Appendices to Contrastive Conditional Neural Processes

This contains three sections for describing experimental datasets, additional experimental results and detailed model implementation.

A. Dataset Description

The experiments in this work are three folds, including 1D time-series function, 2D predator-prey dynamics and high-dimensional time-series datasets.

A.1. 1D time-series functions

We run few-shot regression on four function families, see Tab. 1 for the specific definition of each. Each function instantiation used in training/validation/testing is generated with respect to the specified range therein. By following conventional meta-learning setting, instantiations within a batch are randomly sampled from the dataset of sinusoid functions without replacement, *i.e.* $f_i \neq f_{i'}$, $\forall f_i \in f_B$, with $f_B = \{f_i\}_{i=1:|B|}$.

We generate 500 samples for each function family with splitting training/validation/testing in the 9 : 1 : 1 ratio, without function overlap. For each instantiation, context size is sampled from (0, N] (N refers to N-shot) with extra target size is drawn from (0, 10] in the training phase, whilst in validation and testing phase the context size is fixed at N and prediction is evaluated on the whole sequence. The reported results are acquired on the test set. The dimension of observation space is $\mathcal{Y} \subseteq \mathbb{R}$.

A.2. 1D GP-generated functions

In addition to the above four function families, we also conduct 1D regression experiments where the datagenerating functions are Gaussian Processes with different

1D Synthetic Function					
Family	Form	α	β	\boldsymbol{x}	
Sinusoid	$y = \alpha \sin(x - \beta)$	(-1, 1)	(-0.5, 0.5)	$(-\pi, \pi)$	
Exponentials	$y = \alpha \times \exp(x - \beta)$	(-1, 1)	(-0.5, 0.5)	(-1, 4)	
Oscillators	$y = \alpha \sin(x - \beta) \exp(-0.5t)$	(-1, 1)	(-0.5, 0.5)	(0, 5)	
Straight lines	$y = \alpha x + \beta$	(-1, 1)	(-0.5, 0.5)	(0, 5)	

Table 1. Details of 1D times-series functions. Columns of α , β , x corresponds to the range where a function instantiation is randomly sampled therein (e.g., $f_1 = -0.5 \sin(x - 0.3)$ and $f_2 = 0.4 \sin(x + 0.1)$).

kernels. We use the kernels as in AttnCNP and ConvCNP, including RBF, Periodic and Noisy Matérn:

$$k(x_i, x_j)_{\text{RBF}} = \exp\left(-\frac{d(x_i, x_j)^2}{2l^2}\right)$$

$$k(x_i, x_j)_{\text{PER}} = \exp\left(-\frac{2\sin^2\left(\pi d\left(x_i, x_j\right)/p\right)}{l^2}\right)$$

$$k(x_i, x_j)_{\text{MAT}} = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left(\frac{\sqrt{2\nu}}{l}d\left(x_i, x_j\right)\right)^{\nu} \qquad (1)$$

$$K_{\nu}\left(\frac{\sqrt{2\nu}}{l}d\left(x_i, x_j\right)\right) + \epsilon$$

The number of instantiations we generate for each kernel is 4096, of which 256 are used for testing and 256 for validation, respectively. Other than that, we use the same sampling strategy as time-series functions, and perform N-shot regression.

A.3. 2D predator-prey dynamics

To model the 2D population dynamics, we fit the Lotka-Volterra equations (LV) with CNPs. Given values of y_1 and y_2 at initial time index x = 0, the poluations of both species vary after each time increment based on interaction coefficients α , β , δ , γ . The growth rates are defined by

$$\nabla_x y_1 = \alpha y_1 - \beta y_1 y_2$$

$$\nabla_x y_2 = \delta y_1 y_2 - \gamma y_2$$
(2)

with respect to time increment $x : 0 \rightarrow x_{\max}$. Thus, in each simulated trajectory, the populations at each time index are determined by the initial values of $y_1, y_2, \alpha, \beta, \gamma, \delta$. We consider two modes of simulation. One is *Greek* mode where $\alpha, \beta, \gamma, \delta$ are set to fixed values with y_1 and y_2 are randomly sampled from the given range, while in *Population* mode y_1 and y_2 are assigned with fixed initial numbers with $\alpha, \beta, \gamma, \delta$ become random variables (Tab. 2). For both modes, we run 200 trials as the meta-dataset with accumulating 150 timesteps, *i.e.* $x_{\max} = 150$ for every trial. We generate 200 samples for each function family with splitting training/validation/testing in the 9 : 1 : 1 ratio, without function overlap. For each trajectory, context size is sampled from (0, 80] with extra target size is drawn from (0, 20] in the training phase, whilst in validation and testing phase

2D Population Dynamics						
Mode	y_1	y_2	α	β	γ	δ
Greek	(0.5, 2.0)	(0.5, 2.0)	4/3	2/3	1	1
Population	1.6	0.8	(0.9, 1.1)	(0.05, 0.15)	(1.25, 1.75)	(0.5, 1.0

Table 2. Details of 2D LV systems. y_1 and y_2 denote the initial population of predator and prey, resepctively. $\alpha, \beta, \delta, \gamma$ refer to the interaciton coefficients. Either initial population or interaction coefficients are set to fixed in a specific mode, while another group is randomly initialized. All the values have been normalized to avoid the impact caused by magnitude.

the context size is fixed at 80 and prediction is evaluated on the whole trajectory. The reported results are acquired on the test set. In this case, the dimension of observation space is $\mathcal{Y} \subseteq \mathbb{R}^2$.

A.4. Higher-dimensional time-series

We experiment with two higher-dimensional time-series dataset, where observations are depicted as images, including a BouncingBall dataset and a RotMNIST dataset. For Bouncing Ball, each trajectory contains the movements of three interacting balls within a rectangular box, with the length of 20 steps, where each timestep is framed as a 32*32image. We randomly grab 10000 abd 500 trajectories for training and testing, respectively. For RotMNIST, each trajectory contains the rotation of a handwritten digit "3" presenting 16 angles (so as with the length of 16 steps), where each timestep is frame as a 28 * 28 image. Randomly drawn 400 sequences are used for experiment with a 9:1 ratio for splitting training/testing set. During the training phase for both datasets, the context size and extra target size are randomly drawn in the range (0, 5] and (0, 5]. In this case, the dimension of observation space are $\mathcal{Y} \subseteq \mathbb{R}^{784}$ (RotMNIST) and $\mathcal{Y} \subseteq \mathbb{R}^{1024}$ (Bouncing Ball).



Figure 1. Examples of 5 consecutive steps of Bouncing Ball data and RotMNIST data.

B. Additional Experimental Results

B.1. Few-shot Regression on GP-generated Function

We supplement the results for 5-shot and 20-shot regression for GP-generated functions in three different ker-



Figure 2. Comparison on Computational Efficiency

nels, with reconstruction error (MSE) used as the evaluation metrics. It is noticeable that ConvCNP shows particular merits when predicting periodic data, while CCNP performs better with noise contained (Noisy Matérn). Also, the *translation-equivariance* assumption baked into ConvCNP may not hold for every case.

	5-shot regression			20-shot regression		
	RBF	Periodic	Matérn	RBF	Periodic	Matérn
CNP	0.723 ± 0.008	0.588 ± 0.001	0.972 ± 0.005	0.456 ± 0.019	0.535 ± 0.001	0.923 ± 0.044
AttnCNP	0.587 ± 0.004	0.552 ± 0.007	0.908 ± 0.007	0.119 ± 0.003	0.497 ± 0.022	0.569 ± 0.283
ConvCNP	0.581 ± 0.004	$\textbf{0.278} \pm \textbf{0.007}$	0.785 ± 0.013	0.102 ± 0.001	$\textbf{0.065} \pm \textbf{0.128}$	0.468 ± 0.151
CCNP	$\overline{\textbf{0.509}\pm\textbf{0.032}}$	$\underline{0.443 \pm 0.138}$	$\overline{\textbf{0.635}\pm\textbf{0.004}}$	$\textbf{0.100} \pm \textbf{0.015}$	$\underline{0.183 \pm 0.029}$	$\overline{0.412\pm0.002}$

Table 3. Additional 1D Few-Shot Regression with GP in three kernels, running over 3 different seeds.

B.2. Running Efficiency

We have discussed that one of the limitations of CCNP is running efficiency due to the negative sampling steps when performing contrastive learning. We provide a comparison of initial relative elapsed time between CNP, AttnCNP, and CCNP running on the same epochs (see Fig. 2). CCNP primarily emphasizes predictive efficacy, so its efficiency may be sacrificed somewhat. Meanwhile, AttnCNP and CCNP are more comparable since both involve calculating attention, while CCNP costs more time to train contrastive targets. It might be possible to optimize CCNP's efficiency by using ideas like MoCo.

B.3. Ablation studies for Projection Heads

As it is a commonsense that the projection head plays a significant role in determining the performance of contrastive learning, we also examine the effects of setting different sizes for projection head through studies performed within high-dimensional sequences.

Dim.	RotMNIST MSE ($\times 10^{-2}$)	BouncingBall MSE ($\times 10^{-1}$)
8	0.751 ± 0.126	0.537 ± 0.002
16	0.687 ± 0.078	0.511 ± 0.001
32	0.654 ± 0.072	0.483 ± 0.008
64	$\underline{0.648 \pm 0.044}$	$\underline{0.470 \pm 0.005}$
128	$\textbf{0.646} \pm \textbf{0.079}$	$\textbf{0.458} \pm \textbf{0.004}$

Table 4. Ablation studies of Projection Head with pred_size=10

C. Additional Implementation Details

We implement the model with PyTorch 1.8.0 on a Nvidia GTX Titan XP GPU. See below for the details of CCNP's components.

C.1. CCNP for 1D

Input Encoder architecture for 1D.¹

$$(x, \boldsymbol{y}) \in (\mathbb{R}, \mathbb{R}) \to \text{FCL}(64) \to \text{ReLU} \\ \to \text{FCL}(64) \to \text{ReLU} \\ \to \text{FCL}(64) \to \text{ReLU} \\ \to \text{FCL}(64) \Rightarrow \mathbf{r}_C/\mathbf{r}_T/\mathbf{r}_F$$

Position-Aware Self-Attention architecture for 1D.

$$\label{eq:MultiHeadAttention} \begin{split} \text{MultiHeadAttention} &= \{K, Q, V, A, H\} \\ \text{with } K = \text{FCL}(64), \text{key_transformation} \end{split}$$

Q = FCL(64), query_transformation V = FCL(64), value_transformation A = DotProductAttention(K, Q, V)H = FCL(64), head_fusion

where

 $\begin{aligned} \text{DotProductAttention} &= \{K, Q, V\} \\ & \text{with } K = \text{FCL}(64), \text{key_transformation} \\ & Q = \text{FCL}(64), \text{query_transformation} \\ & V = \text{FCL}(64), \text{value_transformation} \end{aligned}$

Temporal Contrastive Component for 1D.

$$\begin{aligned} (x_t, \mathbf{r}_T) \in (\mathbb{R}, \mathbb{R}^{64}) &\to \mathrm{FCL}(64) \to \mathrm{ReLU} \\ &\to \mathrm{FCL}(64) \Rightarrow \varphi(x_t, \mathbf{r}_T) \\ \varphi(x_t, \mathbf{r}_T) &\to \mathrm{FCL}(8) \Rightarrow \hat{\boldsymbol{z}}_t \\ \boldsymbol{y}_t &\to \mathrm{FCL}(8) \Rightarrow \boldsymbol{z}_t \\ (\hat{\boldsymbol{z}}_t, \boldsymbol{z}_t) &\to \mathrm{InfoNCE} \\ \mathrm{with} \ \tau = 0.5 \end{aligned}$$

¹FCL(d) = Fully Connected Layer(output_dimension)

Function Contrastive Component for 1D. Taking 2 instantiations for illustration. We use instantiations within a whole batch in practice.

$$\begin{split} \mathbf{x}_{C}^{f_{1}}, \mathbf{x}_{C}^{f_{2}} \rightarrow \mathbf{x}_{C_{1}}^{f_{1}}, \mathbf{x}_{C_{2}}^{f_{2}}, \mathbf{x}_{C_{1}}^{f_{2}}, \mathbf{x}_{C_{2}}^{f_{2}} \\ \mathbf{x}_{C_{1}}^{f_{1}}, \mathbf{x}_{C_{2}}^{f_{1}}, \mathbf{x}_{C_{1}}^{f_{2}}, \mathbf{x}_{C_{2}}^{f_{2}} \rightarrow \mathbf{r}_{C_{1}}^{f_{1}}, \mathbf{r}_{C_{2}}^{f_{1}}, \mathbf{r}_{C_{2}}^{f_{2}}, \mathbf{r}_{C_{2}}^{f_{2}} \\ \mathbf{r}_{C_{1}}^{f_{1}} \rightarrow \text{FCL}(8) \Rightarrow \mathbf{q}_{i}^{f_{1}} \\ \mathbf{r}_{C_{2}}^{f_{2}} \rightarrow \text{FCL}(8) \Rightarrow \mathbf{q}_{j}^{f_{1}} \\ \mathbf{r}_{C_{1}}^{f_{2}} \rightarrow \text{FCL}(8) \Rightarrow \mathbf{q}_{i}^{f_{2}} \\ \mathbf{r}_{C_{1}}^{f_{2}} \rightarrow \text{FCL}(8) \Rightarrow \mathbf{q}_{j}^{f_{2}} \\ \mathbf{q}_{i}^{f_{1}}, \mathbf{q}_{j}^{f_{1}}, \mathbf{q}_{i}^{f_{2}}, \mathbf{q}_{j}^{f_{2}} \rightarrow \text{InfoNCE} \\ \text{with } \tau = 0.5 \end{split}$$

Output Decoder architecture for 1D.

$$(x_t, \mathbf{r}_C, \mathbf{r}_T, \mathbf{r}_F) \in (\mathbb{R}, \mathbb{R}^{64}, \mathbb{R}^{64}, \mathbb{R}^{64}) \to \operatorname{concat}(\cdot, \cdot, \cdot, \cdot)$$

$$\to \operatorname{FCL}(64) \to \operatorname{ReLU}$$

$$\to \operatorname{FCL}(64) \to \operatorname{ReLU}$$

$$\to \operatorname{FCL}(64) \to \operatorname{ReLU}$$

$$\to \operatorname{FCL}(64) \to \operatorname{ReLU}$$

$$\Rightarrow \mathbf{r}_g$$

$$\mathbf{r}_g \in \mathbb{R}^{64} \to \operatorname{FCL}(1) \Rightarrow \hat{\mu}_t$$

$$\mathbf{r}_g \in \mathbb{R}^{64} \to \operatorname{FCL}(1) \Rightarrow \hat{\sigma}_t$$

C.2. CCNP for 2D

(x

Similar to 1D, except for $\boldsymbol{y} \in \mathbb{R}^2$, thus in decoder

$$\mathbf{r}_g \in \mathbb{R}^{64} \to \mathrm{FCL}(2) \Rightarrow \hat{\boldsymbol{\mu}}_t$$
$$\mathbf{r}_g \in \mathbb{R}^{64} \to \mathrm{FCL}(2) \Rightarrow \hat{\boldsymbol{\sigma}}_t$$

C.3. CCNP for High-dimensional data

We replace $\psi(\cdot)$ for encoding observations y with Convolutional Blocks², with BN³. For decoding we use ConvT⁴ Input Encoder architecture for RotMNIST.

$$\begin{aligned} \boldsymbol{y} \in \mathbb{R}^{784} \to \operatorname{Conv}(16, 5, 2, 2) \to \operatorname{BN}(16) \to \operatorname{ReLU} \\ \to \operatorname{Conv}(32, 5, 2, 2) \to \operatorname{BN}(32) \to \operatorname{ReLU} \\ \to \operatorname{Conv}(64, 5, 2, 2) \to \operatorname{BN}(64) \to \operatorname{ReLU} \\ \to \operatorname{Conv}(128, 5, 2, 2) \to \operatorname{BN}(128) \to \operatorname{ReLU} \\ \in \mathbb{R}, \psi(\boldsymbol{y})) \to \operatorname{FCL}(128) \Rightarrow \mathbf{r}_C/\mathbf{r}_T/\mathbf{r}_F \end{aligned}$$

Output Decoder architecture for RotMNIST.

$$(x_t, \mathbf{r}_C, \mathbf{r}_T, \mathbf{r}_F) \in (\mathbb{R}, \mathbb{R}^{128}, \mathbb{R}^{128}, \mathbb{R}^{128}) \to \operatorname{concat}(\cdot, \cdot, \cdot, \cdot) \\ \to \operatorname{FCL}(72) \Rightarrow \mathbf{r}_q$$

 $^{{}^{2}}$ Conv(f, k, s, p) = Convolution2D(feat_maps, kernel, stride, pad)

 $^{^{3}}$ BN(d) = BatchNormalization2D(dim)

⁴ConvT(f, k, s, p) = ConvTranspose2D(feat_map, kernel, stride, pad)

$$\begin{split} \mathbf{r}_g &\rightarrow \text{reshape}(\text{batch_size}, 8, 28, 28) \\ &\rightarrow \text{ConvT}(128, 3, 1, 0) \rightarrow \text{BN}(128) \rightarrow \text{ReLU} \\ &\rightarrow \text{ConvT}(64, 5, 2, 0) \rightarrow \text{BN}(64) \rightarrow \text{ReLU} \\ &\rightarrow \text{ConvT}(32, 5, 2, 1) \rightarrow \text{BN}(32) \rightarrow \text{ReLU} \\ &\rightarrow \text{ConvT}(1, 5, 1, 2) \Rightarrow \hat{\boldsymbol{\mu}} \in \mathbb{R}^{28 \times 28} \end{split}$$

Input Encoder architecture for Bouncing Ball.

$$\begin{aligned} \boldsymbol{y} \in \mathbb{R}^{1024} \to \operatorname{Conv}(16, 5, 2, 2) \to \operatorname{BN}(16) \to \operatorname{ReLU} \\ \to \operatorname{Conv}(32, 5, 2, 2) \to \operatorname{BN}(32) \to \operatorname{ReLU} \\ \to \operatorname{Conv}(64, 5, 2, 2) \to \operatorname{BN}(64) \to \operatorname{ReLU} \\ \to \operatorname{Conv}(128, 5, 2, 2) \to \operatorname{BN}(128) \to \operatorname{ReLU} \\ (x \in \mathbb{R}, \psi(\boldsymbol{y})) \to \operatorname{FCL}(128) \Rightarrow \mathbf{r}_C/\mathbf{r}_T/\mathbf{r}_F \end{aligned}$$

Output Decoder architecture for Bouncing Ball.

$$(x_t, \mathbf{r}_C, \mathbf{r}_T, \mathbf{r}_F) \in (\mathbb{R}, \mathbb{R}^{128}, \mathbb{R}^{128}, \mathbb{R}^{128}) \to \operatorname{concat}(\cdot, \cdot, \cdot, \cdot) \\ \to \operatorname{FCL}(72) \Rightarrow \mathbf{r}_g$$

$$\begin{split} \mathbf{r}_g &\to \mathrm{reshape}(\mathrm{batch_size}, 8, 32, 32) \\ &\to \mathrm{ConvT}(128, 3, 2, 1) \to \mathrm{BN}(128) \to \mathrm{ReLU} \\ &\to \mathrm{ConvT}(64, 5, 2, 1) \to \mathrm{BN}(64) \to \mathrm{ReLU} \\ &\to \mathrm{ConvT}(32, 5, 2, 1) \to \mathrm{BN}(32) \to \mathrm{ReLU} \\ &\to \mathrm{ConvT}(1, 5, 1, 2) \Rightarrow \hat{\boldsymbol{\mu}} \in \mathbb{R}^{32 \times 32} \end{split}$$

References