

# Unsupervised Salient Object Detection with Spectral Cluster Voting

Gyungin Shin<sup>1</sup>, Samuel Albanie<sup>2</sup>, Weidi Xie<sup>1,3</sup>

<sup>1</sup>Visual Geometry Group, Department of Engineering Science University of Oxford, UK

<sup>2</sup>Department of Engineering, University of Cambridge, UK

<sup>3</sup>Cooperative Medianet Innovation Center, Shanghai Jiao Tong University, China

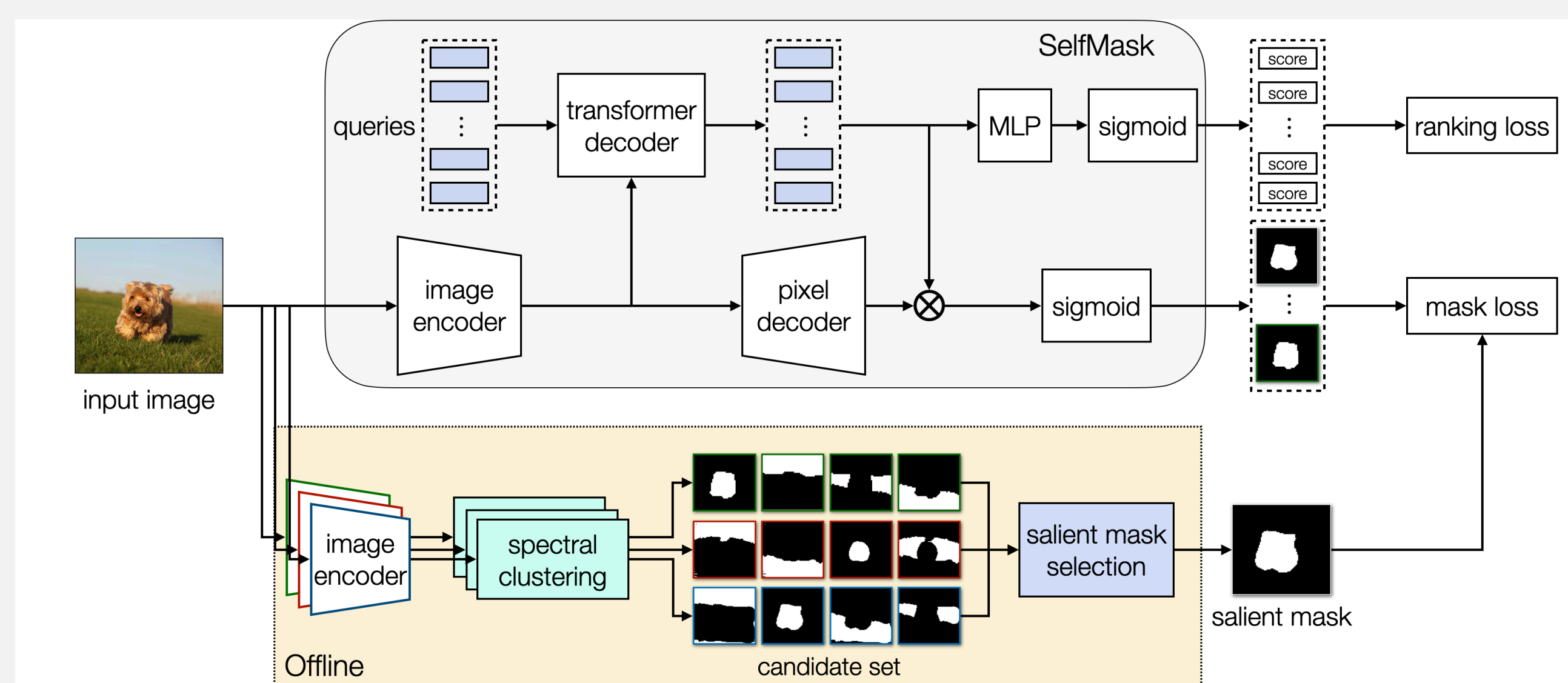
gyungin@robots.ox.ac.uk, sma71@cam.ac.uk, weidi@sjtu.edu.cn code: <https://github.com/NoelShin/selfmask>



## Overview & contribution

- We revisit spectral clustering and demonstrate its potential for discovering salient objects across various self-supervised features, e.g., MoCov2, SwAV, and DINO.
- Given mask proposals from multiple applications of spectral clustering on different self-supervised features, we pick the most salient mask with a proposed winner-takes-all voting which leverages framing and distinctiveness priors for filtering non-salient masks.
- Using the selected object segmentation as pseudo groundtruth masks, we train a salient object detector, termed SELFmask, and show that the model outperforms prior approaches on three unsupervised SOD benchmarks.

## Proposed method



Given several different self-supervised encoders,

- We first generate a set of pseudo-mask candidates per image using spectral clustering. In the figure, we show 12 masks from clusterings ( $k=4$ ) on three different encoder features.
- We select the most salient mask among them via the proposed winner-takes-all voting strategy and use it as a pseudo-mask for the image.
- Then we train SELFmask on the pseudo-masks with two loss functions: a **Mask loss** and a **Ranking loss**. **Mask loss** encourages multiple predictions made by the model to be similar to the pseudo-mask. At the same time, the model is tasked to output an objectness score for each predicted mask such that a score of a prediction closer to the salient mask to be higher via the **Ranking loss**.
- At inference time, we only select a prediction with the highest objectness score.

## Experiments

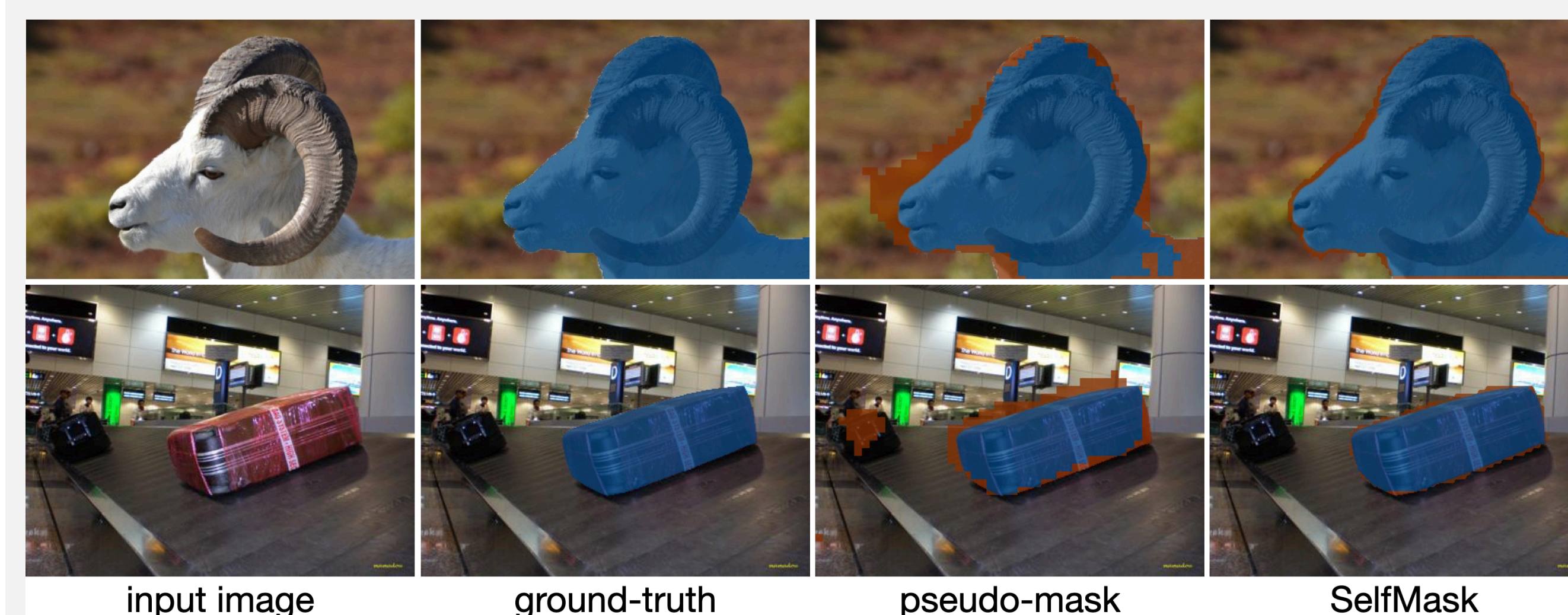


Figure 2: **Sample visualisations of the pseudo-masks and predictions from our model.** Blue and orange coloured regions denote the intersection and difference between a ground-truth and a predicted mask.

Model	Arch.	Cluster.	DUT-OMRON $k=\{2, 3, 4\}$	DUTS-TE $k=\{2, 3, 4\}$	ECSSD $k=\{2, 3, 4\}$
<b>Convolutional Nets</b>					
MoCov2	ResNet50	$k$ -means	.375	.415	.500
MoCov2	ResNet50	spectral	<b>.387</b>	<b>.454</b>	<b>.627</b>
SwAV	ResNet50	$k$ -means	.399	.444	.542
SwAV	ResNet50	spectral	<b>.401</b>	<b>.458</b>	<b>.590</b>
<b>Vision Transformer</b>					
DINO	ViT-S/16	$k$ -means	.377	.392	.541
DINO	ViT-S/16	spectral	<b>.394</b>	<b>.417</b>	<b>.577</b>
DINO	ViT-S/8	$k$ -means	.369	.377	.551
DINO	ViT-S/8	spectral	<b>.398</b>	<b>.411</b>	<b>.587</b>

Table 1: **Spectral clustering dominates  $k$ -means for self-supervised features.** We report upper bound IoUs to compare the quality of masks produced by  $k$ -means and spectral clustering on *self*-supervised features with two different encoder architectures. We report the average of the results from  $k=\{2, 3, 4\}$ .

Features	$k$	Pseudo-mask UB
DINO MoCov2 SwAV		
✗ ✓ ✓	2	.508 .562
	2, 3	.561 .626
	2, 3, 4	.580 .658
✓ ✗ ✓	2	.473 .553
	2, 3	.538 .644
	2, 3, 4	.559 .682
✓ ✓ ✗	2	.459 .546
	2, 3	.536 .648
	2, 3, 4	.566 .688
✓ ✓ ✓	2	.511 .584
	2, 3	.567 .664
	2, 3, 4	<b>.590 .698</b>

Table 2: **Forming a candidate set with various self-supervised features and multiple  $k$  values improves IoU of both pseudo-masks and upper bound masks (UB).** We compare cases with different combinations of self-supervised features and cluster numbers of  $k=2$ ,  $\{2, 3\}$  or  $\{2, 3, 4\}$  on HKU-IS.

Selection	Framing prior	HKU-IS SOD
random	✗	.206 .197
	✓	.464 .277
center	✗	.362 .122
	✓	.442 .392
voting (ours)	✗	.081 .200
	✓	<b>.590 .447</b>

Table 3: **Winner-takes-all voting and the framing prior both significantly improve mask quality.** We compare our voting strategy to different selection strategies along with the effect of framing prior under the IoU metric. Selection is performed from a candidate set including DINO, MoCov2 and SwAV features with  $k = \{2, 3, 4\}$ .

Model	DUT-OMRON			DUTS-TE			ECSSD		
	Acc	IoU	$\max F_\beta$	Acc	IoU	$\max F_\beta$	Acc	IoU	$\max F_\beta$
HS	.843	.433	.561	.826	.369	.504	.847	.508	.673
wCtr	.838	.416	.541	.835	.392	.522	.862	.517	.684
WSC	.865	.387	.523	.862	.384	.528	.852	.498	.683
DeepUSPS	.779	.305	.414	.773	.305	.425	.795	.440	.584
BigBiGAN	.856	.453	.549	.878	.498	.608	.899	.672	.782
E-BigBiGAN	.860	.464	.563	.882	.511	.624	.906	.684	.797
Melas-Kyriazi et al.	.883	.509	-	.893	.528	-	.915	.713	-
LOST	.797	.410	.473	.871	.518	.611	.895	.654	.758
LOST <sup>†</sup>	.818	.489	.578	.887	.572	.697	.916	.723	.837
TokenCut	.880	.533	.600	.903	.576	.672	.918	.712	.803
TokenCut <sup>†</sup>	.897	.618	.697	.914	.624	.755	.934	.772	.874
pseudo-masks (Ours)	.811	.403	-	.845	.466	-	.893	.646	-
SELFmask (Ours)	.901	.582	.680	.923	.626	.750	.944	.781	.889
SELFmask <sup>†</sup> (Ours)	<b>.919</b>	<b>.655</b>	<b>.852</b>	<b>.933</b>	<b>.660</b>	<b>.882</b>	<b>.955</b>	<b>.818</b>	<b>.956</b>

Table 4: **Comparison to state-of-the-art unsupervised saliency detection methods on three salient object detection benchmarks.** We observe that SELFmask yields improved performance over prior state-of-the-art approaches across all benchmarks when trained with the pseudo-masks for the DUTS training images. The best score per column is highlighted in bold. <sup>†</sup> applies Bilateral solver for post-processing.

## Conclusion

- In this work, we address the challenging problem of unsupervised salient object detection (SOD).
- For this, we first observe that self-supervised features exhibit significantly greater object segmentation potential with spectral clustering than with  $k$ -means.
- Inspired by this observation, we extract foreground regions among multiple masks generated from different self-supervised features, and varying cluster numbers based on the winner-takes-all voting.
- By using the selected masks as pseudo-masks, we train a saliency detection network and show promising results compared to previous unsupervised methods on various SOD benchmarks.