

Poster: THU-AM-286

JUNE 18-22, 2023

CVPR



VANCOUVER, CANADA



Highlight

Interactive Segmentation as Gaussian Process Classification

Minghao Zhou^{1,2}, Hong Wang² (✉), Qian Zhao¹, Yuexiang Li²,
Yawen Huang², Deyu Meng^{1,3,4}, Yefeng Zheng²

¹Xi'an Jiaotong University ²Tencent Jarvis Lab ³Pengcheng Laboratory

⁴Macau University of Science and Technology



腾讯天衍实验室
TENCENT JARVIS

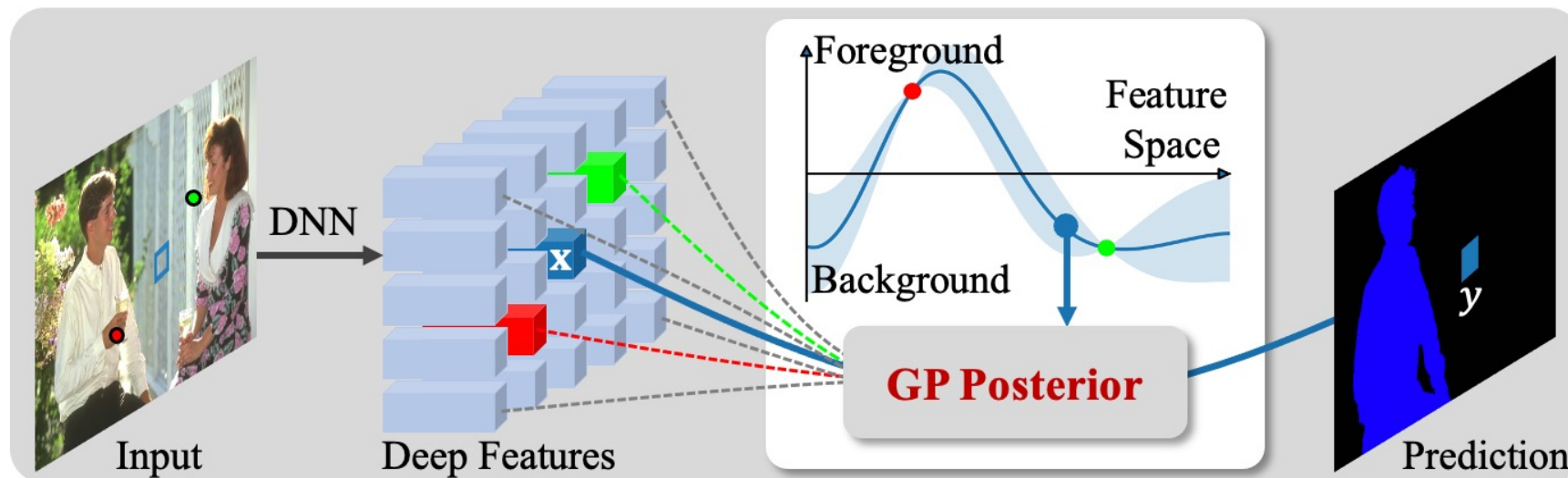


Overview

Overview

Brand-new Perspective for Interactive Segmentation

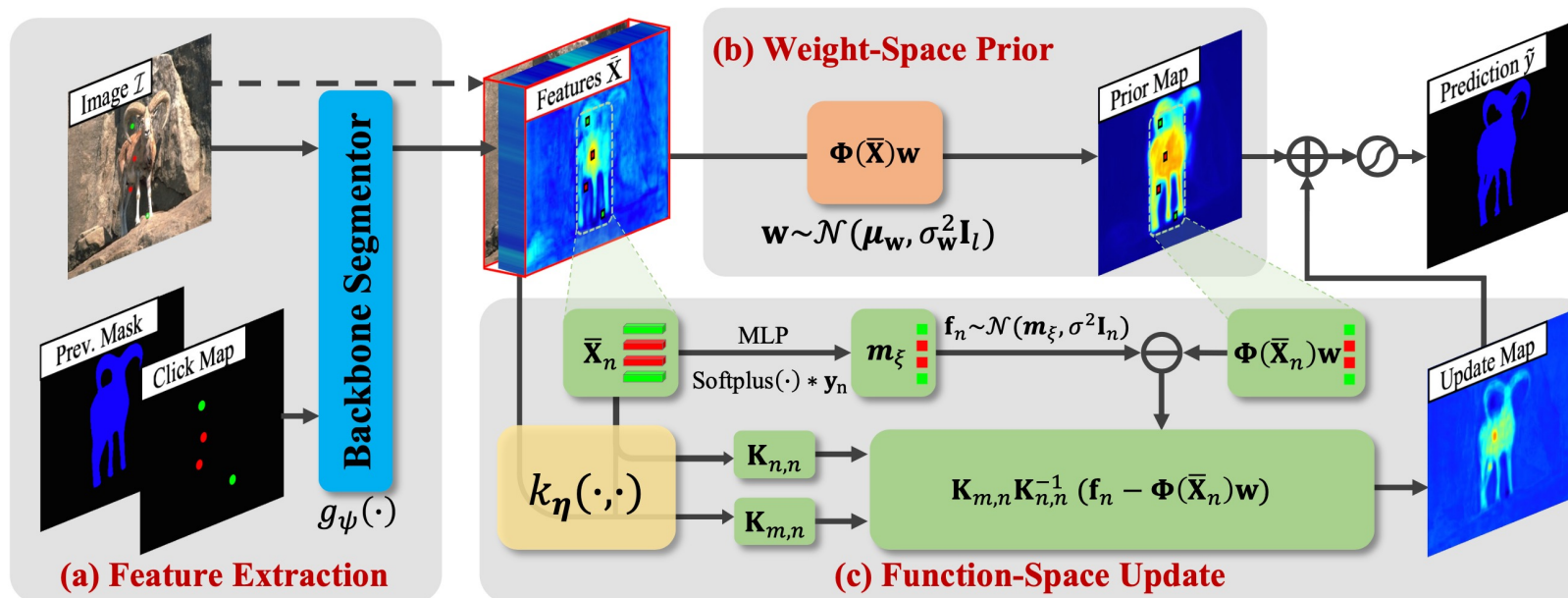
For each image, we model interactive segmentation as a Gaussian process classification task, with clicked/unclicked pixels as training/testing data.



Overview

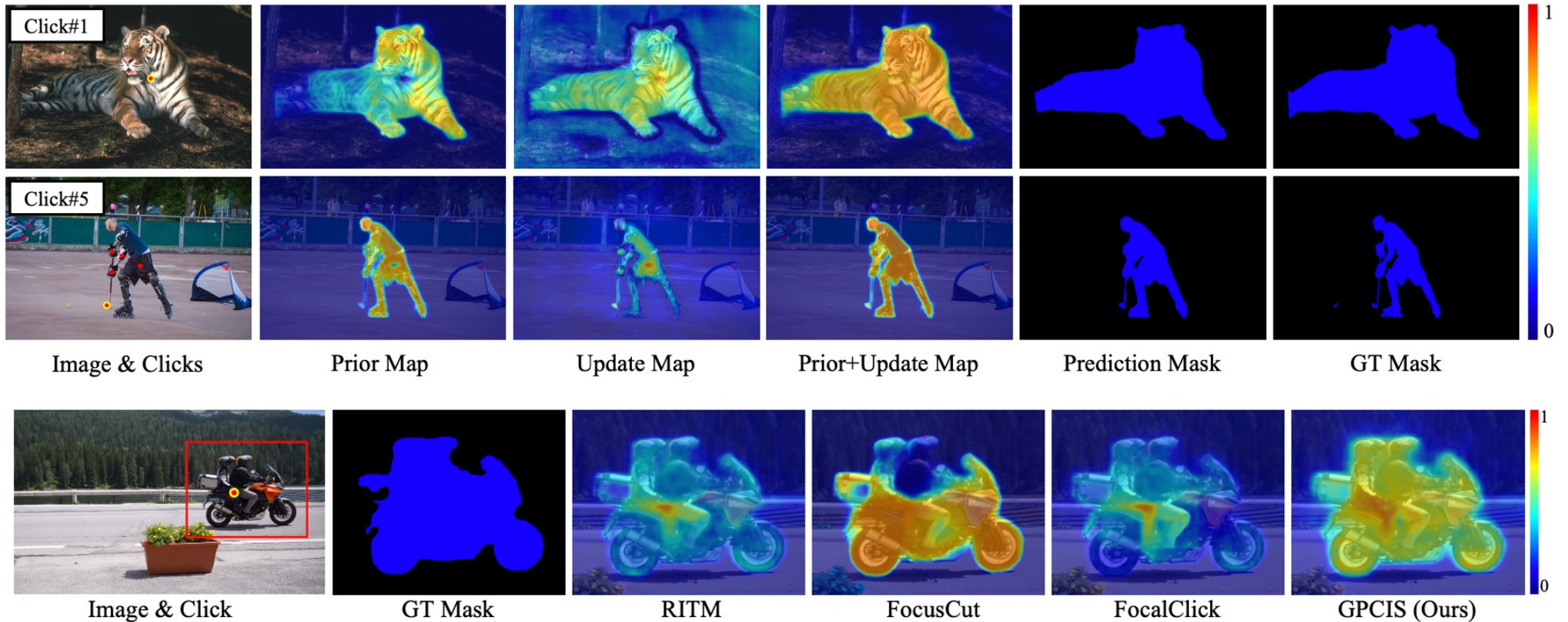
Gaussian Process Classification-based Interactive Segmentation (GPCIS) Model

- Explicitly guide the information propagation procedure;
- Provide theoretical support for accurate predictions at clicks;
- Concise framework and clear working mechanism.



Overview

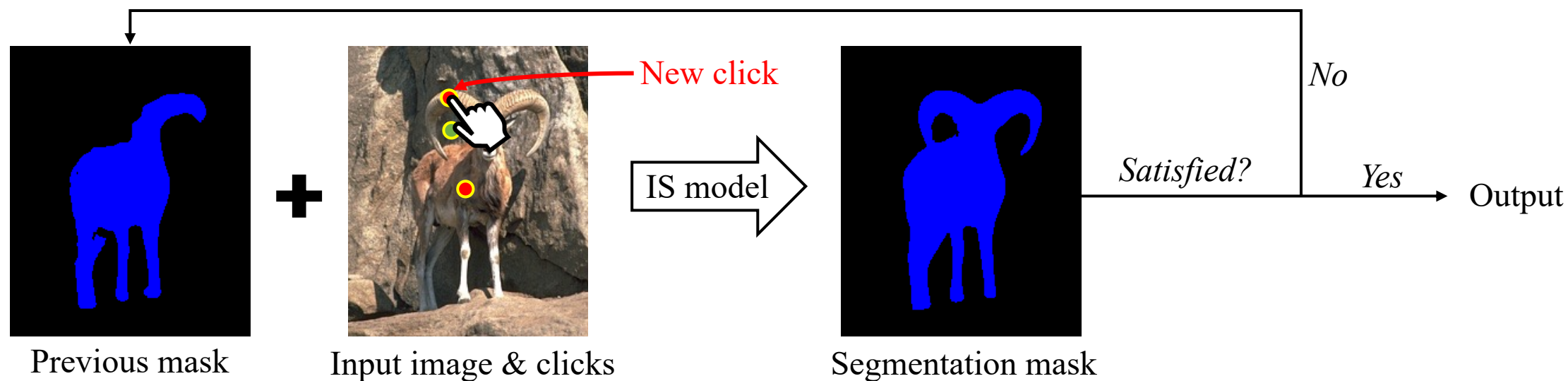
Experiments for Model Verification and Performance Evaluation



Introduction

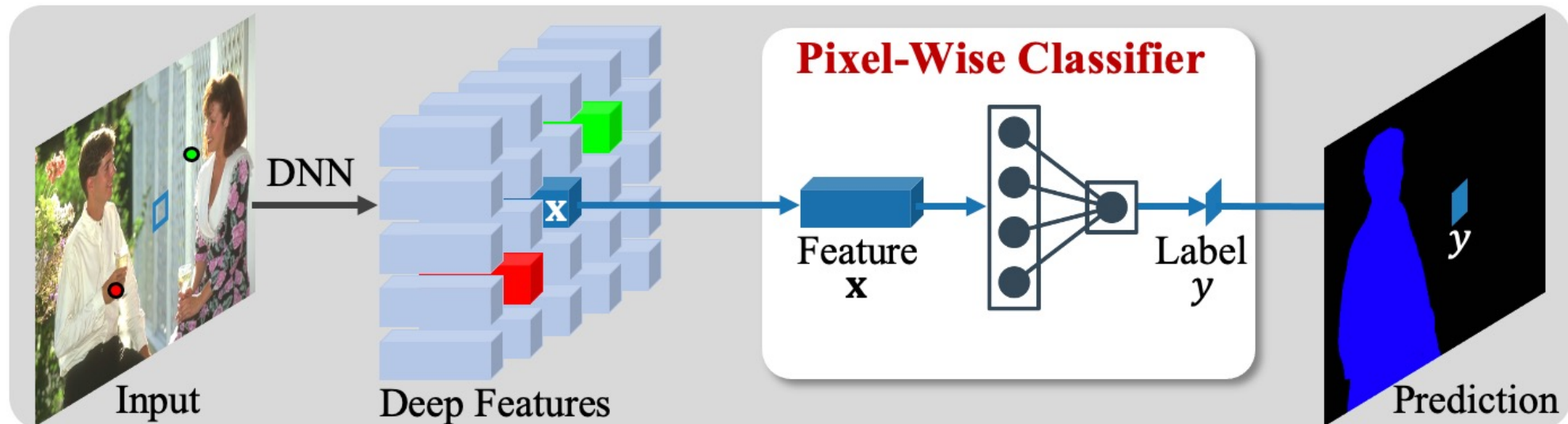
Introduction

Interactive Segmentation (IS) Task



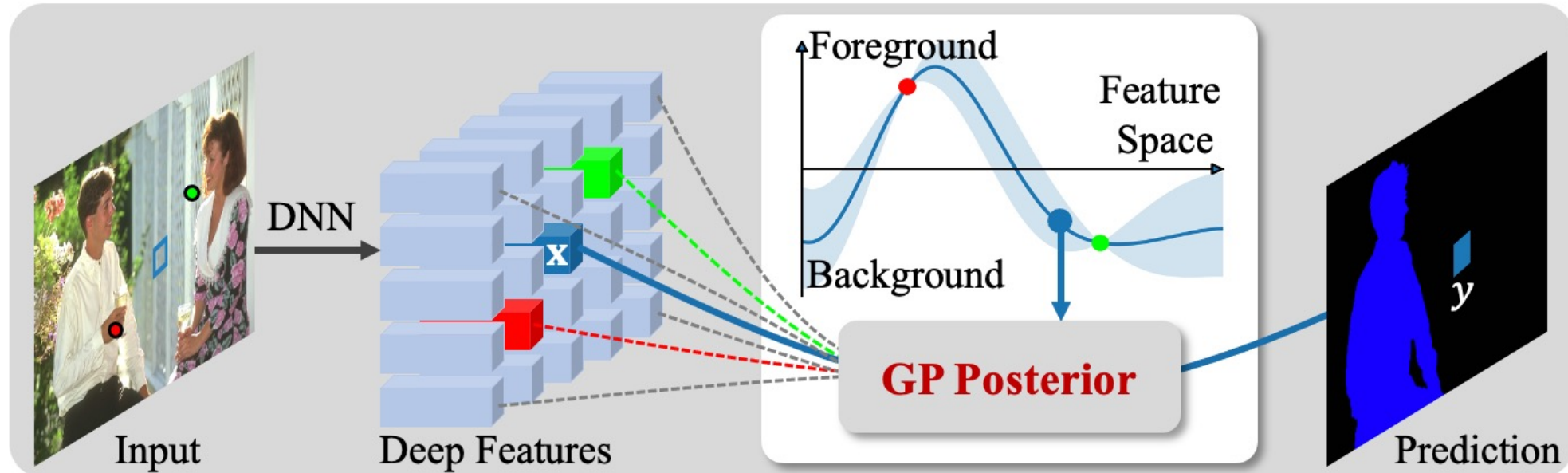
Current Deep Learning-based Interactive Segmentation Models

- Generally perform pixel-wise classification without specific designs;
- No explicit theoretical support that the clicked regions can be properly activated and correctly classified.



Gaussian Process(GP)-based Framework for Interactive Segmentation

- Explicitly measure the relations between data points with a kernel function;
- Promote accurate predictions for training data;
- Can be flexibly integrated with deep networks via deep kernel learning.



Methodology

Methodology

Model Formulation

$\mathbf{X}_n/\mathbf{X}_*$: features of clicked/unclicked pixels;

$\mathbf{y}_n/\mathbf{y}_*$: labels of clicked/unclicked pixels;

\mathbf{f}_* : latent classification scores of unclicked pixels.

Goal: posterior probability of labels \mathbf{y}_* at unclicked pixels

$$p(\mathbf{y}_* | \mathbf{X}_*, \mathbf{X}_n, \mathbf{y}_n) = \int p(\mathbf{y}_* | \mathbf{f}_*) p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}_n, \mathbf{y}_n) d\mathbf{f}_*$$

Monte Carlo Approximate
 $p(\mathbf{y}_* | \mathbf{X}_*, \mathbf{X}_n, \mathbf{y}_n) \approx \text{sig}(\tilde{\mathbf{f}}_*)$

$$p(\mathbf{y}_* | \mathbf{f}_*) = \prod_{c=1}^* \text{sig}(y_c f_c)$$

GP posterior

key step

Sample $\tilde{\mathbf{f}}_*$

Overcome the Intractability of the GP Posterior

GP Posterior: $p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}_n, \mathbf{y}_n) = \int \underbrace{p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}_n, \mathbf{f}_n)}_{\text{Gaussian}} \underbrace{p(\mathbf{f}_n | \mathbf{X}_n, \mathbf{y}_n)}_{\text{non-Gaussian}} d\mathbf{f}_n \leftarrow \text{intractable}$

$p(\mathbf{f}_n | \mathbf{X}_n, \mathbf{y}_n) \propto \underbrace{p(\mathbf{y}_n | \mathbf{X}_n, \mathbf{f}_n)}_{\text{non-Gaussian}} \underbrace{p(\mathbf{f}_n | \mathbf{X}_n)}_{\text{Gaussian}}$

Approximate with $q(\mathbf{f}_n | \mathbf{X}_n, \mathbf{y}_n) = \mathcal{N}(\mathbf{m}_\xi(\mathbf{X}_n, \mathbf{y}_n), \sigma^2 \mathbf{I}_n)$

where $\mathbf{m}_\xi(\mathbf{X}_n, \mathbf{y}_n) = \text{Softplus}(\text{MLP}_\xi(\mathbf{X}_n)) * \mathbf{y}_n$

$$\min_q D_{KL}(q(\mathbf{f}_n | \mathbf{X}_n, \mathbf{y}_n) || p(\mathbf{f}_n | \mathbf{X}_n, \mathbf{y}_n))$$

Achieve Efficient Sampling of the GP Posterior

Approximated tractable GP posterior: $p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}_n, \mathbf{y}_n) \sim \mathcal{N}(\boldsymbol{\mu}_{*|n}, \mathbf{K}_{*,*|n})$



Standard approach to draw samples

$$\tilde{\mathbf{f}}_* = \boldsymbol{\mu}_{*|n} + \underbrace{\mathbf{K}_{*,*|n}^{1/2}}_{O(*^3)} \boldsymbol{\zeta} \text{ with } \boldsymbol{\zeta} \sim \mathcal{N}(0, \mathbf{I}_n)$$

We adopt a decoupled sampling framework[1] with the linear complexity $O(*)$ as:

$$\tilde{\mathbf{f}}_* = \underbrace{\boldsymbol{\Phi}(\mathbf{X}_*) \mathbf{w}}_{\text{weight-space prior}} + \underbrace{\mathbf{K}_{*,n} \mathbf{K}_{n,n}^{-1} (\mathbf{f}_n - \boldsymbol{\Phi}(\mathbf{X}_n) \mathbf{w})}_{\text{function-space update}},$$

Remark: Theoretical Support for Accurate Predictions at Clicks

$$\tilde{\mathbf{f}}_* = \underbrace{\Phi(\mathbf{X}_*)\mathbf{w}}_{\text{weight-space prior}} + \underbrace{\mathbf{K}_{*,n}\mathbf{K}_{n,n}^{-1}(\mathbf{f}_n - \Phi(\mathbf{X}_n)\mathbf{w})}_{\text{function-space update}},$$



Replacing * with n

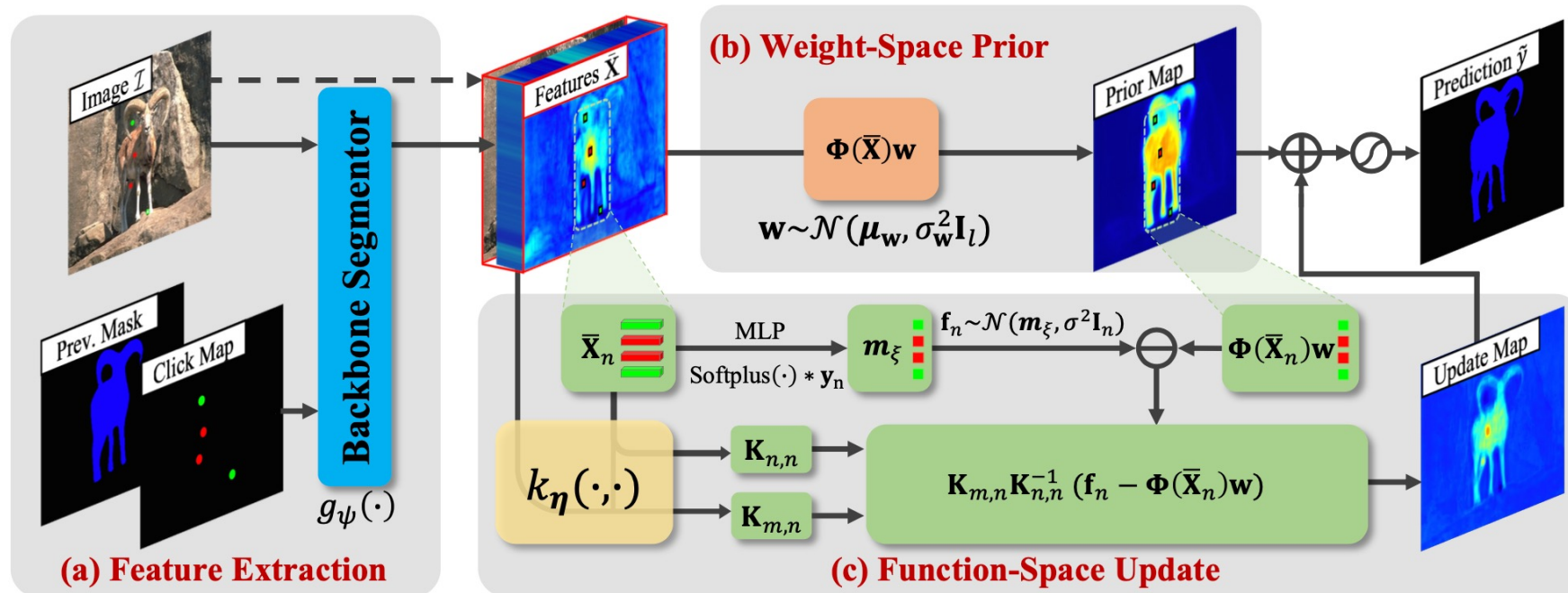
$$\tilde{\mathbf{f}}_n \approx \mathbf{f}_n \approx \mathbf{m}_\xi = \boxed{\text{Softplus}(\text{MLP}_\xi(\mathbf{X}_n))} * \mathbf{y}_n$$

> 0

Methodology

GPCIS Framework

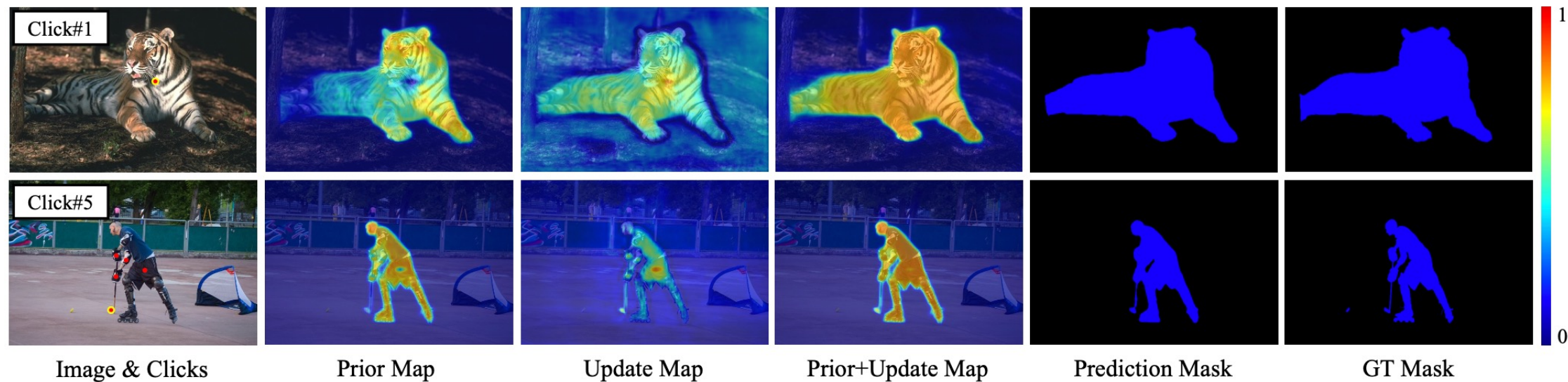
$$\tilde{\mathbf{f}}_* = \underbrace{\Phi(\mathbf{X}_*)\mathbf{w}}_{\text{weight-space prior}} + \underbrace{\mathbf{K}_{*,n}\mathbf{K}_{n,n}^{-1}(\mathbf{f}_n - \Phi(\mathbf{X}_n)\mathbf{w})}_{\text{function-space update}}$$



Experiments

Experiments

Visualization of the Decoupled GP Posterior

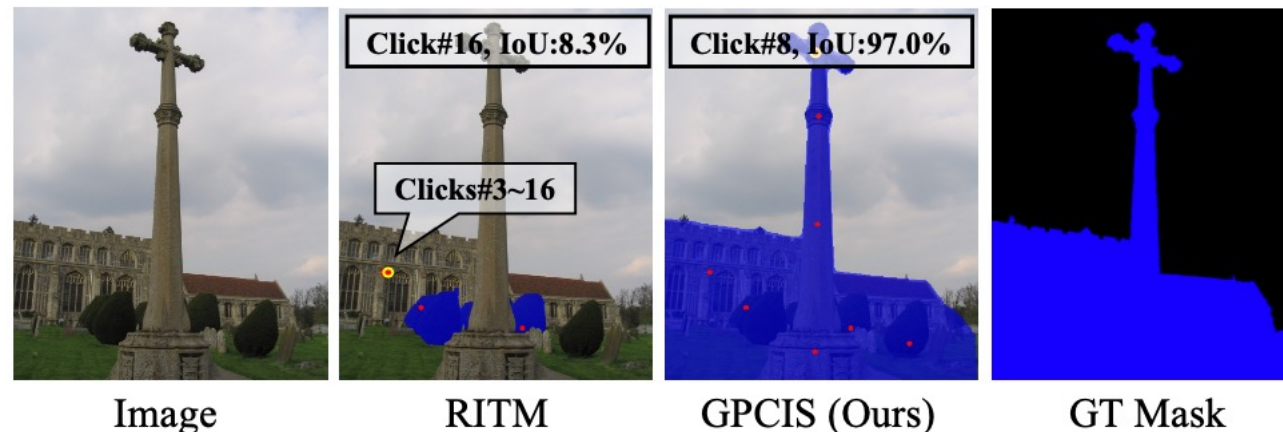


Accuracy at Clicked Pixels

Table 1. The effect of ϵ^2 on the NoIC of our proposed GPCIS with the backbone segmentor ResNet50 on the DAVIS dataset [36].

ϵ^2	10^{-1}	10^{-2}	10^{-3}	10^{-4}	10^{-5}	10^{-6}	10^{-7}
NoIC	36	30	21	15	15	8	2

NoIC : the Number of Incorrectly classified Clicks over a testing dataset



Experiments

NoC Performance over Four Evaluation Datasets

Table 2. NoC@85 and NoC@90 of different competing methods on four datasets, *i.e.*, GrabCut, Berkeley, SBD, and DAVIS. ‘*’ denotes the models trained on the Augmented PASCAL VOC dataset [9, 13]. Bold and underlined results indicate the top 1st and 2nd rank, respectively.

Backbone	Method	GrabCut [38]		Berkeley [32]		SBD [13]		DAVIS [36]		Average	
		NoC@85	NoC@90	NoC@85	NoC@90	NoC@85	NoC@90	NoC@85	NoC@90	NoC@85	NoC@90
DeepLab-LargeFOV [4]	* RIS-Net [24] ('17)	-	5.00	-	6.03	-	-	-	-	-	-
CAN [53]	LD [23] ('18)	3.20	4.79	-	-	7.41	10.78	5.95	9.57	-	-
FCN [29]	*DOS [51] ('16)	5.08	6.08	-	-	9.22	12.80	9.03	12.58	-	-
	*CMG [31] ('19)	-	3.58	-	5.60	-	-	-	-	-	-
DenseNet [17]	BRS [18] ('19)	2.60	3.60	-	5.08	6.59	9.78	5.58	8.24	-	6.68
Xception-65 [8]	*CA [22] ('20)	-	3.07	-	4.94	-	-	5.16	-	-	-
SegFormerB0-S2 [7, 50]	RITM [41] ('21)	<u>1.62</u>	<u>1.82</u>	1.84	<u>2.92</u>	<u>4.26</u>	<u>6.38</u>	<u>4.65</u>	<u>6.13</u>	<u>3.09</u>	<u>4.31</u>
	FocalClick [7] ('22)	1.66	1.90	-	3.14	4.34	6.51	5.02	7.06	-	4.65
	GPCIS (Ours)	1.60	1.76	1.84	2.70	4.16	6.28	4.45	6.04	3.01	4.20
HRNet18s-S2 [7, 43]	RITM [41] ('21)	2.00	2.24	<u>2.13</u>	3.19	<u>4.29</u>	<u>6.36</u>	<u>4.89</u>	<u>6.54</u>	<u>3.33</u>	4.58
	FocalClick [7] ('22)	<u>1.86</u>	<u>2.06</u>	-	<u>3.14</u>	4.30	6.52	4.92	<u>6.48</u>	-	<u>4.55</u>
	GPCIS (Ours)	1.74	1.94	1.83	2.65	4.28	6.25	4.62	6.16	3.12	4.25
ResNet50 [15]	*FCANet [26] ('20)	2.18	2.62	-	4.66	-	-	5.54	8.83	-	-
	f-BRS-B [40] ('20)	2.20	2.64	2.17	4.22	4.55	7.45	5.44	7.81	3.59	5.53
	CDNet [6] ('21)	2.22	2.64	-	3.69	4.37	7.87	5.17	6.66	-	5.22
	RITM [41] ('21)	2.16	2.30	1.90	2.95	3.97	5.92	<u>4.56</u>	6.05	3.15	4.31
	FocusCut [25] ('22)	1.60	1.78	<u>1.86</u>	3.44	3.62	5.66	5.00	6.38	<u>3.02</u>	4.32
	FocalClick [7] ('22)	1.92	2.14	<u>1.87</u>	<u>2.86</u>	3.84	5.82	4.61	<u>6.01</u>	<u>3.06</u>	<u>4.21</u>
	GPCIS (Ours)	<u>1.64</u>	<u>1.82</u>	1.60	2.60	<u>3.80</u>	<u>5.71</u>	4.37	5.89	2.85	4.00

Experiments

Performance on Different Metrics

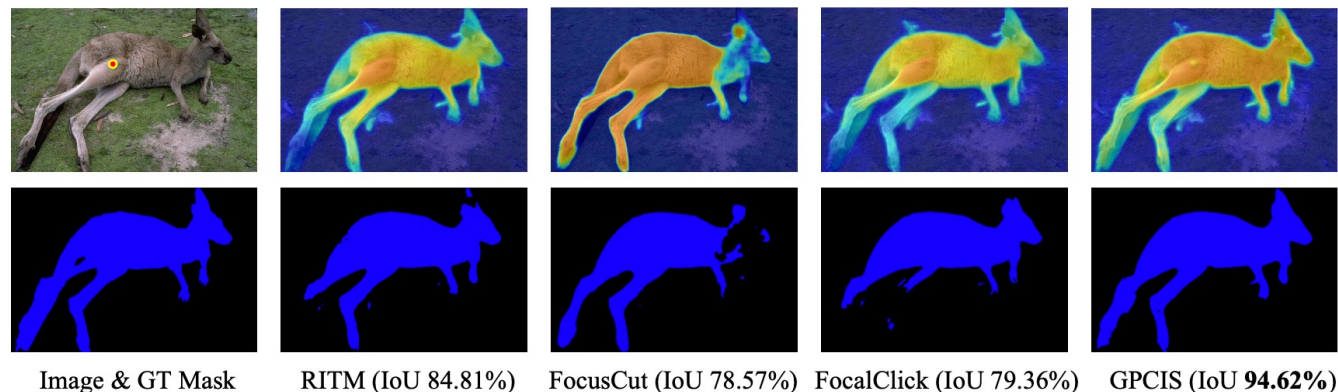
Table 3. Quantitative evaluation on different metrics, and comparisons on parameters and inference time. Here the backbone segmentor is ResNet50, and Second Per Click (SPC) is averagely computed over DAVIS with the testing image size of 384×384 on an NVIDIA V100 GPU. Lower NoC₁₀₀@90, NoF₁₀₀@90, NoIC, #Params, SPC and higher IoU&1, IoU&5 indicate better performance.

Method	Berkeley [32]					DAVIS [36]					#Params (MB)	SPC (ms)
	NoC ₁₀₀ @90	NoF ₁₀₀ @90	IoU&1	IoU&5	NoIC	NoC ₁₀₀ @90	NoF ₁₀₀ @90	IoU&1	IoU&5	NoIC		
f-BRS-B [40]	6.21	2	77.06%	85.00%	1	22.62	57	70.97%	83.87%	0	39.44	116.53
CDNet [6]	-	-	-	-	-	18.59	48	-	-	-	39.90	57.76
RITM [41]	<u>3.75</u>	1	76.88%	94.66%	2	18.09	51	<u>72.89%</u>	<u>89.14%</u>	74	39.48	34.24
FocusCut [25]	4.63	1	<u>78.89%</u>	92.89%	1	19.00	<u>45</u>	<u>72.71%</u>	87.58%	6	40.36	950.68
FocalClick [7]	4.46	2	<u>75.59%</u>	<u>94.90%</u>	0	<u>17.74</u>	49	70.76%	88.90%	42	39.50	41.80
GPCIS (<i>Ours</i>)	3.36	1	79.43%	95.11%	0	17.03	44	75.67%	89.60%	<u>2</u>	39.39	<u>38.82</u>

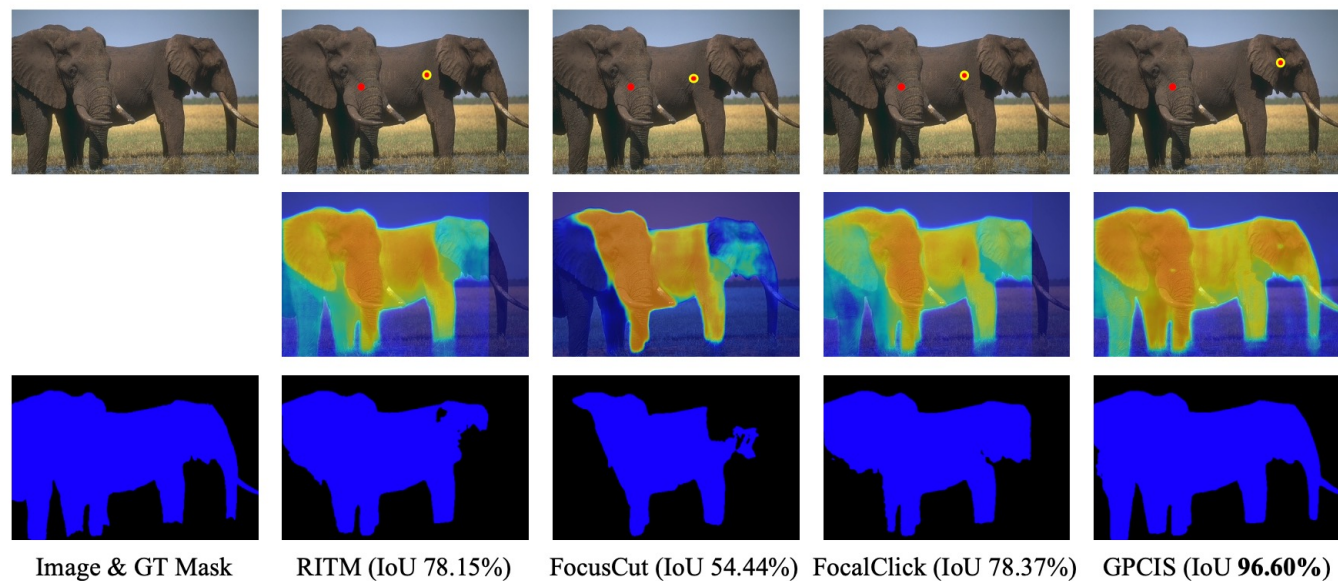
Experiments

Quality Performance

1 Click



2 Clicks



Thank you!



Github Page



Virtual Site



woshizhouminghao@stu.xjtu.edu.cn