) candidates 106 for objects and b) interaction with transformer encoder features in the decoder layers for generat-107 ing outputs. We propose a gating network on each decoder layer to mix multiple groups of queries



Figure 2: The Effectiveness of Mixture-of-Queries Decoder (MoQ Decoder). According to the comparison between the test image and ground truth, hierarchical prediction with MoQ Decoder can refine the output instance masks, including more accurate details than those without MoQ Decoder.

(several queries comprise an expert), deciding which experts are selected for the next layer forwarding. The gating network accepts the encoded features as input, and the final output of each layer
is the weighted recombination of the various experts and the output of the gating network. This
mechanism can refine outputs hierarchically to a fine-grained level via a mixture of experts, which
can generate more accurate instance masks, as shown in Figure 2. In favor of these two designs, our
method can utilize multiple expert queries to integrate effective clues of camouflaged objects in both
spatial and frequency domains, which can achieve outstanding segmentation performance.

- 132 In summary, our main contributions are three-folds.
  - We propose a Mixture-of-Queries Transformer (MoQT) trained in an end-to-end manner for CIS, which takes advantage of the priori of camouflage principles, and refers to the human habit of segmenting camouflaged instances.
  - We proposed a Frequency Enhancement Feature Extractor (FEFE) and a Mixture-of-Queries Decoder (MoQ Decoder) for our MoQT, where FEFE is used for color removal and contour enhancement. The MoQ Decoder aims to mix multiple groups of queries hierarchically to provide more accurate predictions.
  - Extensive experimental results on COD10K and NC4K show consistent performance gains compared with 18 baseline methods and verify the superiority of our method.
- 142 143 144

145

147

121

122

123 124

133

134

135

136

137

138

139

140

141

2 RELATED WORK

## 146 2.1 CAMOUFLAGED OBJECT DETECTION

Camouflaged Object Detection is usually considered as one of the most important origins of CIS and 148 aims to identify the camouflaged objects from the background and has witnessed the development 149 of art and biology Fan et al. (2020a); Le et al. (2019). Early research Huerta et al. (2007); Pan et al. 150 (2011); Sengottuvelan et al. (2008) in COD mainly uses handcrafted features (e.g., gradient, texture, 151 and intensity features) to tell the camouflaged objects from their surroundings. Later, deep learn-152 ing (DL) improves COD's performance in an end-to-end manner, and plenty of DL-based methods 153 Pang et al. (2022); Yang et al. (2021); Piotrowski & Campbell (1982); Xu et al. (2021); Zhong 154 et al. (2022); Mei et al. (2021); Ren et al. (2021) have been proposed. For example, ZoomNet Pang et al. (2022) discusses how to capture camouflaged objects in complex surroundings in a multi-scale 156 manner. Moreover, UGTR Yang et al. (2021) combines the benefits of both Bayesian learning and 157 transformer-based reasoning to handle camouflaged object detection with probabilistic and deter-158 ministic information. Some works Piotrowski & Campbell (1982); Xu et al. (2021); Zhong et al. (2022) even go beyond the RGB domain and explore frequency clues for better performance. In this 159 paper, a Frequency Enhancement Feature Extractor, which refines frequency clues with contour en-160 hancement and color removal, is adopted and allows full rein to both the camouflaged characteristics 161 and the surrounding textures.

#### 162 2.2 CAMOUFLAGED INSTANCE SEGMENTATION

163

164 Camouflage Instance Segmentation (CIS) learns most lessons from traditional instance segmenta-165 tion. The purpose of instance segmentation is to assign pixel-level mask prediction for various 166 instances. Nowadays, instance segmentation methods can be roughly divided into two parts: One-167 stage approaches Bolya et al. (2019); Wang et al. (2020a;b) and two-stage approaches Ren et al. 168 (2015); He et al. (2017a); Cai & Vasconcelos (2018); Chen et al. (2019). Two-stage methods apply mask segmentation after proposal region detection, such as Faster R-CNN Ren et al. (2015), Mask 170 R-CNN He et al. (2017a), Cascade R-CNN Cai & Vasconcelos (2018), and HTC Chen et al. (2019). CFL Le et al. (2021), the first attempt in CIS, also applies two-stage instance segmentation methods. 171 However, one-stage methods show faster inference than two-stage methods and achieve compara-172 ble performance. For example, YOLACT Bolya et al. (2019) adopts two parallel tasks to produce 173 non-local prototype masks with adaptive coefficients. Furthermore, SOLO Wang et al. (2020a) and 174 SOLO-v2 Wang et al. (2020b) predict the instances' center and then decouple the instance masks 175 with kernel feature learning. Recently, researchers have found transformers Cheng et al. (2021a; 176 2022b) show excellent performance on instance segmentation with the assistance of attention mech-177 anisms and instance-specific prototypes. Therefore, transformer-based methods like OSFormer Pei 178 et al. (2022), DCNet Luo et al. (2023) and UOFormer Dong et al. (2023) utilize transformers in 179 CIS and achieve great progress. Inspired by Pei et al. (2022); Luo et al. (2023); Dong et al. (2023), 180 our Mixture-of-Queries Transformer (MoQT) introduces a Mixture-of-Queries Decoder (MoQ Decoder) in the transformer decoder to combine the capabilities of multi-experts hierarchically, which 181 enhances camouflage semantics and refines details of instance masks. 182

183

#### 3 METHOD

185 187

188

#### 3.1 ARCHITECTURE OVERVIEW

189 The overall framework of our proposed model is presented in Figure 3. The whole architecture of 190 our method is a typical MaskFormer-like Cheng et al. (2021b) model, composed of a Frequency 191 Enhancement Feature Extractor (FEFE), a Pixel Decoder, and a Mixture-of-Queries decoder (MoQ 192 Decoder). In the FEFE, we get valuable multi-scale features enhanced by the Fourier transform 193 for revealing the camouflaged clues, where the phase component and amplitude can be used for 194 modeling the information of contours and colors, respectively. We use a contour Enhancement 195 Module (CEM) and a Color Remove Module (CRM) to mine the potential information of contours 196 and eliminate the interference of colors for capturing clues of camouflaged instances. Then, the Pixel Decoder (based on FPN Lin et al. (2017)) gradually upsamples low-resolution features from 197 the output of the backbone to generate high-resolution per-pixel embeddings. The MoQ Decoder 198 computes from per-pixel embeddings and some initialized experts (a series of queries) to get the 199 output prediction. Specifically, in MoQ Decoder, we propose a Mixture-of-Queries Layer (MoQ 200 Laver) after each decoder layer and transform M experts (each expert includes N queries) via self 201 and cross attention mechanisms, where the MoQ Layer is used to combine the M experts of queries 202 hierarchically. Finally, following previous work, we use a mask head and a matching algorithm to 203 output the CIS prediction.

- 204
- 205 206

207

#### 32 FREQUENCY ENHANCEMENT FEATURE EXTRACTOR (FEFE)

208 As mentioned in our Introduction, the camouflage clues are mainly comprised of high-level seman-209 tics (e.g., contours and semantics) and low-level statistics (e.g., colors and styles), which can be 210 reflected by the phrase and amplitude components of Fourier spectrum, respectively. As shown in 211 Figure 1, it is believed that enhancing the influence of contours and eliminating the interference of 212 colors would certainly benefit the performance of segmenting the camouflaged instances. Thus, to 213 explore the camouflaged clues, we design FEFE to model the colors and contours of camouflaged objects, where the phrase and amplitude components are applied to identify the camouflaged se-214 mantics from surroundings in FEFE. Specifically, suppose H and W are the height and width of the 215 input, and the Fourier transformation  $\mathcal{F}(x)$  performed on each channel with a given camouflaged



Figure 3: The Architecture of Our Proposed Model. Our method mainly consists of a Frequency Enhancement Feature Extractor (FEFE), a Pixel Decoder, and a Mixture-of-Queries Decoder (MoQ 235 Decoder). (1) The FEFE captures suitable camouflaged clues with the contour enhancement and 236 color remove modules in the frequency domain. (2) The Pixel Decoder is the same as previous 237 works, based on the FPN architecture, which is used to gradually upsample low-resolution features 238 from the output of the FEFE to generate high-resolution per-pixel embeddings. The detailed ex-239 planation of Pixel Decoder is presented in the Appendix. (3) The MoQ Decoder determines object 240 candidates by multiple cooperation expert queries and hierarchically refines the instance masks with 241 encoded features. 242

image  $x \in \mathcal{R}^{3 \times H \times W}$  can be denoted as:

244 245 246

247

260 261 262

243

$$\mathcal{F}(x) = \sum_{i=0}^{H-1} \sum_{j=0}^{W-1} x[i,j] e^{-J2\pi(\frac{i}{H}u + \frac{j}{W}v)} = \mathcal{A}(x) e^{J\mathcal{P}(x)},\tag{1}$$

where J represents the imaginary unit,  $\mathcal{A}(x)$  (modeling colors) and  $\mathcal{P}(x)$  (modeling contours) are the amplitude and phrase components.

Then, we can get multi-scale image features  $\mathbf{F}^k \in \mathcal{R}^{H \times W \times C}, k \in \{2, 3, 4, 5\}$  extracted from a 251 backbone network with the origin image x. Besides, we feed  $\mathcal{A}(x)$  into a lightweight CNN and a  $1 \times 1$  convolution to obtain the global camouflaged color information  $\mathbf{F}_{color} \in \mathcal{R}^{C}$ , and eliminate its 253 interference via Color Remove Module (CRM). While for the phrase component  $\mathcal{P}(x)$ , due to that 254  $\mathcal{P}(x)$  includes some information on contours and textures, extracting multi-scale features of  $\mathcal{P}(x)$ 255 can present unique advantages to mining camouflaged clues. Thus, we feed  $\mathcal{P}(x)$  into the backbone 256 and obtain hierarchical features  $\mathbf{F}_{contour}^k, k \in \{2, 3, 4, 5\}$  to explore the effects of contours as much 257 as possible by the Contour Enhancement Module (CEM). Formally, the process of FEFE, including 258 CRM and CEM for each scale feature  $\mathbf{F} \in {\{\mathbf{F}^2, \mathbf{F}^3, \mathbf{F}^4, \mathbf{F}^5\}}$ , can be expressed as: 259

$$\mathbf{F}_{refine} = \lambda \mathbf{F} \odot \mathbf{M}_{color} + (1 - \lambda) \mathbf{F} \odot \mathbf{M}_{contour},$$
  

$$\mathbf{M}_{color} = \delta \operatorname{Conv}(\operatorname{avg}_{-}\operatorname{c}((\mathbf{F} - \mathbf{F}_{color})^{2})),$$
  

$$\mathbf{M}_{contour} = \delta \operatorname{Conv}\left(\delta(\operatorname{MLP}(\operatorname{avg}_{-}\operatorname{s}(\mathbf{F}_{contour}))) \odot \mathbf{F}\right),$$
(2)

where avg\_c and avg\_s indicate average pooling along spatial and channel axis, and  $\delta$  is an activation function. CRM and CEM are designd to generate  $\mathbf{M}_{color}$  and  $\mathbf{M}_{contour}$ , respectively. With the above module, we can get multi-scale refined features  $\mathbf{F}_{refine}^{k}, k \in \{2, 3, 4, 5\}$ . Further, to acquire more fine-grained features for more accurate segmentation, we fuse  $\mathbf{F}_{refine}^{k}, k \in \{2, 3, 4, 5\}$  by feeding these features into the pixel decoder based on FPN Lin et al. (2017) architecture, which is used to gradually upsample low-resolution features from the output of the FEFE to generate highresolution per-pixel embeddings  $\mathcal{X} \in \mathcal{R}^{\frac{H}{4} \times \frac{W}{4} \times C}$ .



Figure 4: Illustration of the proposed Mixture-of-Queries Layer with a gated network to output the weights of each expert. The output is the weighted mixture of each input experts.

### 3.3 MIXTURE-OF-QUERIES DECODER (MOQ DECODER)

In order to capture camouflaged instances, the popular transformer-based architecture like Mask-Former Cheng et al. (2021b) and Mask2Former Cheng et al. (2022a), proposes a set of queries to identify whether each pixel belongs to a camouflaged instance. Meanwhile, as discussed in our Introduction, humans may segment camouflaged instances by gradually searching and multi-person collaboration. Inspired by the discussion, we propose a Mixture-of-Queries Decoder (MoQ Decoder) for hierarchically segmenting camouflaged instances.

## 293 3.3.1 MIXTURE-OF-QUERIES MECHANISM

281

282

283 284

287

288

289

290

291

292

303

Different from the standard Transformer decoder architecture, in each layer, we introduce a Mixtureof-Quries Layer (MoQ Layer) after the original decoder layer, and initialize M experts  $E_i$ ,  $i \in [1, M]$ , where each expert contains N queries  $q_i$ ,  $i \in [0, N-1]$ ,  $q_i \in \mathcal{R}^d$ . Each query is responsible for an object candidate. So, we have  $E_i = \{q_0, q_1, \cdots, q_{N-1}\}, q_i \in \mathcal{R}^d$ . Further, the detailed architecture of the designed MoQ Layer is illustrated in Figure 4, where the gated network G outputs a sparse (M+1)-dimensional vector  $G(x) \in \mathcal{R}^{M+1}$  to indicate the weights of each expert. Therefore, given the input of each decoder layer E and the M initialized experts  $[E_1, E_2, \cdots, E_M]$ , the output y of the MoQ Layer can be written as follows:

$$y = \mathbf{E} \cdot \operatorname{softmax}(G(\mathcal{X})), \tag{3}$$

where  $\mathcal{X}$  is the output of pixel decoder, and  $G(\mathcal{X})$  indicates the output of the gated network and E =  $[E', E_1, E_2, \dots, E_M]$ . E' is the output of the standard transformer decoder layer fed with E. Besides, the forward process of each MoQ Decoder Layer (including a original decoder layer and a MoQ layer) can be formulated as:

$$Q = W^{Q} \cdot E, \quad K = W^{K} \cdot \mathcal{X}, \quad V = W^{V} \cdot \mathcal{X},$$

$$E' = \text{LN} \left( E + \text{crossattention}(Q, K, V) \right),$$

$$E' = [E', E_1, E_2, \cdots, E_M] \cdot \text{softmax}(G(\mathcal{X})),$$

$$O = \text{LN}(E' + \text{MLP}(E')),$$

$$(4)$$

where LN is layer normalization and MLP denotes the multi-layer perception network. During the
 training process, to provide deep supervision, we follow Mask2Former to adopt auxiliary losses with
 additional mask prediction heads and Hungarian match loss after each MoQ Layer.

317 318 3.3.2 DISCUSSIONS

To elaborate further on our proposed MoQ Decoder, we provide detailed discussions about the differences between the MoQ Decoder and the vanilla transformer decoder. **1.** Compared with the vanilla transformer decoder, our MoQ Decoder does not contain just one group of queries for capturing various instances but multiple groups of queries in each MoQ Layer, (named as "expert" for each group of queries, noted as E,  $E_1$  and  $E_2$  in Figure 3). Benefiting from the success of MoE, there is a MoQ Layer for combining multi-experts of queries hierarchically to get a more accurate

324	
325	Table 1: Performance Comparison of Various Methods. We display performance of state-of-the-art
326	methods and our MoQT. The best results are in <b>bold</b> .

327	Methods		COD10K-Test			C4K-Te	Doromo(M)	
328			$AP_{50}$	$AP_{75}$	AP	$AP_{50}$	$AP_{75}$	Faranis(WI)
200	Mask R-CNN He et al. (2017b)	25.0	55.5	20.4	27.7	58.6	22.7	43.9
329	MS R-CNN Huang et al. (2019)	30.1	57.2	28.7	31.0	58.7	29.4	60.0
330	Cascade R-CNN Cai & Vasconcelos (2019)	25.3	56.1	21.3	29.5	60.8	24.8	71.7
331	HTC Chen et al. (2019)	28.1	56.3	25.1	29.8	59.0	26.6	76.9
220	BlendMask Chen et al. (2020)	28.2	56.4	25.2	27.7	56.7	24.2	35.8
332	Mask Transfiner Ke et al. (2022)	28.7	56.3	26.4	29.4	56.7	27.2	44.3
333	YOLACT Bolya et al. (2019)	24.3	53.3	19.7	32.1	65.3	27.9	-
334	CondInst Tian et al. (2020)	30.6	63.6	26.1	33.4	67.4	29.4	34.1
335	QueryInst Fang et al. (2021)	28.5	60.1	23.1	33.0	66.7	29.4	-
000	SOTR Guo et al. (2021)	27.9	58.7	24.1	29.3	61.0	25.6	63.1
336	SOLOv2 Wang et al. (2020b)	32.5	63.2	29.9	34.4	65.9	31.9	46.2
337	MaskFormer Cheng et al. (2021a)	38.2	65.1	37.9	44.6	71.9	45.8	45.0
338	Mask2Former Cheng et al. (2022b)	39.4	67.7	38.5	45.8	73.6	47.5	43.9
220	OSFormer Pei et al. (2022)	41.0	71.1	40.8	42.5	72.5	42.3	46.6
339	DCNet Luo et al. (2023)	45.3	70.7	47.5	52.8	77.1	56.5	53.4
340	UQFormer Dong et al. (2023)	45.2	71.6	46.6	47.2	74.2	49.2	37.5
341	CamoFourier Le et al. (2023)	43.52	74.84	42.65	44.95	75.67	44.28	-
342	MSPNet Li et al. (2024)	39.7	69.8	39.8	41.8	71.8	42.3	48.09
343	Ours	47.99	73.01	51.77	54.73	78.45	58.97	61.68

prediction. 2. Inspired by the human habit of segmenting camouflaged instances by gradual searches
and multi-person cooperation. Each layer of MoQ Decoder always has *M* newly initialized experts,
indicating that different depths include different multi-experts responsible for capturing instances.
While the vanilla transformer decoder just initializes one set of queries (*i.e.* one expert) in the first
layer. Due to these two differences between the architecture of the standard transformer decoder and
our proposed MoQ Decoder, our method can obtain a more accurate segmentation mask for each
instance.

### 3.4 OBJECTIVE FUNCTION

As shown in Figure 3, with the fused feature  $\mathcal{X} \in \mathcal{R}^{\frac{H}{4} \times \frac{W}{4} \times C}$  output by Pixel Decoder and the instance candidates  $\hat{E} \in \mathcal{R}^{N \times C}$  generated by the MoQ Decoder, we can finally obtain the segmentation map, which can be formulated as:

$$Mask = \mathcal{X} \times \hat{E}.$$
 (5)

To train the whole network, following DETR Carion et al. (2020), we adopt a Hungarian matching algorithm to match a ground truth label with each predicted segment instance. If no suitable label exists, a special label ("no object") is assigned. Therefore, including the instances and mask supervision, the objective function contains three terms: Cross-entropy Loss  $\mathcal{L}_{CE}$  for the instance score, Focal Loss  $\mathcal{L}_{focal}$  and Dice Loss  $\mathcal{L}_{dice}$  for the mask predictions after each MoQ Layer, written as:

$$\mathcal{L}_{total} = \sum_{l=0}^{L} \mathcal{L}_{CE} + \alpha \cdot \mathcal{L}_{focal} + \beta \cdot \mathcal{L}_{dice}, \tag{6}$$

where L means the amount of decoder layers. By default, we set  $\alpha = 20$  and  $\beta = 1$ .

### 4 EXPERIMENTS

#### 4.1 EXPERIMENTAL SETUPS

Following the mainstream works of CIS Dong et al. (2023); Luo et al. (2023), we evaluate our
method in two datasets: COD10K and NC4K. COD10K includes 3040 training images and 2026
testing images, while NC4K contains 4121 test images for evaluating the generalization of proposed
models. To provide a fair comparison, we train models in the training set in COD10K, and meanwhile test models in both test sets of COD10K and NC4K, which is a standard setting proposed in

Table 2: Performance Comparison of Proposed
Modules. We perform an ablation study on
COD10K and NC4K to validate our proposed modules' effectiveness. "FEFE" and "MoQ" represent Frequency Enhancement Feature Extractor and
Mixture-of-Queries Decoder, respectively.

COD10K-Test

 $AP_{50}$ 

71.33

72.76

73.01

 $AP_{75}$ 

49.31

50.61

51.77

AP

53.22

53.87

54.73

Table 3: Performance Comparison of Various Backbones. We evaluate multiple methods' performance with various backbones on COD10K and NC4K.

Method	Backbone	COD10K	NC4K
OSFormer		41.0	42.5
DCNet	ResNet-50	45.3	52.8
Ours		48.0	54.7
OSFormer		42.0	44.4
DCNet	ResNet-101	46.8	53.5
Ours		48.6	55.4
OSFormer		47.7	50.2
DCNet	Swin-Tiny	50.3	56.3
Ours		51.4	58.1
OSFormer		52.1	56.7
DCNet	Swin-Small	52.3	58.4
Ours		53.2	59.2

previous works Luo et al. (2023); Pei et al. (2022). In order to comprehensively evaluate the models, We use  $AP_{50}$ ,  $AP_{75}$ , and AP scores as evaluation metrics to quantify the performance of our method and baselines Luo et al. (2023); Dong et al. (2023); Pei et al. (2022). Besides, the implementation details of our experiments are presented in our *Appendix*.

NC4K-Test

 $AP_{50}$ 

77.84

78.02

78.45

 $AP_{75}$ 

57.29

58.12

58.97

### 4.2 Comparison with State-of-the-Art Methods

The CIS task is a relatively novel task that has been proposed in recent years, and only a few previous works are involved in this task, such as OSFormer Pei et al. (2022), DCNet Luo et al. (2023), and UQFormer Dong et al. (2023). Consequently, we also adopt several popular generic instance segmentation methods as baselines on the CIS task for a more comprehensive test. And for a fair comparison, the backbone of these methods is configured as ResNet-50. The performance comparison results are shown in Table 1. It is easy to observe that our proposed model can consistently outperform the state-of-the-art methods by a large margin on both COD10K and NC4K test sets.

405 (1) Results on COD10K. As shown in Table 1, we compare our proposed model with 5 CIS models 406 (i.e., OSFormer Pei et al. (2022), DCNet Luo et al. (2023), UQFormer Dong et al. (2023), Camo-407 Fourier Le et al. (2023), and MSPNet Li et al. (2024)) and 13 generic instance segmentation models. 408 Our model can achieve 51.77% in  $AP_{75}$ , which outperforms the second best method DCNet Luo 409 et al. (2023) by 4.27% in  $AP_{75}$ . In AP, our model also gets a performance improvement of 2.69%. Notice that our method does not achieve the highest value in  $AP_{50}$ , instead of a comparable per-410 formance of 73.01% in  $AP_{50}$ . These results indicate that our method can acquire more accurate 411 segmentation masks of camouflaged objects. 412

413 (2) Results on NC4K. Likewise, we evaluate these methods on NC4K dataset, and the results on this 414 test set reflect the generalization ability of these models. Our model yields 58.97% in  $AP_{75}$ , while 415 the previous best method DCNet is 56.5%, which demonstrates that our method gets an obvious gain 416 of 2.47% in  $AP_{75}$ , suggesting a great generalization ability of our model as well. In AP, our model 417 achieves the highest performance metrics of 54.73%, surpassing the second best method (DCNet) 418 by 1.93%. Besides, our model also obtains a 1.35% improvement in  $AP_{50}$ . The overall metrics of 419 various AP values reflect our method's obvious superiority over other baselines.

420 421

384

386

387

389 390 391

392

393

394

395

397

FEFE

~

1

MoQ

V

1

AP

45.76

47.12

47.99

4.3 Ablation Studies and Visualizations

To look deeper into our proposed method, in this section, we present a series of ablation studies to demonstrate the effectiveness of each proposed module.

**Effectiveness of proposed modules.** To explore the effectiveness of the proposed FEFE and MoQ Decoder, we validate the importance of each component by removing them one at a time. As shown in Table 2, the performance without MoQ Decoder drops by 2.23 % in AP, 1.68% in  $AP_{AP_{50}}$  and 2.46% in  $AP_{AP_{75}}$  on COD10K-Test. On NC4K-Test, the metrics of AP,  $AP_{50}$  and  $AP_{75}$  are also reduced by 1.51%, 0.66% and 1.68%, respectively. Similarly, if the components of FEFE are ablated, there is a drop in segmentation performance as well. For example, on COD10K-Test, the performance just achieves 47.12% in AP, 72.76% in  $AP_{50}$  and 50.61% in  $AP_{75}$ , which are consistently lower than that without any modules ablated (as shown in the last row of Table 2). The

447 448 449

450



Figure 5: Performance Comparison of Various Numbers of Queries in Each Expert. AP and  $AP_{75}$  of our MoQT with various numbers of queries on COD10K-Test (a) and NC4K-Test (b) are shown.

Table 4: Performance Comparison of Various Number of Decoder Layers. We apply various numbers of decoder layers, and the performance is shown as follows. The best results are in bold.

Decoder Lovers	COD10K-Test			N	C4K-Te	Darame(M)	
Decouer Layers	AP	$AP_{50}$	$AP_{75}$	AP	$AP_{50}$	$AP_{75}$	T at at its(IVI)
2	46.50	71.99	50.23	53.47	78.14	57.69	54.62
4	47.02	72.34	51.01	53.61	78.32	57.80	57.93
6	47.99	73.01	51.77	54.73	78.45	58.97	61.68
8	47.45	72.36	51.28	53.73	78.36	57.88	65.43
10	46.89	71.53	50.07	53.35	78.01	57.10	68.36
12	47.20	72.06	50.15	53.82	78.22	58.10	71.18

reduced performance demonstrates that these two proposed modules can capture clues of camou flaged instances and provide accurate segmentation. With both modules, our method can lead to
 huge performance gains in evaluation metrics.

Various backbones. To further explore the potential of our model, we equip it with different 464 feature extractor backbones, such as ResNet-50 He et al. (2016), ResNet-101 He et al. (2016), 465 SwinTransformer-Tiny (Swin-Tiny) Liu et al. (2021), and SwinTransformer-Small (Swin-Small) 466 Liu et al. (2021). For a fair comparison with baselines, all these models are pretrained on ImageNet-467 1k Deng et al. (2009). The results are presented in Table 3. With the same backbone, our method 468 shows the best performance among compared baselines, which indicates our method outperforms 469 the state-of-the-art methods. For example, when ResNet-101 is the backbone, the metrics of AP of 470 our method are 48.6% and 55.4% on COD10K and NC4K datasets, respectively, while the second best method just reaches 46.8% and 53.5%. With a larger backbone, the results also prove that our 471 method has the potential for further improvement. 472

473 Ablation on the number of queries. Object queries are essential in the transformer architecture for 474 prediction. Therefore, we study the performance with different numbers of queries in each expert 475 group. As shown in Figure 5, we change the number of queries from 10 to 40 and evaluate the 476 performance metrics of AP and  $AP_{75}$  in both COD10K and NC4K test sets. In fact, the number 477 of queries in each group should be larger than the actual count of objects to avoid instance fusion, which is determined by the dataset distribution. Moreover, it can be seen that when the number is 478 set as 10, our model obtains the best performance on both datasets. For example, when the number 479 is 10, AP and AP<sub>75</sub> is 47.99% and 51.77% on COD10K, respectively. Meanwhile, AP and AP<sub>75</sub> 480 reach 54.73% and 58.97%. 481

**482 Analysis of the number of decoder layers.** We apply auxiliary losses after each decoding layer, as 483 formulated in Equation (6). Hence, the number of decoder layers *L* is important for the segmentation 484 performance. As presented in Table 4, we vary the number of decoder layers, picked from the set 485  $\{2, 4, 6, 8, 10, 12\}$ . We find that the overall performance of the model improves with the increase of *L*. And when L = 6, the model can get the best performance. There is no additional performance







Figure 7: Visualizations of Various Methods. Different colored masks indicate different instances.
gain when the *L* continues to increase, which may be caused by limited data to train the model for

514 further improvement.

**Impacts about Hyper-parameters.** We study the impacts of the Hyper-parameters  $\alpha$  and  $\beta$  in Equation (6). On both COD10K and NC4K datasets, when  $\alpha = 20$ , the best performance is achieved, proved by the metrics of AP = 47.99% and 54.73%, respectively, shown in Figure 6(a). Therefore, we choose  $\alpha = 20$  in our method by default. For the hyper-parameter  $\beta$ , we change the value of  $\beta$  from 0.1 to 2, and the results as presented in Figure 6 (b). It can be seen that the model gets the best performance when  $\beta = 1$ . Therefore, to get the best performance, we set  $\alpha$  as 20, and  $\beta$  as 1.

521 Visualization Results. To comprehensively evaluate our method, we also present some qualita-522 tive analysis, as shown in Figure 7, referring to *appendix* for more visualization. We visualize the 523 segmentation masks of various methods, including OSFormer Pei et al. (2022), DCNet Luo et al. 524 (2023), and our method, to demonstrate the performance with qualitative results. It can be seen that our method performs better than previous methods, which can be proved by the clear boundaries 525 and accurate masks of our method (shown in the last row of Figure 7). In short, our method not 526 only improves the evaluation metrics on two datasets but also gains in visual results of segmentation 527 masks. 528

529

498 499

500 501

509 510 511

## 530 5 CONCLUSION

531 532

In this paper, we propose a novel Mixture-of-Queries Transformer (MoQT) for camouflaged instance segmentation. MoQT applies a Frequency Enhancement Feature Extractor for feature extraction in the frequency domain, with the assistance of a contour enhancement module and a color removal module. Besides, a Mixture-of-Queries Decoder uses multiple expert groups of queries as candidates and shares semantic information with transformer encoder features. Multi-scale features enable MoQT to refine prediction hierarchically and get fine-grained instance masks with collaboration of multiple groups of queries. Compared with plenty of state-of-the-art baselines, our proposed MoQT shows outstanding performance on two benchmark datasets, demonstrating the proposed method's effectiveness.

## 540 REFERENCES

558

572

576

577

- Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation.
  In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9157–9166, 2019.
- Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In
   *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6154–6162,
   2018.
- 548
  549
  549
  550
  551
  2019.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and
   Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pp. 213–229. Springer, 2020.
- Hao Chen, Kunyang Sun, Zhi Tian, Chunhua Shen, Yongming Huang, and Youliang Yan. Blend mask: Top-down meets bottom-up for instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8573–8581, 2020.
- Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4974–4983, 2019.
- Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need
   for semantic segmentation. volume 34, pp. 17864–17875, 2021a.
- Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. 2021b.
- Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked attention mask transformer for universal image segmentation. 2022a.
- Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked attention mask transformer for universal image segmentation. 2022b.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
  - Bo Dong, Jialun Pei, Rongrong Gao, Tian-Zhu Xiang, Shuo Wang, and Huan Xiong. A unified query-based paradigm for camouflaged instance segmentation, 2023.
- Deng-Ping Fan, Ge-Peng Ji, Guolei Sun, Ming-Ming Cheng, Jianbing Shen, and Ling Shao. Camouflaged object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2777–2787, 2020a.
- Deng-Ping Fan, Ge-Peng Ji, Tao Zhou, Geng Chen, Huazhu Fu, Jianbing Shen, and Ling Shao.
   Pranet: Parallel reverse attention network for polyp segmentation. In *International conference on medical image computing and computer-assisted intervention*, pp. 263–273. Springer, 2020b.
- Deng-Ping Fan, Tao Zhou, Ge-Peng Ji, Yi Zhou, Geng Chen, Huazhu Fu, Jianbing Shen, and Ling
   Shao. Inf-net: Automatic covid-19 lung infection segmentation from ct images. *IEEE transac- tions on medical imaging*, 39(8):2626–2637, 2020c.
- Yuxin Fang, Shusheng Yang, Xinggang Wang, Yu Li, Chen Fang, Ying Shan, Bin Feng, and Wenyu Liu. Instances as queries. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6910–6919, 2021.
- Huazhu Fu, Fei Li, José Ignacio Orlando, Hrvoje Bogunović, Xu Sun, Jingan Liao, Yanwu Xu,
   Shaochong Zhang, and Xiulan Zhang. Refuge: Retinal fundus glaucoma challenge, 2019. URL https://dx.doi.org/10.21227/tz6e-r977.

609

594	Ruohao Guo, Dantong Niu, Liao Ou, and Zhenbo Li. Sotr: Segmenting objects with transformers.
595	In Proceedings of the IEEE/CVF international conference on computer vision, pp. 7157–7166,
596	2021.
597	

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017a.
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017b.
- Zhaojin Huang, Lichao Huang, Yongchao Gong, Chang Huang, and Xinggang Wang. Mask scoring
   r-cnn. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
   pp. 6409–6418, 2019.
- Iván Huerta, Daniel Rowe, Mikhail Mozerov, and Jordi Gonzàlez. Improving background subtraction based on a casuistry of colour-motion segmentation problems. In *Iberian Conference on Pattern Recognition and Image Analysis*, pp. 475–482. Springer, 2007.
- Lei Ke, Martin Danelljan, Xia Li, Yu-Wing Tai, Chi-Keung Tang, and Fisher Yu. Mask transfiner for
   high-quality instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4412–4421, 2022.
- 616
   617
   618
   618
   618
   619
   619
   610
   610
   6110
   6111
   6112
   6112
   612
   612
   613
   614
   614
   615
   616
   616
   617
   618
   618
   619
   619
   610
   610
   6112
   612
   612
   613
   614
   614
   615
   616
   617
   618
   617
   618
   618
   619
   618
   619
   618
   619
   618
   619
   614
   614
   615
   616
   617
   618
   616
   617
   618
   618
   619
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618
   618</
- Minh-Quan Le, Minh-Triet Tran, Trung-Nghia Le, Tam V. Nguyen, and Thanh-Toan Do. Unveiling camouflage: A learnable fourier-based augmentation for camouflaged object detection and instance segmentation, 2023.
- Trung-Nghia Le, Tam V Nguyen, Zhongliang Nie, Minh-Triet Tran, and Akihiro Sugimoto.
   Anabranch network for camouflaged object segmentation. *Computer vision and image under*standing, 184:45–56, 2019.
- Trung-Nghia Le, Yubo Cao, Tan-Cong Nguyen, Minh-Quan Le, Khanh-Duy Nguyen, Thanh-Toan Do, Minh-Triet Tran, and Tam V Nguyen. Camouflaged instance segmentation in-the-wild:
   Dataset, method, and benchmark suite. *IEEE Transactions on Image Processing*, 31:287–300, 2021.
- Chen Li, Ge Jiao, Guowen Yue, Rong He, and Jiayu Huang. Multi-scale pooling learning for camouflaged instance segmentation. *Applied Intelligence*, pp. 1–15, 2024.
- Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie.
   Feature pyramid networks for object detection. In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 936–944, 2017.
- <sup>636</sup>
   <sup>637</sup> Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- Naisong Luo, Yuwen Pan, Rui Sun, Tianzhu Zhang, Zhiwei Xiong, and Feng Wu. Camouflaged
   instance segmentation via explicit de-camouflaging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17918–17927, 2023.
- Haiyang Mei, Ge-Peng Ji, Ziqi Wei, Xin Yang, Xiaopeng Wei, and Deng-Ping Fan. Camouflaged object segmentation with distraction mining. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8772–8781, 2021.
- 647 Yuxin Pan, Yiwang Chen, Qiang Fu, Ping Zhang, Xin Xu, et al. Study on the camouflaged target detection method based on 3d convexity. *Modern Applied Science*, 5(4):152, 2011.

648	Youwei Pang, Xiaogi Zhao, Tian-Zhu Xiang, Lihe Zhang, and Huchuan Lu. Zoom in and out: A
649	mixed-scale triplet network for camouflaged object detection. In <i>Proceedings of the IEEE/CVF</i>
650	Conference on computer vision and pattern recognition, pp. 2160–2170, 2022.
651	
652	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen Zaming Lin, Natelia Cimelshein, Luca Antiga, et al. Buterahi An imperativa etula high
653	nerformance deep learning library. Advances in neural information processing systems 32, 2010
654	performance deep rearring notary. Advances in neural information processing systems, 52, 2019.
655	Jialun Pei, Tianyang Cheng, Deng-Ping Fan, He Tang, Chuanbo Chen, and Luc Van Gool. Osformer:
656	One-stage camouflaged instance segmentation with transformers. In European conference on
657	computer vision. Springer, 2022.
658	Leon N Piotrowski and Fergus W Campbell. A demonstration of the visual importance and flexibility
660	of spatial-frequency amplitude and phase. <i>Perception</i> , 11(3):337–346, 1982.
661	Lingthe Day Viennei Her Lei 7hr Vernies V. Vernies V. Weining Weng 7iller Days
662	and Pheng Ann Heng. Deep texture aware features for camouflaged object detection. <i>IEEE</i>
663	Transactions on Circuits and Systems for Video Technology, 33(3):1157–1167, 2021.
664	
665	Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object
666	detection with region proposal networks. Advances in neural information processing systems, 28,
667	2015.
668	P Sengottuvelan, Amitabh Wahi, and A Shanmugam. Performance of decamouflaging through ex-
669	ploratory image analysis. In 2008 First International Conference on Emerging Trends in Engi-
670	neering and Technology, pp. 6–10. IEEE, 2008.
671	7bi Tian Chunhua Shan and Hao Chan. Conditional convolutions for instance segmentation. In
672	Computer Vision-ECCV 2020: 16th European Conference Glassow IIK August 23–28 2020
673	Proceedings. Part I 16. pp. 282–298. Springer, 2020.
674	
675	Hrvoje Turić, Hrvoje Dujmić, and Vladan Papić. Two-stage segmentation of aerial images for search
676	and rescue. Information Technology and Control, 39(2), 2010.
677	Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by
678	locations. In Computer Vision-ECCV 2020: 16th European Conference, Glasgow, UK, August
679	23-28, 2020, Proceedings, Part XVIII 16, pp. 649-665. Springer, 2020a.
680	Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen, Solov2: Dynamic and fast
690	instance segmentation. Advances in Neural information processing systems, 33:17721–17732,
683	2020b.
684	Cinuci Vu Duinang Zhang, Va Zhang, Vanfang Wang, and Oi Tian. A fourier based framework
685	for domain generalization. In Proceedings of the IFFF/CVF conference on computer vision and
686	pattern recognition, pp. 14383–14392, 2021.
687	
688	Fan Yang, Qiang Zhai, Xin Li, Rui Huang, Ao Luo, Hong Cheng, and Deng-Ping Fan. Uncertainty-
689	guided transformer reasoning for camouflaged object detection. In <i>Proceedings of the IEEE/CVF</i>
690	International Conference on Computer Vision, pp. 4140–4155, 2021.
691	Yijie Zhong, Bo Li, Lv Tang, Senyun Kuang, Shuang Wu, and Shouhong Ding. Detecting camou-
692	flaged object in frequency domain. In Proceedings of the IEEE/CVF Conference on Computer
693	Vision and Pattern Recognition, pp. 4504–4513, 2022.
694	Peng Zhou, Bor-Chun Chen, Xintong Han, Mahyar Najibi, Abhinay Shriyastaya, Ser-Nam Lim
695	and Larry Davis. Generate, segment, and refine: Towards generic manipulation segmentation.
696	In Proceedings of the AAAI conference on artificial intelligence, volume 34, pp. 13058–13065,
697	2020.
698	
699	
700	
701	