Scale-Invariant Continuous Implicit Neural Representa tions For Remote Sensing Object Counting

Anonymous authors Paper under double-blind review

Abstract

Many object counting methods rely on density map estimation (DME) using convolutional neural networks (CNNs) on discrete grid image representations. However, these methods struggle with large variations in object size or input image resolution, typically due to different imaging conditions and perspective effects. Worse yet, discrete grid representations of density maps result in information loss with blurred or vanished details for lowresolution inputs. To overcome these limitations, we design Scale-Invariant Implicit neural representations for counting (SI-INR) to map arbitrary-scale input signals into a continuous function space, where each function produces density values over continuous spatial coordinates. SI-INR achieves robust counting performances with respect to changing object sizes, extensive experiments on commonly used diverse datasets have validated the proposed method.

024

004 005

006

008 009

010 011

012

013

014

015

016

017

018

019

021

022

1 Introduction

025 026

027 Understanding the distribution and abundance of people as well as geographic entities such 028 as buildings and cars within these environments becomes crucial for various "smart city" 029 applications, such as urban planning, traffic management and beyond. Object counting holds promising potential for such tasks, which have also been studied in other domains too, including crowd counting for security (Ma et al., 2019; Li et al., 2018), animal crowd estimations (Ma et al., 2015), and cell counting for biomedicine (Paul Cohen et al., 2017). 032 Successful counting methods have been developed by introducing deep learning (Fu et al., 033 2015; Wen et al., 2021) and self-attention (Gao et al., 2020; Lin et al., 2022). In recent years, 034 the best-performing methods are mostly based on density map estimation (DME) (Ma et al., 2019; Gao et al., 2022; Lin et al., 2022; Wan et al., 2021), training convolutional neural 036 networks (CNNs) to generate discrete density maps. 037

However, several challenges persist in applying current deep learning methods for reliable 038 counting: 1) Scale-Dependence: CNNs (LeCun et al., 1998) lack intrinsic scale equivari-039 ance, leading to degraded performance when object sizes deviate from those seen during 040 training. This issue is particularly pronounced for inputs at resolutions differing from the 041 training set, as CNNs rely on fixed receptive fields that cannot dynamically adapt to scale 042 variations; 2) Expressiveness-Bottleneck: Traditional grid-based density maps approximate 043 object distributions with Gaussian kernels, imposing a fixed spatial structure that misaligns 044 with irregular object arrangements (Wan & Chan, 2019). Gaussian smoothing reduces noise but blurs local details, degrading fidelity in dense and sparse regions, and limiting counting accuracy in complex scenes. 046

047To address these issues, we design a new object counting framework named Scale-Invariant048Implicit Neural Representations (SI-INR) mapping arbitrary-scale discrete images into 2D049continuous functions which are invariant to the object or structure scales. This allows the050model to preserve the fine details and reduce potential information loss for better counting051accuracy and generalizability. Moreover, the scale-invariance, an important property for052the mapping between input images and output density maps, is explicitly introduced as the053inductive bias of model itself to potentially improve data efficiency and model robustness.
Our main contributions can be summarized as follows:

- 1. We propose Scale-Invariant Continuous Implicit Neural Representations (SI-INR), an object counting framework mapping discrete grid signals into continuous 2D functions which are invariant to image scaling.
- 2. As a proof of concept, we give a realization of SI-INR using existing Scale-Equivariant Steerable Network (SESN) (Sosnovik et al., 2019) and a novel deep neural operator based INR module. A sampling-based optimization objective is then derived for efficient model training.
- 3. We conduct extensive experiments to evaluate the effectiveness of our SI-INR on object counting, demonstrating notable advancements in performance, especially on the remote sensing counting dataset.
- 063 064 065

054

055

057

060

061

062

2 Related Work

067 Object counting: Object counting, e.g. well-studied crowd counting (Lin et al., 2001), has 068 been mostly developed by detecting or segmenting individual objects in the scene. End-to-069 end learning to directly map image features to object counts has been the most success-070 ful counting strategy with rapid advancements in deep learning Krizhevsky et al. (2012); 071 Wang et al. (2015), especially for object counting in densely populated scenes (Chan et al., 2008). Counting based on Density Map Estimation (DME) using convolutional neural net-072 works (CNNs) (Lempitsky & Zisserman, 2010; Fu et al., 2015; Ranjan et al., 2018; Sindagi 073 & Patel, 2019) to preserve translation-invariant multi-scale image features has shown supe-074 rior performance over conventional object counting techniques. More recent ASPDNet (Gao 075 et al., 2020) and PSGCNet (Gao et al., 2022) have integrated attention, deformable convo-076 lution, pyramidal scale modules (PSM) to address challenges in counting such as complex 077 cluttered backgrounds, viewing perspective, object appearance, and size variability. Besides, 078 Huang et al. (2023) propose an optimized global regression model EfreeNet which is more 079 more annotation-efficient. Yi et al. (2023) introduce a lightweight multiscale context fusion module (LMCFM) and a lightweight counting scale pooling module (LCSPM) to reduce the 081 numbers of network parameters and computing cost. These models have achieved state-082 of-the-art (SOTA) counting performance on the RSOC (Remote Sensing Object Counting) 083 dataset (Gao et al., 2020).

Scale-equivariance and invariance: The concept of scale-equivariance and invariance was 085 first proposed in image processing and computer vision (Lowe, 1999; 2004). To handle vari-086 ations in scale effectively, multi-scale features can be learned by applying the convolutions 087 to several rescaled versions of the images or feature maps in every layer (Kanazawa et al., 088 2014; Marcos et al., 2018) or by rescaling trainable filters (Xu et al., 2014). Cai et al. (2016) proposed a pyramidal structure to learn scale-dependent features, which is widely 089 used in the object detection field nowadays. Later, Gaussian scale-space theory (Lindeberg, 090 1994) and group theory (Cohen & Welling, 2016) have been widely used for achieving scale-091 equivariance and invariance. Yang et al. (2023); Lindeberg (2022) parameterized convolu-092 tional filters as a linear combination of Gaussian derivative filters with different scales, and achieved scale-equivariance in image classification and segmentation tasks. Unlike models 094 rooted in Gaussian scale-space theory, Sosnovik et al. (2019) proposed a Scale-Equivariant Steerable Network (SESN), which utilizes steerable filters parameterized by a trainable lin-096 ear combination of pre-calculated Hermite basis functions. These models all first build a scale-equivariant model, and use a simple pooling layer or rescale the outputs to con-098 vert the model into a scale-invariant one if needed (Sosnovik et al., 2019). However, such 099 methods have significant demands on memory and computational resources and can lose information when doing the equivariance to invariance conversion. In SI-INR, we adopt 100 a scale-equivariant model to learn a deep representation that is equivariant to input scale 101 variations, and introduce a scale-invariant model to convert the equivariant mapping to an 102 invariant one. 103

Implicit neural representations: Implicit Neural Representations (INRs) allow for continuous
flexible representations of complex objects and scenes in computer vision (Mescheder et al.,
2019; Oechsle et al., 2019; Barrowclough et al., 2021; Wang et al., 2022). Together with
recently developed positional encoding strategies (Tancik et al., 2020; Sitzmann et al., 2020)
and end-to-end hypernetwork-based learning (Dupont et al., 2022; Lee et al., 2024; Kim

et al., 2023), efficient training to capture high-frequency details has been achieved to better
 model complex natural signals.

However, hypernetwork-based INR models are not translation-equivariant, making them
unsuitable for handling scale variations in input signals. As a result, these models typically require the input to be of a fixed size. More recent Hierarchical Neural Operator
Transformer (HiNOTE) (Luo et al., 2024) integrates neural operators in implicit neural representation, which can preserve more local information and improve the generalizability of INR models.

In SI-INR, we adopt a lightweight INR implementation and replace the traditional coordinate input by a continuous latent. In this case, our INR model transforms continuous latent to the continuous representation of targets, which can be seen as a deep neural operator (Kovachki et al., 2023). This approach improves stability, accelerates model training, and offers greater flexibility when incorporating images of varying sizes during training.

121 122 123

124

3 Method

We propose a novel object counting framework, Scale-Invariant Implicit Neural Representation (SI-INR), that explicitly models scale-invariance in the adopted continuous INR for robust object counting. We start the discussion by first presenting the problems of existing methods in Section 3.1. Next, we describe our solution and concept of SI-INR with the corresponding analysis in Section 3.2.1. We then describe the detailed construct of SI-INR using SESN and a deep neural operator based INR architectures in Section 3.2.2. Lastly in Section 3.3, we provide our sampling based model training for SI-INR on object counting.

132 133

134

155

156

3.1 Problem statement

In many real-world computer vision applications, input images can vary widely in dimension and size, which requires the corresponding object detection and counting models to be able to robustly identify objects ranging from a few pixels to thousands of pixels in scale. While existing off-the-shelf CNNs often take the input images of the fixed size, the flexibility to handle arbitrary resolution inputs and generate arbitrary resolution outputs is highly desirable to take the best advantage of heterogeneous data available.

For a given image I, existing methods aim to establish a mapping $f : \mathbb{R}^{d_{\mathbf{I}}} \Rightarrow \mathbb{R}^{d_{\mathbf{D}^{g_t}}}$ from input image $\mathbf{I} \in \mathbb{R}^{d_{\mathbf{I}}}$ to corresponding outputs $\mathbf{D}^{g_t} \in \mathbb{R}^{d_{\mathbf{D}^{g_t}}}$, where both input images and outputs are in typical discrete-grid representations and the outputs' resolution $d_{\mathbf{D}^{g_t}}$ is typically lower than inputs' resolution reflected by the input dimension $d_{\mathbf{I}}$.

However, since traditional convolution models are not scale-invariant and always have a fixed downsampling ratio, these methods usually struggle with handling varying resolution inputs and objects in different scales. where $f(p_1(h)(I_a)) \neq f(I_a)$. Here *h* denotes an element of the Scale-Translation group *H* and represents one scale-translation operator, $p_1(\cdot)$ denotes the corresponding group actions of *h* acting on the image domain. We provide a more detailed explanation of scale equivariance in Appendix A.1.

To address this limitation, we propose SI-INR to model the mapping for continuous signal representations. SI-INR learns a mapping $\Psi : \mathcal{I} \to \mathcal{F}$ from input image space \mathcal{I} to the continuous function space \mathcal{F} .

$$\Psi(p_1(h)(\mathbf{I}_a))(\mathbf{x}) = \Psi(\mathbf{I}_a)(\mathbf{x}) = \mathbf{D}^{gt}(\mathbf{x}),\tag{1}$$

157 where we emphasize that the input image space \mathcal{I} here is more flexible considering arbitrary 158 normalized spatial coordinates, $\mathbf{x} \in [0, 1]^2$, sampling from continuous image space. $\Psi(\mathbf{I}_a)$ 159 denotes the predicted continuous representation of density maps for \mathbf{I}_a . \mathbf{D}^{gt} denotes a 160 continuous ground truth. Due to this formulation, SI-INR allows the learned model to 161 handle varying resolution inputs, detect as many target objects at different scales as possible, and generate arbitrary resolution outputs.



Figure 1: Schematic diagram of Scale-Invariant Implicit Neural Representation (SI-INR) 174 and existing Density Map Estimation (DME) methods. SI-INR learns scale-invariant con-175 tinuous representations in three steps: first, a scale-equivariant backbone is designed to 176 extract deterministic scale-equivariant features; then, a scale-invariant encoder is adopted 177 to aggregate scale-equivariant features in different scales; and finally, an INR decoder con-178 verts extracted features into an invariant output, a continuous representation of task targets. 179 Visualization of continuous and discrete representations demonstrates that continuous rep-180 resentations preserve more information, leading to better reconstruction of the continuous 181 output.

184

3.2 Scale-Invariant Implicit Neural Representations

Following the above formulation, to achieve a scale-invariant mapping from the image space It to the function space \mathcal{F} , we propose the SI-INR modular framework consisting of three components: a Scale-Equivariant backbone (SE-Backbone), a Hybrid Pyramidal Scale Module (HPSM), and an INR decoder. The SE backbone is designed to extract deterministic scale-equivariant features resilient to different resolutions of inputs and sizes of objects; then the HPSM merges scale-equivariant features in different scales; and finally, the INR decoder converts extracted features into a scale-invariant output, which is a continuous representation of density maps.

192 193 194

3.2.1 Model Components

Traditional computer vision backbones struggle with scale variations of objects where different sizes of the same objects can have different appearances in the feature maps. This results in inconsistent counting estimates and requires the model to have higher capacity to memorize the same objects in different scales (Zhan et al., 2022).

199 To construct a model which is data efficient and capable of handling unseen resolution in-200 puts, we want our model to be scale-invariant so that consistent outputs can be generated 201 with respect to varying scales of objects in input images. As the combination of scale-202 equivariant and scale-invariant models is scale-invariant (Sosnovik et al., 2019), we can 203 construct the model by ensuring that either each component in the framework is invariant, or the framework appropriately combines scale-equivariant and scale-invariant components. 204 It is more reasonable to implement a scale-equivariant mapping first to preserve fine de-205 tails (Sosnovik et al., 2019), and then convert the equivariant mapping into a scale-invariant 206 one for consequent prediction tasks. 207

SI-INR adopts this strategy and takes three steps to achieve scale-invariant mapping for arbitrary-scale signals as shown in Figure 1: first, a scale-equivariant model $B(\cdot)$ is adopted to extract deterministic features so that scale changes in objects will only affect the scale of feature maps while keeping the appearance:

$$B(p_1(h)(\mathbf{I}_a)) = p_B(h)(B)(\mathbf{I}_a), \tag{2}$$

where $p_B(\cdot)$ denotes the corresponding group actions of h acting on the feature domain. We emphasize that any scale-equivariant method can be adopted into SI-INR to handle different tasks. 216 Second, a scale-invariant converter transforms the equivariant features into scale-invariant 217 ones for consequent prediction tasks. At this step, a scale-invariant encoder $E(\cdot)$ is adopted 218 to integrate extracted equivariant features: 219

$$E(p_B(h)(B(\mathbf{I}_a))) = E(B(\mathbf{I}_a)).$$
(3)

Equation (3) demonstrates that the output of $E(\cdot)$ is invariant to scale changes in the input 222 signals. While $E(\cdot)$ can map the same signal at different scales into one latent space, it 223 cannot map different signals at varying scales into a unified space. To address this problem, 224 we introduce a scale-invariant continuous representation mapping \mathcal{H} . For any query image 225 I_a, \mathcal{H} maps the output of $E(\cdot)$ into a continuous function $u_a: [0,1]^2 \to \mathbb{R}_{>0}$, which means the 226 corresponding density value for arbitrary normalized query position \mathbf{x} can be predicted by 227 evaluating the mapped continues function at \mathbf{x} : $u_a(\mathbf{x})$. In SI-INR, we extend the function to 228 $u(x; z, \theta_{\text{INR}})$, where $z \in \mathbb{R}^m$ represents the latent features extracted by the encoder. Thus, 229 SI-INR learns a conditional continuous representation of the input by optimizing over θ_{INR} 230 and z. This allows for task-specific predictions such as continuous density estimation.

Training a continuous representation learner not only achieves a uniform format of output for different scales input but also enables flexible training sample-based algorithms to enhance the training as we discuss in Section 3.3.

We now prove the scale-invariance of our SI-INR for any scale-translation action on I_a in Theorem 1.

Theorem 1 Given a scale-translation operation h and an input image I_a , SI-INR is scaleinvariant:

240 Proof:

241

242 243

244

245 246

247

258 259 260

220 221

$$\Psi(p_1(h)(\mathbf{I}_a))(\mathbf{x}) = \mathcal{H}(E(B(p_1(h)(\mathbf{I}_a)))(\mathbf{x}) = \mathcal{H}(E(B(\mathbf{I}_a)))(\mathbf{x}) = \Psi(\mathbf{I}_a)(\mathbf{x}), \tag{4}$$

SI-INR not only is invariant to the change of scales, but also maps signals into one latent continuous function space.

3.2.2 SESN and Deep neural operator based Realization

248 In order to achieve the scale-invariant model, we first map the input signals into an equiv-249 ariant space through a scale-equivariant backbone $B(\cdot)$. Thanks to the efficient INR-based 250 equivariant-to-invariant conversion architecture, any scale-equivariant models can be applied in SI-INR, we choose SESN (Sosnovik et al., 2019) to build our scale-equivariant backbone 251 $B(\cdot)$ as an example. The introduction of SESN is given in Appendix A.2. Since the sum-252 mation of two scale-equivariant models is still scale-equivariant by Lemma 1 in Appendix 253 A.2, we build our backbone based on the residual architecture, where the input is added to 254 the output of each equivariant layer. 255

In SI-INR, the scale-invariant converter encoder $E(\cdot)$ is designed by a combination of a scaling operator $S(\cdot)$ and a CNN-based model G:

$$E(B(p_1(h_1)(\mathbf{I}_a))) = G(S(B(p_1(h_1)(\mathbf{I}_a)))) = G(p_B(h' \cdot h_1^{-1})p_B(h_1)(B(\mathbf{I}_a)))$$

= $G(p_B(h')(B(\mathbf{I}_a))) = E(B((\mathbf{I}_a))),$ (5)

where the scaling operator rescales any feature map $p_B(h_1)(B(\mathbf{I}_a))$ into a consistent scale: $p_B(h')(B(\mathbf{I}_a)), S(\cdot) = p_B(h' \cdot h_1^{-1})$. This scaling operator ensures the invariant property of our encoder $E(\cdot)$. Note that h' is a hyperparameter which can rescale the derived equivariant features into a more reasonable size $\mathbb{R}^{l_w \times l_h \times C_S^L}$, where l_w , l_h , and C_S^L denote the size and number of channels.

Prior knowledge about the original scales of test images can further improve the prediction accuracy of SI-INR. In the cases where the images in the dataset are unscaled and mutually independent, we can maintain the scale-invariance property by setting the scaling factor to 1.

5

In the CNN-based model G, we propose a Hybrid Pyramidal Scale Module (HPSM) consisting of a Pyramidal Scale Module (Gao et al., 2022) that implements convolutions with increasing kernel size to detect equivariant features in different scales, and scale-invariant convolutions (SESC with maximum scale projection) to convert extracted features into invariant ones.

Finally, SI-INR efficiently utilizes extracted information by introducing an implicit neural representation model *H* which maps different signals into a specific continuous function space. This mapping fully utilizes its continuous property to predict and recover fine details. Representing these deterministic features in a continuous function can capture more information compared with a discrete one. In our experiments in Section 4, SI-INR achieves better counting performance and we show that increasing the number of samples increases counting accuracy for small target objects such as compact cars.

282 \mathcal{H} outputs a continuous representation u_a of the target output for image \mathbf{I}_a , Here we give 283 the expression of the INR decoder: 284

$$\mathcal{H}(E(B(\mathbf{I}_a)), \boldsymbol{\theta}_{INR})(\mathbf{x}) = u_a(\mathbf{x} | \mathbf{z}^a, \boldsymbol{\theta}_{INR}), \tag{6}$$

where θ_{INR} denotes the trainable parameters and the INR model consists of L_{INR} linear layers, $\theta_{INR} = [W_1, b_1, \dots, W_{L_{INR}}, b_{L_{INR}}].$

To build a more flexible INR model that can handle complex scenarios, we generate a continuous latent \mathbf{z}^a , for any query position x, we compute corresponding $\mathbf{z}^a_{\mathbf{x}}$ by sampling from \mathbf{z}^a and set it as the INR model's input.

In this way, our INR model treats input features as continuous functions instead of 2D discrete arrays, which can fully utilize local details in downstream analyses (Luo et al., 2024). We call this continuous-to-continuous INR module a deep neural operator based INR.

With all these, the predicted value at position x can be estimated as:

$$u_a(\mathbf{x}) = \boldsymbol{\phi}_{L_{INR}}^{INR}(\boldsymbol{W}_{L_{INR}}^T \dots \boldsymbol{\phi}_1^{INR}(\boldsymbol{W}_1^T(\mathbf{z}_{\mathbf{x}}^a) + \boldsymbol{b}_1) \dots + \boldsymbol{b}_{L_{INR}}),$$
(7)

where $\phi^{INR}(\cdot)$ denotes the activation function.

285

296

297 298 299

300

307

317318319320

We now prove the scale-invariance of our SI-INR's realization. Given a scale-translation operation h and an input image I_a , we have:

$$u_{p_{\mathbf{I}_{a}}(h)(\mathbf{I}_{a})}(\cdot) = \mathcal{H}(E(p_{B}(h' \cdot h^{-1})(B(p_{\mathbf{I}_{a}}(h)(\mathbf{I}_{a})))), \boldsymbol{\theta}_{INR})(\cdot)$$

$$= \mathcal{H}(E(p_{B}(h' \cdot h^{-1})(p_{B}(h)(B(\mathbf{I}_{a})))), \boldsymbol{\theta}_{INR})(\cdot)$$

$$= \mathcal{H}(E(p_{B}(h')(B(\mathbf{I}_{a}))), \boldsymbol{\theta}_{INR})(\cdot) = u_{a}(\cdot),$$

(8)

where $p_{I_a}(h)$ and $p_B(h)$ denote the group actions of h on the input image domain and $B(\cdot)$'s output domain. For any scale-translation action $p_{I_a}(h)$, the output for image I_a is always $u_a(\cdot)$.

3113.3 Training with Regional Sampling

To train a continuous representation of the density map, a continuous ground truth is needed. We achieve this by directly constructing the likelihood function of any position **x** given label y_n as $p(\mathbf{x} | y_n) = \mathcal{N}(\mathbf{x}; \mathbf{m}_n, \sigma^2 \mathbf{1}_{2 \times 2})$. The density map \mathbf{D}^{gt} is then modeled as 2D stochastic processes in the continuous spatial domain:

$$\mathbf{D}^{gt}\left(\mathbf{x}\right) \stackrel{\text{def}}{=} \sum_{n=1}^{N} \mathcal{N}\left(\mathbf{x}; \mathbf{m}_{n}, \sigma^{2} \mathbf{1}_{2 \times 2}\right) = \sum_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{\|\mathbf{x} - \mathbf{m}_{n}\|_{2}^{2}}{2\sigma^{2}}\right), \quad (9)$$

where **x** denotes the normalized spatial coordinates, $\mathbf{x} \in [0, 1]^2$, $\mathcal{N}(\mathbf{x}; \mathbf{m}_n, \sigma^2 \mathbf{1}_{2 \times 2})$ denotes the 2D Gaussian distribution with the mean \mathbf{m}_n and isotropic covariance matrix $\sigma^2 \mathbf{1}_{2 \times 2}$. Similar to Ma et al. (2019), we incorporate a Bayesian counting loss function in SI-INR, which is robust to noise and object occlusion. Together with the MAE counting loss, the final minimization objective is:

326 327

328 329

330

331

332

333

334

335 336

338 339 340

341

342

343 344

345 346

347

$$\mathcal{L} = \frac{1}{A} \sum_{a=1}^{A} \Big[\sum_{n=0}^{N} (\mathbb{E}_{\mathbf{x} \sim p_s(\mathbf{x})}(c_{n,a}) - \mathbb{E}_{\mathbf{x} \sim p_s(\mathbf{x})}(c_{n,a}^{\text{gt}})) + \kappa (\sum_{n=0}^{N} (\mathbb{E}_{\mathbf{x} \sim p_s(\mathbf{x})}(c_{n,a})) - N) \Big], \quad (10)$$

where $c_{n,a}^{gt}(\cdot)$ and $c_{n,a}(\cdot)$ denote the contribution of \mathbf{D}_a^{gt} to the *n*-th object label, $p_s(\mathbf{x})$ is any probability distribution of \mathbf{x} which enables our model to be trained using any existing stochastic optimization algorithm, and κ is a hyperparameter that balances the two loss function terms. To efficiently compute the loss, we introduce a regional sampling approach, where \mathbf{x} is uniformly sampled from a pre-defined grid. The loss function is then rewritten as:

$$\mathcal{L} = \frac{1}{A} \sum_{a=1}^{A} \Big[\sum_{n=0}^{N} (\mathbb{E}_{\mathbf{x} \sim \text{Unif}[0,1]^2}(c_{n,a}) - \mathbb{E}_{\mathbf{x} \sim \text{Unif}[0,1]^2}(c_{n,a}^{\text{gt}})) + \kappa (\sum_{n=0}^{N} (\mathbb{E}_{\mathbf{x} \sim \text{Unif}[0,1]^2}(c_{n,a})) - N) \Big].$$
(11)

The detailed derivation is provided in Appendix A.3. The count predictions can be hereby obtained by sampling uniformly over the normalized image domain and computing the summation of $\mathbf{D}_{a}(\mathbf{x})$.

4 Experiments

4.1 Experimental Setup

Datasets. We evaluate the model's performance on three challenging datasets: (1) the Remote Sensing Object Counting (RSOC) dataset (Gao et al., 2020); (2) the Car Parking
Lot Dataset (CARPK) (Hsieh et al., 2017); and (3) the Pontifical Catholic University of Parana+ Dataset (PUCPR+) (Hsieh et al., 2017). Details on the datasets and the train-test split are provided in Appendix B.1.

Baselines. We compare our SI-INR with three baseline methods: ASPDNet (Gao et al., 2020), an attention-based network with scale pyramid and deformable convolutions; PSGCNet (Gao et al., 2022), which integrates pyramidal scale and global context modules; and eFreeNet (Huang et al., 2023), an ensemble of first-rank-then-estimate networks. Network architecture details are available in Appendix B.2.

358 Implementation. We implement our scale-equivariant backbone by stacking four SESC 359 layers (Sosnovik et al., 2019) with residual connections where the kernel size of SESC is set to 360 3 to reduce computational costs. Our scale-invariant encoder $E(\cdot)$ consists of two parts. The 361 first part involves a scaling operator, the second part involves a VGG-19 network (Simonyan 362 & Zisserman, 2014) to learn deep features, following a pyramidal scale module that has 363 increasing kernel sizes from 3×3 to 11×11 , and applies SESC layers with maximum scale projection. Our INR-based decoder network consists of four fully connected layers with 364 residual connections, and one fully connected layer with learnable parameters to output raw 365 density maps. We use the Adam (Kingma & Ba, 2014) optimizer for both our SI-INR and 366 baseline models, and set the learning rate to be 1e - 4. We initialize parameters in SI-INR 367 by random sampling from a Gaussian distribution $\mathcal{N}(0, 0.01^2)$. We set $\sigma = 8$ in our loss 368 function and $S_{INR} = 64$ when generating density maps for different inputs unless specified. 369 We evaluate our SI-INR with baseline models on the RSOC dataset, CARPK dataset, and 370 PUCPR+ dataset, following Huang et al. (2023)'s setup with RSOC images resized into 371 512×512 . Data augmentation has been implemented during model training by randomly 372 flipping the input images horizontally. We select the models with the lowest RMSE and 373 proper density maps in the first 300 training epochs and report the results. We run all 374 our experiments with the fixed random seed 64 on a workstation with a NVIDIA V100 375 32GB GPU. We adopt two widely used metrics in object counting tasks following previous work (Gao et al., 2020; Ma et al., 2019) to evaluate baselines and our model: the Mean 376 Absolute Error (MAE) and the Root Mean Squared Error (RMSE). We additionally compare 377 SI-INR with SOTA crowd-counting methods on the UCF-QNRF dataset in Appendix B.5.



Figure 2: Predicted density maps by SI-INR and other baselines for four test images from RSOC. The test images (Test Images) and corresponding density maps (GT) are randomly sampled. The illustrated density maps are predicted by PSGCNet with MSE loss (PS-GCNet+MSE), PSGCNet with Bayesian counting loss (PSGCNet+Bayes), ASPDNet with MSE loss (ASPD+MSE), ASPDNet with Bayesian counting loss (ASPD+Bayes), and SI-INR. Warmer colors denote higher values while cooler colors denote lower values.

4.2 Main Results

- We visualize the predicted density maps generated by SI-INR and other SOTA methods in Figure 2, excluding eFreeNet (Huang et al., 2023), as it is a regression-based counting method. All four images are randomly sampled from the RSOC dataset, with the first three from the RSOC ship dataset and the last one from the RSOC small-vehicle dataset. As shown, SI-INR delivers more accurate counting performance, particularly when the objects' appearance and scales are complex. In the first three images, SI-INR generates clearer density maps. In the last image from the RSOC small-vehicle dataset, where the cars are too small for the other SOTA methods to detect, SI-INR's scale-invariant property enables it to produce higher-quality density maps and achieve better counting accuracy.
- Inference Efficiency. In the inference stage of the RSOC building dataset, ASPD-Net requires approximately 15.13 seconds, PSGC-Net takes around 2.47 seconds, eFreeNet takes around 3.84 seconds, and our SI-INR model requires about 3.87 seconds. SI-INR does take longer during the training phase compared to PSGC-Net due to the integration of scale-equivariant models and the use of stacks of linear layers in the INR. However, the inference cost remains acceptable thanks to the simple design of our INR decoder, which consists of only 4 linear layers, and our lightweight scale-equivariant backbone.

Quantitative Results. Our performance evaluation on different benchmark datasets with the reported experimental results on RSOC in Table 1, and CARPK and PUCPR+ in Ta-ble 2, respectively. It shows hat SI-INR significantly improves the MAE on all the datasets compared with our baseline model PSGCNet. Also, SI-INR achieves comparable results, es-pecially on PUCPR+ datasets. Note that the RSOC small-vehicle and RSOC ship datasets exhibit the largest scale variations and the smallest target objects within the RSOC dataset as we show in Appendix B.1. For the RSOC ship and small-vehicle datasets, SI-INR achieves superior counting performance primarily due to its flexibility in generating outputs at ar-bitrary resolutions. Compared with other methods with fixed downsample ratio and can only generate 64×64 density maps, SI-INR improves the counting performance by directly generating larger and clearer density maps as we showed in Table 5. These results highlight the significant advantages of SI-INR in handling targets across varying scales.

Table 1: Comparison of counting performances on the RSOC datasets.

Mothod	L	Loss		Building		Small-vehicle		Large-vehicle		Ship	
Method	MSE.	Bayes.	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	
MCNN	\checkmark		13.65	16.56	488.65	1317.44	36.56	56.55	263.91	412.30	
eFreeNet	\checkmark		6.99	9.61	195.86	463.62	14.55	19.77	65.34	85.45	
PSGCNet	\checkmark		7.33	11.02	346.78	952.64	21.54	32.75	75.27	94.79	
PSGCNet		\checkmark	7.18	10.98	196.25	360.15	14.47	26.19	72.07	98.06	
ASPDNet	\checkmark		7.40	11.06	378.23	978.93	18.76	31.06	63.32	84.85	
ASPDNet		\checkmark	7.59	11.07	365.69	1101.25	16.61	29.26	64.82	89.24	
SI-INR		\checkmark	6.54	9.80	157.18	306.43	12.61	21.78	59.76	81.79	

Table 2: Comparison of counting performances on the CARPK and PUCPR+ datasets.

Method	Loss		CARPK		PUCPR+	
Method	MSE.	Bayes.	MAE	RMSE	MAE	RMSE
MCNN (Zhang et al., 2016)	\checkmark		24.95	39.63	21.86	29.53
eFreeNet (Huang et al., 2023)	\checkmark		46.42	52.34	18.98	23.03
PSGCNet (Gao et al., 2022)		\checkmark	11.07	14.55	3.87	4.86
PSGCNet (Gao et al., 2022)		\checkmark	7.71	10.28	3.17	5.27
ASPDNet (Gao et al., 2020)	\checkmark		10.01	12.84	4.21	5.02
ASPDNet (Gao et al., 2020)		\checkmark	9.98	13.19	4.48	5.93
SAFECount (You et al., 2023)	\checkmark		5.33	7.04	2.24	3.44
SI-INR		\checkmark	5.54	7.43	2.09	2.70

Generalization Results. We evaluate the robustness of SI-INR and baseline methods to the scale variation in testing data by testing models using images with resolutions different from training images. To further increase the scale variance inside the RSOC dataset, we keep the training images in a fixed scale 512×512 , rescale test images into 5 different resolutions: 205×205 , 307×307 , 410×410 , 512×512 , 614×614 , and test our SI-INR as well as baselines on the rescaled test images. The varying scales in the test set better reflect real-world conditions, where objects may appear at different sizes due to changes in altitude, camera settings, or image cropping. SI-INR significantly outperforms all baseline models when the variation of resolution in test image is presented. The performance advantages over baseline models illustrate that our SI-INR is not only more robust compared to other baselines when processing images with unseen resolutions, but also more data efficient. Especially on PUCPR+ dataset, SI-INR reduces MAE by 70.9% and RMSE by 71.64% compared with effREENet. Furthermore, PSGCNet employs a traditional pyramidal architecture to address object scale variance in object counting. However, as illustrated in Figure 3 and Figure 5, SI-INR achieves superior counting accuracy and produces higher-quality density maps when facing different resolution inputs even with extremely low resolution like $104 \times$ 104, demonstrating its enhanced ability to handle scale variance compared to traditional methods.

- As shown in Figure 3 and Figure 5, SI-INR achieves a relative scale-invariance model while
 truly scale-invariance model is impractical for SESN, SESN relies on group convolutions
 to approximate scale-equivariance. However, the finite set of sampled scales used during
 training and inference means that scale-equivariance is not exact but rather approximate
 within a certain range of scales. Incorporating exact scale-equivariance for all scales would
 require infinite representations, which is computationally infeasible.
- 482 Showcases. We further visualize these results in Figure 3. Although the performances of
 483 all the models degrade with small-scale inputs, our SI-INR can still produce density maps
 484 with the objects well separated. Moreover, the underestimation of object counts is much
 485 less severe in SI-INR compared to all the baselines, which demonstrates the robustness of
 SI-INR under scale variation. We provide additional results in Appendix B.3.



Figure 3: Predicted density maps by SI-INR and other baselines for two test images from RSOC. Two test images are rescaled to 205×205 , 307×307 , 410×410 , 512×512 , 614×614 before fed into the models.

Table 3: Counting performance of handling unseen scales images on the CARPK, PUCPR+, RSOC building datasets and RSOC large-vehicle datasets.

ſ	Mothod	Loss		CARPK		PUCPR+		Building		Large-vehicle	
	Method	MSE.	Bayes.	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
ſ	eFreeNet	\checkmark		50.37	56.83	30.45	35.59	16.97	19.72	49.39	57.26
	PSGCNet	\checkmark		39.36	54.13	49.51	77.05	11.56	16.43	28.74	46.22
	PSGCNet		\checkmark	37.63	52.46	32.58	52.38	12.09	16.96	22.47	39.99
	ASPDNet	\checkmark		41.28	53.62	46.31	75.18	11.31	15.60	26.86	40.17
	ASPDNet		\checkmark	37.25	52.26	35.02	63.90	11.37	16.11	22.11	39.37
	SI-INR		\checkmark	24.30	28.09	8.85	11.91	7.96	11.29	21.89	30.47

4.3 Ablation studies

In this section, we evaluate the effect of each constituting component, test the sensitivity of the counting performance of our SI-INR model with respect to Sampling Rate S_{INR} . We further discuss the effect of different scale-equivariant models in Appendix B.4.

Effect of constituting components. We conduct ablation experiments to study the effect of
each SI-INR component using the RSOC large-vehicle dataset. The results are reported in
Table 4. We observe progressive counting performance improvement by introducing each
of our model components, which shows that all of the scale-equivariant backbone, hybrid
pyramidal scale module, and INR-decoder help improve the counting accuracy.

Table 4: Contributions of each component (SE-backbone, hybrid pyramidal scale module,INR-decoder) in SI-INR.

533	Mathad	RSOC-Large-vehicle		
534	Method	MAE	RMSE	
535	VGG19	20.26	32.75	
536	VGG19+HPSM+INR	14.70	25.74	
537	SE-Backbone+HPSM	15.65	25.43	
538	SE-Backbone+INR	16.28	28.87	
539	SE-Backbone+HPSM+INR	12.61	21.78	

Method	RSO	C-ship	RSOC-small-vehicle			
Method	MAE	RMSE	MAE	RMSE		
$S_{INR}=8$	127.14	173.21	277.50	1017.56		
$S_{INR}=16$	65.49	93.03	273.78	1016.11		
$S_{INR}=32$	62.97	86.81	255.12	821.21		
$S_{INR}=64$	62.26	83.98	243.66	731.51		
$S_{INR}=128$	59.76	81.79	157.18	306.43		

Table 5: Effect of sampling rate S_{INR} on object counting in SI-INR.

Effect of sampling rate S_{INR} . We also report the prediction accuracy and include the 552 predicted density maps by our SI-INR trained with the loss estimated by sampling from the grids of different size $S_{INR} \times S_{INR}$ in Table 5. In this section, we evaluate the counting 554 performance of SI-INR on the RSOC ship and RSOC small-vehicle datasets with S_{INR} = 8, 16, 32, 64, 128. The results show that SI-INR achieves better counting performance as S_{INR} increases from 8 to 128. This effect is particularly pronounced for the RSOC smallvehicle dataset, where the sampling ratio significantly impacts counting accuracy. Since the targets are very small, training with a higher sampling ratio helps the model more accurately locate the vehicles. Besides, this continuous property helps make it easy to balance the computation costs and the counting accuracy requirement. We further visualize several results in Appndix B.7.

5 Conclusions

565 In this paper, we introduce SI-INR, a novel scale-invariant INR implementation that maps discrete grid image signals into continuous 2D function space, maintaining invariance to scal-566 ing variation of the input signals. For object counting, SI-INR achieves SOTA performance, 567 our experiments demonstrate that SI-INR is exceptionally robust and flexible compared to 568 existing methods, capable of processing images of unseen resolutions during testing and ef-569 fectively handling images of various scales during training. This flexibility allows SI-INR to 570 learn and capture more detailed features from different input training images. SI-INR can 571 be easily applied to other image analysis tasks to achieve arbitrary-scale SOTA performance 572 robustly with respect to input image size/resolution. Future work will focus on applying 573 SI-INR to multi-task scenarios, integrating object detection, image segmentation, and depth 574 estimation alongside counting.

575 576 577

578

580

581

583

585

586

588

589

590

553

555

556

557

558

559

560

561 562 563

Reproducibility Statement 6

We have ensured the reproducibility of our experiments by providing detailed descriptions of 579 the model architectures, data augmentation steps, training procedures and hyperparameters setup in Section 4.1. The datasets are introduced in Appendix B.1. For baseline models, we also give training details and hyperparameter configurations in Appendix B.2. Additionally, 582 the code is included in the Supplementary Material and will be made publicly available in a repository to support further verification and replication by other researchers. 584

References

Oliver JD Barrowclough, Georg Muntingh, Varatharajan Nainamalai, and Ivar Stangeby. Binary segmentation of medical images using implicit spline representations and deep learning. Computer Aided Geometric Design, 85:101972, 2021.

591 Zhaowei Cai, Quanfu Fan, Rogerio S Feris, and Nuno Vasconcelos. A unified multi-scale deep convolutional neural network for fast object detection. In Computer Vision–ECCV 592 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14, pp. 354–370. Springer, 2016.

- Antoni B Chan, Zhang-Sheng John Liang, and Nuno Vasconcelos. Privacy preserving crowd 595 monitoring: Counting people without people models or tracking. In 2008 IEEE Conference 596 on Computer Vision and Pattern Recognition, pp. 1–7. IEEE, 2008. 597 I Chen, Wei-Ting Chen, Yu-Wei Liu, Ming-Hsuan Yang, Sy-Yen Kuo, et al. Improving point-598 based crowd counting and localization based on auxiliary point guidance. In European 599 Conference on Computer Vision, pp. 428–444. Springer, 2025. 600 Zhi-Qi Cheng, Qi Dai, Hong Li, Jingkuan Song, Xiao Wu, and Alexander G Hauptmann. 602 Rethinking spatial invariance of convolutional networks for object counting. In Proceed-603 ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 604 19638-19648, 2022. 605 Taco Cohen and Max Welling. Group equivariant convolutional networks. In International 606 conference on machine learning, pp. 2990–2999. PMLR, 2016. 607 608 Li Dong, Haijun Zhang, Yuzhu Ji, and Yuxin Ding. Crowd counting by using multi-level 609 density-based spatial information: A multi-scale cnn framework. Information Sciences, 610 528:79-91, 2020. 611 612 Emilien Dupont, Yee Whye Teh, and Arnaud Doucet. Generative models as distributions 613 of functions. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.), Proceedings of The 25th International Conference on Artificial Intelligence and Statistics, 614 volume 151 of Proceedings of Machine Learning Research, pp. 2989–3015. PMLR, 28–30 615 Mar 2022. URL https://proceedings.mlr.press/v151/dupont22a.html. 616 617 Min Fu, Pei Xu, Xudong Li, Qihe Liu, Mao Ye, and Ce Zhu. Fast crowd density estimation 618 with convolutional neural networks. Engineering Applications of Artificial Intelligence, 619 43:81-88, 2015. 620 621 Guangshuai Gao, Qingjie Liu, and Yunhong Wang. Counting from sky: A large-scale data 622 set for remote sensing object counting and a benchmark method. IEEE Transactions on 623 Geoscience and Remote Sensing, 59(5):3642-3655, 2020. 624 Guangshuai Gao, Qingjie Liu, Zhenghui Hu, Lu Li, Qi Wen, and Yunhong Wang. PSGCNet: 625 A pyramidal scale and global context guided network for dense object counting in remote-626 sensing images. IEEE Transactions on Geoscience and Remote Sensing, 60:1–12, 2022. 627 628 Meng-Ru Hsieh, Yen-Liang Lin, and Winston H Hsu. Drone-based object counting by 629 spatially regularized regional proposal network. In Proceedings of the IEEE international conference on computer vision, pp. 4145–4153, 2017. 630 631 Yongbo Huang, Yuanpei Jin, Liqiang Zhang, and Yishu Liu. Remote sensing object counting 632 through regression ensembles and learning to rank. IEEE Transactions on Geoscience and 633 Remote Sensing, 61:1–17, 2023. 634 635 Haroon Idrees, Muhmmad Tayyab, Kishan Athrey, Dong Zhang, Somaya Al-Maadeed, Nasir 636 Rajpoot, and Mubarak Shah. Composition loss for counting, density map estimation and 637 localization in dense crowds. In Proceedings of the European conference on computer 638 vision (ECCV), pp. 532–546, 2018. 639 Angjoo Kanazawa, Abhishek Sharma, and David Jacobs. Locally scale-invariant convolu-640 tional neural networks. arXiv preprint arXiv:1412.5104, 2014. 641 642 Chiheon Kim, Dovup Lee, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Generalizable 643 implicit neural representations via instance pattern composers. In Proceedings of the 644 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11808–11817, 645 2023.646 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv
- ⁶⁴⁷ Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.

676

677

678

682

683

684

- Nikola Kovachki, Zongyi Li, Burigede Liu, Kamyar Azizzadenesheli, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Neural operator: Learning maps between function spaces with applications to pdes. Journal of Machine Learning Research, 24(89): 1–97, 2023.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger (eds.), Advances in Neural Information Processing Systems, volume 25. Curran Associates, Inc., 2012. URL https://proceedings.neurips.cc/paper_files/paper/2012/file/ c399862d3b9d6b76c8436e924a68c45b-Paper.pdf.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning
 applied to document recognition. Proceedings of the IEEE, 86(11):2278–2324, 1998.
- ⁶⁶⁰ Doyup Lee, Chiheon Kim, Minsu Cho, and WOOK SHIN HAN. Locality-aware generalizable
 ⁶⁶¹ implicit neural representation. Advances in Neural Information Processing Systems, 36, 2024.
 ⁶⁶³
- Victor Lempitsky and Andrew Zisserman. Learning to count objects in images. Advances in Neural Information Processing Systems, 23, 2010.
- Yuhong Li, Xiaofan Zhang, and Deming Chen. CSRNet: Dilated convolutional neural networks for understanding the highly congested scenes. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1091–1100, 2018.
- ⁶⁶⁹
 ⁶⁷⁰ Zhaoxin Li, Shuhua Lu, Yishan Dong, and Jingyuan Guo. Msffa: a multi-scale feature fusion and attention mechanism network for crowd counting. The Visual Computer, 39 (3):1045–1056, 2023.
- Dingkang Liang, Wei Xu, and Xiang Bai. An end-to-end transformer model for crowd localization. In European Conference on Computer Vision, pp. 38–54. Springer, 2022.
 - Hui Lin, Zhiheng Ma, Rongrong Ji, Yaowei Wang, and Xiaopeng Hong. Boosting crowd counting via multifaceted attention. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 19628–19637, 2022.
- Sheng-Fuu Lin, Jaw-Yeh Chen, and Hung-Xin Chao. Estimation of number of people in crowded scenes using perspective transformation. IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans, 31(6):645–654, 2001.
 - Tony Lindeberg. Scale-space theory: A basic tool for analyzing structures at different scales. Journal of applied statistics, 21(1-2):225–270, 1994.
- Tony Lindeberg. Scale-covariant and scale-invariant gaussian derivative networks. Journal of Mathematical Imaging and Vision, 64(3):223-242, 2022.
- ⁶⁸⁷ Chengxin Liu, Hao Lu, Zhiguo Cao, and Tongliang Liu. Point-query quadtree for crowd counting, localization, and more. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1676–1685, 2023.
- ⁶⁹⁰
 ⁶⁹¹ David G Lowe. Object recognition from local scale-invariant features. In Proceedings of the seventh IEEE international conference on computer vision, volume 2, pp. 1150–1157. Ieee, 1999.
- David G Lowe. Distinctive image features from scale-invariant keypoints. International
 journal of computer vision, 60:91–110, 2004.
- Kihaier Luo, Xiaoning Qian, and Byung-Jun Yoon. Hierarchical neural operator transformer
 with learnable frequency-aware loss prior for arbitrary-scale super-resolution. arXiv
 preprint arXiv:2405.12202, 2024.
- Zheng Ma, Lei Yu, and Antoni B Chan. Small instance detection by integer programming on object density maps. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3689–3697, 2015.

- Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong. Bayesian loss for crowd count estimation with point supervision. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 6142–6151, 2019.
- Diego Marcos, Benjamin Kellenberger, Sylvain Lobry, and Devis Tuia. Scale equivariance in cnns with vector fields. arXiv preprint arXiv:1807.11783, 2018.
- Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4460–4470, 2019.
- Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger.
 Texture fields: Learning texture representations in function space. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 4531–4540, 2019.
- Joseph Paul Cohen, Genevieve Boucher, Craig A Glastonbury, Henry Z Lo, and Yoshua
 Bengio. Count-ception: Counting by fully convolutional redundant counting. In Proceedings of the IEEE International Conference on Computer Vision workshops, pp. 18–26, 2017.
- Md Ashiqur Rahman and Raymond A Yeh. Truly scale-equivariant deep nets with fourier
 layers. Advances in Neural Information Processing Systems, 36:6092–6104, 2023.
- Viresh Ranjan, Hieu Le, and Minh Hoai. Iterative crowd counting. In Proceedings of the European conference on computer vision (ECCV), pp. 270–285, 2018.
- Jihye Ryu and Kwangho Song. Crowd counting and individual localization using pseudo square label. IEEE Access, 2024.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
 image recognition. arXiv preprint arXiv:1409.1556, 2014.
- Vishwanath A Sindagi and Vishal M Patel. Ha-ccn: Hierarchical attention-based crowd counting network. IEEE Transactions on Image Processing, 29:323–335, 2019.
- Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. Advances in Neural Information Processing Systems, 33:7462–7473, 2020.
- Qingyu Song, Changan Wang, Zhengkai Jiang, Yabiao Wang, Ying Tai, Chengjie Wang,
 Jilin Li, Feiyue Huang, and Yang Wu. Rethinking counting and localization in crowds: A
 purely point-based framework. In Proceedings of the IEEE/CVF International Conference
 on Computer Vision, pp. 3365–3374, 2021.
- 740 Ivan Sosnovik, Michał Szmaja, and Arnold Smeulders. Scale-equivariant steerable networks.
 741 arXiv preprint arXiv:1910.11093, 2019.
- Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan,
 Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let
 networks learn high frequency functions in low dimensional domains. Advances in Neural
 Information Processing Systems, 33:7537–7547, 2020.
- Jia Wan and Antoni Chan. Adaptive density map generation for crowd counting. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1130–1139, 2019.
- Jia Wan, Ziquan Liu, and Antoni B Chan. A generalized loss function for crowd counting and localization. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 1974–1983, 2021.
- Chuan Wang, Hua Zhang, Liang Yang, Si Liu, and Xiaochun Cao. Deep people counting in extremely dense crowds. In Proceedings of the 23rd ACM international conference on Multimedia, pp. 1299–1302, 2015.

- Yuehao Wang, Yonghao Long, Siu Hin Fan, and Qi Dou. Neural rendering for stereo 3d reconstruction of deformable tissues in robotic surgery. In International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 431–441. Springer, 2022.
- Longyin Wen, Dawei Du, Pengfei Zhu, Qinghua Hu, Qilong Wang, Liefeng Bo, and Siwei Lyu. Detection, tracking, and counting meets drones in crowds: A benchmark. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7812–7821, 2021.
- Yichong Xu, Tianjun Xiao, Jiaxing Zhang, Kuiyuan Yang, and Zheng Zhang. Scale-invariant convolutional neural networks. arXiv preprint arXiv:1411.6369, 2014.
- Yilong Yang, Srinandan Dasmahapatra, and Sasan Mahmoodi. Scale-equivariant unet for
 histopathology image segmentation. arXiv preprint arXiv:2304.04595, 2023.
- Jun Yi, Zhilong Shen, Fan Chen, Yiheng Zhao, Shan Xiao, and Wei Zhou. A lightweight
 multiscale feature fusion network for remote sensing object counting. IEEE Transactions
 on Geoscience and Remote Sensing, 61:1–13, 2023.
- Zhiyuan You, Kai Yang, Wenhan Luo, Xin Lu, Lei Cui, and Xinyi Le. Few-shot object counting with similarity-aware feature enhancement. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 6315–6324, 2023.
- Weiwei Zhan, Guangmin Sun, and Yu Li. Scale-equivariant steerable networks for crowd counting. In 2022 7th International Conference on Control and Robotics Engineering (ICCRE), pp. 174–179. IEEE, 2022.
- Yingying Zhang, Desen Zhou, Siqin Chen, Shenghua Gao, and Yi Ma. Single-image crowd counting via multi-column convolutional neural network. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 589–597, 2016.
- Fushun Zhu, Hua Yan, Xinyue Chen, Tong Li, and Zhengyu Zhang. A multi-scale and multi-level feature aggregation network for crowd counting. Neurocomputing, 423:46–56, 2021.
- 786 787 788

A Supplementary information of model construction

790 A.1 Scale equivariance and invariance

791 Consider a Scale-Translation Group consisting of a Scaling Group G_S and a Translation 792 Group G_T , $H = \{h = (s,t) | s \in G_S, t \in G_T\}$, where h denotes an element of H and rep-793 resents one scale-translation operator, G_S denotes the Scaling Group, which accounts for 794 transformations that scale an object or function, and G_T denotes the Translation Group, 795 which handles shifting the object or function within its domain. Besides, s is the scaling pa-796 rameter, indicating how the input is stretched or compressed; t is the translation parameter, 797 specifying the shifting in the domain.

From the group theory, given an image $I_a \in V_1$, a mapping $\Phi: V_1 \to V_2$ is scale-equivariant if:

801 802

805

$$\Phi(p_1(h)(\mathbf{I}_a)) = p_2(h)(\Phi(\mathbf{I}_a)),$$
(12)

where V_2 denotes the output domain, $p_1(\cdot)$ and $p_2(\cdot)$ denote the corresponding group actions of h acting on V_1 , V_2 . If $p_2(h)$ is the identity mapping, the mapping Φ is scale-invariant.

A.2 Scale-translation equivariance of SESN

Consider a steerable convolution filter (Sosnovik et al., 2019), $\varphi_m(x) = m^{-1}\varphi(m^{-1}x)$, which has the following property:

$$p(s^{-1})(\varphi_m)(x) = \varphi_m(sx) = s^{-1}\varphi_{s^{-1}m}(x),$$
(13)

. .

820

843

849

850

854

859

where p(s) denotes the group action of s on convolution filters. The scaling of this filter is the transformation of its parameters.

With scale-translation group H and steerable convolution filters $\varphi_m(\cdot)$, the groupequivariant convolution on f can be defined as:

$$[f \star_{H} \psi_{m}](s,t) = \int_{S} \int_{T} f(s',t') p(s,t) [\psi_{m}](s',t') d\mu(s') d\mu(t')$$

$$= \sum_{s'} \int_{T} f(s',t') \psi_{sm} (s^{-1}s',t'-t) dt' = \sum_{s'} [f(s',\cdot) \star \psi_{sm} (s^{-1}s',\cdot)](t).$$
(14)

The proof of translation-equivariance of Equation (14) is as follows:

$$[p(\hat{t})[f] \star_{H} \psi_{m}] (s,t) = \sum_{s'} [p(\hat{t})[f] (s', \cdot) \star \psi_{sm} (s^{-1}s', \cdot)] (t)$$

$$= \sum_{s'} p(\hat{t}) [f (s', \cdot) \star \psi_{sm} (s^{-1}s', \cdot)] (t)$$

$$= p(\hat{t}) \left\{ \sum_{s'} [f (s', \cdot) \star \psi_{sm} (s^{-1}s', \cdot)] \right\} (t)$$

$$= p(\hat{t}) [f \star_{H} \psi_{m}] (s, t).$$

$$(15)$$

The proof of scale-equivariance of Equation (14) can be obtained:

$$[p(\hat{s})[f] \star_{H} \psi_{m}](s,t) = \sum_{s'} [p(\hat{s})[f](s',\cdot) \star \psi_{sm}(s^{-1}s',\cdot)](t)$$

$$= \sum_{s'} p(\hat{s}) [f(\hat{s}^{-1}s',\cdot) \star \psi_{\hat{s}^{-1}sm}(s^{-1}s',\cdot)](t)$$

$$= \sum_{s''} [f(s'',\cdot) \star \psi_{\hat{s}^{-1}sm}(\hat{s}s^{-1}s'',\cdot)](\hat{s}^{-1}t)$$

$$= [f \star_{H} \psi_{m}](\hat{s}^{-1}s,\hat{s}^{-1}t)$$

$$= p(\hat{s}) [f \star_{H} \psi_{m}](s,t).$$
(16)

Finally, we have the proof of scale-translation equivariance of Equation (14):

$$p(\hat{s}\hat{t})[f] \star_{H} \psi_{m} = p(\hat{s})p(\hat{t})[f] \star_{H} \psi_{m} = p(\hat{s}) \left[p(\hat{t})[f] \star_{H} \psi_{m}\right]$$

= $p(\hat{s})p(\hat{t}) \left[f \star_{H} \psi_{m}\right] = p(\hat{s}\hat{t}) \left[f \star_{H} \psi_{m}\right].$ (17)

The summation of two scale-equivariant models is still scale-equivariant by the following Lemma.

Lemma 1 The summation of two equivariant mappings $\Phi_1 : V_1 \to V_2, \Phi_2 : V_1 \to V_2$ is still equivariant.

Proof: For any scale-translation operator h we have:

$$(\Phi_1 + \Phi_2)(p_1(h)(\mathbf{I}_a)) = \Phi_1(p_1(h)(\mathbf{I}_a)) + \Phi_2(p_1(h)(\mathbf{I}_a)) = p_2(h)(\Phi_1(\mathbf{I}_a)) + p_2(h)(\Phi_2(\mathbf{I}_a)) = p_2(h)((\Phi_1 + \Phi_2)(\mathbf{I}_a)).$$
(18)

A.3 Derivation of the minimization objective \mathcal{L}_{ELBO}

Here we detail the process of deriving the minimization objective in Section 3.3 of the Main Text.

For a given image **I**, we define the counting annotation map as $\mathcal{D}_{\mathbf{I}} = \{(\mathbf{m}_n, y_n)\}_1^N$. Where $\mathbf{m}_n \in [0, 1]^2$ denotes the normalized image-coordinate position of the *n*-th object, $y_n = n$ is

the corresponding label for each object, and N denotes the total number of labeled objects in I. Typically, people generate ground truth density maps \mathbf{D}^{gt} by convolving this annotation 866 map with a Gaussian kernel (Fu et al., 2015; Paul Cohen et al., 2017; Gao et al., 2020).

867 Following Ma et al. (2019), we define $y(\cdot): \mathbb{R}^2 \to \{1, \cdots, N\}$, with $y(\mathbf{x})$ as whether the 868 location **x** belongs to one of the N objects computed by a prior distribution $p(y(\cdot))$. The 869 posterior distribution can be expressed as: 870

> $p(y(\cdot) = n \mid \mathcal{D}_{\mathbf{I}_a}) = \frac{\mathcal{N}\left(\cdot; \mathbf{m}_n, \sigma^2 \mathbf{1}_{2 \times 2}\right)}{\sum_{i=1}^N \mathcal{N}\left(\cdot; \mathbf{m}_i, \sigma^2 \mathbf{1}_{2 \times 2}\right)}.$ (19)

Based on the definition of $c_n^{pt}(\cdot)$ in Ma et al. (2019), the contribution in the likelihood can be expressed as follows:

$$c_n^{gt}(\cdot) = p\left(y(\cdot) = n \mid \mathcal{D}_{\mathbf{I}_a}\right) \times \mathbf{D}^{gt} \left(\cdot\right) = \frac{\mathcal{N}\left(\cdot; \mathbf{m}_n, \sigma^2 \mathbf{1}_{2 \times 2}\right)}{\sum_{i=1}^N \mathcal{N}\left(\cdot; \mathbf{m}_i, \sigma^2 \mathbf{1}_{2 \times 2}\right)} \times \sum_{i=1}^N \mathcal{N}\left(\cdot; \mathbf{m}_i, \sigma^2 \mathbf{1}_{2 \times 2}\right)$$
$$= \mathcal{N}\left(\cdot; \mathbf{m}_n, \sigma^2 \mathbf{1}_{2 \times 2}\right), \tag{20}$$

We can find that $c_n^{pt}(\cdot)$ is a Gaussian distribution centered at m_n , and the summation of 882 $c_n^{gt}(\cdot)$ equals one. 883

Thanks to the continuous property of SI-INR. Bayesian counting loss can be represented as:

$$\mathcal{L}_{BAY} = \frac{1}{A} \sum_{a=1}^{A} \Big[\sum_{n=0}^{N} (\mathbb{E}_{\mathbf{x} \sim p_s(\mathbf{x})}(c_{n,a}) - \mathbb{E}_{\mathbf{x} \sim p_s(\mathbf{x})}(c_{n,a}^{\text{gt}})) \Big].$$
(21)

where $p_s(\mathbf{x})$ is any probability distribution of \mathbf{x} . This optimization objective enables our 889 model to be trained using any existing stochastic optimization algorithm. Without loss of generality, we can select $p(\mathbf{x})$ as a uniform distribution, $\mathbf{x} \sim \text{Uniform}[0, 1]^2$.

Besides, we add MAE counting loss into our loss function, the final minimization objective is:

$$\mathcal{L}_{SI-INR} = \frac{1}{A} \sum_{a=1}^{A} \left[\sum_{n=0}^{N} (\mathbb{E}_{\mathbf{x} \sim \text{Unif}[0,1]^2}(c_{n,a}) - \mathbb{E}_{\mathbf{x} \sim \text{Unif}[0,1]^2}(c_{n,a}^{\text{gt}})) + \kappa (\sum_{n=0}^{N} (\mathbb{E}_{\mathbf{x} \sim \text{Unif}[0,1]^2}(c_{n,a})) - N) \right]$$
(22)

899 900

901

902 903

904

871

872 873 874

875

881

884

885 886 887

890

891 892

893

> В Sumplementary information of experiments

B.1 Data

(Gao et al., 2020): The Remote Sensing Object Counting (RSOC) dataset is a RSOC 905 large-scale benchmark specifically designed for counting objects in satellite imagery. It 906 includes a total of 3,057 images with 286,539 annotated object instances. The dataset is 907 divided into four distinct subdatasets, each focused on a different object type: Buildings, 908 Small Vehicles, Large Vehicles, and Ships. The RSOC Buildings dataset contains 1,205 909 training images and 1,263 test images, where the image resolution is 512×512 . The RSOC 910 Small Vehicles dataset has 222 training images and 58 test images. The image resolution 911 ranges from 421×799 to 12029×5014 . Large Vehicles consists of 108 training images 912 and 64 test images, The image resolution ranges from 731×596 to 6327×5662 . And the 913 Ships subset has 97 training images and 60 test images. The image resolution ranges from 914 606×1065 to 6335×3591 .

915 (Hsieh et al., 2017): The Car Parking Lot Dataset (CARPK) is a benchmark CARPK 916 for car counting tasks, consisting of 1,148 images taken from drone perspectives over four 917 parking lots, containing 89,777 annotated cars. These images capture real-world scenarios

Figure 4: Example images from the RSOC dataset.

with dense vehicle arrangements, making the dataset challenging for object detection and counting tasks. The average resolution of the images is 1280×720 pixels, providing detailed aerial views. Each image is annotated with bounding boxes around individual cars, making the dataset suitable for both object counting and detection. The dataset is split into 989 training images and 459 testing images.

PUCPR+ (Hsieh et al., 2017): The Pontifical Catholic University of Parana+ Dataset (PUCPR+) is a specialized car counting resource where all images are captured from the 10th floor of a building. PUCOR+ contains 125 images with 16,456 cars, where 100 images are set for training, while the remaining images are utilized for testing the models.

Visualization We further provide several exemplar images from RSOC datasets in Figure 948 4. It can be found that the objects within the same image naturally appear of similar 949 size. However, remote sensing datasets, including RSOC, encompass images with a wide 950 range of resolutions. As a result, object sizes vary significantly across different images, even 951 if they appear uniform within a single image. Furthermore, we resize images to various 952 resolutions to evaluate robustness to scale variability which further increases the range of 953 scale differences across different images in our experiments. In Figure 4, the top-left two 954 images are both from the RSOC large-vehicle dataset, clearly showing that the cars in the 955 second image are three times larger than those in the first image. Similarly, the bottom-left 956 two images, from the RSOC small-vehicle dataset, highlight the differences in visibility: cars 957 are clearly seen in the first image but are almost invisible in the second.

B.2 Baselines

920 921

923 924 925

933 934

935 936

937

938

939

940 941

942

943

944

945

946 947

958 959

960

963

961 In this section, we delve into the training specifics for all baseline models utilized in our 962 experiments.

964 (Gao et al., 2020): ASPDNet is an advanced attention-based network that in-ASPDNet 965 tegrates scale pyramids and deformable convolutions to effectively utilize attention mecha-966 nisms. This architecture captures extensive contextual and high-level semantic information, 967 which aids in reducing the impact of cluttered backgrounds while emphasizing the regions of 968 interest. In our study, we follow its original network design, set the batch size as 16, replace 969 the original Stochastic Gradient Descent (SGD) optimizer with ADAM (Kingma & Ba, 2014) optimizer, and set the learning rate as 1e - 4 to enhance the training results. Our training 970 spans 200 epochs. ASPDNet is trained under MSE counting and Bayes counting (Ma et al., 971 2019) respectively to get the best counting performance.

972 Input scale:102 Input scale:153 973 974 975 976 977 978 979 980 981 982 983 GT COUNT: 50 984 985 986 987

988 Figure 5: Predicted density maps by SI-INR and other baselines for two test images from 989 RSOC. Two test images are rescaled to 102×102 , 153×153 , 205×205 , 256×256 , 307×307 990 before fed into the models.

PSGCNet (Gao et al., 2022): PSGCNet integrates pyramidal scale and global context 992 modules to handle scale variations of remote sensing images. We follow the network setup, setting the learning rate as 1e - 4, and using a batch size of 16. We trained PSGCNet with 994 original Bayesian-based counting loss and MSE respectively. Our training spans 200 epochs. 995

(Huang et al., 2023): The eFreeNet is an ensemble of first-rank-then-estimate eFreeNet 997 networks that tailors a ranking metric optimization scheme to fit object counting. The study 998 employs the default network architecture. In the optimization setup, we set the backbone's 999 learning rate as 1e - 5 while 1e - 5 for other components following the original setup. The 1000 ensemble number is set as 8 to get the best counting performance. Our training spans 3000 1001 epochs with a batch size of 8.

1002

991

993

1003 B.3 Additional qualitative results 1004

1005 We provide additional qualitative results here. For the figures from the RSOC building dataset, two test images are rescaled to 102×102 , 153×153 , 205×205 , 256×256 , and 307×307 before being input into the models. The results in Figure 5 highlight not only 1007 the counting accuracy of our model but also its robustness and ability to handle inputs of 1008 varying resolutions. 1009

- 1010
- B.4 Effect of different Scale-equivariant models 1011

1012 In our experiments, we test SESC (Sosnovik et al., 2019) and scale-equivariant Fourier lay-1013 ers (Rahman & Yeh, 2023) in SI-INR, Compared with SESC, scale-equivariant Fourier layers 1014 demand more computational resources, especially processing images over 512 resolutions. 1015 On the RSOC building dataset, training one epoch for SI-INR with SESC takes 161 seconds, 1016 compared with ASPDNet's 167 seconds and PSGCNet's 121 seconds. However, SI-INR with 1017 scale-equivariant Fourier layers takes more than 900 seconds, which is close to one order of 1018 magnitude more costly.

1019

1020 Comparison on UCF-QNRF dataset B.51021

1022 We compare our SI-INR with state-of-the-art methods as well as our baselines on the UCF-1023 QNRF (University of Central Florida - Qatar National Research Fund) dataset (Idrees et al., 2018), which is a highly diverse dataset consisting of 1,535 images with over 1.2 million 1024 annotated individuals, spanning a wide range of crowd densities and changing object sizes. 1025 We report the results in Table 6.

1027	Table 6: Performance Comparison on	the UCF	-QNRF Da	ataset
1028	Model	MAE	RMSE	
1029	MMNet (Dong et al., 2020)	104.00	178.00	
1030	MSFFA (Li et al., 2023)	94.60	170.60	
1031	MFANet (Zhu et al., 2021)	97.7	166.00	
1032	CLTR (Liang et al., 2022)	85.80	141.30	
1033	Bayesian+ (Ma et al., 2019)	88.70	154.80	
1034	P2PNet (Song et al., 2021)	85.32	154.50	
1035	GauNet (Cheng et al., 2022)	81.60	153.71	
1035	APGCC (Chen et al., 2025)	80.10	136.60	
1030	PSL-Net (Ryu & Song, 2024)	85.50	144.40	
1037	PET (Liu et al., 2023)	79.53	144.32	
1038	PSGCNet (Baseline)	86.30	149.50	
1039	SI-INR (Ours)	80.89	134.73	
1040				

1048

1026

The reported results indicate that SI-INR achieves competitive performance, with a Mean Absolute Error (MAE) of 80.89 and a Root Mean Squared Error (RMSE) of 134.73. Additionally, Compared with existing density map based methods, such as Bayesian+ (Ma et al., 2019), GauNet (Cheng et al., 2022), and our baseline PSGCNet (Gao et al., 2022), SI-INR demonstrates consistent improvements in both metrics. These results highlight the effectiveness of our proposed approach.

1049 B.6 Effect of different methods for handling multi-scale challenges

1050 For efficiency and adaptability, scale-equivariant methods adjust to different scales without 1051 having separate filters for each scale, unlike traditional methods that may rely on resizing 1052 inputs or using multiple filters for different scales. Many traditional multi-scale methods, 1053 such as image pyramids or multi-resolution networks, may struggle with high computational 1054 costs because they process the same image at multiple resolutions, leading to increased 1055 complexity especially with large images or when handling many scales. Besides, traditional multi-scale approaches do not focus on deriving scale-invariant outputs compared with our 1056 SI-INR. 1057

In our experiment, we have compared our SI-INR with four SOTA methods (PSGCNet (Gao et al., 2022), MMNet (Dong et al., 2020), MFANet (Zhu et al., 2021), MSFFA (Li et al., 2023)) that mainly focus on handling multi-scale challenges. PSGCNet applies a pyramidal network to handle multi-scale challenges, MMNet leverages multi-level density-based spatial information, MFANet introduces multi-level feature aggregation, and MSFFA integrates multi-scale feature fusion and attention mechanisms. As we report in Table 6, our SI-INR outperforms these methods on the UCF-QNRF crowd-counting dataset, which demonstrates its superior performance in handling objects of different sizes.

1066 1067 B.7 Visualization of different sampling rate S_{INR} of SI-INR

1068
1069To further demonstrate the effect of sample rate S_{INR} of SI-INR, we visualize SI-INR's
outputs when setting the sample rate from 8 to 128 in the Figure 6. In this ablation
experiment, we let the well-trained SI-INR model directly generate 5 different resolution
density maps, we can find that SI-INR can generate high-quality density maps when the
sample rate sampling rate S_{INR} increases.

- 1073 1074
- 1075
- 1076
- 1077
- 1078
- 1079

