

Neural Point Processes for Pixel-wise Regression

Chengzhi Shi, Gözde Özcan, Miquel Sirera Perelló, Yuanyuan Li, Nina Iftikhar Shamsi, Stratis Ioannidis
Northeastern University, Boston, MA, USA

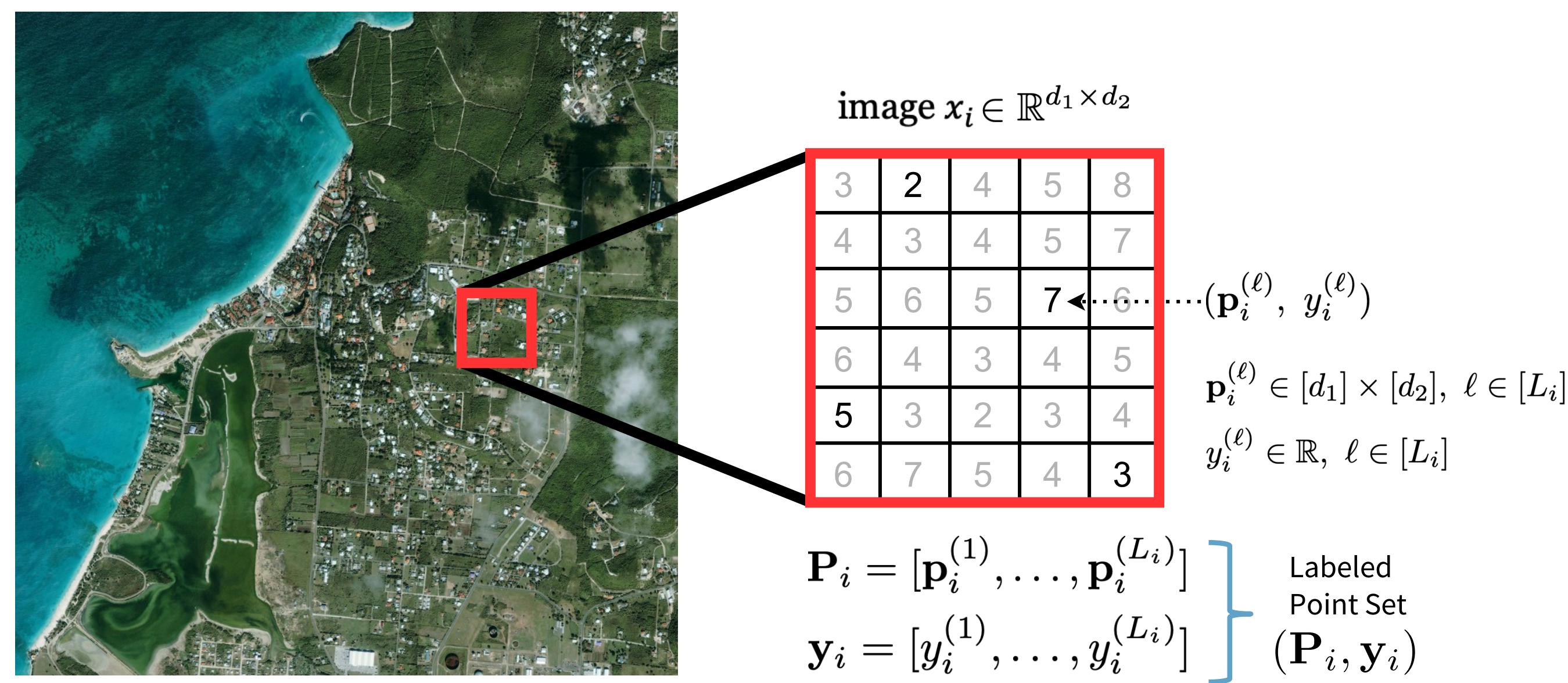


Motivation

- Many real-world tasks only provide **labels at a small subset of pixels**. Some examples are:
 - Medical Imaging
 - Satellite Imaging
 - Remote Sensing
 - Air Pollution Data
 - Challenge:** how to **generalize from sparse supervision**?
 - The goal is **pixel-wise regression**: to learn a model that predicts the continuous label value at **every pixel location**, based only on the image and the sparse supervision
-

Problem Formulation

- Let the **dataset** $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{P}_i, \mathbf{y}_i)\}_{i=1}^n$ be defined such that each sample consists of an image $\mathbf{x}_i \in \mathbb{R}^{d_1 \times d_2}$, where $i \in [n]$, and an associated set of labeled points $(\mathbf{P}_i, \mathbf{y}_i)$, detailed below:



- We fit a neural network (NN) $f(\mathbf{x}; \theta) \in \mathbb{R}^{d_1 \times d_2}$, parametrized by $\theta \in \mathbb{R}^{m_\theta}$, that takes an input image and predicts values at all possible pixel locations $\mathbf{p} \in [d_1] \times [d_2]$

- Base approach:** minimize the Euclidean error only at labeled points

$$\mathcal{L}_{\text{MSE}}(\theta, \mathcal{D}) = \sum_{i=1}^n \sum_{\ell=1}^{L_i} \left(f_{\mathbf{p}_i^{(\ell)}}(\mathbf{x}_i; \theta) - y_i^{(\ell)} \right)^2 = \sum_{i=1}^n \|f_{\mathbf{P}_i}(\mathbf{x}_i; \theta) - \mathbf{y}_i\|^2$$

Problem: Ignores spatial relationships between nearby pixels

Acknowledgement

Research was sponsored by the United States Army Core of Engineers (USACE) Engineer Research and Development Center (ERDC) Geospatial Research Laboratory (GRL) and was accomplished under Cooperative Agreement Federal Award Identification Number (FAIN) W9132V-22-2-0001. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of USACE EDRC GRL or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

Our Method: Neural Point Processes (NPPs)

Bringing structure to sparse pixel-wise regression

- NPP approach: we introduce spatial correlations by modeling labels are modeled as a Gaussian Process (GP) over the DNN output

- We assume that for every point \mathbf{p}_i , the corresponding label y_i is given by:

$$y_i = g_i(\mathbf{p}_i) + \varepsilon$$

where $\varepsilon \sim N(0, \sigma_0^2)$ is i.i.d noise, and $g_i(\cdot)$ is a GP:

$$g_i(\cdot) \sim \mathcal{GP}(m_i(\cdot), k_i(\cdot, \cdot))$$

with mean function $m_i(\mathbf{p}) = f_{\mathbf{p}}(\mathbf{x}_i; \theta) \in \mathbb{R}$ (output of the NN) and kernel function $k_i(\mathbf{p}, \mathbf{p}') = k(\mathbf{p}, \mathbf{p}'; \zeta_i)$ (parametric PSD kernel)

The NPP method:

- Encourage Smoothness
- Accounts for proximity

- With this setup, we can obtain our MLE of θ by minimizing:

$$\mathcal{L}_{\text{NPP}}(\theta, \zeta; \mathcal{D}) = \frac{1}{2} \sum_{i=1}^n \left[\log |k(\mathbf{p}_i, \mathbf{p}_i; \zeta) + 2\sigma_0^2 \mathbf{I}| + (\mathbf{y}_i - f_{\mathbf{p}_i}(\mathbf{x}_i; \theta))^T (k(\mathbf{p}_i, \mathbf{p}_i; \zeta) + \sigma_0^2 \mathbf{I})^{-1} (\mathbf{y}_i - f_{\mathbf{p}_i}(\mathbf{x}_i; \theta)) \right]$$

Kernel Regularization Term Squared Mahalanobis distance Additional noise term

enforces correlations between the values of points that are proximal

Key insight: since NPPs model outputs as a **Gaussian Process**, we enable **test-time updates**

- When partial labels ($\mathbf{P}^\dagger, \mathbf{Y}^\dagger$) are available, we can compute the **posterior distribution** over predictions for the rest of the points we want to predict, namely $(\mathbf{P}^*, \mathbf{Y}^*)$

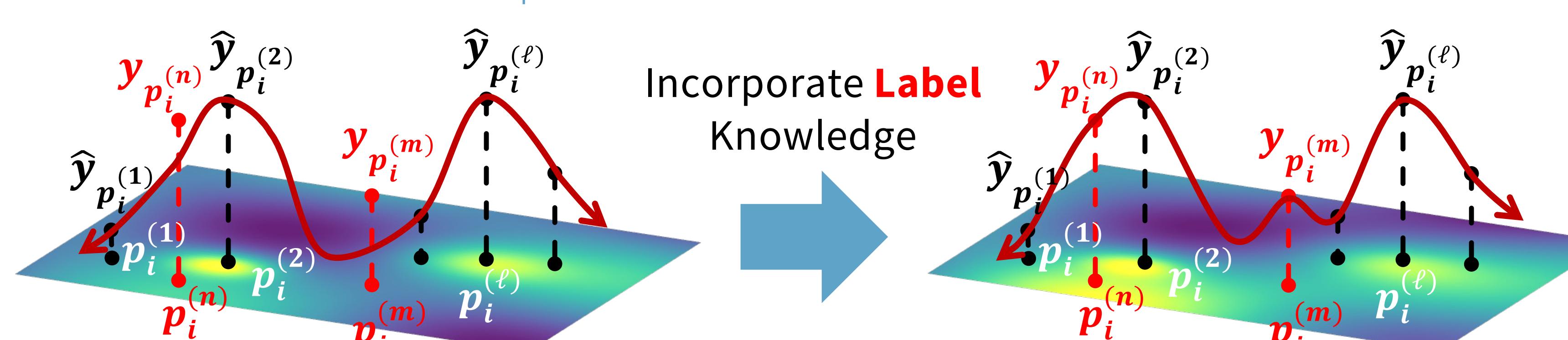
- Conditioned on observations \mathbf{Y}^\dagger , the posterior distribution of \mathbf{Y}^* is Gaussian with the following mean and covariance:

$$\begin{aligned} \mathbb{E}[\mathbf{Y}^*|\mathbf{Y}^\dagger] &= m(\mathbf{P}^*) + k(\mathbf{P}^*, \mathbf{P}^\dagger) \mathbf{K}_{\dagger, \dagger}^{-1} (\mathbf{y}^\dagger - m(\mathbf{P}^\dagger)) \\ \text{cov}(\mathbf{Y}^*) &= k(\mathbf{P}^*; \mathbf{P}^*) - k(\mathbf{P}^*, \mathbf{P}^\dagger) \mathbf{K}_{\dagger, \dagger}^{-1} k(\mathbf{P}^\dagger; \mathbf{P}^*) \end{aligned}$$

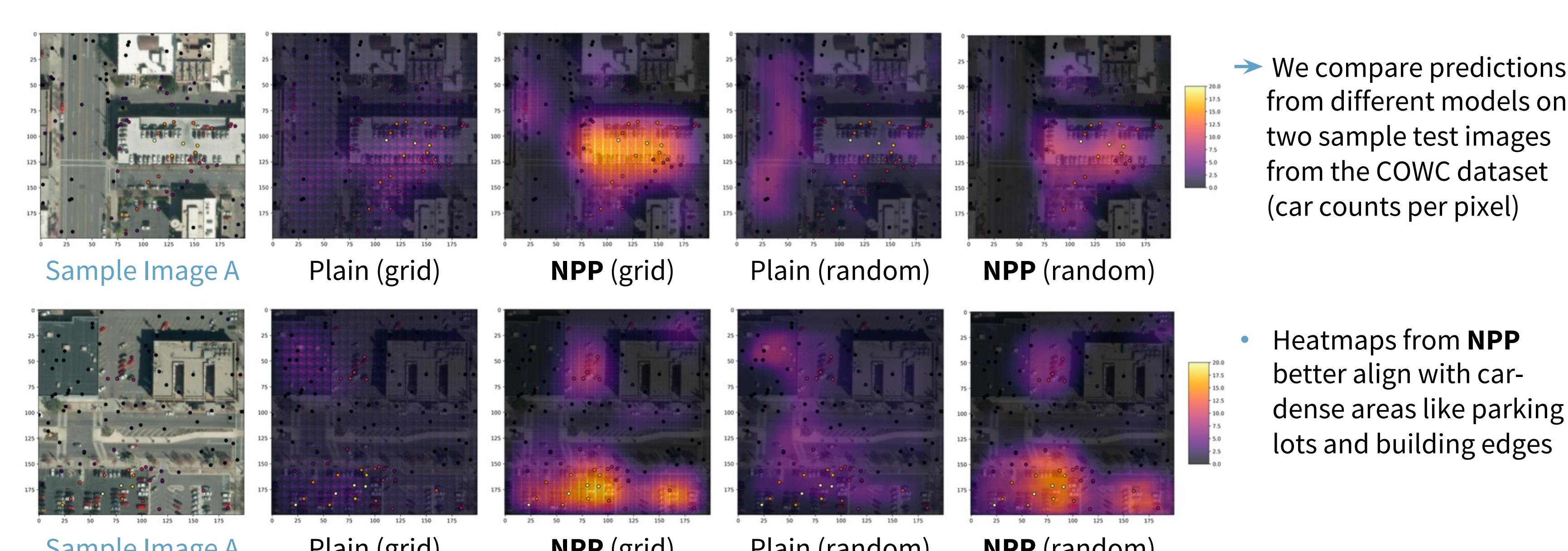
where $\mathbf{K}_{\dagger, \dagger} \equiv k(\mathbf{P}^\dagger; \mathbf{P}^\dagger) + \sigma_0^2 \mathbf{I}$

- This provides **refined estimates** and quantifies **uncertainty** — a full **Bayesian posterior**
- Check our **Partial Label Revelation Experiments** in the results

A quick note on complexity
• Our method's complexity will scale with the number of labeled points, not with the image size! ($L \ll d_1 \times d_2$)



A Visual Comparison



Experiments

Baselines Compared

- Plain:** Standard MSE-trained network
- NP (Neural Processes)** [Garnelo et al., 2018]
- ConvNP (Convolutional Neural Processes)** [Gordon et al., 2020]
- NPP (Ours):** Gaussian Process regularized training
- NPP-GP (Ours):** NPP + posterior update with partial labels

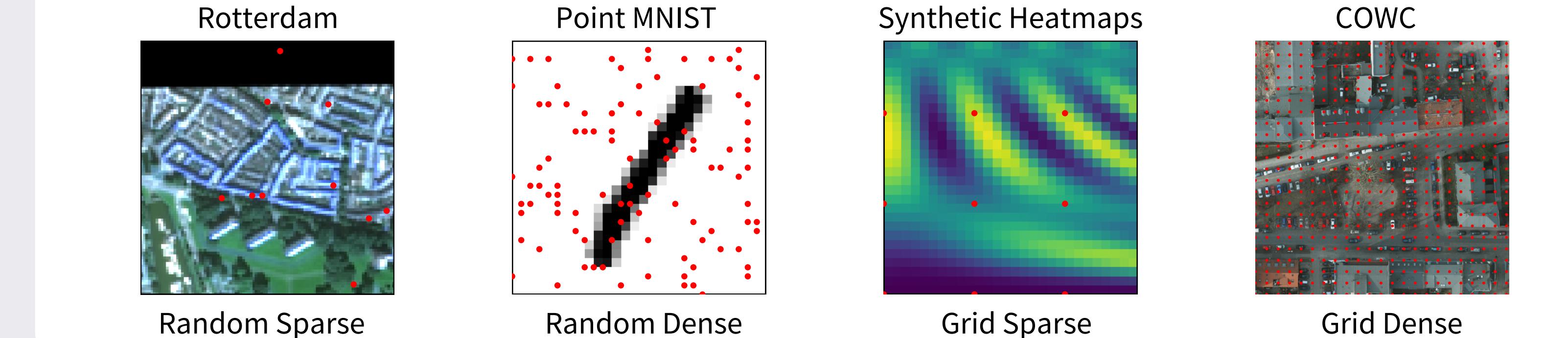
Architectures (Plain, NPP, NPP-GP methods)

- AE (Autoencoder)**
- DDPM + AE (Denoising Diffusion Probabilistic Model + AE)**

Kernel Strategies (NPP, NPP-GP methods)

- Static:** Fixed kernel (e.g., RBF), tuned via validation
- Learnable:** Kernel params optimized during training
- Context-aware:** parameters regressed from the input

Dataset	Point Distribution	# Samples	Width x Height	# Labels Sparse	# Labels Dense
Synthetic Heatmaps	grid	1000	28 x 28	9	100
	random	1000	28 x 28	10	100
Point MNIST	grid	1000	28 x 28	9	100
	random	1000	28 x 28	10	100
Rotterdam	grid	1000	100 x 100	16	121
	random	1000	100 x 100	10	100
COWC	grid	1000	200 x 200	81	529
	random	1000	200 x 200	100	500



Results

Metrics

- MSE ↓:** Measures the average squared difference between predicted and true value
- R² ↑:** Measures how well predictions explain the variance in the true data
- NPPs consistently outperform standard MSE baselines, NP, and ConvNP, especially in sparse label settings and when partial labels are available at inference time — a setup we refer to as **partial label revelation**, where some ground truth labels are revealed during inference to improve prediction accuracy across methods

	Datasets	Point pattern	Rotterdam		COWC		Partial label revelation
			Grid	Random	Grid	Random	
Real-world	Sparse	Plain	1.79	-0.058	1.14	0.297	17.77
		NPP (ours)	1.76	-0.046	0.834	0.437	4.89
	Dense	NP	1.44	0.047	1.64	-0.027	14.00
		ConvNP	0.729	-4.44	1.12	-8.10	16.6
Synthetic	Sparse	Plain	1.55	0.100	0.486	0.678	10.9
		NPP (ours)	1.25	0.286	0.453	0.700	27.1
	Dense	NP	1.25	0.287	0.443	5.60	5.02
		ConvNP	0.581	-7.27	0.344	0.379	13.2
Datasets	Point pattern	PMNIST	67.5	0.087	0.456	0.992	14.91
		Synthetic Heatmaps	56.3	0.239	0.451	0.993	27.1
	Point pattern	Grid	72.0	0.026	0.451	0.993	0.666
		Random	78.2	-0.072	44.8	0.404	104

Partial Label Revelation Experiments: We evaluate model performance when a few point labels are revealed at inference time, enabling refined predictions

