

A Fixed-Point Approach to Unified Prompt-Based Counting

Wei Lin, Antoni B. Chan

Department of Computer Science, City University of Hong Kong
Tat Chee 83, Kowloon Tong, Hong Kong SAR, China
elonlin24@gmail.com, abchan@cityu.edu.hk

Overview

This supplementary material is organized as follows:

1. Results on cross-prompt and multi-prompt adaptation is discussed to show the advantages of the unified prompt-based counting;
2. More details about the three components in the proposed prompt-based counting model are presented;
3. The pseudo-code is presented to show the whole computation process of fixed-point loss in contrastive training.
4. The detailed results in ablation studies.
5. Additional visualizations of the proposed method are presented to showcase its advantages and disadvantages.

Cross/Multi-Prompt Adaptation

Converting different types of prompts into the same representation, prompt masks, offers two additional advantages: (a) When new types of prompts are introduced, we can seamlessly transform them into prompts and directly apply them in our model without requiring training; (b) This approach enables us to combine multiple prompts for counting the object of interest.

To showcase the first advantage, we incorporate an additional prompt type, namely instance mask prompts, for prompt-based counting without requiring extra training. The instance mask is generated by converting the box annotations into instance masks using the SAM method (Kirillov et al. 2023). The results are presented in Table S1. The instance prompts yield comparable outcomes to box-guided training, as the prompts originate from box annotations. However, when using instance prompts, the estimation errors are slightly lower than the results obtained from box annotations, owing to the more precise representation provided by instance masks.

Regarding multiple prompts, we present the results in Table S2. Overall, the performance remains consistent with the results obtained using a single type of prompt. However, when the text prompt is combined with other prompts, such as box or point prompts, its performance improves. The original MAE and MSE of the text prompt are 16.81 and 105.83, respectively. These values decrease to 16.71 and 105.80 when the box prompt is added. Meanwhile, the MAE

prompts	train	validation		test	
		MAE	MSE	MAE	MSE
box		16.87	59.45	16.68	105.08
text	w/	16.92	58.92	16.81	105.83
point		17.16	59.38	15.86	103.27
instance	w/o	16.74	59.28	16.58	104.79

Table S1: Cross-prompt adaptation.

box	text	point	validation		test	
			MAE	MSE	MAE	MSE
✓	✓		16.89	58.91	16.71	105.80
✓		✓	16.76	59.36	16.50	104.79
	✓	✓	16.96	59.01	16.78	105.79
✓	✓	✓	16.73	58.88	16.64	105.75

Table S2: Counting with Multiple types of prompts.

and MSE decrease to 16.78 and 105.79 when combined with point prompts.

Model Structure

Our prompt-based counting model consists of three main components: an image backbone, the cross-attention module, and the density predictor, as illustrated in Figure 2 of the main paper.

Backbone & Image Feature. The backbone utilizes the first four blocks of ResNet-101 (He et al. 2016) as the backbone, the image feature can be represented as:

$$\mathbf{F} = \mathcal{G} \left(\left[\mathbf{F}_{\frac{1}{8}}, \mathcal{U}_2 \left(\mathbf{F}_{\frac{1}{16}} \right) \right], \Theta_{\mathcal{G}} \right), \quad (\text{S1})$$

where \mathcal{G} is a `conv_3x3` layer with parameters $\Theta_{\mathcal{G}}$, and \mathcal{U}_2 represents bilinear interpolation with a scale factor of 2. The inputs to \mathcal{G} , denoted as $\mathbf{F}_{\frac{1}{8}}$ and $\mathbf{F}_{\frac{1}{16}}$, are feature maps extracted from the last two blocks in the backbone. Their subscripts indicate the downsampling rate compared to the input. \mathbf{F} has the same resolution with $\mathbf{F}_{\frac{1}{8}}$.

Cross Attention & Density Feature. The aim of cross attention is to let the prompt token $\mathbf{t} \in \mathbb{R}^C$ communicate with the flattened feature $\mathbf{F} \in \mathbb{R}^{N \times C}$, and then generate density feature accordingly. Specifically, we use \mathcal{F} to generate *query*

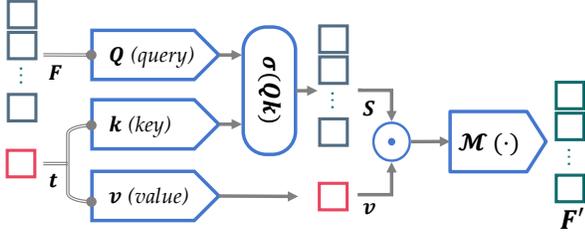


Figure S1: The pipeline of cross-attention module.

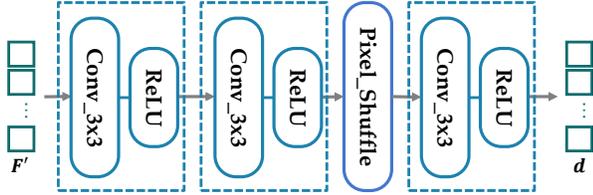


Figure S2: The pipeline of density predictor.

embedding, and use \mathbf{t} to generate *key* and *value* embeddings:

$$\mathbf{Q} = \mathbf{F}\mathbf{W}_q, \quad \mathbf{k} = \mathbf{W}_k\mathbf{t}, \quad \mathbf{v} = \mathbf{W}_v\mathbf{t} \quad (\text{S2})$$

in which $\mathbf{W}_q \in \mathbb{R}^{C \times C'}$, $\mathbf{W}_k \in \mathbb{R}^{C' \times C}$, and $\mathbf{W}_v \in \mathbb{R}^{C' \times C}$ are learnable parameters to perform cross attention. Prior attention models use `softmax` to compute the weights for each *value* embedding. However, there is only one token so that the result of `softmax` would always be one, contrast to the motivation of weight. Thus, we replace the `softmax` with Sigmoid function $\sigma(x) = \frac{1}{1+e^{-x}}$, resulting the following way to generate density features:

$$\mathbf{F}' = \mathcal{M}(\mathbf{S}^\top \mathbf{v}), \quad \mathbf{S} = \sigma(\mathbf{Q}\mathbf{k}) \in \mathbb{R}_+^N, \quad (\text{S3})$$

where \mathcal{M} is a linear projection to transform the attention results into density features, denoted as \mathbf{F}' . The visualized pipeline is displayed in Figure S1.

Density Predictor. We utilize the identical density predictor released in SPDCN (Lin et al. 2022), benefiting from its straightforward yet effective performance. As illustrated in Figure S2, the density predictor is composed of three combined layers of `conv_3x3` and `ReLU`. Before reaching the final convolutional layer, a `pixel_shuffle` layer with scale factor of 8 is employed to upscale the density features, aligning them with the resolution of the input image.

Pseudo Code of Contrastive Training

Algorithm S1 summarizes the procedure for computing the fixed-point loss under the contrastive training scheme. At the beginning, $\mathbf{d}_i^{(0)}$ and $\mathbf{d}_j^{(0)}$ are initialized with \mathbf{m}_i and $\mathbf{0}_{\mathbf{m}_j}$ respectively. Subsequently, the fixed-point estimation is performed iteratively.

In the t -th iteration, we first compute the prompt token $\mathbf{t}^{(t+1)}$ based on both the positive and negative prompt masks, $\mathbf{d}_i^{(t)}$ and $\mathbf{d}_j^{(t)}$. Following this, the positive and negative density map are updated via the function \mathcal{D} with parameters θ ,

Algorithm S1: Contrastive Training

Input: Positive sample $(\mathbf{F}_i, \mathbf{m}_i, \mathbf{d}_i')$, negative sample $(\mathbf{F}_j, \mathbf{m}_j, \mathbf{d}_j')$, and the number of iteration T .

Output: Contrastive fixed-point loss $\hat{\mathcal{L}}_{ij}$.

- 1: Initialize $\mathbf{d}_i^{(0)} \leftarrow \mathbf{m}_i$ and $\mathbf{d}_j^{(0)} \leftarrow \mathbf{0}_{\mathbf{m}_j}$.
 - 2: **for** $t \leftarrow 0$ to $T-1$ **do**
 - 3: {Detach $\mathbf{d}_i^{(T-1)}$ and $\mathbf{d}_j^{(T-1)}$ from the computation graph according to (17) in the main paper.}
 - 4: **if** $t = T-1$ **then**
 - 5: $\mathbf{d}_i^{(t)} \leftarrow \text{DETACH}(\mathbf{d}_i^{(t)})$.
 - 6: $\mathbf{d}_j^{(t)} \leftarrow \text{DETACH}(\mathbf{d}_j^{(t)})$.
 - 7: **end if**
 - 8: Comprehensive token $\mathbf{t}^{(t+1)} \leftarrow \frac{\mathbf{F}_i^\top \mathbf{d}_i^{(t)} + \mathbf{F}_j^\top \mathbf{d}_j^{(t)}}{\|\mathbf{d}_i^{(t)}\|_1 + \|\mathbf{d}_j^{(t)}\|_1}$.
 - 9: The positive density map $\mathbf{d}_i^{(t+1)} \leftarrow \mathcal{D}_\theta(\mathbf{F}_i, \mathbf{t}^{(t+1)})$.
 - 10: The negative density map $\mathbf{d}_j^{(t+1)} \leftarrow \mathcal{D}_\theta(\mathbf{F}_j, \mathbf{t}^{(t+1)})$.
 - 11: **end for**
 - 12: The token at infinity $\mathbf{t}^{(\infty)} \leftarrow \frac{\mathbf{F}_i^\top \mathbf{d}_i'}{\|\mathbf{d}_i'\|_1}$.
 - 13: The positive density at infinity $\mathbf{d}_i^{(\infty)} \leftarrow \mathcal{D}_\theta(\mathbf{F}_i, \mathbf{t}^{(\infty)})$.
 - 14: The negative density at infinity $\mathbf{d}_j^{(\infty)} \leftarrow \mathcal{D}_\theta(\mathbf{F}_j, \mathbf{t}^{(\infty)})$.
 - 15: The positive fixed-point loss $\hat{\mathcal{L}}_i \leftarrow \left\| \mathbf{d}_i^{(\infty)} - \mathbf{d}_i' \right\| + \left\| \mathbf{d}_i^{(T)} - \text{DETACH}(\mathbf{d}_i^{(\infty)}) \right\|^2$.
 - 16: The negative fixed-point loss $\hat{\mathcal{L}}_j \leftarrow \left\| \mathbf{d}_j^{(\infty)} \right\| + \left\| \mathbf{d}_j^{(T)} - \text{DETACH}(\mathbf{d}_j^{(\infty)}) \right\|^2$.
 - 17: **return** $\hat{\mathcal{L}}_{ij} \leftarrow \hat{\mathcal{L}}_i + \hat{\mathcal{L}}_j$.
-

as described in (9) of the main paper. Here \mathcal{D} comprises the cross-attention module in Figure S1 and the density predictor in Figure S2. Note that both density maps are detached from the computation graph before the final iteration.

To compute the fixed-point loss, we also predict the corresponding density map with the token at infinity, $\mathbf{t}^{(\infty)}$. The final loss function is computed for both the positive and negative samples. The distinction lies in that the prediction for the positive component should align with the ground truth, while the prediction for the negative component should approach an all-zero density map. It is worth noting that $\mathbf{d}^{(\infty)}$ is detached in the finite loss calculation to prevent degenerate solutions where both $\mathbf{d}^{(T)}$ and $\mathbf{d}^{(\infty)}$ converge towards zero (Chen and He 2021).

Detailed Results in Ablation Studies

In the main paper, we only present the average MAE/MSE across three types of prompts to conserve space. Here, we provide the detailed corresponding data. Table S3 compares the methods for generating text prompt masks; Table S4 presents the results while comparing fixed-point loss with L2 loss; Table S5 discusses the enhancements achieved by different components in fixed-points; and Table S6 provides

Mask	box (val)		text (val)		point (val)	
	MAE	MSE	MAE	MSE	MAE	MSE
cosine	18.55	69.32	20.59	79.30	18.46	69.70
softmax	16.87	59.45	16.92	58.92	17.16	59.38
Mask	box (test)		text (test)		point (test)	
	MAE	MSE	MAE	MSE	MAE	MSE
cosine	16.97	108.13	17.84	108.77	16.12	105.82
softmax	16.68	105.08	16.81	105.83	15.86	103.27

Table S3: Comparison of text prompt mask generation.

loss function	prompt	validation		test	
		MAE	MSE	MAE	MSE
L2	box	18.84	66.96	17.26	107.62
	text	18.37	65.49	16.83	105.83
	point	20.07	69.52	17.55	106.03
fixed-point	box	16.87	59.45	16.68	105.08
	text	16.92	58.92	16.81	105.83
	point	17.16	59.38	15.86	103.27

Table S4: Comparison between MSE and fixed-point loss.

the N-MAE and N-MSE values for different prompts.

More Visualization

In this section, we also present more visualizations of the proposed method. In Figure S3, we provide some counting examples with text prompts. It demonstrates that our model can generate accurate density maps even when the given image contains multiple types of objects. For convenience, we use $d^{(t)}$ to denote the prediction in the t -th iteration. The visualization shows that $d^{(1)}$ consistently predicts a lower count than $d^{(2)}$. When compared with $d^{(1)}$, the count of $d^{(2)}$ is closer to the ground truth.

In Figure S4, we illustrate some failure cases. In the first row of Figure S4, the minimum repetitive unit is a single *lens*. However, the object of interest, *sunglasses*, naturally has two lenses, doubling the repetitive unit. This discrepancy causes the model to predict a count that is twice the ground truth. In the second row of Figure S4, two types of pills with completely different appearances and shapes are presented. The model only counts one of them, while the other one is

Prompt	infinity	finity	validation set		test set	
			MAE	MSE	MAE	MSE
box	\mathcal{L}_∞	-	27.08	72.56	30.06	111.42
	-	$\mathcal{L}_{T'}$	20.67	68.34	17.29	107.67
	\mathcal{L}_∞	$\mathcal{L}_{T'}$	16.87	61.86	16.46	106.99
	\mathcal{L}_∞	\mathcal{L}_T	16.87	59.45	16.68	105.08
text	\mathcal{L}_∞	-	27.93	74.35	30.82	110.88
	-	$\mathcal{L}_{T'}$	21.07	70.71	17.31	106.90
	\mathcal{L}_∞	$\mathcal{L}_{T'}$	16.99	61.47	16.71	107.32
	\mathcal{L}_∞	\mathcal{L}_T	16.92	58.92	16.81	105.83
point	\mathcal{L}_∞	-	19.73	68.76	20.03	108.32
	-	$\mathcal{L}_{T'}$	20.64	67.13	17.53	107.67
	\mathcal{L}_∞	$\mathcal{L}_{T'}$	17.54	63.10	15.69	103.42
	\mathcal{L}_∞	\mathcal{L}_T	17.16	59.38	15.88	103.27

Table S5: Ablation study on different combinations of infinity and finity part in fixed-point loss.

contrastive training	prompt	validation set		test set	
		N-MAE	N-MSE	N-MAE	N-MSE
w/o	box	38.27	73.37	50.91	88.28
	text	37.89	70.08	49.85	88.34
	point	38.12	75.81	50.68	85.67
w/	box	1.90	17.57	3.85	14.84
	text	1.83	17.29	3.90	15.03
	point	1.97	16.62	3.92	15.70

Table S6: Ablation study on contrastive training.

recognized as irrelevant object.

Figure S5 showcases two examples of counting with different prompts. These diverse prompts also exhibit varying levels of performance, but the estimation errors remain consistent across the entire dataset.

References

- Chen, X.; and He, K. 2021. Exploring Simple Siamese Representation Learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 15750–15758.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.
- Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A. C.; Lo, W.-Y.; Dollár, P.; and Girshick, R. 2023. Segment Anything. In *International Conference on Computer Vision (ICCV)*.
- Lin, W.; Yang, K.; Ma, X.; Gao, J.; Liu, L.; Liu, S.; Hou, J.; Yi, S.; and Chan, A. B. 2022. Scale-Prior Deformable Convolution for Exemplar-Guided Class-Agnostic Counting. In *British Machine Vision Conference (BMVC)*, 313.

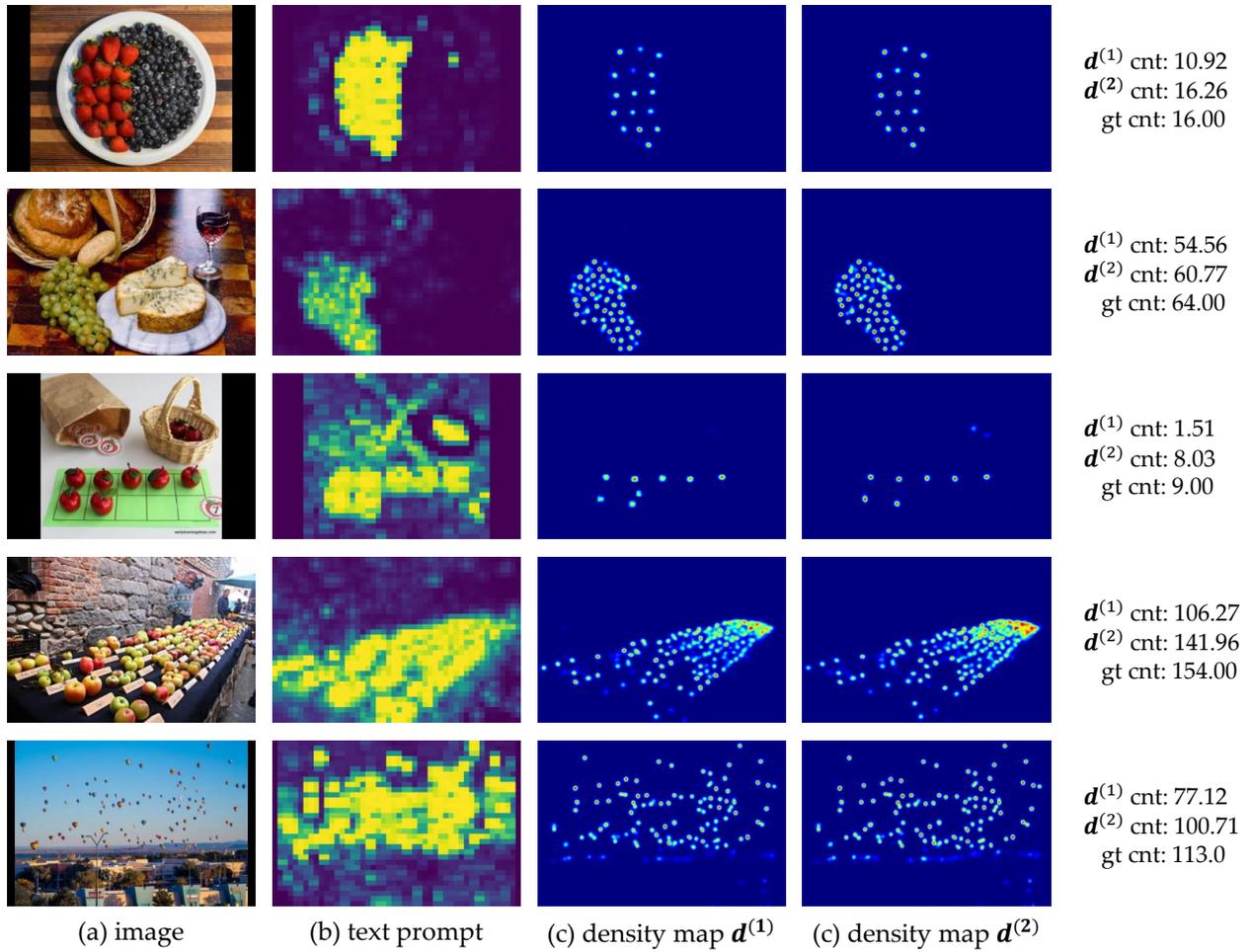


Figure S3: Examples of counting with text prompts.

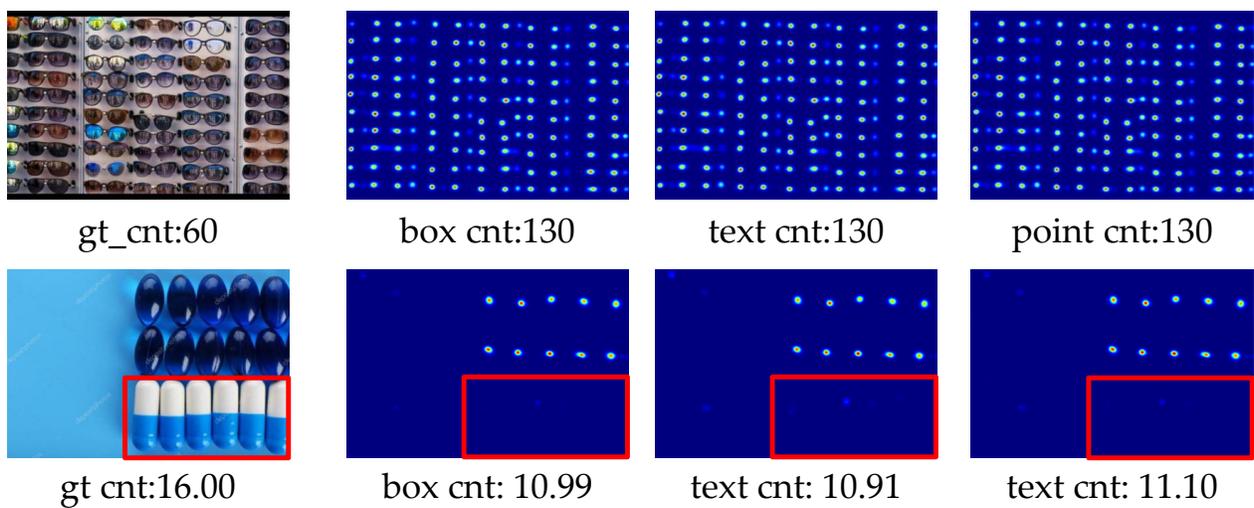


Figure S4: Typical failure cases.

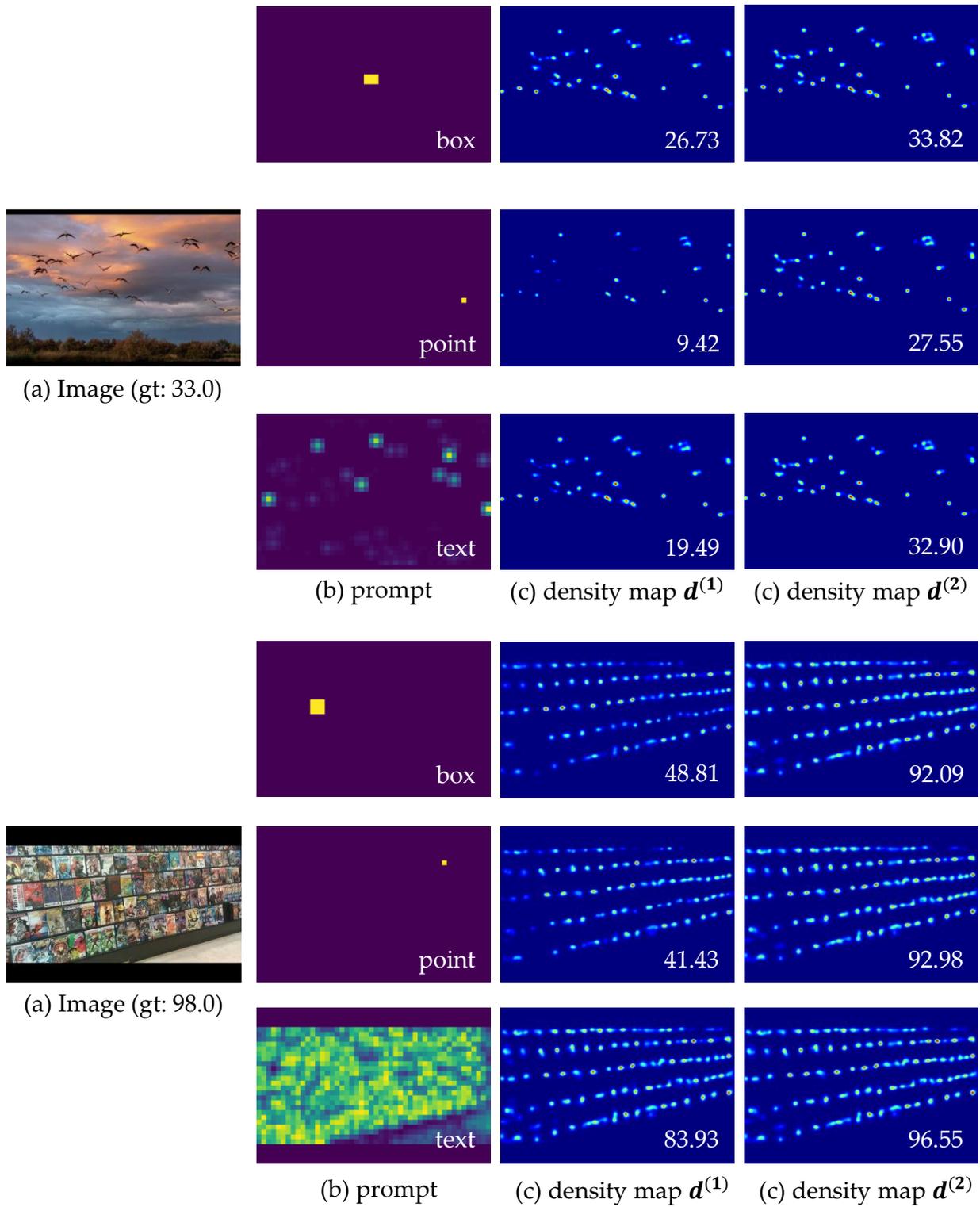


Figure S5: Examples of counting with different types of prompts.