

Supplementary material

A. Hyperparameters and training settings

For all our models we used the following settings in the UCI experiments:

Kernels. All kernels were the RBF, using a lengthscale parameter per input dimension, initialized to the squareroot of the dimension.

Inducing points. For data with more than 128 training data points the inducing point locations were chosen using the `kmeans2` from the `scipy` package, with 128 points. Otherwise they were set to the data. For latent variable layers, the GP layer above has extra dimensions. The inducing points for the extra dimensions were padded with random draws from a standard normal variable.

Linear projections between layers. We implemented the linear mean functions and multioutput structure using a linear projection of 5 independent GPs concatenated with their inputs. We initialized the projection matrix to the first 5 principle components of the data concatenated with the identity matrix.

Amortization of the variational parameters. We used three layer fully connected network with skip connections between each layer. We used the `tanh` non-linearity with 10 units for the inner layers. We used the weight initialization from (Glorot and Bengio, 2010) and the exponential function to enforce positivity of the standard deviation parameters. We added a bias of -5 in the final layer to ensure the standard deviations were small at the start of optimization.

Likelihood. The likelihood variance was initialized 0.01.

Parameterizations. All positive model parameters were constrained to be positive using the `softplus` function, clipped to 10^{-6} . The variational parameters for the sparse GP layers were parameterized by the mean and the square root of the covariance.

Optimization. The final GP layers were optimized using natural gradients on the natural parameters, with an initial step size of 0.01. All other parameters were optimized using the Adam optimizer (Kingma and Ba, 2015) with all parameters set to their tensorflow default values except the initial step size of 0.005. Optimizing just final layer using natural gradients and the inner layers with the Adam optimizer in this way is show by Hebbal et al. (2019) to be an effective strategy in practice. We used a batch size of 512 and trained for 100K iterations, annealing the learning rate of both the adam and natural gradient steps by a factor of 0.98 per 1000 iterations. The mean functions and kernel projection matrices (the P matrices that correlate the outputs) were *not* optimized.

Importance weights. We used 5 importance weights for all the latent variables models with importance-weighted variational inference.

Predictions. We used 2000 samples for prediction, sampling from the prior for w for the latent variables models. We then used a kernel density estimator (with Silverman’s rule to set the bandwidth) to estimate the density.

Splits and preprocessing. We used 90% splits on the shuffled data, and rescaled the inputs and outputs to unit standard deviation and zero mean. For dataset with more than 10K test points we used a random sample of 10K test points (the same sample for each split). NB we have reported the test log-likelihoods without restoring the scaling, to facilitate comparisons between data. The splits were the same for all the models. We implemented our models using `gpflow` (Matthews et al., 2017).

Comparison of methods. The error bars in Table 1 are standard errors over the five splits. These error bars are likely to *overestimate* the uncertainty for comparisons between methods, however, as there may be correlations between the performance of methods over the splits (as the splits are the same for each method). The average ranks were computed for each split separately, and then averaged over the splits and method. To mitigate the effect of correlations between splits we report the median difference in test log-likelihood from the single layer GP.

VAE models The conditional VAE models used `tanh` activations and weight initializations from (Glorot and Bengio, 2010). Optimization was performed used the adam optimizer, with the same batch size and learning rate schedule as the GP models.

B. Variational posterior derivation

Here we derive the sparse GP posterior $q(f)$ and its KL divergence from the prior, following [Hensman et al. \(2013\)](#). This derivation applies to all the GP layers in the DGP model.

We begin with the approach of [Titsias \(2009\)](#) and form a variational distribution over the function f by conditioning the prior on a set of inducing points $\{f(\tilde{x}_i)\}_{i=1}^M$. We write $\tilde{\mathbf{f}}$ for the vector with i th component $f(\tilde{x}_i)$. In [Titsias \(2009\)](#), $q(f)$ is defined as

$$q(f) = p(f|\tilde{\mathbf{f}})q(\tilde{\mathbf{f}}), \quad (12)$$

where $q(\tilde{\mathbf{f}})$ is unspecified. In the single layer case, analytic results can be used to show that $q(\tilde{\mathbf{f}})$ has an optimal form that is Gaussian. Instead, we follow [Hensman et al. \(2013\)](#) and *assume* that

$$q(\tilde{\mathbf{f}}) = \mathcal{N}(\tilde{\mathbf{f}}|\tilde{\mathbf{m}}, \tilde{\mathbf{S}}). \quad (13)$$

As this choice is conjugate to the prior conditional $p(f|\tilde{\mathbf{f}})$, we can use standard Gaussians identities to show that

$$q(f) = \mathcal{GP}(\mu, \Sigma), \quad (14)$$

where

$$\mu(x) = \mathbf{k}(x)^\top \tilde{\mathbf{K}}^{-1} \mathbf{m} \quad (15)$$

$$\Sigma(x, x') = k(x, x') - \mathbf{k}(x)^\top \tilde{\mathbf{K}}^{-1} (\tilde{\mathbf{K}} - \mathbf{S}) \tilde{\mathbf{K}}^{-1} \mathbf{k}(x') \quad (16)$$

with $[\mathbf{k}(x)]_i = k(x, \tilde{x}_i)$ and $[\tilde{\mathbf{K}}]_{ij} = k(\tilde{x}_i, \tilde{x}_j)$. The KL divergence is given by

$$\text{KL}(q(f)||p(f)) = -\mathbb{E}_{q(f)} \log \frac{p(f)}{q(f)} \quad (17)$$

$$= -\mathbb{E}_{q(f)} \log \frac{p(f|\tilde{\mathbf{f}})p(\tilde{\mathbf{f}})}{p(f|\tilde{\mathbf{f}})q(\tilde{\mathbf{f}})} \quad (18)$$

$$= -\mathbb{E}_{q(\tilde{\mathbf{f}})} \log \frac{p(\tilde{\mathbf{f}})}{q(\tilde{\mathbf{f}})} \quad (19)$$

$$= \text{KL}(q(\tilde{\mathbf{f}})||p(\tilde{\mathbf{f}})), \quad (20)$$

which is finite.

The parameters $\tilde{\mathbf{m}}$ and $\tilde{\mathbf{S}}$ together with $\{\tilde{x}_i\}$ are variational parameters.

For multiple outputs, use the same matrix P for correlating the outputs as the prior. In practice, we implement the correlated output model by multiplying independent GPs by the matrix P , so we use independent outputs for the variational distribution.

C. Further tables and figures

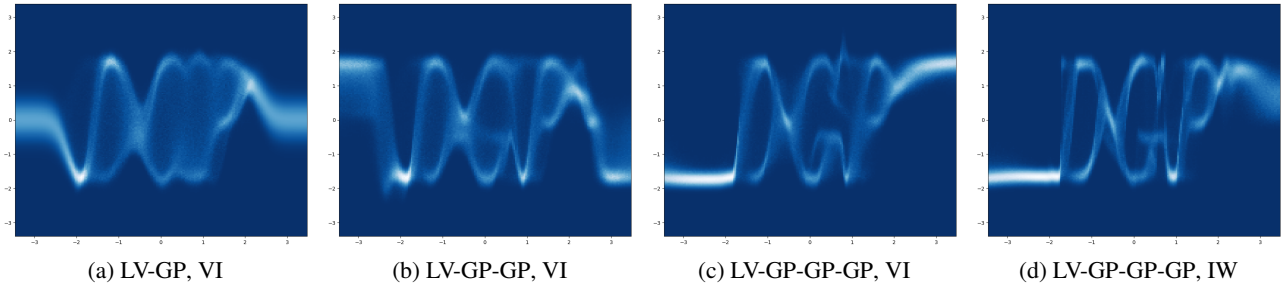


Figure 5: Posteriors for further models for the ‘DGP’ data

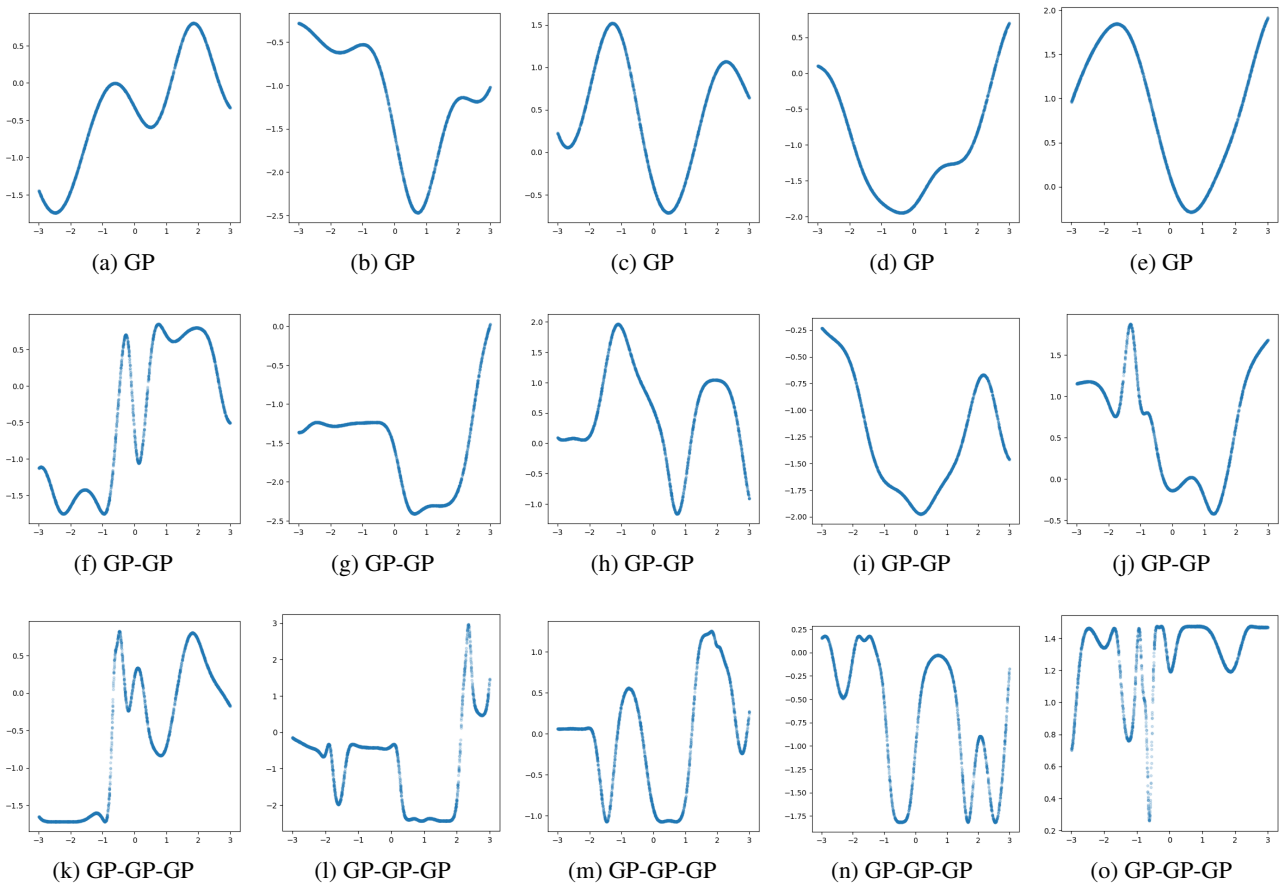


Figure 6: Prior draws from one, two and three layer models without latent variables

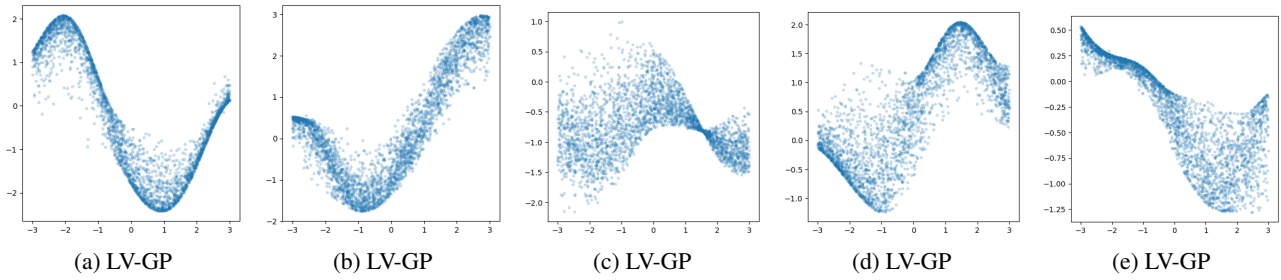


Figure 7: Prior draws from the LV-GP model

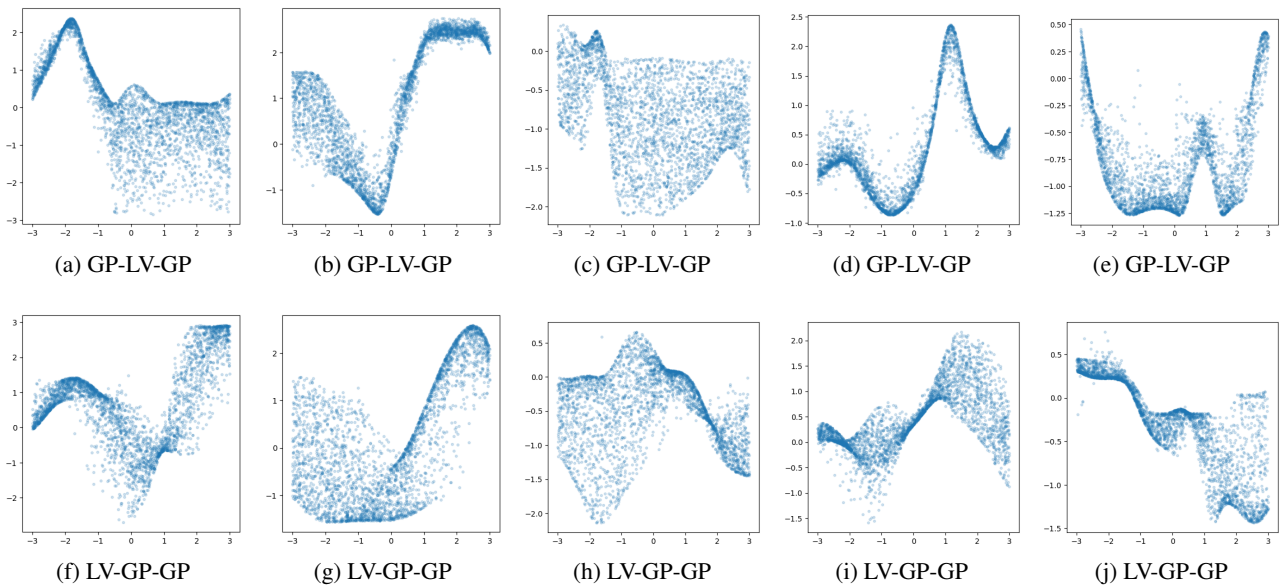


Figure 8: Prior draws from the LV-GP-GP and GP-LV-GP models

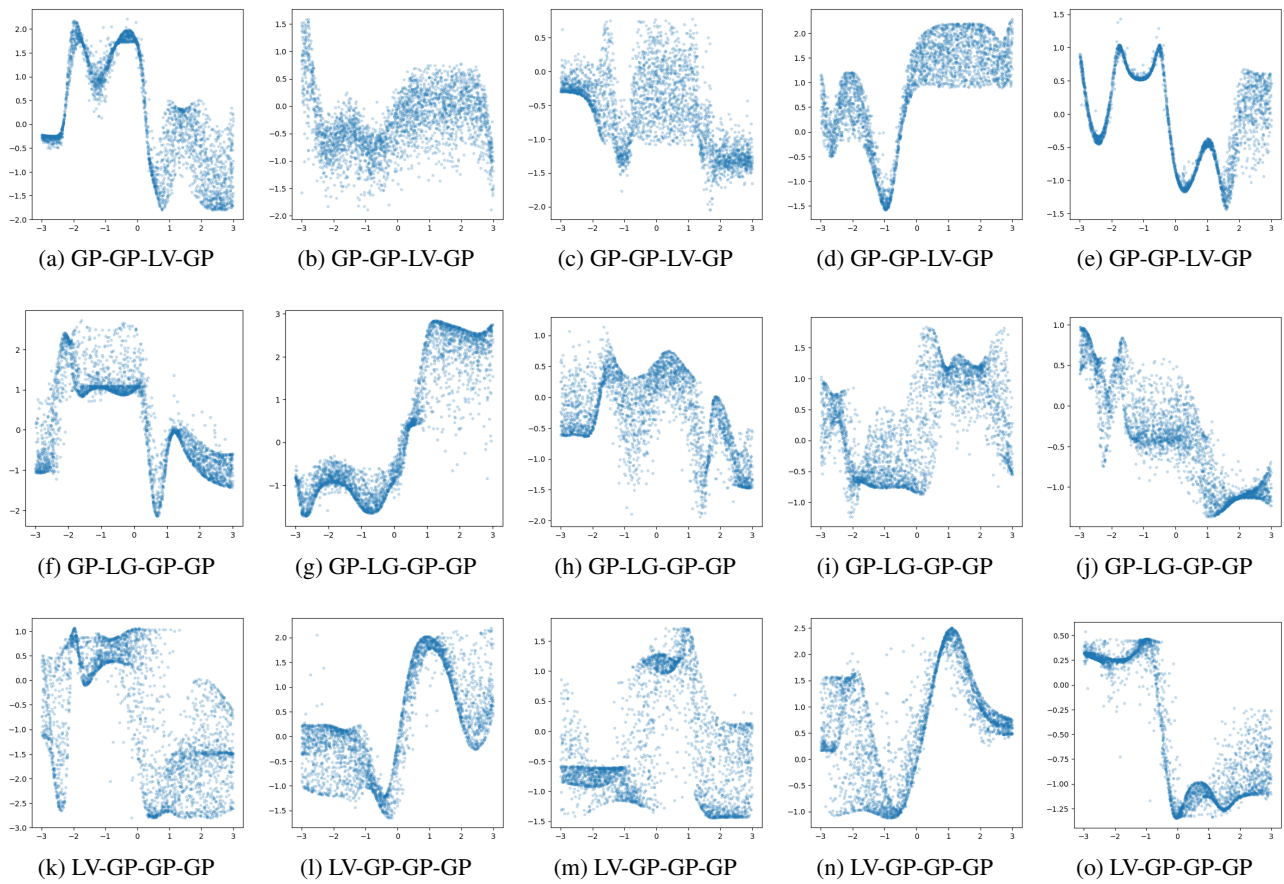


Figure 9: Prior draws from LV-GP-GP-GP, GP-LG-GP-GP and GP-GP-LV-GP models

Dataset	N	D	Linear		CVAE 50	CVAE 100, 100	GP		GP-GP		LV-GP		LV-GP-GP		LV-GP-GP-GP	
			VI	VI			VI	VI	VI	VI	VI	VI	VI	VI	VI	VI
challenger	23	4			0.5191	0.9990	0.8870	0.7739		0.9099	0.7001	0.7115	0.8466	0.8484	0.7187	
fertility	100	9			0.9408				0.9598	0.9341			0.9928	0.9969		
concreteslump	103	7							0.9970	0.9922						
autos	159	25			0.9823	0.8798			0.9903	0.9874			0.9990	0.9990	0.9990	
servo	167	4			0.8358				0.9989	0.9975	0.9990	0.9990			0.9988	
breastcancer	194	33							0.9787	0.9814			0.9911	0.9943		
machine	209	7			0.9181	0.8967			0.9974	0.9989						
yacht	308	6			0.9987				0.9989	0.9987	0.9988	0.9987	0.9985	0.9990	0.9988	
automp	392	7			0.9097	0.9948			0.9989	0.9979	0.9979	0.9989	0.9980	0.9981	0.9981	
boston	506	13			0.9916				0.9964	0.7034	0.7074	0.7470	0.7398	0.7608	0.7642	
forest*	517	12			0.8453	0.7145			0.9981							
stock	536	11														
pendulum	630	9													0.9988	
energy	768	8									0.9987				0.9990	
concrete	1030	8													0.9989	
solar*	1066	10			0.9986	0.9782				0.3380	0.3509	0.2728	0.2945	0.2557	0.2883	
airfoil	1503	5			0.3028	0.3228							0.9988	0.9988	0.9988	
winered	1599	11			0.9989	0.9986							0.9234	0.9751	0.8473	
gas	2565	128			0.9790	0.6980				0.9983	0.9935	0.9989	0.9983	0.9990	0.9986	
skillercraft	3338	19			0.9984					0.9989	0.9983	0.9989	0.9988	0.9990	0.9989	
sm1†	4137	26														
winewhite	4898	11												0.9990	0.9990	
parkinsons	5875	20														
kin8nm†	8192	8			0.9989											
pumadyn32nm	8192	32														
power*	9568	4			0.9964								0.9964	0.9963	0.9855	
naval	11934	14												0.9770		
pol†	15000	26			0.9970											
elevators	16599	18			0.9979								0.9839	0.9508	0.9980	
bike‡	17379	17			0.9815								0.9987	0.9976	0.9979	
kin40k‡	40000	8			0.9990								0.9990	0.9990		
protein	45730	9			0.8839	0.8800							0.8904	0.8761	0.8691	
tamielctric	45781	3			0.9551	0.9502							0.9550	0.9549	0.9551	
keggdirected†	48827	20			0.8279	0.5350							0.6839	0.5187	0.7039	
slice†	53500	385												0.9988	0.9988	
keggundirected*	63608	27			0.4804	0.9882							0.5511	0.2049	0.9295	
3droad	434874	3			0.9421	0.9479							0.9650	0.9467	0.9744	
song	515345	90			0.9936	0.9957							0.9969	0.9893	0.9963	
buzz	583250	77			0.9952	0.9909							0.9988	0.9938	0.9984	
nytaxi	1420068	8			0.9608	0.9422							0.9515	0.9019	0.9467	
houseelectric‡	2049280	11											0.9962	0.9671	0.9933	

Table 3: Median Shapiro–Wilk test statistic for the test points. Blanks have values of one (perfectly Gaussian). Smaller values indicate more evidence for non-Gaussian marginal distributions.

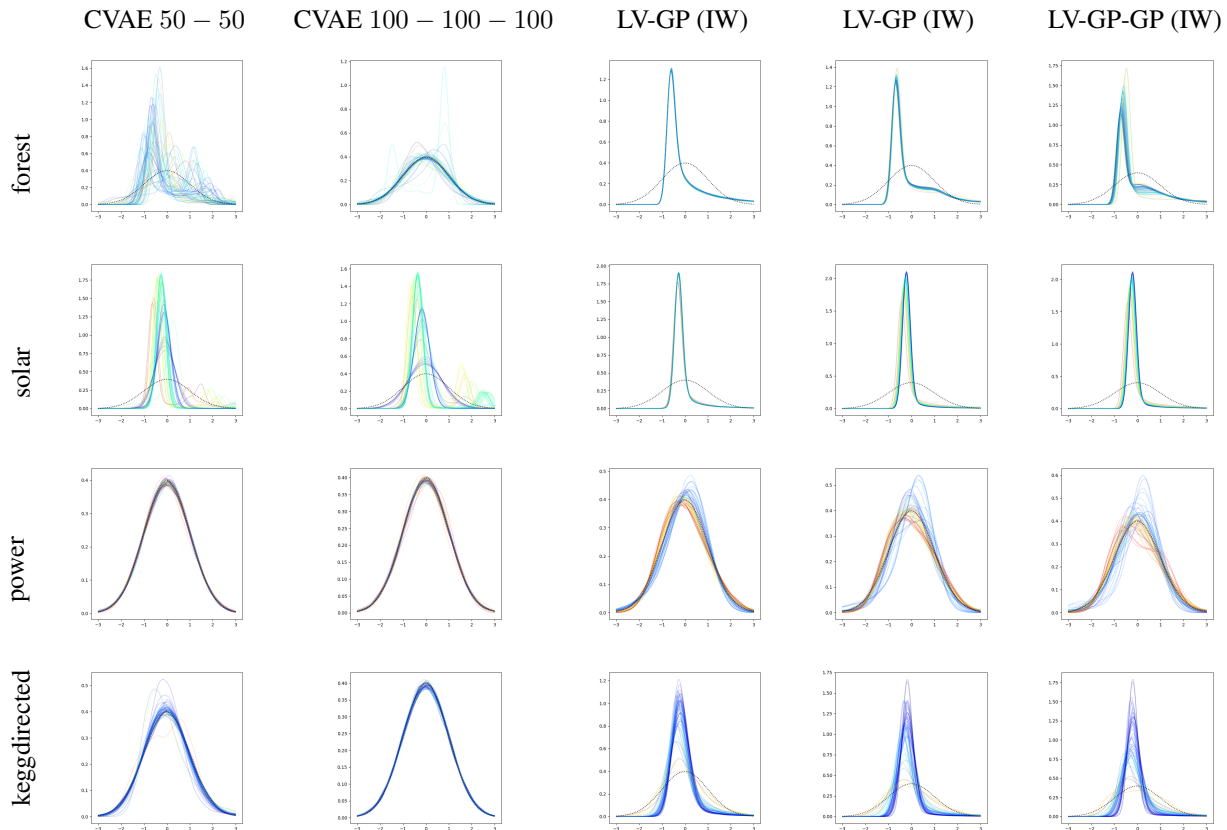


Figure 10: Posterior marginals, re-scaled to zero mean and unit standard deviation, with the Gaussian density marked as a dotted curve. Colors are according to the principal component of the inputs. These are the four datasets highlighted in the main text for being prominent examples benefiting from latent variables. In these cases the additional depth of model leads to more non-Gaussian marginals.

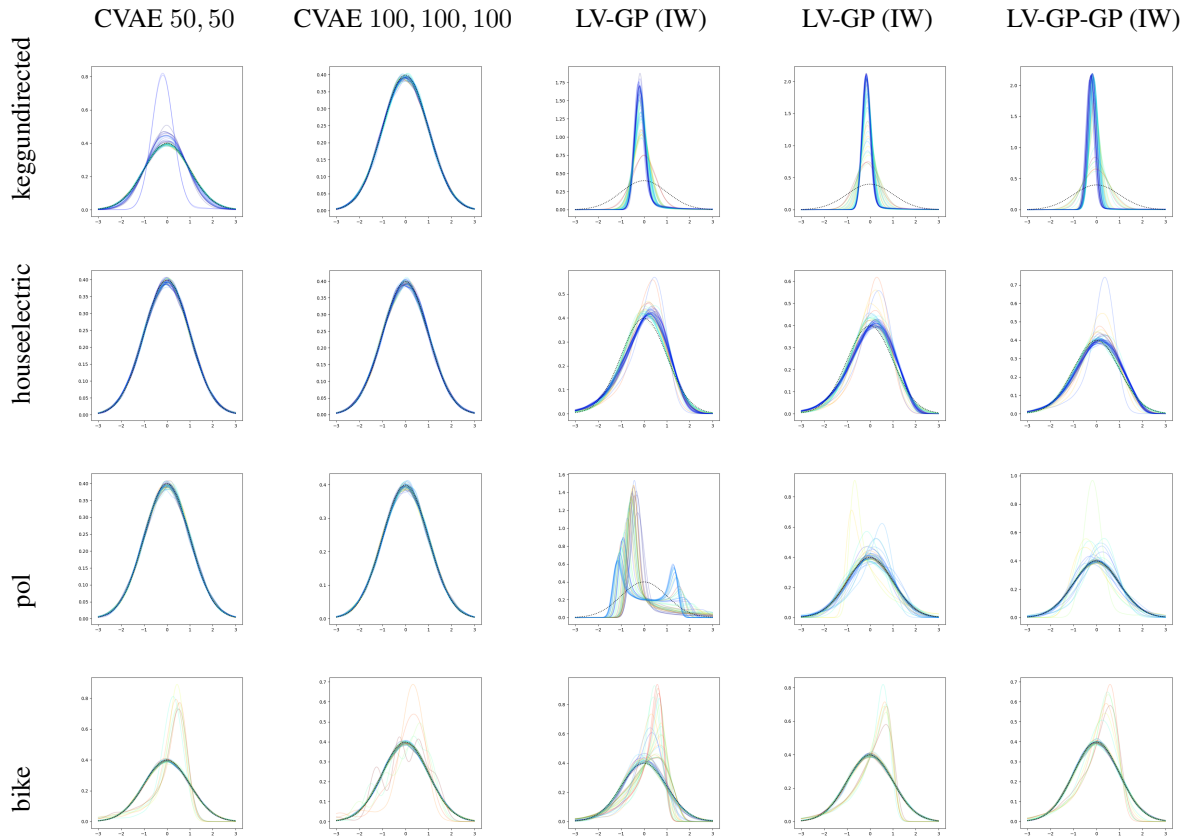


Figure 11: As Figure 10, but for the four datasets highlighted in the main text for benefiting from both latent variables and depth. Note that, unlike in Figure 10, for the ‘pol’ and ‘bike’ data the deeper models actually have *simpler* marginals. This is because the deep model has more complexity in the mapping, so can more closely fit the data, whereas the simple model must explain the data with a complex noise distribution.