

Visual Consensus Prompting for Co-Salient Object Detection

Jie Wang¹, Nana Yu¹, Zihao Zhang¹, Yahong Han^{1*}

¹College of Intelligence and Computing, Tianjin University, Tianjin, China

wangjiexy@tju.edu.cn, yunana@tju.edu.cn, zhangzihao2490@tju.edu.cn, yahong@tju.edu.cn

Abstract

Existing co-salient object detection (CoSOD) methods generally employ a three-stage architecture (i.e., encoding, consensus extraction & dispersion, and prediction) along with a typical full fine-tuning paradigm. Although they yield certain benefits, they exhibit two notable limitations: 1) This architecture relies on encoded features to facilitate consensus extraction, but the meticulously extracted consensus does not provide timely guidance to the encoding stage. 2) This paradigm involves globally updating all parameters of the model, which is parameter-inefficient and hinders the effective representation of knowledge within the foundation model for this task. Therefore, in this paper, we propose an interaction-effective and parameter-efficient concise architecture for the CoSOD task, addressing two key limitations. It introduces, for the first time, a parameter-efficient prompt tuning paradigm and seamlessly embeds consensus into the prompts to formulate task-specific Visual Consensus Prompts (VCP). Our VCP aims to induce the frozen foundation model to perform better on CoSOD tasks by formulating task-specific visual consensus prompts with minimized tunable parameters. Concretely, the primary insight of the purposeful Consensus Prompt Generator (CPG) is to enforce limited tunable parameters to focus on co-salient representations and generate consensus prompts. The formulated Consensus Prompt Disperser (CPD) leverages consensus prompts to form task-specific visual consensus prompts, thereby arousing the powerful potential of pre-trained models in addressing CoSOD tasks. Extensive experiments demonstrate that our concise VCP outperforms 13 cutting-edge full fine-tuning models, achieving the new state of the art (with 6.8% improvement in F_m metrics on the most challenging CoCA dataset). Source code has been available at <https://github.com/WJ-CV/VCP>.

1. Introduction

Co-Salient Object Detection (CoSOD) is a group-based image understanding task aimed at detecting salient objects that commonly appear among a group of relevant images.

*Corresponding author.

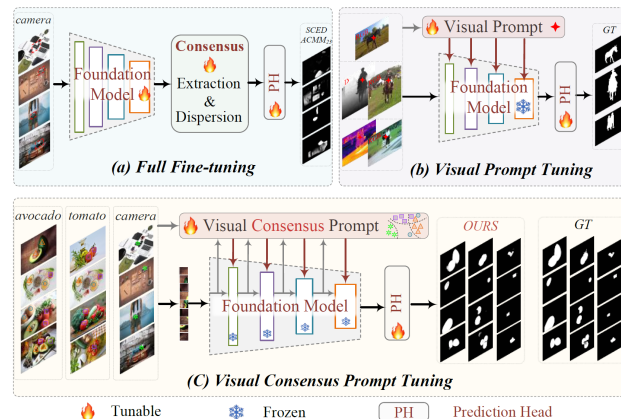


Figure 1. Existing relevant methods VS. our VCP. (a) Existing CoSOD methods based on typical architectural patterns and full fine-tuning paradigms. (b) Introducing simple tunable parameters as visual prompts to address foreground segmentation tasks in single-scene images. (c) Our proposed VCP and some visualization results. The frozen foundation model is mined to generate task-specific visual consensus prompts (with minimized tunable parameters), thereby inducing it to effectively perform CoSOD.

CoSOD methods have demonstrated their effectiveness as a key pre-processing step for various computer vision tasks, such as co-segmentation [39], co-localization [9], and object tracking [34]. Additionally, with the effectiveness of CoSOD methods continue to improve, they are increasingly being utilized in numerous practical applications [8, 30, 37].

Numerous impressive CoSOD methods have been proposed and continually advancing the performance of this task. As illustrated in Fig.1 (a), most of these methods [33, 35, 40, 42, 44, 46] employ a common three-stage architecture: first, encoding multi-scale features using a foundation model pretrained on large-scale datasets; next, designing unique consensus mining schemes to focus on and mine co-salient object representations within the group of relevant images to complete consensus extraction and further disperse the consensus into multi-scale features; finally, obtaining co-saliency predictions using a prediction head. In this architecture, the effectiveness of encoded features greatly facilitates the extraction and dispersion of consensus, but the meticulously extracted consensus struggles to

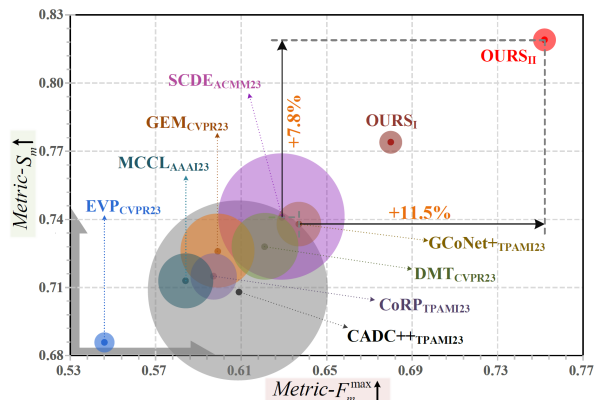


Figure 2. Quantitative comparison of our VCP with 8 representative methods on the CoCA [42] dataset regarding S_m , F_m^{\max} metrics, and tunable parameters. The bubble area represents the tunable parameters (M). SCED [33], GEM [31], MCCL [44], GCoNet+ [43], CoPR [46], DMT [16], and CADC++ [41] are all full fine-tuning CoSOD methods. EVP [20] is based on prompt learning for SOD tasks, and we retrain it using the CoSOD dataset.

guide feature encoding in a timely manner, as the encoder is only fine-tuned at the end of parameter optimization. In a nutshell, this architecture lacks efficient interaction between encoding and consensus, resulting in neither fully realizing their potential. More critically, existing methods typically employ the full fine-tuning paradigm, which involves adapting the model to the specific CoSOD task by tuning all parameters (including the large-scale pre-trained foundation model) using existing CoSOD datasets [19, 28, 42]. The full fine-tuning paradigm has two main limitations: 1) It is parameter-inefficient and requires many repetitions of tuning and storage of the entire pre-trained model, resulting in significant computational and storage overheads. 2) The quality and quantity of fine-tuning data also restrict the effectiveness of knowledge representation of the pre-trained foundation model in this task. These limitations make it difficult for existing methods to achieve more effective and efficient CoSOD performance (Fig.1 (a) and Fig.2). Additionally, they increase the difficulty of CoSOD methods in practical applications, and this difficulty further increases with the inherent trend of increasing scale of foundation models (e.g., Transformer series or large models).

The parameter-efficient prompt tuning paradigm [15, 18, 21] has been proposed and promoted in Natural Language Processing (NLP) [2, 25, 26]. Inspired by that, recent works [3, 12, 23, 45] have introduced the prompt tuning paradigm into visual tasks, and guided the model to solve visual recognition tasks by freezing the foundation model and adding a few learnable parameters as task-specific visual prompts. Additionally, some promising works utilize the idea of prompt tuning to address foreground segmentation tasks (Fig.1 (b)). These methods [20, 22] preset simple tunable embeddings as foreground prompts to address SOD-related tasks. As shown in Fig. 2, these methods per-

form poorly when tackling specific CoSOD tasks, primarily due to two main limitations: (1) These methods solely focus on single-image foregrounds regardless of their category, while the CoSOD not only expands a group dimension but also involves numerous non-co-salient objects within the images as interference. (2) Introducing simple tunable parameters as visual prompts struggles to model intra-group co-salient representations and to achieve effective CoSOD performance. Hence, we ponder: How to formulate task-specific prompts for CoSOD to adapt the foundation model to effectively and efficiently facilitate this task?

At this juncture, it is surprising to find that if the extraction and dispersion of consensus (the key representations specific to the CoSOD task) are embedded into visual prompts, the architecture we seek for efficient interaction between consensus and encoding can perfectly align with the prompt tuning paradigm. In this spirit, we propose a concise and parameter-efficient Visual Consensus Prompting (VCP) for co-salient object detection. Our VCP aims to induce the frozen pre-trained model to effectively and efficiently execute CoSOD tasks by constructing task-specific visual consensus prompts (with minimized tunable parameters). Given the critical importance of visual consensus for CoSOD tasks, we design two key components to support the implementation of VCP: the Consensus Prompt Generator (CPG) and the Consensus Prompt Dispenser (CPD). The primary insight of CPG is to enforce limited tunable parameters to focus on intra-group co-salient representations from frozen embedding features, thereby generating consensus prompts. CPD leverages consensus prompts to form task-specific visual consensus prompts and adaptively harnesses the powerful potential of frozen Transformer layers in addressing CoSOD. We summarize the contribution of our work as follows:

- An interaction-effective and parameter-efficient visual consensus prompts architecture for CoSOD is proposed, which can serve as a powerful alternative to existing methods based on the common architecture and the full fine-tuning paradigm.
- The proposed Consensus Prompt Generator systematically aggregates task-specific consensus prompts, while the Consensus Prompt Dispenser adaptively induces the foundation model to perform better in this task.
- Extensive experiments demonstrate the comprehensive performance advantage of our VCP. Particularly on the COCA dataset, which best reflects the model’s robustness, our VCP outperforms the state-of-the-art method by **5.6%** and **6.8%** in terms of S_m and F_m , respectively.

It’s our belief that the CoSOD task can be conducted with a concise architecture in an interaction-effective and parameter-efficient manner. We hope the impressive performance of VCP can demonstrate the potential of the consensus prompt tuning paradigm to the CoSOD community.

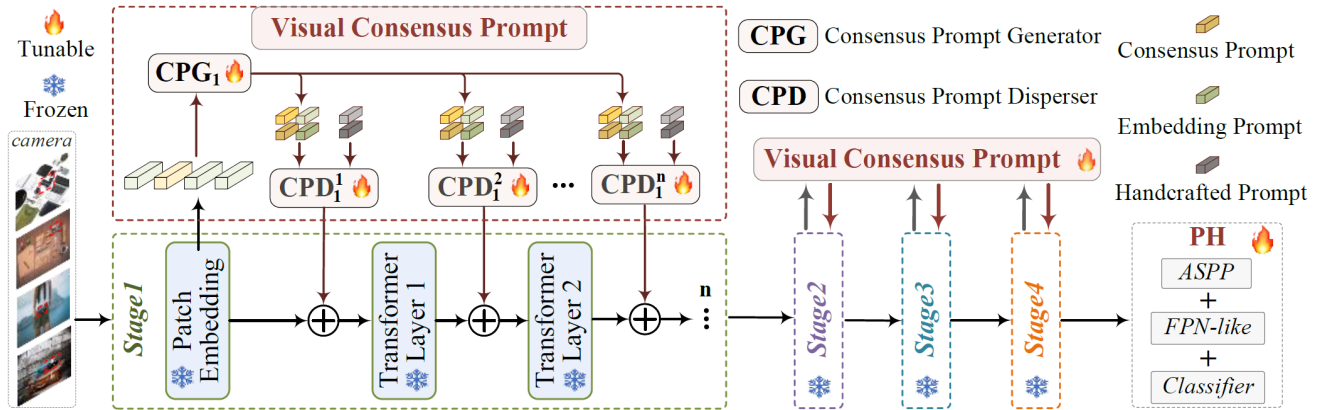


Figure 3. Overall framework pipeline of our proposed concise and parameter-efficient VCP model. We induce the frozen foundation model to perform better on the CoSOD task by formulating Visual Consensus Prompts with minimal tunable parameters. The proposed Consensus Prompt Generator (CPG) and Consensus Prompt Disperser (CPD) support the implementation of VCP. The CPG mines intra-group co-salient representations of the frozen embeddings to generate consensus prompts P_{Co} . The CPD utilizes P_{Co} to form Visual Consensus Prompts and induce the frozen transformer layers to perform the CoSOD task.

2. Related works

2.1. Co-Salient Object Detection

CoSOD aims to discover and segment salient objects that appear commonly within a group of images. Recent CoSOD methods [8, 29, 33, 35, 40, 42, 44] have achieved impressive performance, and their implementation paradigms can be broadly summarized as consensus extraction and consensus dispersion. Consensus extraction aims to focus on the salient objects common to the group of images and form a consensus representation. Consensus dispersion utilizes the extracted consensus to guide the model’s decoding and generate predictions. Xu *et al.* [33] employ a hierarchical Transformer to extract semantic consensus and propagate consensus based on the Transformer. Some promising methods have attempted to use saliency predictions to filter out non-salient backgrounds or to further refine consensus representations [14, 43, 46]. However, existing methods lack effective interaction between encoding and consensus, and they are all based on the full fine-tuning strategy, which involves adapting the model to the specific CoSOD task by adjusting all parameters (including the large-scale pre-trained foundation model) using existing CoSOD datasets. This parameter-inefficient tuning paradigm incurs significant computational and storage overheads, making it not only challenging to achieve better performance but also unfriendly to practical applications. In contrast, an interaction-effective and parameter-efficient architecture for CoSOD is proposed, which can serve as a powerful alternative to existing prevailing methods.

2.2. Visual Prompting Tuning

Prompt-based learning is pioneered and driven by the GPT series [2, 25, 26] in the field of NLP. It initially referred

to adding task-specific descriptions to downstream inputs to assist language models in handling downstream tasks, rather than solely adapting pre-trained models to fit downstream tasks through full fine-tuning. In addition to subsequent work on how to construct better prompt texts [13, 27], recent efforts have also proposed treating prompts as continuous vectors specific to tasks and directly optimizing them through gradients during the fine-tuning process, known as Prompt Tuning [15, 18, 21]. Besides, prompt learning also shows its effectiveness in many computer vision tasks [3, 12, 20, 22, 23, 45]. VPT [12] provides a set of learnable parameters pre-prepared for Transformer encoders, significantly outperforming full fine-tuning on 20 downstream identification tasks. Unlike simply adding some linear layers or tunable parameters to address visual recognition tasks [3, 12, 23], some promising methods incorporate prompt learning into foreground segmentation tasks. EVP [20] introduces handcrafted features to create explicit visual prompts for effective single-image SOD. VSCoDe [22] incorporates task-specific and domain-specific prompts to address single-modal or multi-modal SOD tasks through full fine-tuning of the model. These methods primarily focus on foreground objects in a single scene and only perform simple prompt tuning, making it difficult to model intra-group consensus and achieve satisfactory co-salient object segmentation results.

3. Proposed Methodology

3.1. Overview

In this section, we propose Visual Consensus Prompt tuning, i.e., VCP, for CoSOD, it is a concise and parameter-efficient alternative to existing full fine-tuning methods [16, 31, 33, 41, 46]. The proposed VCP is a predefined

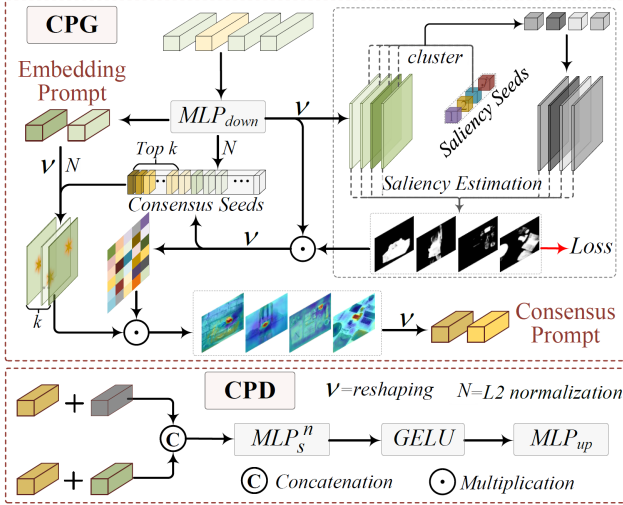


Figure 4. Overall pipeline of the proposed CPG and CPD. The CPG utilizes predefined saliency seeds to generate saliency estimation maps through clustering, thereby obtaining consensus seeds. By selecting top-k representative consensus seeds, consensus prompts P_{Co} are obtained. The CPD utilizes P_{Co} to generate visual consensus prompts P_{Visual}^{Co} and induce the frozen transformer layers to address the CoSOD task.

prompt specific to the CoSOD task, which reformulates this task to closely resemble those addressed during the original pre-training, thereby inducing the foundation model to perform well on this task. Concisely, the core of our VCP is tuning only a few task-specific parameters to adapt large-scale pre-trained Transformer models for effectively addressing CoSOD tasks. For CoSOD tasks, given a set of related images $\{I_i \in \mathbb{R}^{C \times H \times W}\}_{i=1}^N$, the model is required to search for and segment intra-group co-salient objects in order to generate co-salient prediction maps $\{M_i\}_{i=1}^N$. The main insight of VCP is to enforce limited tunable parameters to focus on and generate task-specific consensus prompts, and utilize these prompts to induce the foundation model to perform effectively. As shown in Fig. 3, the proposed VCP primarily consists of two key components: the CPG and the CPD. The CPG is responsible for mining intra-group co-salient representations of the frozen embeddings to generate consensus prompts. The CPD utilizes these consensus prompts to form task-specific visual consensus prompts and induce the frozen transformer layers to perform the CoSOD task.

Segformer [32] is employed as our foundation model, which is a hierarchical Transformer architecture pretrained on ImageNet, utilized for semantic segmentation tasks. Segformer extracts multi-scale features through four stages, each comprising an overlapping patch embedding layer and several visual transformer blocks. The extracted embedding features $E = \{E_s \in \mathbb{R}^{N \times L_s \times C_s}, L_s = H_s W_s\}_{s=1}^4$ are utilized by CPG to mine intra-group co-salient representations for generating consensus prompts P_{Co} and embed-

ding prompts P_{Em} . Inspired by traditional methods [11, 24] that employ handcrafted image features to aid segmentation tasks, we utilize fast Fourier transform [20], which generates Handcrafted Prompt P_{Hand} through several predefined tunable embedding layers. CPD leverages consensus prompts to integrate three types of prompts (P_{Co} , P_{Em} and P_{Hand}) to form task-specific visual consensus prompts P_{Visual}^{Co} , thereby adaptively inducing different depths of transformer layers. We modify the original decoding part of Segformer, reconstructing a concise Prediction Head (PH) and integrating a classifier.

3.2. Consensus Prompt Generator

Effectively mining intra-group consensus representations and suppressing irrelevant feature components is crucial for specific CoSOD tasks. Therefore, we propose a CPG, which can efficiently extract co-salient representations from frozen embedding features within the group and generate consensus prompts P_{Co} . As depicted in Fig. 4, initially, the CPG initializes j learnable saliency seed tensors, which are utilized to generate prototype representations of intra-image salient objects via clustering. These prototype representations guide the model in performing the initial saliency estimation. Subsequently, leveraging the obtained saliency estimation map, we generate pixel embeddings for all salient objects, i.e., consensus seeds. From these consensus seeds, k highly relevant pixel embeddings are selected to form consensus representations, which are then mapped back to the original embeddings to generate consensus prompts P_{Co} .

The embedding features obtained by the foundation model have a high dimensionality. A tunable linear layer (with a scale r to control the tunable parameters) is used to reduce the dimensionality of the input embedding features to form the initial embedding prompt $P_{Em} \in \mathbb{R}^{N \times L_s \times C_r}$, $C_r = C_s/r$. Inspired by [7, 14, 46], we leverage saliency priors to assist in consensus extraction. However, unlike previous approaches, we no longer rely on additional large-scale saliency detection datasets to aid in training the prediction head. Instead, we leverage the concept of prototype learning by learning saliency seeds (cluster centers) within the embedding to accomplish salient object representation. Specifically, j learnable saliency seeds $S_{seed} \in \mathbb{R}^{j \times C_r}$ are pre-defined. The reshaped embedding features $P_{em} \in \mathbb{R}^{N \times C_r \times H_s \times W_s}$ are first subjected to normalization and convolution operations to form soft allocation probability scores $S_{soft} \in \mathbb{R}^{N \times j \times L_s}$.

$$S_{soft} = \nu(\text{Softmax}(\text{conv}(L_2(P_{em})))), \quad (1)$$

where L_2 and ν represent L2 normalization and reshaping operations, respectively. Then, based on the $S_{soft} \in \mathbb{R}^{N \times j \times L_s}$, the residuals $Res \in \mathbb{R}^{N \times j \times C_r \times L_s}$ between the embedding features and the saliency seeds $S_{seed} \in \mathbb{R}^{j \times C_r}$

are computed.

$$Res = (\nu(P_{em}) - \nu(S_{seed})) \times \nu(S_{assign}), \quad (2)$$

These residuals $Res \in \mathbb{R}^{N \times j \times C_r \times L_s}$ are then utilized to weightedly sum the embedding features, resulting in updated saliency seeds representations $S_{seed}^{update} \in \mathbb{R}^{N \times j \times C_r}$.

$$S_{seed}^{update} = \sum_{i=1}^{L_s} Res(N, j, C_r, i), \quad (3)$$

Finally, the updated saliency seeds are used to enhance the salient objects within the embedding features, thereby generating the saliency estimation map $\{M^s\}_{s=1}^4$.

$$M^s = conv[MLP(L_2(S_{seed}^{update})), P_{em}]. \quad (4)$$

Where $[\cdot]$ represents channel concatenation. We utilize a prototype learning-based approach to achieve attention on salient objects in complex scenes by learning saliency prototypes present in clustered embeddings. Simultaneously, we utilize CoSOD labels for multi-stage supervision to maintain consistency between the saliency estimation maps and the consensus attention targets. We utilize the obtained saliency estimation map to further filter out non-co-salient object representations in the embedding features, forming consensus prompts $P_{Co} \in \mathbb{R}^{N \times L_s \times C_r}$. Specifically, we utilize the saliency estimation map $\{M^s\}_{s=1}^4$ to reshape all embedding representations $P_{em} \in \mathbb{R}^{N \times C_r \times H_s \times W_s}$ within the group into pixel patch embeddings, namely, consensus seeds $Co_{seed} \in \mathbb{R}^{N \times L_s \times C_r}$. Next, we search within the consensus seeds for the top-k pixel embeddings as representative consensus seeds $Co_{seed}^{rep} \in \mathbb{R}^{k \times C_r}$, which exhibit the highest correlation among them.

$$score = Co_{seed} \otimes \nu(average(P_{em} \times M^s))^T, \quad (5)$$

$$Co_{seed}^{rep} = gather(Co_{seed}, argtopk(score)), \quad (6)$$

where \otimes represents matrix multiplication. Finally, the obtained representative consensus seeds are mapped back to the original embeddings to form effective consensus features, and spatial attention $S_{att} \in \mathbb{R}^{N \times 1 \times H_s \times W_s}$ is utilized to further enhance the consensus representation. The enhanced consensus features are reshaped to form consensus prompts $P_{Co} \in \mathbb{R}^{N \times L_s \times C_r}$.

$$F_{Co} = conv(L_2(P_{em}), weight = \nu(Co_{seed}^{rep})), \quad (7)$$

$$P_{Co} = \nu(conv(F_{Co} \times S_{att})). \quad (8)$$

3.3. Consensus Prompt Disperser

To effectively leverage the obtained consensus prompts to induce the foundation model to perform well on CoSOD tasks, a concise CPD is further proposed. Firstly, the obtained consensus prompts P_{Co} are used to further guide the consensus expression within the embedding prompts

P_{Em} through simple integration, denoted as Embedding Consensus Prompts $P_{Em}^{Co} \in \mathbb{R}^{N \times L_s \times C_r}$. The P_{Em}^{Co} are more effective for tuning the frozen foundation model's embedding components. Additionally, inspired by some methods that incorporate handcrafted features into models and gain benefits. For example, the EVP [20] utilizes handcrafted features derived from fast Fourier transform as Explicit Visual Prompts to achieve effective foreground segmentation. Therefore, we also consider incorporating handcrafted features P_{Hand} as part of the visual prompts to focus on high-frequency components within individual images. However, since the handcrafted prompts used in EVP only focus on individual independent images, their effectiveness is significantly discounted when dealing with CSOD tasks that require emphasis on intra-group relevance (as demonstrated in the experimental section). Taking these considerations into account, we once again utilize the obtained consensus prompts P_{Co} to highlight the consensus high-frequency components in the handcrafted prompts $P_{Hand} \in \mathbb{R}^{N \times L_s \times C_r}$, denoted as Handcrafted Consensus Prompts $P_{Hand}^{Co} \in \mathbb{R}^{N \times L_s \times C_r}$. The handcrafted prompts P_{Hand} introduced in our method are obtained directly from the input image through predefined multi-stage overlapping patch embeddings after fast Fourier transformation.

The P_{Em}^{Co} and P_{Hand}^{Co} are further integrated to obtain Visual Consensus Prompts $P_{Visual}^{Co} \in \mathbb{R}^{N \times L_s \times 2C_r}$.

$$P_{Visual}^{Co} = [P_{Em} + P_{Co}, P_{Hand} + P_{Co}], \quad (9)$$

Utilizing these visual consensus prompts to perform adaptive tunes on different depths of vision transformer layers is also crucial for inducing the foundation model to address CoSOD tasks. Therefore, multiple unshared linear layers are used to achieve adaptive tunes of different transformer layers. Then, an intra-stage shared linear layer is utilized to maintain consistency in dimensions between the prompts $P = \{P_S^n \in \mathbb{R}^{N \times L_s \times C_s}\}_{s=1}^4$ and the transformer features.

$$P_S^n = MLP_{up}(gelu(MLP_S^n(P_{Visual}^{Co}))). \quad (10)$$

Finally, the obtained task-specific prompts P will induce the foundation model to achieve effective and efficient CoSOD performance.

3.4. Prediction Head and Objective Function

The decoder of Segformer [32] utilizes multiple linear layers to up-project multi-stage features to a unified dimension (typically higher, $d=768$). The multi-scale features are concatenated in the high dimension and down-projected through linear layers to generate predictions. The decoder of Segformer possesses a significant number of parameters (3.15M), which greatly increases the number of tunable parameters, contradicting the motivation to minimize tunable parameters. Therefore, a concise prediction head is designed, primarily comprising an ASPP, FPN-like, and an

Table 1. Quantitative comparison between the proposed VCP and 13 SOTA methods on three benchmark datasets regarding six comprehensive quantitative metrics, tunable parameters, and model size. DUT-class [42], COCO-9k [19], and COCO-SEG [28] are widely used training datasets in CoSOD and we denote them as D, C and S, respectively. “↑” means that the higher the numerical value, the better the model performance. Red, blue, and green represent the top three performances, respectively.

Datasets	Metrics	EVP	CADC	GCoNet	DCFM	DMT	CoRP	CADC++	GEM	GCoNet+	MCCL	SCED	UniTR	CONDA	AOURS _I	OOURS _{II}	VS.
		CVPR ₂₃ S+D	CCV ₂₁ C+D	CVPR ₂₁ D	CVPR ₂₂ C	CVPR ₂₃ C+D	TPAMI ₂₃ C+D	TPAMI ₂₃ D	CVPR ₂₃ S+D	TPAMI ₂₃ S+D	AAAI ₂₃ S+D	ACMM ₂₃ S+D	TMM ₂₄ S+D	ECCV ₂₄ S+D	— C+D	— S+D	
CoCA	$S_m \uparrow$	0.686	0.681	0.673	0.710	0.728	0.715	0.708	0.726	0.738	0.713	0.741	0.708	0.763	0.774	0.819	+5.6%
	$E_m^{\max} \uparrow$	0.760	0.744	0.760	0.783	0.800	0.769	0.791	0.808	0.814	0.796	0.828	0.789	0.839	0.829	0.871	+3.2%
	$F_m^{\max} \uparrow$	0.546	0.548	0.544	0.598	0.621	0.597	0.609	0.599	0.637	0.584	0.629	0.574	0.685	0.680	0.752	+6.7%
	MAE \downarrow	0.126	0.132	0.105	0.085	0.107	0.103	0.107	0.095	0.081	0.097	0.084	0.089	0.089	0.069	0.054	-2.7%
	$E_m \uparrow$	0.708	0.690	0.739	0.778	0.754	0.747	—	0.767	0.783	0.764	0.804	0.766	0.790	0.813	0.830	+2.6%
	$F_m \uparrow$	0.510	0.503	0.531	0.593	0.590	0.584	—	0.566	0.612	0.560	0.610	0.559	0.640	0.660	0.708	+6.8%
CoSOD3k	$S_m \uparrow$	0.839	0.801	0.802	0.810	0.852	0.838	0.823	0.853	0.843	0.854	0.865	0.852	0.862	0.874	0.895	+3.0%
	$E_m^{\max} \uparrow$	0.854	0.840	0.860	0.874	0.895	0.890	0.876	0.911	0.901	0.902	0.923	0.903	0.911	0.918	0.938	+1.5%
	$F_m^{\max} \uparrow$	0.813	0.759	0.778	0.805	0.836	0.827	0.808	0.829	0.834	0.832	0.859	0.834	0.853	0.868	0.893	+3.4%
	MAE \downarrow	0.076	0.096	0.071	0.067	0.064	0.060	0.070	0.062	0.062	0.062	0.053	0.058	0.056	0.049	0.043	-1.0%
	$E_m \uparrow$	0.852	0.824	0.857	0.871	0.882	0.886	—	0.885	0.872	0.879	0.900	0.890	0.893	0.901	0.909	+0.9%
	$F_m \uparrow$	0.774	0.743	0.770	0.800	0.815	0.815	—	0.803	0.813	0.809	0.842	0.821	0.827	0.850	0.860	+1.8%
CoSal2015	$S_m \uparrow$	0.876	0.866	0.845	0.838	0.896	0.867	0.875	0.885	0.881	0.887	0.894	0.896	0.900	0.911	0.927	+2.7%
	$E_m^{\max} \uparrow$	0.914	0.906	0.887	0.892	0.933	0.912	0.922	0.933	0.924	0.923	0.938	0.940	0.944	0.944	0.962	+1.8%
	$F_m^{\max} \uparrow$	0.874	0.862	0.847	0.856	0.903	0.882	0.889	0.882	0.891	0.887	0.908	0.902	0.908	0.920	0.941	+3.3%
	MAE \downarrow	0.068	0.064	0.068	0.067	0.047	0.055	0.047	0.053	0.056	0.054	0.045	0.041	0.045	0.037	0.030	-1.1%
	$E_m \uparrow$	0.881	0.874	0.884	0.888	0.922	0.907	—	0.913	0.902	0.907	0.924	0.928	0.923	0.933	0.944	+1.6%
	$F_m \uparrow$	0.836	0.826	0.838	0.850	0.880	0.869	—	0.856	0.870	0.868	0.892	0.887	0.887	0.902	0.915	+2.3%
Tunable Param. (M)	3.7	392.85	280.36	142.3	40.4	20	393.21	52.3	18.4	27.1	156.7	146.6	24.1	4.94	4.94		
Model Size (MB)	14.1	1498.7	541.7	542.9	154.4	228.3	—	199.7	70.4	104.5	1750	541	139.5	19	19		

Table 2. Ablation analysis on the main components of our VCP architecture. “BL” represents the baseline, which utilizes solely the frozen foundation model and the tunable FPN-like decoder. P_{Hand} represents $P_{Em}+P_{Hand}$ and P_{Co} represents $P_{Em}+P_{Co}$.

Combination				Tunable	Model Size	CoCA			CoSOD3k			CoSal2015		
BL	PH	P_{Hand}	P_{Co}	Param. (M)	(MB)	$S_m \uparrow$	$F_m^{\max} \uparrow$	MAE \downarrow	$S_m \uparrow$	$F_m^{\max} \uparrow$	MAE \downarrow	$S_m \uparrow$	$F_m^{\max} \uparrow$	MAE \downarrow
✓				0.14	0.52	0.635	0.470	0.138	0.777	0.749	0.099	0.813	0.814	0.099
✓	✓			1.49	5.73	0.659	0.498	0.125	0.832	0.795	0.073	0.866	0.862	0.063
✓	✓	✓		2.04	7.87	0.704	0.561	0.111	0.852	0.831	0.067	0.894	0.893	0.054
✓	✓		✓	4.59	17.67	0.763	0.662	0.073	0.873	0.870	0.052	0.911	0.923	0.043
✓	✓	Addition		4.72	18.14	0.755	0.646	0.080	0.873	0.861	0.048	0.910	0.916	0.036
✓	✓	Concatenation		4.94	19.02	0.774	0.680	0.069	0.874	0.868	0.049	0.911	0.920	0.037

additional linear classifier. This prediction head accepts a smaller unified projection dimension and has fewer tunable parameters (1.49M).

The prediction head generates a final prediction map M and classification predictions V_{cla} . We adopt a weighted combination L_w of commonly used binary cross-entropy loss (BCE) and IOU loss to constrain the prediction map. The predicted classes are supervised using cross-entropy loss L_{ce} . Simultaneously, we utilize the saliency estimation maps generated by the CPG in four stages as initial co-salient object predictions $\{M^s\}_{s=1}^4$ to constrain them towards the final prediction targets. The overall loss can be expressed as:

$$L = \alpha L_w(M, GT) + \beta \sum_{s=1}^4 L_w(M^s, GT) + \lambda L_{ce}(V_{cla}, V_{lab}). \quad (11)$$

Where α , β , and λ are set to 10, 2, and 0.1, respectively.

4. Experiments

4.1. Datasets and Evaluation Metrics

Datasets. All experiments are evaluated based on the most commonly used three datasets for CoSOD tasks: CoCA [42], CoSOD3k [7], and CoSal2015 [36]. The three primary training datasets widely used in CoSOD tasks are the DUT-class [42], COCO-9k [19], and COCO-SEG [28] datasets, and we denote them as D, C, and S, respectively. The combinations of training datasets primarily used by existing methods include D+C [14, 16, 40, 46], as well as D+S [17, 31, 33, 44]. In the comparative experiments, we conduct experiments using both combinations of training data to ensure fairness of comparison. In all ablation experiments, we use the DUT-class and COCO-9k (D+C) datasets as the training set.

Evaluation Metrics. Six evaluation metrics are used for comparison, including MAE [4], Mean-Em (E_m), E_m^{\max}

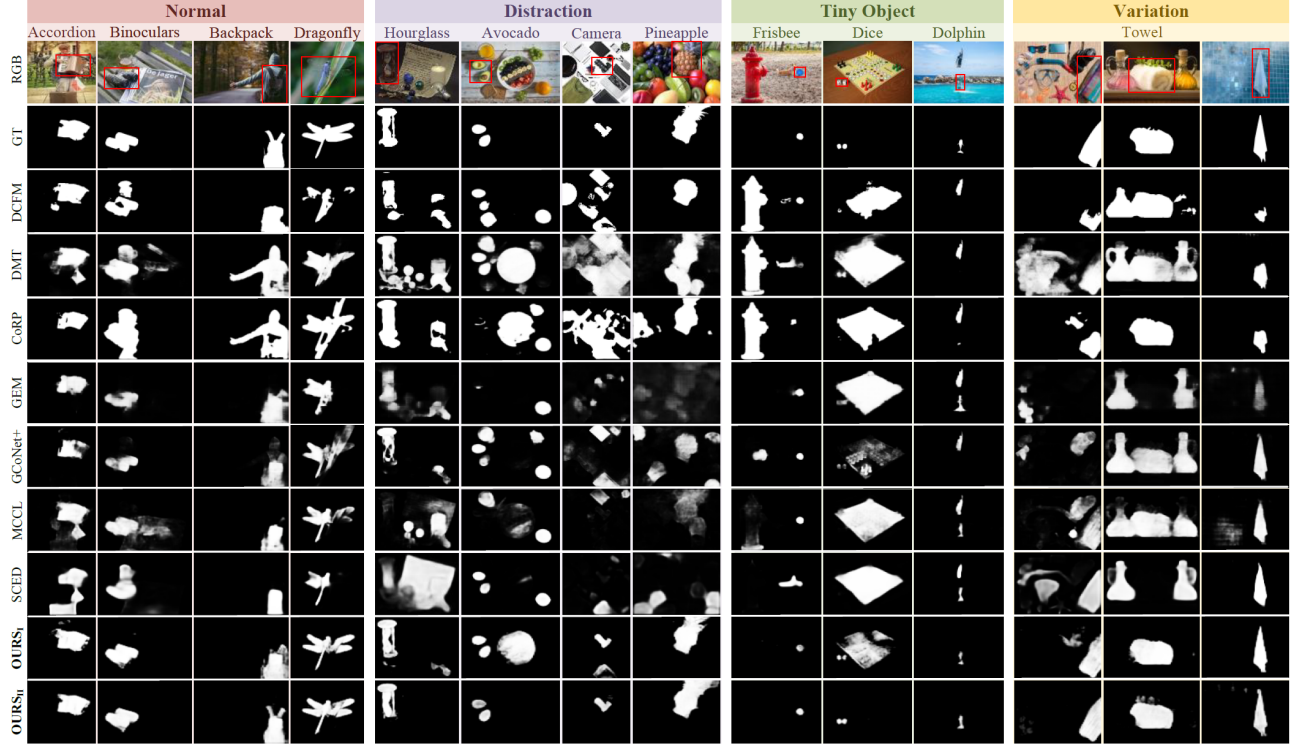


Figure 5. Visual comparison between our VCP and the most representative seven methods across four scenarios.

Table 3. Ablation on the CPD details and some auxiliary components in the architecture.

Settings	Tunable Param. (M)	Model Size (MB)	COCA			CoSOD3k			CoSal2015		
			$S_m \uparrow$	$F_m^{\max} \uparrow$	MAE \downarrow	$S_m \uparrow$	$F_m^{\max} \uparrow$	MAE \downarrow	$S_m \uparrow$	$F_m^{\max} \uparrow$	MAE \downarrow
w/o $L_{ce}(V_{cla}, V_{lab})$	4.92	18.93	0.761	0.653	0.078	0.872	0.864	0.051	0.911	0.919	0.038
w/o $L_w(M^s, GT)$	4.94	19.02	0.742	0.634	0.096	0.870	0.853	0.055	0.911	0.915	0.037
w/ Segformer Head	3.84	14.73	0.747	0.656	0.101	0.872	0.861	0.054	0.909	0.916	0.038
w/ Share MLP	4.53	17.41	0.768	0.669	0.067	0.871	0.865	0.051	0.905	0.914	0.042
w/ Unshare MLP	5.78	22.24	0.763	0.655	0.076	0.873	0.865	0.051	0.911	0.918	0.038
OURS(Adaptive MLP)	4.94	19.02	0.774	0.680	0.069	0.874	0.868	0.049	0.911	0.920	0.037

[6], S_m [5], Mean-Fm (F_m) and F_m^{\max} [1], which are to assess the average pixel-wise absolute difference, local and global similarity, structural similarity between the predictions and the ground truths, and the weighted harmonic mean of precision and recall, respectively.

4.2. Implementation Details

The proposed method is based on the Pytorch framework while using a pre-trained model of Segformer [32] on the ImageNet dataset, and all experiments are conducted using a single NVIDIA 3090 GPU. We randomly pick N samples from three different groups in each training batch.

$$N = \min(N_g(A), N_g(B), N_g(C), 16). \quad (12)$$

Where N_g means the number of images in the corresponding group. For inference, all samples in each group are input at one time. We train our network using the AdamW optimizer, the initial learning rate is initially set to $5e-4$, and

cosine decay is applied to the learning rate. When training with D+C, we train the model for 100 epochs with a learning rate decay to $1e-4$, and the total training time is around 9 hours and the inference time is around 65.3 FPS. When training with D+S, we train for 200 epochs with a learning rate decay to $1e-5$. The inputs are resized into 288×288 for both training and inference.

4.3. Comparisons with State-of-the-art Methods

To demonstrate the effectiveness of our VCP, we compare it with 13 SOTA CoSOD methods from the past three years, including SCED [33], MCCL [44], GCoNet+ [43], GEM [31], CAD++ [41], CoPR [46], DMT [16], DCFM [38], GCoNet [8], CAD [40], UniTR [10], CONDA [17] and EVP [20]. All methods except EVP are based on Full fine-tuning for CoSOD task. EVP addresses foreground segmentation by introducing simple visual prompts. Since there are currently no methods that utilize prompt learning for

Table 4. Ablation on the proposed visual consensus prompts P_{Visual}^{Co} . We introduce other visual prompts proposed in some parameter-efficient prompt tuning methods to replace our VCP in performing CoSOD tasks ($prompt = Adaptformer/VPT/EVP$).

Prompt	Tunable Param. (M)	Model Size (MB)	COCA			CoSOD3k			CoSal2015		
			$S_m \uparrow$	$F_m^{max} \uparrow$	MAE \downarrow	$S_m \uparrow$	$F_m^{max} \uparrow$	MAE \downarrow	$S_m \uparrow$	$F_m^{max} \uparrow$	MAE \downarrow
Ours(Full-tuning)	65.79	251.43	0.740	0.624	0.093	0.866	0.854	0.055	0.905	0.911	0.040
Ours(Only PH)	1.49	5.73	0.659	0.498	0.125	0.832	0.795	0.073	0.866	0.862	0.063
Prompt=Adaptformer [3]	1.54	5.99	0.695	0.547	0.111	0.852	0.831	0.066	0.883	0.881	0.061
Prompt=VPT-Deep [12]	1.60	6.20	0.700	0.554	0.106	0.853	0.832	0.067	0.889	0.889	0.060
Prompt=EVP [20]	2.04	7.87	0.704	0.561	0.111	0.852	0.831	0.067	0.894	0.893	0.054
EVP(S+D) [20]	3.70	14.10	0.686	0.546	0.126	0.839	0.813	0.076	0.876	0.874	0.068
VCP($r = 8$)	3.34	12.89	0.769	0.665	0.069	0.872	0.864	0.051	0.913	0.918	0.035
VCP($r = 4$)	4.94	19.02	0.774	0.680	0.069	0.874	0.868	0.049	0.911	0.920	0.037

CoSOD tasks, we retrain EVP [20] to attempt solving the CoSOD problem. This serves as compelling evidence to validate the effectiveness of our visual consensus prompts for CoSOD tasks. **Quantitative Results:** Table 1 presents a quantitative comparison of our VCP with the most representative works in the past three years across multiple metrics. Compared to 12 full fine-tuning CoSOD methods, our VCP demonstrates a significant performance advantage on three commonly used benchmark test sets. Compared to the prompt-based tuning method EVP [20], our VCP exhibits overwhelmingly superior performance. **Qualitative Results:** Fig. 5 illustrates the visual comparisons between our VCP and seven representative works across four selected scenarios (“Normal”, “Distraction”, “Tiny Object” and “Variation”). Through the visual comparisons in Fig. 5, it can be observed that our VCP achieves more effective localization and segmentation performance even when facing challenging cases such as multiple salient object interferences, small or diverse co-salient objects.

4.4. Ablation Study

We conduct ablation experiments to demonstrate the effectiveness of each component of the proposed VCP. All ablation experiments are conducted using the C+D dataset as the training set, tested on three mainstream benchmark datasets.

Effectiveness of the architecture designs. Table 2 shows the performance contributions of the main components in our architecture. “BL” represents the baseline, which utilizes solely the frozen foundation model and the tunable FPN-like decoder. P_{Hand} represents $P_{Em} + P_{Hand}$ and P_{Co} represents $P_{Em} + P_{Co}$. “Addition” represents the addition to obtain visual consensus prompts, i.e., $P_{Visual}^{Co} = P_{Em}^{Co} + P_{Hand}^{Co}$. It can be observed that the introduced hand-crafted prompts show certain improvement but the enhancement is limited, while our proposed consensus prompts significantly enhance the effectiveness of the model.

Effectiveness of CPD and auxiliary components. Table 3 presents the implementation details of CPD and some ablations of auxiliary components in the architecture. “ $L_{ce}(V_{cla}, V_{lab})$ ” denotes the removal of the classifier loss, and “ $L_w(M^s, GT)$ ” denotes the removal of supervision for

the multiscale predictions generated in CPG. Furthermore, we further execute two schemes to validate the effectiveness of CPD in performing adaptive guidance after obtaining visual consensus prompts. For different depths (n) of Transformer layers within a specific stage, our CPD employs an adaptive tuning scheme, i.e., shared MLP_{up} and unshared MLP_S^n , while “Share MLP” and “Unshare MLP” respectively indicate whether MLP_{up} and MLP_S^n are shared or not. We employ an adaptive approach to disperse visual consensus prompts P_{Visual}^{Co} , achieving a balance between parameter efficiency and effective performance.

Effectiveness of the proposed CPG. To validate the effectiveness of the consensus prompts generated by our CPG, we replace our CPG with other visual prompts (*Adaptformer/VPT/EVP*) from state-of-the-art prompt tuning methods. As seen in Table 4, our VCP significantly outperforms existing visual prompt schemes by a large margin in addressing specific CoSOD tasks, and it can comprehensively outperform our full fine-tuning scheme.

Additional experimental details and scalability validations are provided in the **supplementary materials**.

5. Conclusion

In this work, we propose VCP, a interaction-effective and parameter-efficient visual prompt tuning framework for CoSOD, serving as a potent alternative to existing methods based on the common architecture and the full fine-tuning paradigm. Additionally, our formulated visual consensus prompts (with minimized tunable parameters) effectively and efficiently induce the frozen pre-trained foundation model to play a crucial role in CoSOD tasks. Our CPG effectively enforces limited tunable parameters to focus on intra-group co-salient representations and generate key consensus prompts. Similarly, our CPD accomplishes its mission of adaptively inducing the foundation model using consensus prompts. Extensive experiments demonstrate that our VCP outperforms existing state-of-the-art visual prompt schemes by a large margin when addressing specific CoSOD tasks. Compared to state-of-the-art full fine-tuning methods tailored for CoSOD tasks, our VCP still maintains a comprehensive competitive advantage.

Acknowledgment. This work is supported by the National Nature Science Foundation of China (Nos.62376186).

References

- [1] Radhakrishna Achanta, Sheila Hemami, Francisco Estrada, and Sabine Susstrunk. Frequency-tuned salient region detection. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 1597–1604, 2009. 7
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Proc. Adv. Neural Inform. Process. Syst.*, pages 1877–1901, 2020. 2, 3
- [3] Shoufa Chen, Chongjian Ge, Zhan Tong, Jiangliu Wang, Yibing Song, Jue Wang, and Ping Luo. Adaptformer: Adapting vision transformers for scalable visual recognition. In *Proc. Adv. Neural Inform. Process. Syst.*, pages 16664–16678, 2022. 2, 3, 8
- [4] Ming-Ming Cheng, Jonathan Warrell, Wen-Yan Lin, Shuai Zheng, Vibhav Vineet, and Nigel Crook. Efficient salient region detection with soft image abstraction. In *Proc. IEEE/CVF Int. Conf. Comput. Vision*, pages 1529–1536, 2013. 6
- [5] Deng-Ping Fan, Ming-Ming Cheng, Yun Liu, Tao Li, and Ali Borji. Structure-measure: A new way to evaluate foreground maps. In *Proc. IEEE/CVF Int. Conf. Comput. Vision*, pages 4548–4557, 2017. 7
- [6] Deng-Ping Fan, Cheng Gong, Yang Cao, Bo Ren, Ming-Ming Cheng, and Ali Borji. Enhanced-alignment measure for binary foreground map evaluation. In *Proc. Int. Joint Conf. Artif. Intell.*, pages 698–704, 2018. 7
- [7] Deng-Ping Fan, Tengteng Li, Zheng Lin, Ge-Peng Ji, Dingwen Zhang, Ming-Ming Cheng, Huazhu Fu, and Jianbin Shen. Re-thinking co-salient object detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(8):4339–4354, 2021. 4, 6
- [8] Qi Fan, Deng-Ping Fan, Huazhu Fu, Chi-Keung Tang, Ling Shao, and Yu-Wing Tai. Group collaborative learning for co-salient object detection. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 12288–12298, 2021. 1, 3, 7
- [9] Guoqiang Gong, Xinghan Wang, Yadong Mu, and Qi Tian. Learning temporal co-attention models for unsupervised video action localization. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 9819–9828, 2020. 1
- [10] Ruohao Guo, Xianghua Ying, Yanyu Qi, and Liao Qu. Unitr: A unified transformer-based framework for co-object and multi-modal saliency detection. *IEEE Trans. Multimedia*, 26:7622–7635, 2024. 7
- [11] Hailing Huang, Weiqiang Guo, and Yu Zhang. Detection of copy-move forgery in digital images using sift algorithm. In *Proc. IEEE Pacific-Asia Workshop Comput. Intell. Ind. Appl.*, pages 272–276, 2008. 4
- [12] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *Proc. Eur. Conf. Comput. Vision*, pages 709–727. Springer, 2022. 2, 3, 8
- [13] Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *Trans. Assoc. Comput. Linguistics*, 8:423–438, 2020. 3
- [14] Wen-Da Jin, Jun Xu, Ming-Ming Cheng, Yi Zhang, and Wei Guo. Icnnet: Intra-saliency correlation network for co-saliency detection. In *Proc. Adv. Neural Inform. Process. Syst.*, pages 18749–18759, 2020. 3, 4, 6
- [15] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proc. Conf. Empirical Methods Nat. Lang. Process.*, pages 3045–3059, 2021. 2, 3
- [16] Long Li, Junwei Han, Ni Zhang, Nian Liu, Salman Khan, Hisham Cholakkal, Rao Muhammad Anwer, and Fahad Shahbaz Khan. Discriminative co-saliency and background mining transformer for co-salient object detection. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 7247–7256, 2023. 2, 3, 6, 7
- [17] Long Li, Nian Liu, Dingwen Zhang, Zhongyu Li, Salman Khan, Rao Anwer, Hisham Cholakkal, Junwei Han, and Fahad Shahbaz Khan. Conda: Condensed deep association learning for co-salient object detection. In *Proc. Eur. Conf. Comput. Vision*, 2024. 6, 7
- [18] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proc. Annu. Meet. Assoc. Comput. Linguistics Int. Joint Conf. Nat. Lang. Process.*, pages 4582–4597, 2021. 2, 3
- [19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proc. Eur. Conf. Comput. Vision*, pages 740–755. Springer, 2014. 2, 6
- [20] Weihuang Liu, Xi Shen, Chi-Man Pun, and Xiaodong Cun. Explicit visual prompting for low-level structure segmentations. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 19434–19445, 2023. 2, 3, 4, 5, 7, 8
- [21] Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. *arXiv preprint arXiv:2110.07602*, 2021. 2, 3
- [22] Ziyang Luo, Nian Liu, Wangbo Zhao, Xuguang Yang, Dingwen Zhang, Deng-Ping Fan, Fahad Khan, and Junwei Han. Vscope: General visual salient and camouflaged object detection with 2d prompt learning. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 17169–17180, 2024. 2, 3
- [23] Xing Nie, Bolin Ni, Jianlong Chang, Gaofeng Meng, Chunlei Huo, Shiming Xiang, and Qi Tian. Pro-tuning: Unified prompt tuning for vision tasks. *IEEE Trans. Circuits Syst. Video Technol.*, 2023. 2, 3
- [24] Alin C Popescu and Hany Farid. Exposing digital forgeries by detecting traces of resampling. *IEEE Trans. Signal Process.*, 53(2):758–767, 2005. 4
- [25] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018. 2, 3

- [26] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 2, 3
- [27] Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. In *Proc. Conf. Empirical Methods Nat. Lang. Process.*, pages 4222–4235, 2020. 3
- [28] Chong Wang, Zheng-Jun Zha, Dong Liu, and Hongtao Xie. Robust deep co-saliency detection with group semantic. In *Proc. AAAI Conf. Artif. Intell.*, pages 8917–8924, 2019. 2, 6
- [29] Jie Wang, Nana Yu, Zihao Zhang, and Yahong Han. Single-group generalized rgb and rgb-d co-salient object detection. *IEEE Trans. Circuits Syst. Video Technol.*, 2024. 3
- [30] Wenguan Wang, Qiuxia Lai, Huazhu Fu, Jianbing Shen, Haibin Ling, and Ruigang Yang. Salient object detection in the deep learning era: An in-depth survey. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(6):3239–3259, 2021. 1
- [31] Yang Wu, Huihui Song, Bo Liu, Kaihua Zhang, and Dong Liu. Co-salient object detection with uncertainty-aware group exchange-masking. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 19639–19648, 2023. 2, 3, 6, 7
- [32] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *Proc. Adv. Neural Inform. Process. Syst.*, pages 12077–12090, 2021. 4, 5, 7
- [33] Peiran Xu and Yadong Mu. Co-salient object detection with semantic-level consensus extraction and dispersion. In *Proc. ACM Int. Conf. Multimedia*, pages 2744–2755, 2023. 1, 2, 3, 6, 7
- [34] Xi Yang, Shaoyi Li, Jun Ma, Jun-yan Yang, and Jie Yan. Co-saliency-regularized correlation filter for object tracking. *Signal Processing: Image Communication*, 103:116655, 2022. 1
- [35] Siyue Yu, Jimin Xiao, Bingfeng Zhang, and Eng Gee Lim. Democracy does matter: Comprehensive feature mining for co-salient object detection. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 979–988, 2022. 1, 3
- [36] Dingwen Zhang, Junwei Han, Chao Li, and Jingdong Wang. Co-saliency detection via looking deep and wide. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 2994–3002, 2015. 6
- [37] Dingwen Zhang, Junwei Han, Chao Li, Jingdong Wang, and Xuelong Li. Detection of co-salient objects by looking deep and wide. *Int. J. Comput. Vis.*, 120:215–232, 2016. 1
- [38] Kaihua Zhang, Tengteng Li, Bo Liu, and Qingshan Liu. Co-saliency detection via mask-guided fully convolutional networks with multi-scale label smoothing. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 3095–3104, 2019. 7
- [39] Kaihua Zhang, Yang Wu, Mingliang Dong, Bo Liu, Dong Liu, and Qingshan Liu. Deep object co-segmentation and co-saliency detection via high-order spatial-semantic network modulation. *IEEE Trans. Multimedia*, 25:5733–5746, 2022. 1
- [40] Ni Zhang, Junwei Han, Nian Liu, and Ling Shao. Summarize and search: Learning consensus-aware dynamic convolution for co-saliency detection. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 4167–4176, 2021. 1, 3, 6, 7
- [41] Ni Zhang, Nian Liu, Fang Nan, and Junwei Han. Cadc++: Advanced consensus-aware dynamic convolution for co-salient object detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, 46(5):2741–2757, 2023. 2, 3, 7
- [42] Zhao Zhang, Wenda Jin, Jun Xu, and Ming-Ming Cheng. Gradient-induced co-saliency detection. In *Proc. Eur. Conf. Comput. Vision*, pages 455–472. Springer, 2020. 1, 2, 3, 6
- [43] Peng Zheng, Huazhu Fu, Deng-Ping Fan, Qi Fan, Jie Qin, Yu-Wing Tai, Chi-Keung Tang, and Luc Van Gool. Gconet+: A stronger group collaborative co-salient object detector. *IEEE Trans. Pattern Anal. Mach. Intell.*, 45(9):10929–10946, 2023. 2, 3, 7
- [44] Peng Zheng, Jie Qin, Shuo Wang, Tian-Zhu Xiang, and Huan Xiong. Memory-aided contrastive consensus learning for co-salient object detection. In *Proc. AAAI Conf. Artif. Intell.*, pages 3687–3695, 2023. 1, 2, 3, 6, 7
- [45] Jiawen Zhu, Simiao Lai, Xin Chen, Dong Wang, and Huchuan Lu. Visual prompt multi-modal tracking. In *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pages 9516–9526, 2023. 2, 3
- [46] Ziyue Zhu, Zhao Zhang, Zheng Lin, Xing Sun, and Ming-Ming Cheng. Co-salient object detection with co-representation purification. *IEEE Trans. Pattern Anal. Mach. Intell.*, 45(7):8193–8205, 2023. 1, 2, 3, 4, 6, 7