

# Transformer Neural Process - Kernel Regression

Daniel Jenson<sup>1</sup>Jhonathan Navott<sup>1</sup>Elizaveta Semenova<sup>2,\*</sup>University of Oxford<sup>1</sup>Mengyan Zhang<sup>1</sup>Makkunda Sharma<sup>1</sup>Seth Flaxman<sup>1,\*</sup>Imperial College London<sup>2</sup>

## Abstract

Stochastic processes model various natural phenomena from disease transmission to stock prices, but simulating and quantifying their uncertainty can be computationally challenging. For example, modeling a Gaussian Process with standard statistical methods incurs an  $\mathcal{O}(n^3)$  penalty, and even using state-of-the-art Neural Processes (NPs) incurs an  $\mathcal{O}(n^2)$  penalty due to the attention mechanism. We introduce the Transformer Neural Process - Kernel Regression (TNP-KR), a new architecture that incorporates a novel transformer block we call a Kernel Regression Block (KRBBlock), which reduces the computational complexity of attention in transformer-based Neural Processes (TNPs) from  $\mathcal{O}((n_C + n_T)^2)$  to  $\mathcal{O}(n_C^2 + n_C n_T)$  by eliminating masked computations, where  $n_C$  is the number of context, and  $n_T$  is the number of test points, respectively, and a fast attention variant that further reduces all attention calculations to  $\mathcal{O}(n_C)$  in space and time complexity. In benchmarks spanning such tasks as meta-regression, Bayesian optimization, and image completion, we demonstrate that the full variant matches the performance of state-of-the-art methods while training faster and scaling two orders of magnitude higher in number of test points, and the fast variant nearly matches that performance while scaling to millions of both test and context points on consumer hardware.

## 1 INTRODUCTION

The principle challenge of modern spatiotemporal Bayesian modeling is scale. As the number of observed locations increases from tens to thousands or hundreds of thousands, traditional techniques used to model spatiotemporal phenomena break down. Perhaps the most common method typically employed to model spatiotemporal processes is the Gaussian Process (GP). Gaussian Processes are a particularly well-behaved class of stochastic processes. Specifically, for a finite index set  $\{t \in T\}$ , the collection  $X = (X_1, \dots, X_T)$  follows a multivariate Gaussian distribution. This makes various analytic calculations tractable, facilitating regression, marginalization, and sampling with GPs.

While GPs provide a significant degree of flexibility in modeling, the analytic solutions they yield do not scale well in the number of observed locations. Using a GP to model spatial random effects within a Markov Chain Monte Carlo (MCMC) sampler incurs an  $\mathcal{O}(n^3)$  cost per sample, where  $n$  is the number of observed locations. This is because the covariance matrix must be inverted, or factorized in the case of Cholesky decomposition, at each iteration in order to generate a sample. Unfortunately, this means that for only  $n = 1,000$  locations, nearly a billion operations must be performed to generate a single sample.

In order to accelerate Bayesian inference with spatiotemporal stochastic processes, there have been at least three prominent strains of research. The first is Variational Inference (VI), which aims to recast the inference problem as an optimization problem and maximize the Evidence Lower Bound (ELBO). The second aims to accelerate sampling by using a generative neural network-based approximation. This family tends to leverage Variational Autoencoders (VAEs). The third is a recent family of deep learning models called Neural Processes (NPs). These models use a meta-learning objective, meaning that once trained, the forward pass of the model takes as input “context” or observed points and returns a function. This function

Preliminary work.

Correspondence to: daniel.jenson@worc.ox.ac.uk.

\*These authors jointly supervised this work.

can then be evaluated at any collection of test points and returns both their mean predictions and associated uncertainties.

The NP family of models has grown rapidly over the last few years, but recently Transformer Neural Processes (TNPs), including TNP-D, TNP-ND, and TNP-A have come to dominate the landscape (Nguyen and Grover, 2023). These models, however, suffer from an  $\mathcal{O}(n^2)$  complexity due to the attention mechanism used in transformer encoder blocks. We extend this family with the Transformer Neural Process - Kernel Regression (TNP-KR) model. TNP-KR uses a novel transformer block we call a Kernel Regression Block (KRBlock), which reduces the cost of attention in transformer-based NPs from  $\mathcal{O}((n_C + n_T)^2)$  to  $\mathcal{O}(n_C^2 + n_C n_T)$  where  $n_C$  is the number of context points and  $n_T$  is the number of test points. We also introduce an even faster variant, which uses Performer attention inside the KRBlock (Choromanski et al., 2022). Performer attention uses a kernel approximation to softmax attention and further reduces the complexity to  $\mathcal{O}(n_C)$ , enabling the model to scale to millions of points on consumer hardware. We show that the full variant matches the performance of state-of-the-art methods while training faster and scaling two orders of magnitude higher in number of test points, and the fast variant nearly matches that performance while scaling to millions of both test and context points on consumer hardware.

## 2 BACKGROUND

A number of general techniques have been developed to reduce the computational burden of modeling large spatiotemporal datasets. These include, but are not limited to variational inference (VI) (Blei, Kucukelbir, and McAuliffe, 2017), stochastic process emulation (Mishra et al., 2022; Semenova, Xu, et al., 2022; Semenova, Verma, et al., 2023), and neural processes (NPs) (Garnelo, Schwarz, et al., 2018; Garnelo, Rosenbaum, et al., 2018; Kim et al., 2019; Lee et al., 2020; Gordon et al., 2020; Nguyen and Grover, 2023). There is also a long literature on approximate methods to scale up Gaussian processes in particular which we do not cover in detail, see e.g. Hensman, Fusi, and Lawrence (2013), Rue, Martino, and Chopin (2009), Solin and Särkkä (2020), Wilson, Dann, and Nickisch (2015), and Lindgren, Lindström, and Rue (2010). While using distinct algorithmic approaches, all these methods provide approximations to the posterior distributions.

### 2.1 Variational Inference (VI)

VI (Blei, Kucukelbir, and McAuliffe, 2017; Murphy, 2023) approximates the posterior distribution by framing inference as an optimization problem, aiming to maximize the Evidence Lower Bound (ELBO) by minimizing the Kullback-Leibler (KL) divergence between a variational distribution  $q_\psi(\mathbf{z})$  and the true posterior  $p_\theta(\mathbf{z} | \mathbf{x})$ . Although VI is widely used, its effectiveness depends on selecting an appropriate variational family, and there are no guarantees on how close the ELBO is to the true log-likelihood, making uncertainty estimation challenging when the variational family poorly approximates the true distribution (Yao et al., 2018; Huggins et al., 2020).

### 2.2 Stochastic Process Emulation

Another line of research aims to accelerate sampling by approximating samples from computationally intensive stochastic processes. This is the aim of models like PriorVAE, PriorCVAE, and  $\pi$ -VAE (Semenova, Xu, et al., 2022; Semenova, Verma, et al., 2023; Mishra et al., 2022). Currently these models are all based on Variational Autoencoders (VAEs) (Kingma and Welling, 2022). VAEs consist of an encoder and decoder combined with a latent sampling process. They encode raw data into a vector of latent parameters, which are then used to sample a latent vector. This latent vector is then passed through the decoder, whose objective is to recreate the original data. The advantage of models like these is that if the latent distribution is simple, i.e. a multivariate normal with diagonal covariance, it can be very easy to sample. This means that a fully trained network can generate new samples from the original data distribution by sampling latents and passing them through the decoder. Furthermore, this can often be done in time linear in the number of layers in the network, which can be two orders of magnitude faster than sampling from a real GP (Semenova, Verma, et al., 2023). Because neural networks are inherently differentiable, they can also be transparently integrated into inference frameworks like NumPyro, where gradient-informed samplers like the No-U-Turn Sampler (NUTS) can easily pass gradients through the model. The principle challenge with this class of models is that the number of input and output locations is fixed and ordered, which means a new model must be retrained each time the number of observed locations changes or the location associated with each input changes. These models are also sensitive to the dimensionality of the latent vector, which induces an information bottleneck on autoencoded data and can cause oversmoothing in generated samples.

### 2.3 Neural Processes (NPs)

Neural Processes (NPs) are a family of models that use deep neural networks to represent stochastic processes. NPs are considered “meta-learners” because, instead of modelling a single function  $h : \mathbf{s} \rightarrow \mathbf{f}$ , they take as input context points  $(\mathbf{s}_C, \mathbf{f}_C)$  and return a distribution over functions  $g : (\mathbf{s}_C, \mathbf{f}_C) \rightarrow P_{h|\mathbf{s}_C, \mathbf{f}_C}(h : \mathbf{s}_T \rightarrow \mathbf{f}_T)$ , allowing evaluation at test points  $\mathbf{s}_T$  without retraining. To be valid stochastic processes, NPs must ensure invariance to context/test point ordering and consistency under marginalization.

There are two main classes of Neural Processes, latent neural processes and conditional neural processes. Latent neural processes dedicate part of their architecture to generating latent parameters, which are sampled and passed through the decoder to generate coherent samples, similar to VAEs. The fundamental assumption of latent NPs is that the test points are independent conditional on the latent vector. For instance, if  $(\mathbf{s}_C, \mathbf{f}_C)$  represents a tuple of locations and function values at context (observed) points,  $(\mathbf{s}_T, \mathbf{f}_T)$  represents locations and function values at test (unobserved) points, and  $\mathbf{z}$  represents a sampled latent vector, the likelihood of the function values at test locations can be formulated as follows:

$$\begin{aligned} p(\mathbf{f}_T | \mathbf{s}_T, \mathbf{s}_C, \mathbf{f}_C) &= \int p(\mathbf{z} | \mathbf{s}_C, \mathbf{f}_C) p(\mathbf{f}_T | \mathbf{s}_T, \mathbf{z}) d\mathbf{z} \\ &= \int p(\mathbf{z} | \mathbf{s}_C, \mathbf{f}_C) \prod_{i=1}^{|T|} p(\mathbf{f}_T^{(i)} | \mathbf{s}_T^{(i)}, \mathbf{z}) d\mathbf{z} \end{aligned}$$

Conditional neural processes often have very similar architectures, except they avoid sampling a latent vector and condition on a fixed representation,  $\mathbf{r}$ , of the context points. This implies the following factorization, assuming the encoder has already processed the context points,  $\text{enc}(\mathbf{s}_C, \mathbf{f}_C) = \mathbf{r}$ :

$$p(\mathbf{f}_T | \mathbf{s}_T, \mathbf{s}_C, \mathbf{f}_C) = \prod_{i=1}^{|T|} p(\mathbf{f}_T^{(i)} | \mathbf{s}_T^{(i)}, \mathbf{r})$$

In practice, conditional neural processes tend to perform better. It is unclear whether this is because the latent vectors,  $\mathbf{z}$ , are an insufficient representation of the latent parameters, or if the models are expending some of their finite capacity on producing good estimates of  $\mathbf{z}$  at the cost of final predictive accuracy. Either way, the most performant NPs as measured by log-likelihood scores are conditional, specifically conditional transformer neural processes (TNPs).

The advantage of NPs is that they can be pretrained across a variety of priors and number of context points. And at test time, they can perform GP regression for thousands of paths simultaneously in seconds on a GPU. This can be several orders of magnitude faster than calculating uncertainty for each path within an inference framework. The disadvantage of these networks is that previously they significantly underfit relative to true GPs, particularly at context or observed data points. However, NPs have evolved from simple MLP-based networks to sophisticated convolutional and transformer-based architectures. These newer variants are often just as accurate as baseline GPs and at least an order of magnitude faster.

We detail the essential aspects of modern canonical designs in this section. For brevity, we omit the original NP (Garnelo, Schwarz, et al., 2018) and CNP (Garnelo, Rosenbaum, et al., 2018) models, which used an autoencoder framework with MLPs. Also relevant, but not included here are Bootstrapping Neural Processes (Lee et al., 2020), which represent an extension to all NP models.

#### 2.3.1 Attentive Neural Process (ANP)

The Attentive Neural Process (ANP) was introduced to address several limitations of the NP and CNP architectures (Kim et al., 2019). In particular, NPs and CNPs underfit on observed context points and generally tend to oversmooth the posterior predictive. One of the principle reasons for underfitting is that the decoder is unable to differentially attend to each of the context points when decoding at test points. For instance, if one context point is located at  $\mathbf{s}_c = -2$  and the test point is at  $\mathbf{s}_t = 2$ , i.e. opposite ends of the training region,  $\mathbf{s}_c$  equally influences the context vector used to decode at  $\mathbf{s}_t$  despite having little to no influence on the behavior of function at  $\mathbf{s}_t = 2$ . Accordingly, the authors propose an attentive version of the NP that allows both local and global context to be incorporated when decoding at test locations.

The ANP shares the NP architecture, but replaces MLPs with multihead dot product self-attention in both the latent and deterministic encoding paths. Then, in the deterministic path it adds cross-attention between test points and context points so that the associated context vector,  $\mathbf{r}_t$ , summarizes information relevant to the test location  $\mathbf{s}_t$ , rather than being a global representation. This localized context vector,  $\mathbf{r}_t$ , is passed through the decoder with the global latent representation  $\mathbf{z}$ , and the location,  $\mathbf{s}_t$ , to produce  $(\mu_t, \sigma_t^2)$ . With the ANP and its conditional variant CANP, the NP family started to become a viable replacement to true GPs for inference.

### 2.3.2 Convolutional Conditional Neural Process (ConvCNP)

Convolutional Conditional Neural Processes (ConvCNP) were designed to incorporate translation equivariance into NPs (Gordon et al., 2020). When a model exhibits translation equivariance, it is able to identify a feature or function behavior regardless of how far its input has shifted in the domain. This improves the model’s capacity to generalize beyond the training region.

The authors of ConvCNP define two architectures, one for on-the-grid data and one for off-the-grid data. Here we detail the off-the-grid version since it is more generic and can be used in both cases. First, the domain is partitioned into a fixed, uniform grid. Then, a positive-definite Reproducing Kernel Hilbert Space (RKHS) kernel is evaluated at each of the grid points using the context set. This value is then normalized using a density channel so the point values are invariant to the cardinality of the context set. This grid is then run through a CNN-based architecture, *e.g.* ResNet (He et al., 2016), to create an updated hidden state representation at grid locations. Finally, the decoder uses another RKHS kernel, typically the same one used in the encoder, to decode test points using the underlying hidden state grid values.

ConvCNPs perform very well on the standard NP benchmarks, often exceeding the performance of other models at a fraction of the number of learnable parameters. However, beyond simple examples, ConvCNP requires many more parameters, often on the same order as other models, to perform competitively. Furthermore, ConvCNP is very sensitive to the parameterization of the intermediate grid and the effective receptive field of the CNN layer. For instance, at lower lengthscales, the model performs better when there is a higher grid resolution, but this increase in grid resolution changes the receptive field of the CNN layer, so the CNN’s kernels must be optimized in conjunction. Lastly, due to the fixed intermediate grid, ConvCNP suffers from the curse of dimensionality and is difficult to evaluate in non-contiguous regions of the domain.

### 2.3.3 Transformer Neural Processes (TNP)

Transformer Neural Processes (TNPs) (Nguyen and Grover, 2023) can be considered an extension of the Conditional Attentive Neural Process (CANP) that uses multiple transformer encoder blocks instead of stacked self-attention and cross-attention. A transformer encoder block consists of a self-attention layer followed by a feedforward network with residual connections interspersed. The pseudocode can be viewed in Appendix 5.

Because standard transformers are designed to work with sequences of data and use fixed positional embeddings, the TNP authors had to modify the architecture for NPs. Accordingly, TNPs dispense with the positional embeddings, merge the context and test sequences as input, and introduce a special attention mask. The context sequence consists of location and function value pairs,  $\{(\mathbf{s}_1, \mathbf{f}_1), \dots, (\mathbf{s}_{n_C}, \mathbf{f}_{n_C})\}$ , and the test sequence consists of location and zero-padded function value pairs,  $\{(\mathbf{s}_1, \mathbf{0}), \dots, (\mathbf{s}_{n_T}, \mathbf{0})\}$ . Within each layer, a mask is applied that prevents context points from attending to test points and prevents test points from attending to other test points. This means that context points only attend to other context points and test points also only attend to context points. After the encoding stage, the embeddings for test points are passed through a prediction head that estimates the mean and covariance structures.

There are three original TNP variants: TNP-D, TNP-ND, and TNP-A. TNP-D (Diagonal) assumes test points can be factorized independently conditional on the context points, in line with most CNP variants. TNP-ND (Non-Diagonal) parameterizes a covariance matrix by estimating the elements in the lower triangular matrix of a Cholesky decomposition for test points. Lastly, TNP-A (Autoregressive) assumes that the test points can be factorized autoregressively. This means that each time a test point is predicted, its true value is then added to the context set and used when predicting the next test point.

In practice, we found TNP-D to consistently perform well. On the other hand, we found TNP-ND could become arbitrarily miscalibrated. When there are very few context points, TNP-ND maximizes the log-likelihood by collapsing the standard deviation for test points that lie near context points, which provides very precise estimates. However, this also causes the model to collapse uncertainty estimates for test points that lie far away from context points, leading to both high log-likelihoods and high miscalibration rates. On the other hand, TNP-A assumes that the true test points are observed autoregressively, which is equivalent to running any NP model forward one test point at a time. This method yields log-likelihoods that correspond to sequential observations, rather than joint distributions over all unobserved points. As we are interested in the general case, we do not include experiments on TNP-A.

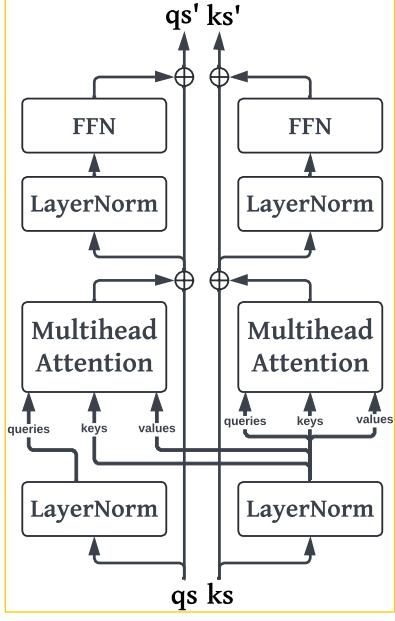


Figure 1: KRBlock architecture. Queries and keys have their own residual pathways, but keys are used to update queries through cross-attention.

### 3 TRANSFORMER NEURAL PROCESS - KERNEL REGRESSION (TNP-KR)

#### 3.1 KRBlock

TNP-KR was inspired by TNP-D, which consistently performs well on NP benchmarks and does not make the same assumptions as TNP-A or suffer from miscalibration like TNP-ND. One of the principal limitations of all TNP variants, however, is the  $O((n_C + n_T)^2)$  attention used by the encoder layers. Recall that  $n_C$  is the number of context points and  $n_T$  is the number of test points. This also is not strictly necessary because the special mask applied after attention discards many of the calculations, namely, those from context points to test points and those from test points to other test points. Accordingly, we introduce the KRBlock, a transformer block that avoids computing  $O(n_T^2 + n_C n_T)$  attention values altogether. Unlike the original transformer encoder block, TNP-KR also uses pre-normalized residual connections (Xiong et al., 2020), which has been shown to stabilize training and improve performance.

We call this a KRBlock because the cross-attention from test to context points can be viewed as a form of Nadaraya-Watson kernel regression where the locations are the feature embeddings and the kernel is a dot product softmax kernel. Specifically, if  $\alpha$  is the dot product softmax kernel,  $\mathbf{q}_i$  is a query embedding

representing a test point,  $\mathbf{k}_j$  is a key embedding representing a context point, and  $\mathbf{v}_j$  is the value embedding associated with the same context point, cross-attention in the KRBlock can be formulated explicitly as Nadaraya-Watson kernel regression:

$$f(\mathbf{q}_i) = \sum_j \frac{\alpha(\mathbf{q}_i, \mathbf{k}_j)}{\sum_j \alpha(\mathbf{q}_i, \mathbf{k}_j)} \mathbf{v}_j \quad (1)$$

Stacking KRBlocks then allows the model to perform iterative kernel regression on increasingly complex internal representations of test and context points. Cross-attention from test to context points costs  $\mathcal{O}(n_C n_T)$  and the self-attention among context points costs  $\mathcal{O}(n_C^2)$ , making the total complexity  $\mathcal{O}(n_C^2 + n_C n_T)$ .

When there are a large number of context points, the  $\mathcal{O}(n_C^2)$  term can still prove computationally prohibitive. For example, in satellite imagery, a comparatively small image of 300x500 results in 150,000 locations or pixels (Heaton et al., 2018). A common application is to inpaint pixels missing due to cloud cover. With 20% of pixels missing, there would still be 120000 context points. This means that even with a KRBlock, the space and time complexity would be on the order of 120000 $^2$ , requiring nearly 60GB of memory and 14.4 billion multiplications per self-attention application. Thus, in order to further scale KRBlocks, we incorporate Performer attention, also known as fast attention (Choromanski et al., 2022).

#### 3.2 Fast Attention

Fast attention is based on a complete algorithm called Fast Attention with Orthogonal Random features (FAVOR+). FAVOR+ allows attention to be calculated in linear space and time complexity without making any assumptions about the sparsity or rank of the attention matrix. It is nearly unbiased and offers uniform convergence and low variance (Choromanski et al., 2022). FAVOR+ constructs an attention matrix  $\mathbf{A}^{L \times L}$  where  $\mathbf{A}(i, j) = K(\mathbf{q}_i, \mathbf{k}_j)$  without ever fully materializing it. For a randomized mapping  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}_+^r$ , entries can be estimated with  $K(\mathbf{q}, \mathbf{k}) = \mathbb{E}[\phi(\mathbf{q})^\top \phi(\mathbf{k})]$ . In fact, most kernels can be modeled using the following mapping:

$$\phi(\mathbf{x}) = \frac{\mathbf{h}(\mathbf{x})}{\sqrt{m}} \left[ f_1(\omega_1^\top \mathbf{x}), \dots, f_l(\omega_l^\top \mathbf{x}), \dots, f_t(\omega_t^\top \mathbf{x}), \dots, f_m(\omega_m^\top \mathbf{x}) \right]$$

where  $f_1, \dots, f_l : \mathbb{R} \rightarrow \mathbb{R}$ ,  $\mathbf{h} : \mathbb{R}^d \rightarrow \mathbb{R}$ , and  $\omega_1, \dots, \omega_m \stackrel{\text{iid}}{\sim} \mathcal{D}$ . When  $\mathcal{D} = \mathcal{N}(0, \mathbf{I}_d)$ , this leads to

the Gaussian kernel,  $K_{\text{gauss}}$ . Ignoring the  $\sqrt{d}$  normalization, the softmax kernel is defined as  $\text{SM}(\mathbf{q}, \mathbf{k}) := \exp(\mathbf{q}^\top \mathbf{k})$ . Furthermore, the softmax kernel can be expressed as a function of  $|K_{\text{gauss}}|$ :

$$\text{SM}(\mathbf{q}, \mathbf{k}) = \exp\left(\frac{\|\mathbf{q}\|^2}{2}\right) K_{\text{gauss}}(\mathbf{q}, \mathbf{k}) \exp\left(\frac{\|\mathbf{k}\|^2}{2}\right)$$

Thus, the softmax kernel can be defined using  $\phi$  where  $h(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}\|^2}{2}\right)$ ,  $f_1 = \sin$ , and  $f_2 = \cos$ . Unfortunately, this representation behaves poorly around zero when the trigonometric functions return negative values, so Choromanski et al. introduce positive random features (PRFs) for softmax. This results in the form:

$$\text{SM}(\mathbf{q}, \mathbf{k}) = \mathbb{E}_{\omega \sim \mathcal{N}(0, \mathbf{I}_d)} \left[ e^{\omega^\top \mathbf{q} - \frac{\|\mathbf{q}\|^2}{2}} e^{\omega^\top \mathbf{k} - \frac{\|\mathbf{k}\|^2}{2}} \right]$$

which implies  $h(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}\|^2}{2}\right)$  and  $f_1 = \exp$ . This approximation is further improved by forcing the  $\omega$  random vectors to be normalized and orthogonal (OR $\oplus$ ).

Mapping the query and key matrices through these random projections,  $\phi(\mathbf{Q}) = \mathbf{Q}' \in \mathbb{R}^{L \times r}$  and  $\phi(\mathbf{K}) = \mathbf{K}' \in \mathbb{R}^{L \times r}$ , attention can be reexpressed as follows:

$$\begin{aligned} \widehat{\text{Attention}}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) &= \hat{\mathbf{D}}^{-1} \underbrace{(\mathbf{Q}'((\mathbf{K}')^\top \mathbf{V}))}_{\mathcal{O}(Lrd)}, \\ \hat{\mathbf{D}} &= \text{diag}(\mathbf{Q}'((\mathbf{K}')^\top \mathbf{1}_L)). \end{aligned}$$

Thus, by carefully controlling the order of matrix multiplications, space complexity can be reduced from  $\mathcal{O}(L^2 + Ld)$  to  $\mathcal{O}(Lr + Ld + rd)$  and time complexity can be reduced from  $\mathcal{O}(L^2 d)$  to  $\mathcal{O}(Lrd)$ . We demonstrate the impact of the KRBlock and fast attention on runtime and memory in Tables 1 and 2.

## 4 EXPERIMENTS

In order to evaluate the performance of TNP-KR, we test its performance across standard benchmarks, which include GPs, image completion, and simple Bayesian Optimization (BO). All experiments were performed on a single 24GB Nvidia RTX 4090 GPU.

### 4.1 1D Gaussian Processes

For one-dimensional GPs, we evaluate the principal NP models on the RBF, periodic, and Matérn  $3/2$  kernels. For each of these kernels,  $\sigma^2$  can be factored out and used to standardize the data. Accordingly, in our tests, we assume the data has been standardized so the models can focus on learning and differentiating lengthscales. While standard benchmarks sample

lengthscales from 0.1 to 0.6 uniformly to model over domain on the domain  $[-2, 2]$ , we sample lengthscales according to  $\ell \sim \text{Beta}(\alpha = 3, \beta = 7)$ , which has both a mean and median of approximately 0.3. This is a more challenging benchmark that allows for greater differentiation among models since more than 50% of lengthscales fall below 0.3 and less than 10% lie above 0.5. (In practice, we found most models could easily learn lengthscales above 0.5). For the periodic kernel, we also sample the period uniformly from 0.5 to 2, which represents between 2 and 8 cycles on the domain  $[-2, 2]$ .

Each model is seeded 5 times and trained on 100,000 batches of size 32. For each seed, the models are evaluated on a final test set consisting of 5,000 batches of size 32. Each sample in each batch consists of 50 randomly sampled locations and 100 linearly spaced points throughout the domain. Between 3 and 50 of the randomly sampled points are used as context points, and the test points consist of all 150 points. Because many of these models use different heuristics to prevent the collapse of standard deviation at or near context points in the test set, we add an observation noise of 0.1. This added noise prevents collapse and allows us to avoid using model-dependent heuristics, which could artificially hamper their performance and prevent honest comparison. Lastly, all models are trained with a single cosine annealing learning rate schedule, gradient norm clipping with a max norm of 3.0, and the YOGI optimizer (Zaheer et al., 2018). Table 3 shows the mean and standard errors across these test sets for each model and each kernel.

Table 1: Microseconds ( $\mu s$ ) per sample by varying number of test points and keeping context points fixed at 100. “OOM” indicates Out of Memory error.

# Test	TNP-D	TNP-KR-Full	TNP-KR-Fast
$10^0$	118	171	194
$10^1$	119	287	327
$10^2$	129	267	282
$10^3$	286	230	226
$10^4$	45,000	474	436
$10^5$	OOM	7,040	4,510
$10^6$	OOM	86,800	64,300

These models can then be used in a simple one-dimensional Bayesian Optimization (BO) setting where each model is given ten observed starting points. BO then uses the model to calculate the expected improvement across the domain and selects the point where that metric is greatest. There is a fixed budget of 50 iterations and the objective is to identify the function’s minimum. Table 4 shows the mean and standard error of regret across 5 seeds for each model

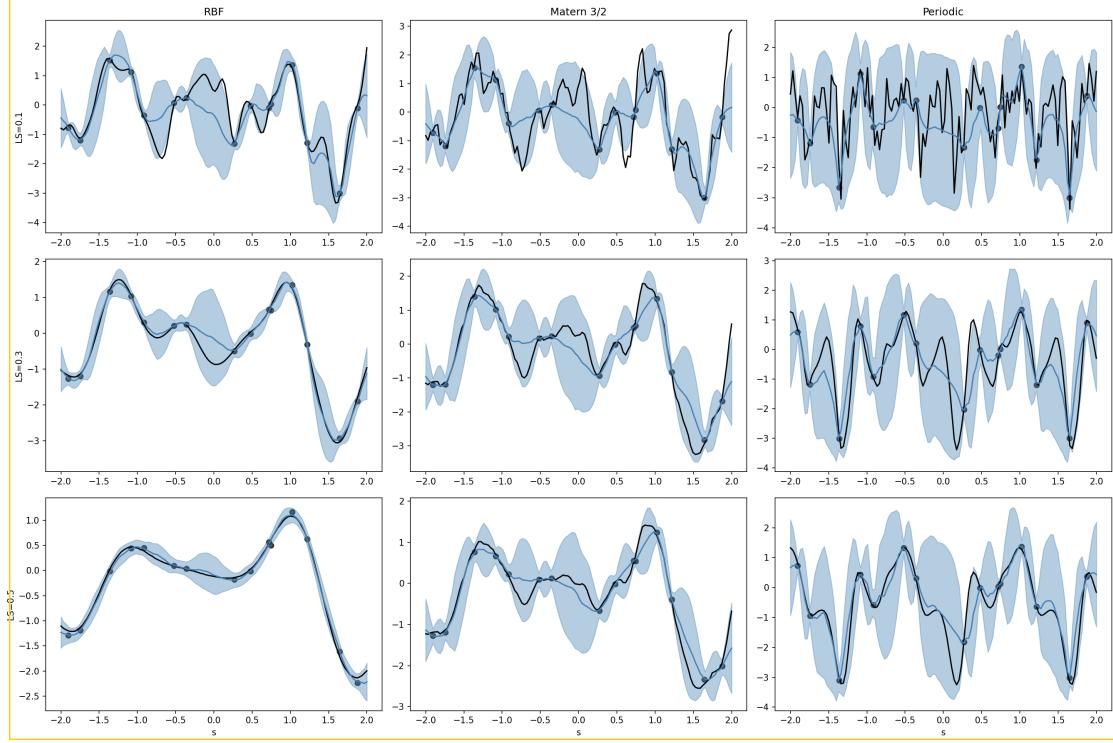


Figure 2: Posterior predictive samples for TNP-KR-Fast. Black line is the true function, black dots are noisy observations, blue line is model mean prediction, and shaded blue area is 95% confidence region.

Table 2: Microseconds ( $\mu\text{s}$ ) per sample by varying number of context points and keeping test points fixed at 100. “OOM” indicates Out of Memory error.

# Context	TNP-D	TNP-KR-Full	TNP-KR-Fast
$10^0$	99	98	138
$10^1$	99	220	277
$10^2$	114	286	299
$10^3$	247	302	234
$10^4$	44,700	36,700	632
$10^5$	OOM	OOM	6,900
$10^6$	OOM	OOM	88,400

Table 3: 1D GP mean and standard error of negative log-likelihoods over 5 runs, where the score for each run consists of the average negative log-likelihood over 5,000 batches of 32. Best model in bold.

Model	Matern 3/2	Periodic	RBF
NP	$0.638 \pm 0.003$	$1.342 \pm 0.001$	$0.338 \pm 0.007$
CNP	$0.547 \pm 0.002$	$1.165 \pm 0.002$	$0.238 \pm 0.006$
BNP	$0.523 \pm 0.003$	$1.165 \pm 0.003$	$0.165 \pm 0.005$
ANP	$0.043 \pm 0.008$	$0.829 \pm 0.007$	$-0.341 \pm 0.003$
CANP	$0.047 \pm 0.014$	$0.818 \pm 0.005$	$-0.310 \pm 0.024$
BANP	$0.023 \pm 0.007$	$0.761 \pm 0.011$	$-0.347 \pm 0.008$
ConvCNP	$-0.018 \pm 0.003$	$0.543 \pm 0.003$	$-0.456 \pm 0.003$
TNP-D	$-0.023 \pm 0.003$	$0.516 \pm 0.003$	$-0.455 \pm 0.004$
TNP-KR-Fast	$-0.006 \pm 0.002$	$0.562 \pm 0.006$	$-0.425 \pm 0.005$
TNP-KR-Full	$-0.024 \pm 0.002$	$0.516 \pm 0.002$	$-0.453 \pm 0.004$

and each kernel. There are 100 tasks per seed. Here, ConvCNP performs best across the board, which is likely due to the fact that the convolutional kernels in this model introduce a strong local inductive bias, so it is able to identify functional patterns that fall within its kernel width with very few context points.

#### 4.2 2D Gaussian Processes

For the two-dimensional Gaussian Processes, we test the RBF kernel on the domain  $[-2, 2]^2$ . The test points consist of 128 randomly selected points on the domain as well as 256 on a uniform grid over the domain. The context points consist of between 12 and

128 of the randomly selected points, which represents between 5% and 50% of the number on the uniform grid. Each model is trained on 100,000 batches of size 16 and tested on 5,000 batches of size 16. The other training settings remain the same as the one-dimensional case. Here TNP-D and TNP-KR with full attention perform best. Although ConvCNP was omitted due to the time required to optimize its kernel and grid for the benchmark, we believe it would also perform competitively here.

Table 4: 1D Bayesian Optimization mean and standard error of 500 samples on domain [-2, 2]. Best model in bold.

Model	Matern 3/2	Periodic	RBF
NP	0.145 ± 0.013	0.131 ± 0.010	0.079 ± 0.010
CNP	0.106 ± 0.011	0.331 ± 0.020	0.061 ± 0.009
BNP	0.126 ± 0.012	0.293 ± 0.019	0.069 ± 0.009
ANP	0.072 ± 0.010	0.123 ± 0.012	0.041 ± 0.009
CANP	0.065 ± 0.010	0.154 ± 0.014	0.042 ± 0.008
BANP	0.059 ± 0.009	0.106 ± 0.011	0.037 ± 0.008
ConvCNP	<b>0.011 ± 0.003</b>	<b>0.043 ± 0.006</b>	<b>0.003 ± 0.001</b>
TNP-D	0.059 ± 0.009	<b>0.044 ± 0.007</b>	0.023 ± 0.007
TNP-KR-Fast	0.062 ± 0.009	0.052 ± 0.007	0.027 ± 0.006
TNP-KR-Full	0.054 ± 0.009	<b>0.039 ± 0.005</b>	0.025 ± 0.006

Table 5: 2D GP mean and standard error of negative log-likelihoods over 5 runs, where the score for each run consists of the average negative log-likelihood over 5,000 batches of 16. Best model in bold.

Model	RBF
NP	1.209 ± 0.002
CNP	1.174 ± 0.002
ANP	0.753 ± 0.012
CANP	0.611 ± 0.011
TNP-D	<b>0.484 ± 0.003</b>
TNP-KR-Fast	0.554 ± 0.002
TNP-KR-Full	<b>0.482 ± 0.005</b>

Table 7: CelebA mean and standard error of negative log-likelihoods over 5 runs, where the score for each run consists of the average negative log-likelihood over 5,000 batches of 32. Best model in bold.

Model	NLL
NP	-0.107 ± 0.002
CNP	-0.127 ± 0.001
ANP	-0.983 ± 0.023
CANP	-0.554 ± 0.096
TNP-D	-1.524 ± 0.005
TNPKR Fast	-1.006 ± 0.025
TNPKR Full	<b>-1.561 ± 0.008</b>

Table 8: Cifar-10 mean and standard error of negative log-likelihoods over 5 runs, where the score for each run consists of the average negative log-likelihood over 5,000 batches of 32. Best model in bold.

Model	NLL
NP	0.054 ± 0.004
CNP	0.009 ± 0.001
ANP	-0.696 ± 0.092
CANP	-0.663 ± 0.096
TNP-D	<b>-1.416 ± 0.004</b>
TNP-KR-Fast	-0.823 ± 0.012
TNP-KR-Full	<b>-1.428 ± 0.007</b>

## 5 CONCLUSION

In this work, we introduce TNP-KR, an extension to TNP-D that incorporates KRBlocks and fast attention. We believe this represents an important step in scaling transformer-based models for large spatiotemporal applications, which include population genetics, epidemiology, and meteorology, among others. In the future, we plan to explore alternative kernels in both full and fast attention, as well as attention bias that explicitly takes into account pairwise distances. We also plan to explore mechanisms for sampling transformer-based models that respect covariance structures in such a fashion that we could embed TNP-KR and its variants in inference frameworks such as NumPyro and Stan.

## Acknowledgements

D.J. acknowledges his Google DeepMind scholarship. E.S. acknowledges support in part by the AI2050 program at Schmidt Sciences (Grant [G-22-64476]). S.F. acknowledges the EPSRC (EP/V002910/2).

Table 6: MNIST mean and standard error of negative log-likelihoods over 5 runs, where the score for each run consists of the average negative log-likelihood over 5,000 batches of 32. Best model in bold.

Model	NLL
NP	-0.716 ± 0.007
CNP	-0.776 ± 0.006
ANP	-0.865 ± 0.021
CANP	-0.855 ± 0.010
TNP-D	<b>-1.189 ± 0.007</b>
TNP-KR-Fast	-1.049 ± 0.006
TNP-KR-Full	<b>-1.172 ± 0.006</b>

## References

- Blei, David M., Alp Kucukelbir, and Jon D. McAuliffe (Apr. 3, 2017). “Variational Inference: A Review for Statisticians”. In: *Journal of the American Statistical Association* 112.518, pp. 859–877. ISSN: 0162-1459, 1537-274X. DOI: 10.1080/01621459.2017.1285773. arXiv: 1601.00670[cs, stat]. URL: <http://arxiv.org/abs/1601.00670> (visited on 07/30/2024).
- Choromanski, Krzysztof et al. (Nov. 19, 2022). “Rethinking Attention with Performers”. In: arXiv:2009.14794. DOI: 10.48550/arXiv.2009.14794. arXiv: 2009.14794[cs, stat]. URL: <http://arxiv.org/abs/2009.14794> (visited on 05/01/2024).
- Garnelo, Marta, Dan Rosenbaum, et al. (July 4, 2018). “Conditional Neural Processes”. In: arXiv:1807.01613. DOI: 10.48550/arXiv.1807.01613. arXiv: 1807.01613[cs, stat]. URL: <http://arxiv.org/abs/1807.01613> (visited on 07/03/2024).
- Garnelo, Marta, Jonathan Schwarz, et al. (2018). *Neural Processes*. arXiv: 1807.01622 [cs.LG]. URL: <https://arxiv.org/abs/1807.01622>.
- Gordon, Jonathan et al. (June 25, 2020). *Convolutional Conditional Neural Processes*. arXiv: 1910.13556[cs, stat]. URL: <http://arxiv.org/abs/1910.13556> (visited on 04/16/2024).
- He, Kaiming et al. (June 2016). “Deep Residual Learning for Image Recognition”. In: pp. 770–778. DOI: 10.1109/CVPR.2016.90.
- Heaton, Matthew J. et al. (Apr. 25, 2018). *A Case Study Competition Among Methods for Analyzing Large Spatial Data*. DOI: 10.48550/arXiv.1710.05013. arXiv: 1710.05013[stat]. URL: <http://arxiv.org/abs/1710.05013> (visited on 12/01/2023).
- Hensman, James, Nicolo Fusi, and Neil D Lawrence (2013). “Gaussian processes for big data”. In: *arXiv preprint arXiv:1309.6835*.
- Huggins, Jonathan et al. (2020). “Validated variational inference via practical posterior error bounds”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR, pp. 1792–1802.
- Kim, Hyunjik et al. (July 9, 2019). “Attentive Neural Processes”. In: arXiv:1901.05761. DOI: 10.48550/arXiv.1901.05761. arXiv: 1901.05761[cs, stat]. URL: <http://arxiv.org/abs/1901.05761> (visited on 04/02/2024).
- Kingma, Diederik P and Max Welling (2022). *Auto-Encoding Variational Bayes*. arXiv: 1312.6114[stat.ML]. URL: <https://arxiv.org/abs/1312.6114>.
- Lee, Juho et al. (Oct. 27, 2020). *Bootstrapping Neural Processes*. DOI: 10.48550/arXiv.2008.02956. arXiv: 2008.02956[cs, stat]. URL: <http://arxiv.org/abs/2008.02956> (visited on 05/19/2024).
- Lindgren, Finn, Johan Lindström, and Håvard Rue (2010). “An explicit link between gaussian fields and gaussian markov random fields; the spde approach”. In.
- Mishra, Swapnil et al. (Sept. 13, 2022).  $\pi$ VAE: a stochastic process prior for Bayesian deep learning with MCMC. DOI: 10.48550/arXiv.2002.06873. arXiv: 2002.06873[cs, stat]. URL: <http://arxiv.org/abs/2002.06873> (visited on 07/30/2024).
- Murphy, Kevin P. (2023). *Probabilistic Machine Learning: Advanced Topics*. MIT Press. URL: <http://probml.github.io/book2>.
- Nguyen, Tung and Aditya Grover (Feb. 7, 2023). *Transformer Neural Processes: Uncertainty-Aware Meta Learning Via Sequence Modeling*. DOI: 10.48550/arXiv.2207.04179. arXiv: 2207.04179[cs]. URL: <http://arxiv.org/abs/2207.04179> (visited on 04/13/2024).
- Rue, Håvard, Sara Martino, and Nicolas Chopin (2009). “Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations”. In: *Journal of the Royal Statistical Society Series B: Statistical Methodology* 71.2, pp. 319–392.
- Semenova, Elizaveta, Prakhar Verma, et al. (2023). “PriorCVAE: scalable MCMC parameter inference with Bayesian deep generative modelling”. In: *arXiv preprint arXiv:2304.04307*.
- Semenova, Elizaveta, Yidan Xu, et al. (June 2022). “PriorVAE: Encoding spatial priors with VAEs for small-area estimation”. In: *Journal of The Royal Society Interface* 19.191, p. 20220094. ISSN: 1742-5662. DOI: 10.1098/rsif.2022.0094. arXiv: 2110.10422[cs, stat]. URL: <http://arxiv.org/abs/2110.10422> (visited on 07/30/2024).
- Solin, Arno and Simo Särkkä (2020). “Hilbert space methods for reduced-rank Gaussian process regression”. In: *Statistics and Computing* 30.2, pp. 419–446.
- Wilson, Andrew Gordon, Christoph Dann, and Hannes Nickisch (2015). “Thoughts on massively scalable Gaussian processes”. In: *arXiv preprint arXiv:1511.01870*.
- Xiong, Ruibin et al. (2020). *On Layer Normalization in the Transformer Architecture*. arXiv: 2002.04745[cs.LG]. URL: <https://arxiv.org/abs/2002.04745>.
- Yao, Yuling et al. (2018). “Yes, but did it work?: Evaluating variational inference”. In: *International*

*Conference on Machine Learning*. PMLR, pp. 5581–5590.

Zaheer, Manzil et al. (2018). “Adaptive Methods for Nonconvex Optimization”. In: *Advances in Neural Information Processing Systems*. Ed. by S. Bengio et al. Vol. 31. Curran Associates, Inc. URL: [https://proceedings.neurips.cc/paper\\_files/paper/2018/file/90365351ccc7437a1309dc64e4db32a3-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/90365351ccc7437a1309dc64e4db32a3-Paper.pdf).

## APPENDIX

### ALGORITHMS

---

**Algorithm 1** Transformer Encoder Block forward pass.

---

```
1: Input:  $x_0$ , mask,  $p_{\text{dropout}}$ 
2:  $x_1 \leftarrow \text{MultiheadAttention}(x, x, x, \text{mask})$ 
3:  $x_2 \leftarrow x_0 + \text{Dropout}(x_1, p_{\text{dropout}})$ 
4:  $x_3 \leftarrow \text{LayerNorm}(x_2)$ 
5:  $x_4 \leftarrow \text{FeedForward}(x_3)$ 
6:  $x_5 \leftarrow x_3 + \text{Dropout}(x_4, p_{\text{dropout}})$ 
7: return  $x_5$ 
```

---

**Algorithm 2** KRBlock forward pass.

---

```
Input:  $qs_0$ ,  $ks_0$ , mask,  $p_{\text{dropout}}$ 
 $qs_1 \leftarrow \text{LayerNorm}_1(qs_0)$ 
 $ks_1 \leftarrow \text{LayerNorm}_1(ks_0)$ 
 $qs_2 \leftarrow \text{MultiheadAttention}(qs_1, ks_1, qs_1, \text{mask})$ 
 $ks_2 \leftarrow \text{MultiheadAttention}(ks_1, ks_1, ks_1, \text{mask})$ 
 $qs_3 \leftarrow qs_0 + \text{Dropout}(qs_2, p_{\text{dropout}})$ 
 $ks_3 \leftarrow ks_0 + \text{Dropout}(ks_2, p_{\text{dropout}})$ 
 $qs_4 \leftarrow \text{LayerNorm}_2(qs_3)$ 
 $ks_4 \leftarrow \text{LayerNorm}_2(ks_3)$ 
 $qs_5 \leftarrow \text{FeedForward}(qs_4)$ 
 $ks_5 \leftarrow \text{FeedForward}(ks_4)$ 
 $qs_6 \leftarrow qs_3 + \text{Dropout}(qs_5, p_{\text{dropout}})$ 
 $ks_6 \leftarrow ks_3 + \text{Dropout}(ks_5, p_{\text{dropout}})$ 
return  $qs_6, ks_6$ 
```

---

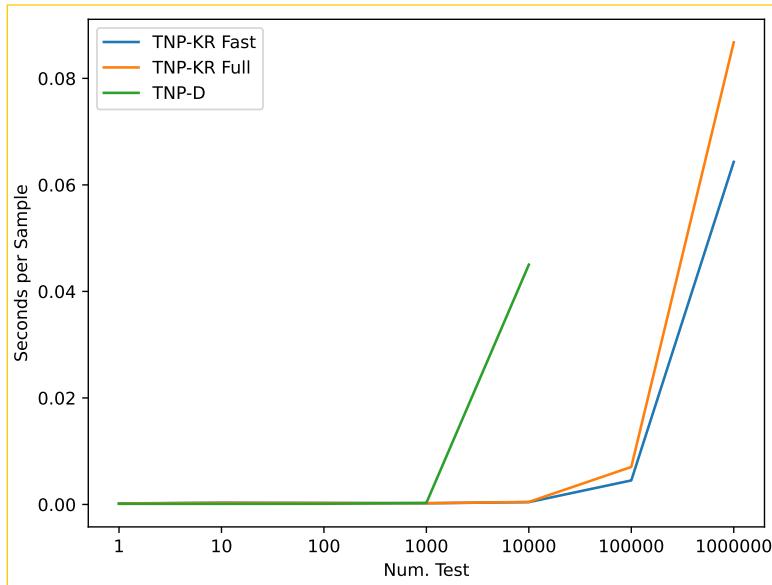
**ADDITIONAL FIGURES**

Figure 3: Scaling test points in TNP models, keeping number of context points fixed at 100.

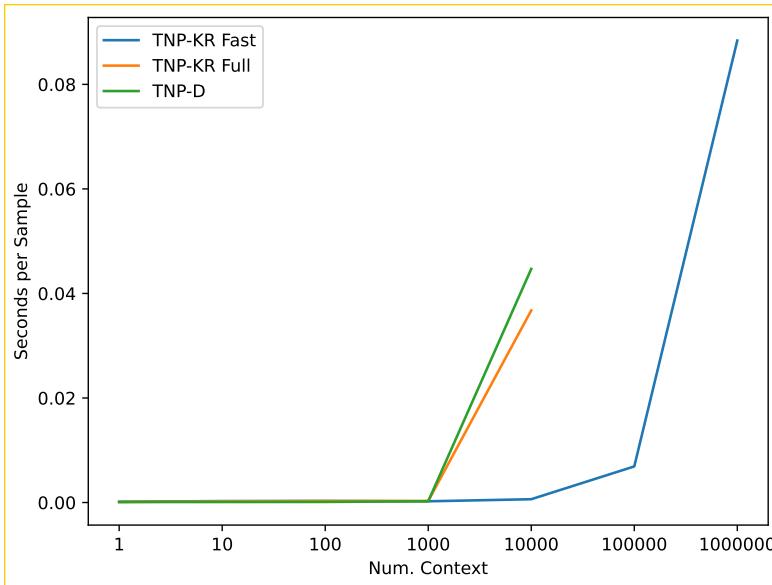


Figure 4: Scaling context points in TNP, keeping number of test points fixed at 100.

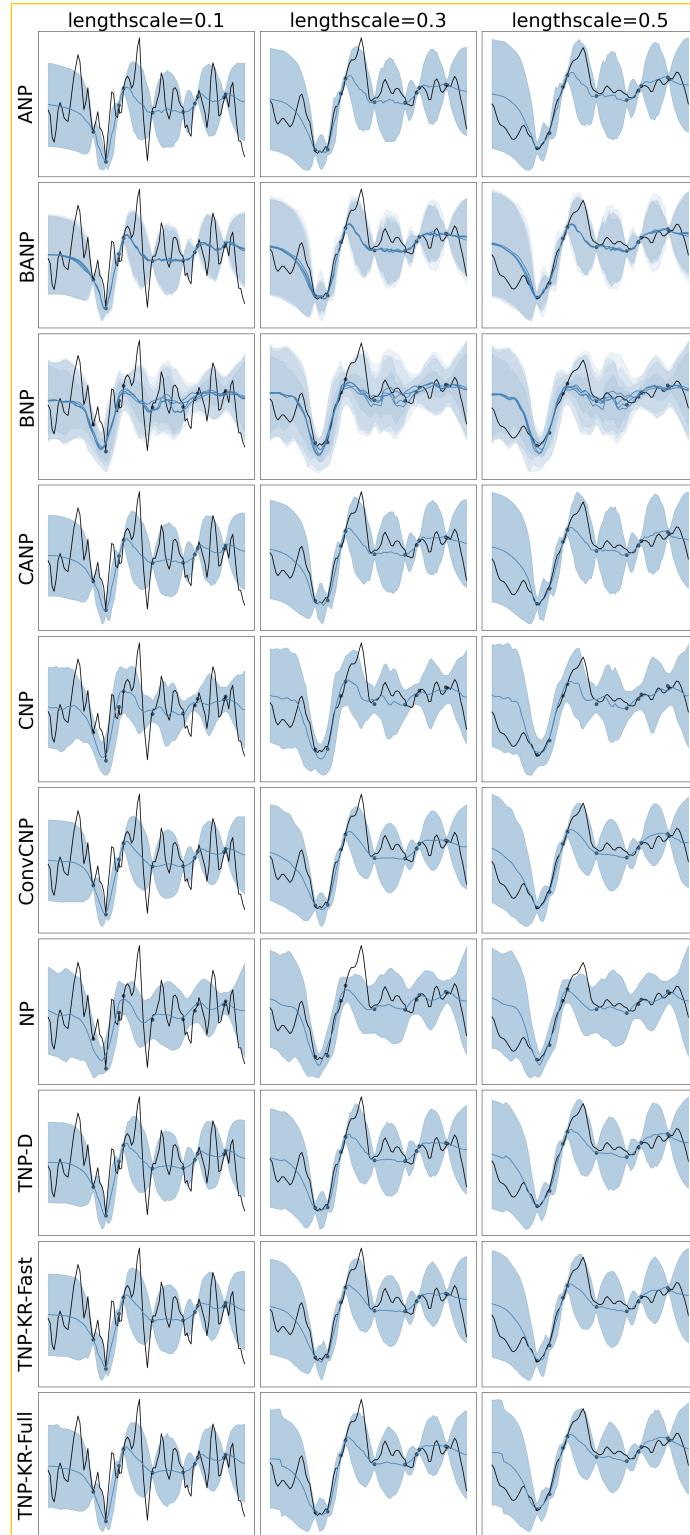


Figure 5: 1D GP Matern 3/2 samples. The black line is the true function, the black dots are noisy observations, the blue line is the model's mean prediction, and the shaded blue is the uncertainty. Models perform similarly, but ConvCNP and the TNP variants have the tightest uncertainty bounds.

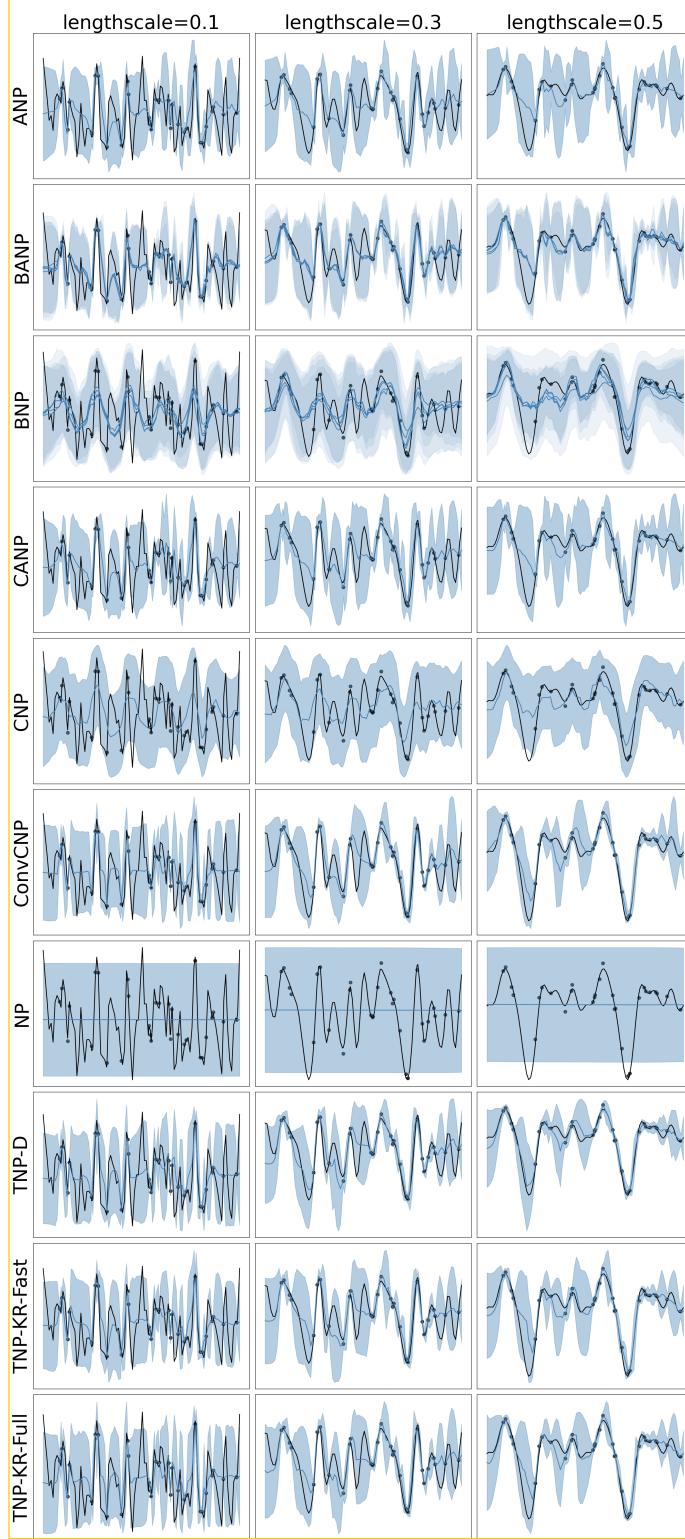


Figure 6: 1D GP periodic samples. The black line is the true function, the black dots are noisy observations, the blue line is the model's mean prediction, and the shaded blue is the uncertainty. Some models struggle to model low period periodic functions, particularly NPs. ConvCNP and TNP variants have the tightest uncertainty bounds.

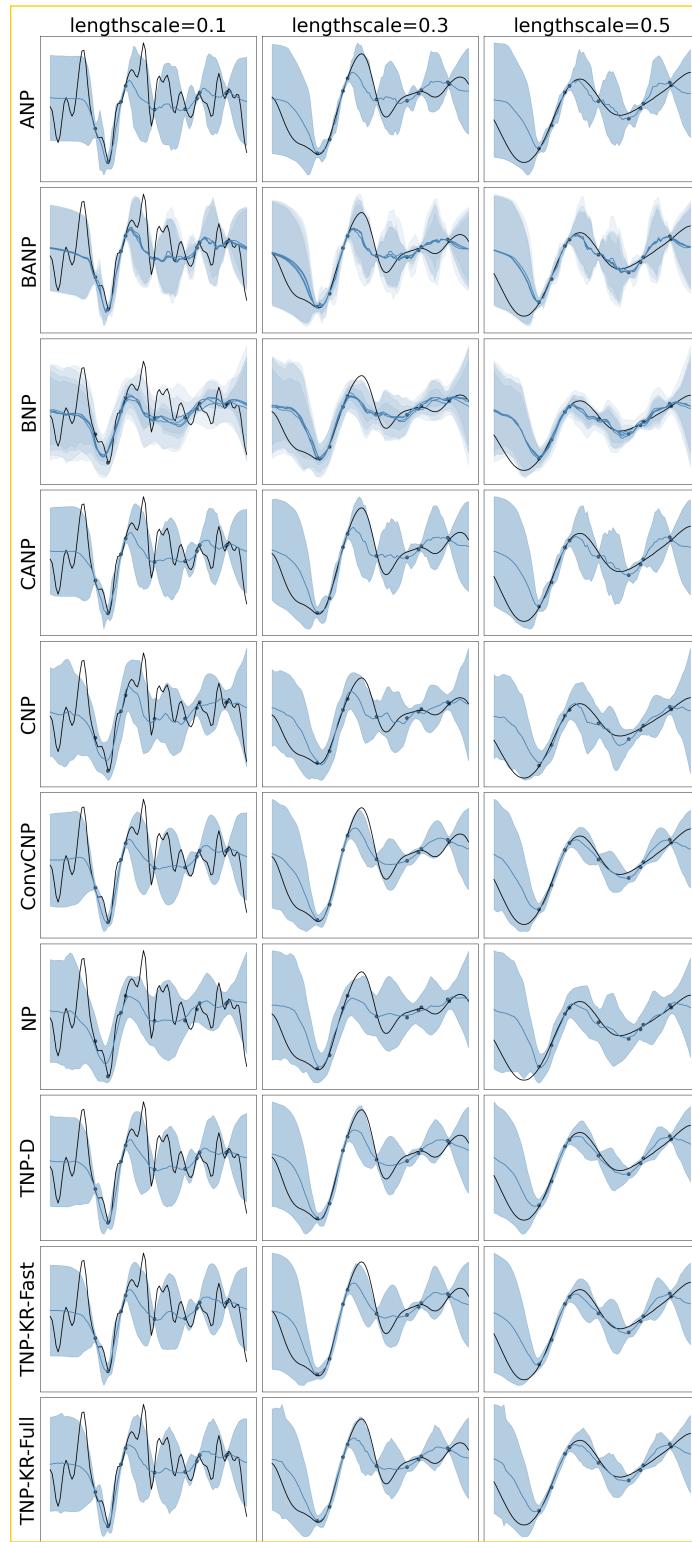


Figure 7: 1D GP RBF samples. The black line is the true function, the black dots are noisy observations, the blue line is the model’s mean prediction, and the shaded blue is the uncertainty. Most models do well on this task, but ConvCNP and the TNP variants have marginally tighter uncertainty bounds.

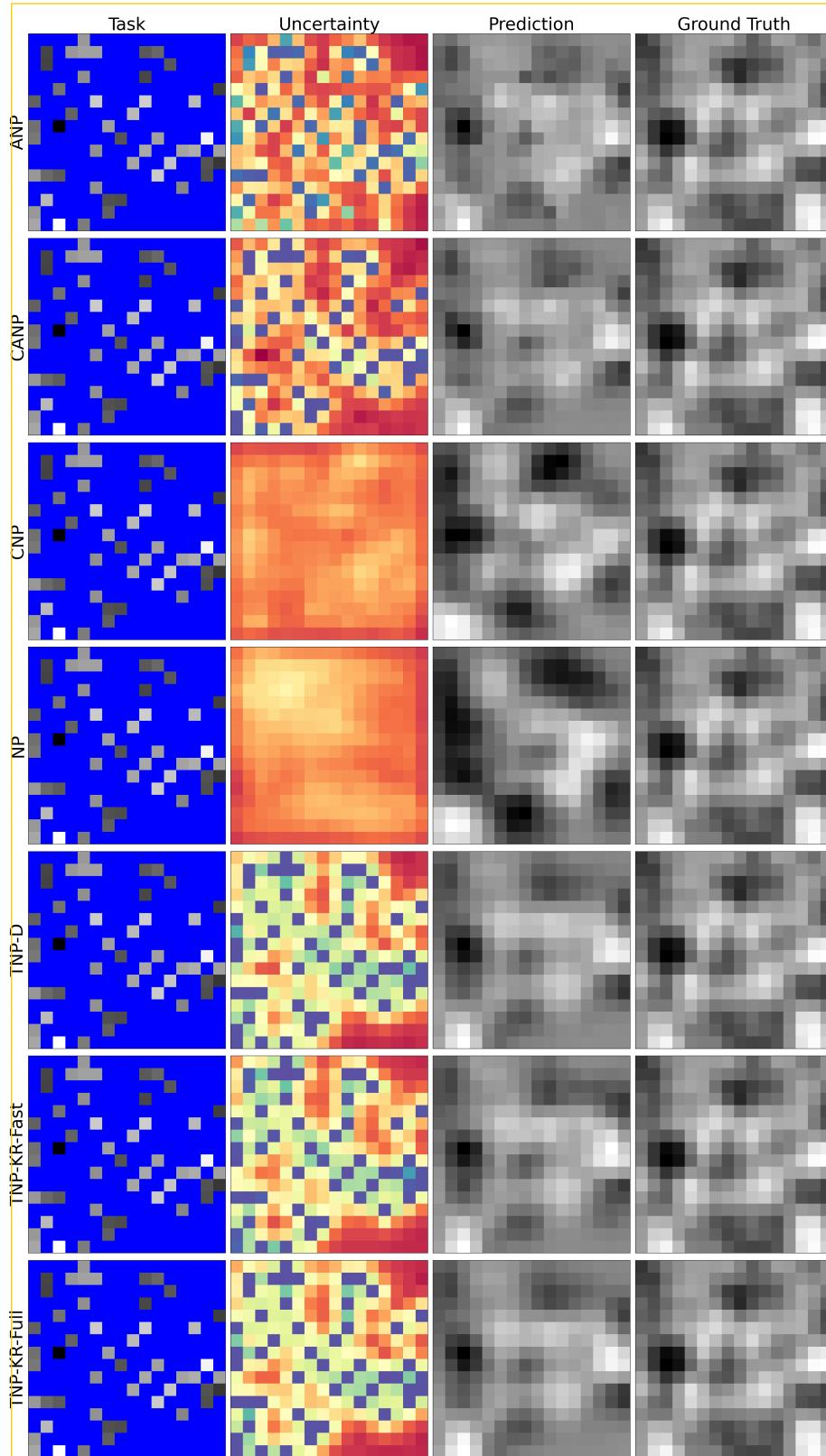


Figure 8: 2D GP RBF sample. In the Task panel, masked blue pixels are unobserved locations. In the Uncertainty panel, warmer colors signify greater uncertainty. NP and CNP struggle to model in two dimensions. TNP variants perform better than ANP variants, largely due to the increased accuracy of transformer blocks.

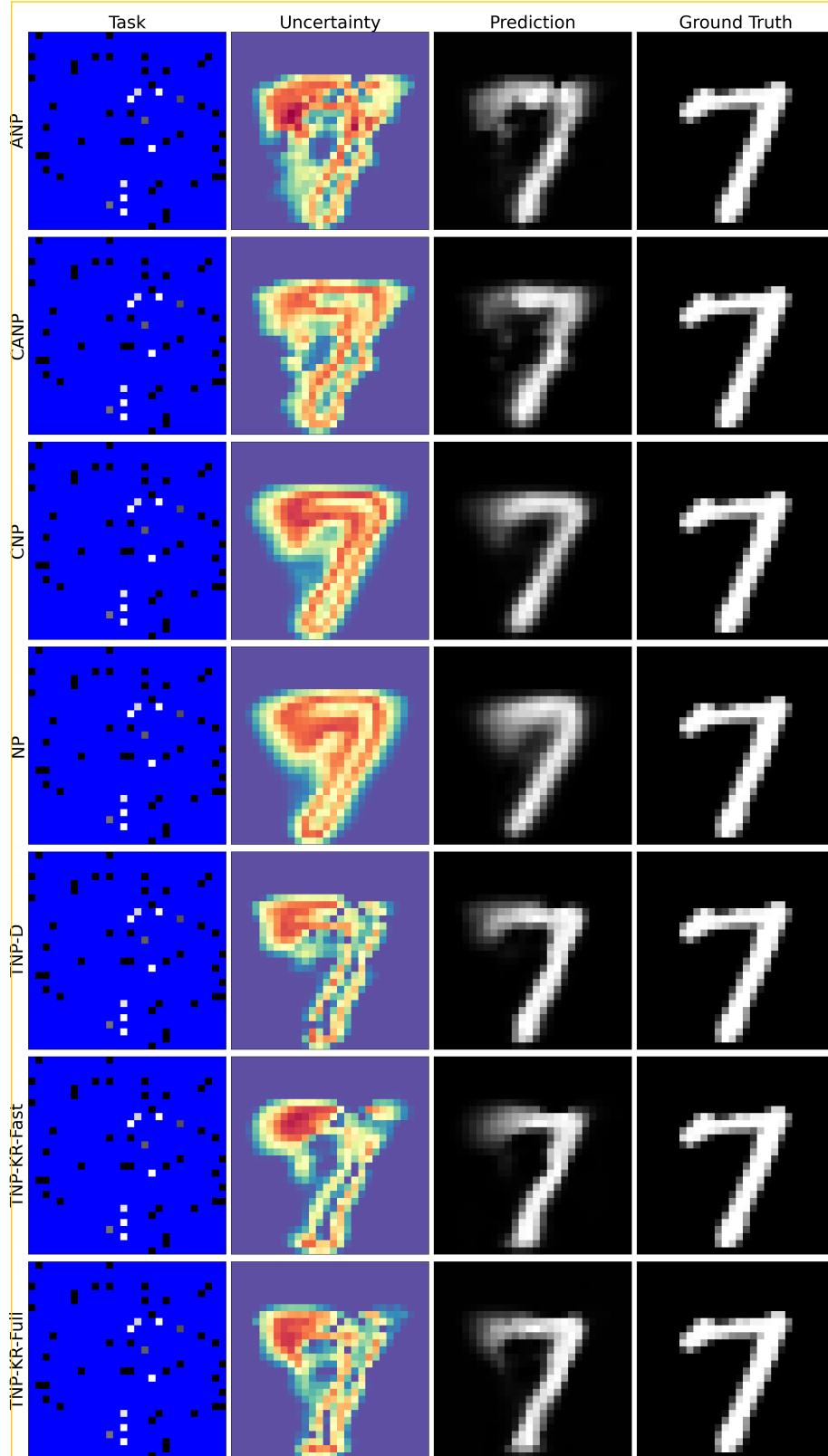


Figure 9: MNIST sample. In the Task panel, masked blue pixels are unobserved locations. In the Uncertainty panel, warmer colors signify greater uncertainty. Predictions are relatively consistent among models, but the TNP variants have better uncertainty bounds, particularly around observed context points.

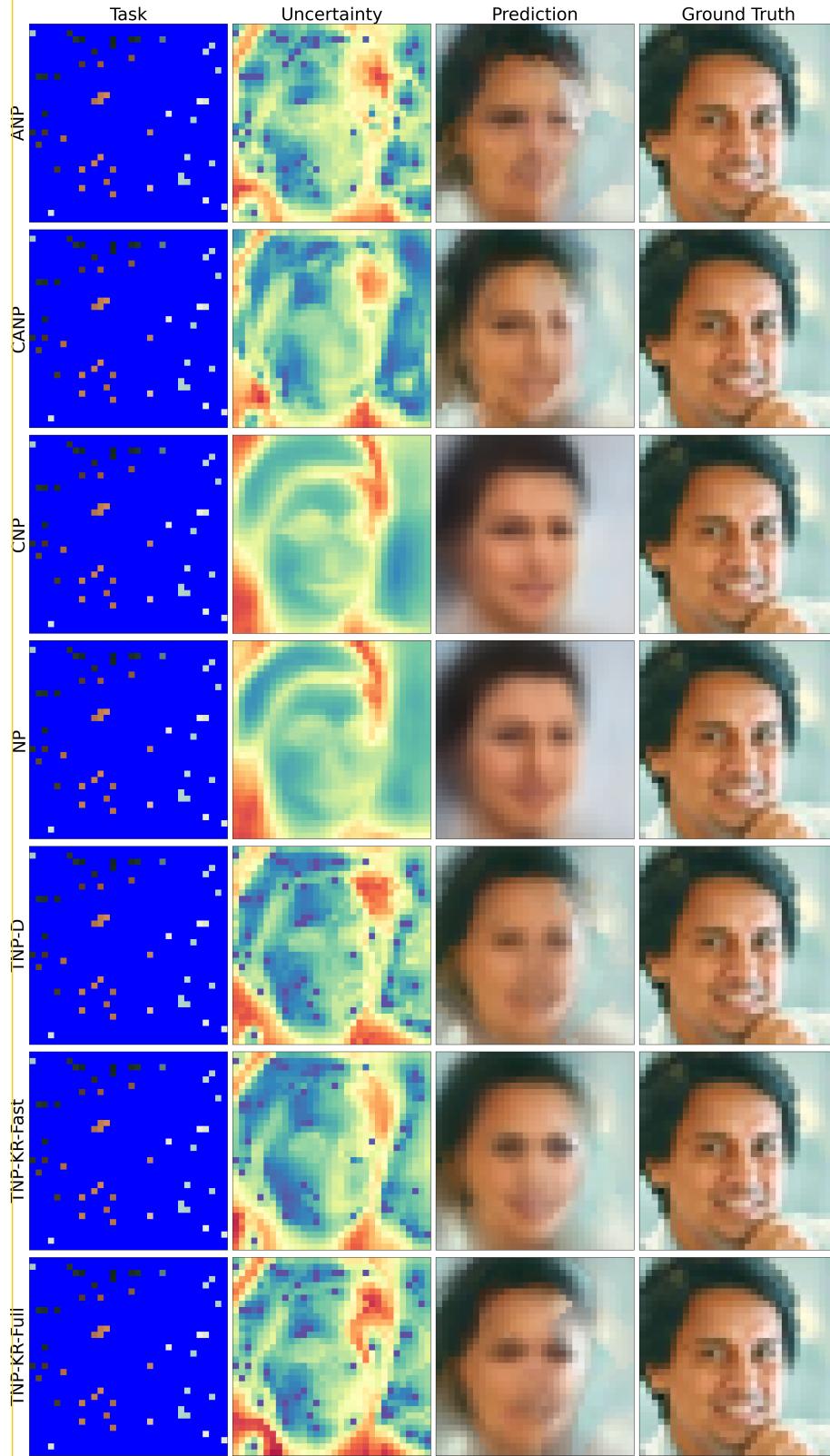


Figure 10: CelebA sample. In the Task panel, masked blue pixels are unobserved locations. In the Uncertainty panel, warmer colors signify greater uncertainty. NP and CNP tend to oversmooth, while TNP and ANP variants are able to identify more details, particularly at observed context points. There are purple dots for TNP and ANP variants, indicating that the model is highly confident in the uncertainty bounds around observed points, but NP and CNP tend to oversmooth at context points.

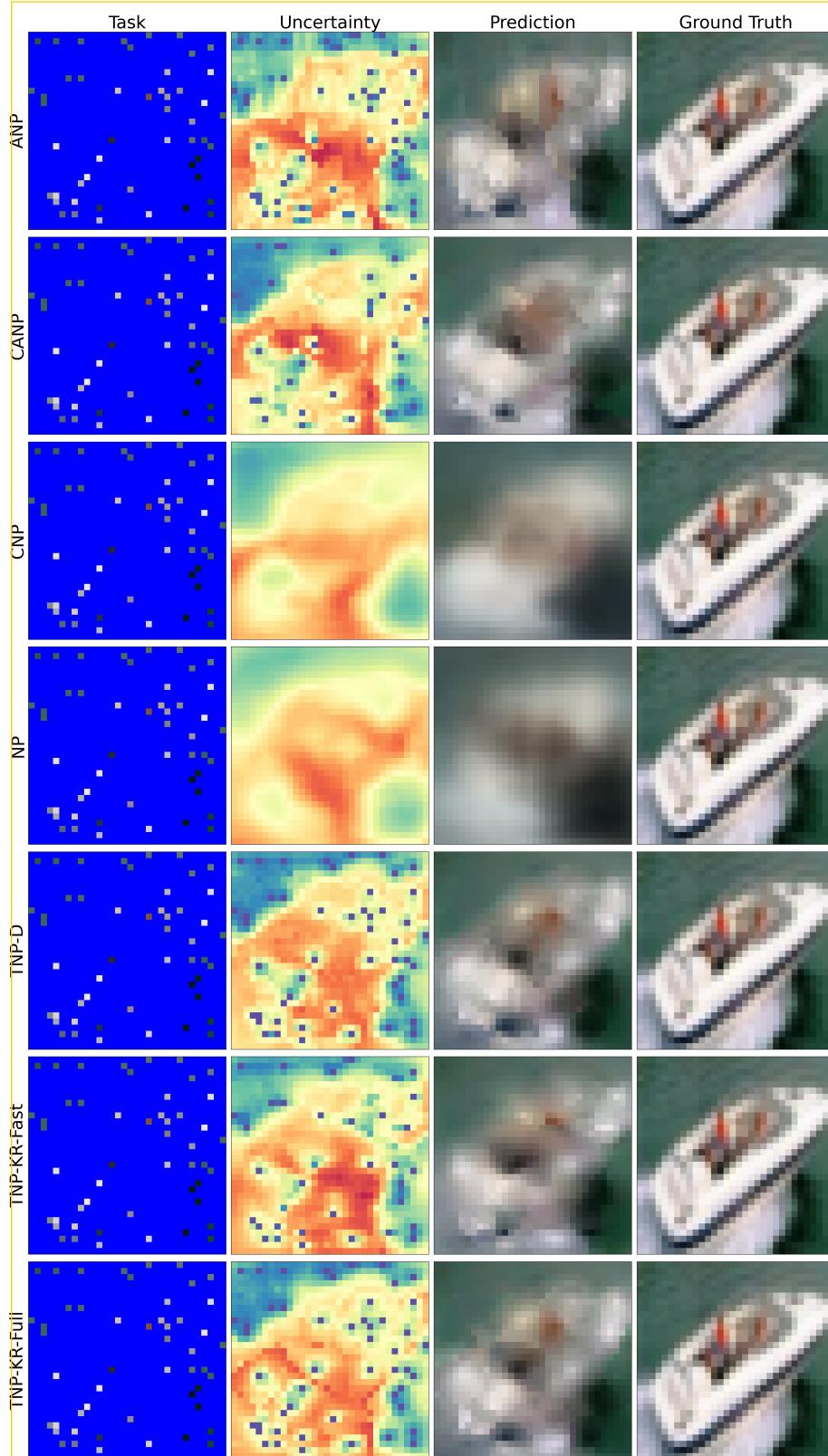


Figure 11: Cifar 10 sample. In the Task panel, masked blue pixels are unobserved locations. In the Uncertainty panel, warmer colors signify greater uncertainty. NP and CNP models oversmooth, while TNP and ANP variants appear grainy due to attention focusing on pixels in different regions of the image.