A Detailed Variational Approximation

In this section, we repeat the derivation of the variational approximation in more detail.

Since exact inference in this model is intractable, we discuss a variational approximation to the model's true marginal likelihood and posterior in this section. Analogously to y, we denote the random vectors which contain the function values of the respective functions and outputs as a and f. The joint probability distribution of the data can then be written as

$$p(\boldsymbol{y}, \boldsymbol{f}, \boldsymbol{a} \mid \boldsymbol{X}) = p(\boldsymbol{f} \mid \boldsymbol{a}) \prod_{d=1}^{D} p(\boldsymbol{y}_{d} \mid \boldsymbol{f}_{d}) p(\boldsymbol{a}_{d} \mid \boldsymbol{X}),$$

$$\boldsymbol{a}_{d} \mid \boldsymbol{X} \sim \mathcal{N}(\boldsymbol{X}, \boldsymbol{K}_{\boldsymbol{a}, \boldsymbol{d}} + \sigma_{\boldsymbol{a}, \boldsymbol{d}}^{2} \mathbf{I}),$$

$$\boldsymbol{f} \mid \boldsymbol{a} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{K}_{\boldsymbol{f}} + \sigma_{\boldsymbol{f}}^{2} \mathbf{I}),$$

$$\boldsymbol{y}_{d} \mid \boldsymbol{f}_{d} \sim \mathcal{N}(\boldsymbol{f}_{d}, \boldsymbol{K}_{\boldsymbol{g}, \boldsymbol{d}} + \sigma_{\boldsymbol{y}, \boldsymbol{d}}^{2} \mathbf{I}).$$
(9)

Here, we use K to refer to the Gram matrix corresponding to the kernel of the respective GP. All but the CPs factorize over both the different levels of the model as well as the different outputs.

To approximate a single deep GP, Hensman and Lawrence [11] proposed nested variational compression in which every GP in the hierarchy is handled independently. While this forces a variational approximation of all intermediate outputs of the stacked processes, it has the appealing property that it allows optimization via stochastic gradient descent [10] and the variational approximation can after training be used independently of the original training data.

A.1 Augmented Model

Nested variational compression focuses on augmenting a full GP model by introducing sets of *inducing variables u* with their *inducing inputs Z*. Those variables are assumed to be latent observations of the same functions and are thus jointly Gaussian with the observed data.

It can be written using its marginals [22] as

$$p(\hat{a}, u) = \mathcal{N}(\hat{a} \mid \mu_{a}, \Sigma_{a}) \mathcal{N}(u \mid Z, K_{uu}), \text{ with}$$

$$\mu_{a} = X + K_{au} K_{uu}^{-1} (u - Z),$$

$$\Sigma_{a} = K_{aa} - K_{au} K_{uu}^{-1} K_{ua},$$
(10)

where, after dropping some indices and explicit conditioning on X and Z for clarity, \hat{a} denotes the function values $a_d(X)$ without noise and we write the Gram matrices as $K_{au} = k_{a,d}(X, Z)$.

While the original model in (9) can be recovered exactly by marginalizing the inducing variables, considering a specific variational approximation of the joint $p(\hat{a}, u)$ gives rise to the desired lower bound in the next subsection. A central assumption of this approximation [22] is that given enough inducing variables at the correct location, they are a sufficient statistic for \hat{a} , implying conditional independence of the entries of \hat{a} given X and u. We introduce such inducing variables for every GP in the model, yielding the set $\{u_{a,d}, u_{f,d}, u_{g,d}\}_{d=1}^{D}$ of inducing variables. Note that for the CP f, we introduce one set of inducing variables $u_{f,d}$ per output f_d . These inducing variables play a crucial role in sharing information between the different outputs.

A.2 Variational Lower Bound

To derive the desired variational lower bound for the log marginal likelihood of the complete model, multiple steps are necessary. First, we will consider the innermost GPs a_d describing the alignment functions. We derive the Scalable Variational GP (SVGP), a lower bound for this model part which can be calculated efficiently and can be used for stochastic optimization, first introduced by Hensman, Fusi, and Lawrence [10]. In order to apply this bound recursively, we will both show how to propagate the uncertainty through the subsequent layers f_d and g_d and how to avoid the inter-layer cross-dependencies using another variational approximation as presented by Hensman and Lawrence [11]. While Hensman and Lawrence considered standard deep GP models, we will show how to apply their results to CPs.

The First Layer Since the inputs X are fully known, we do not need to propagate uncertainty through the GPs a_d . Instead, the uncertainty about the a_d comes from the uncertainty about the correct functions a_d and is introduced by the processes themselves. To derive a lower bound on the marginal log likelihood of a_d , we assume a variational distribution $q(u_{a,d}) \sim \mathcal{N}(m_{a,d}, S_{a,d})$ approximating $p(u_{a,d})$ and additionally assume that $q(\hat{a}_d, u_{a,d}) = p(\hat{a}_d | u_{a,d}) q(u_{a,d})$. After dropping the indices again, using Jensen's inequality we get

$$\log p(\boldsymbol{a} \mid \boldsymbol{X}) = \log \int p(\boldsymbol{a} \mid \boldsymbol{u}) p(\boldsymbol{u}) d\boldsymbol{u}$$

$$= \log \int q(\boldsymbol{u}) \frac{p(\boldsymbol{a} \mid \boldsymbol{u}) p(\boldsymbol{u})}{q(\boldsymbol{u})} d\boldsymbol{u}$$

$$\geq \int q(\boldsymbol{u}) \log \frac{p(\boldsymbol{a} \mid \boldsymbol{u}) p(\boldsymbol{u})}{q(\boldsymbol{u})} d\boldsymbol{u}$$

$$= \int \log p(\boldsymbol{a} \mid \boldsymbol{u}) q(\boldsymbol{u}) d\boldsymbol{u} - \int q(\boldsymbol{u}) \log \frac{q(\boldsymbol{u})}{p(\boldsymbol{u})} d\boldsymbol{u}$$

$$= \mathbb{E}_{q(\boldsymbol{u})} [\log p(\boldsymbol{a} \mid \boldsymbol{u})] - \text{KL}(q(\boldsymbol{u}) \parallel p(\boldsymbol{u})),$$
(11)

where $\mathbb{E}_{q(u)}[\cdot]$ denotes the expected value with respect to the distribution q(u) and $KL(\cdot \| \cdot)$ denotes the KL divergence, which can be evaluated analytically.

To bound the required expectation, we use Jensen's inequality again together with (10) which gives

$$\log p(\boldsymbol{a} \mid \boldsymbol{u}) = \log \int p(\boldsymbol{a} \mid \hat{\boldsymbol{a}}) p(\hat{\boldsymbol{a}} \mid \boldsymbol{u}) d\hat{\boldsymbol{a}}$$

= $\log \int \mathcal{N}(\boldsymbol{a} \mid \hat{\boldsymbol{a}}, \sigma_a^2 \mathbf{I}) \mathcal{N}(\hat{\boldsymbol{a}} \mid \boldsymbol{\mu_a}, \boldsymbol{\Sigma_a}) d\hat{\boldsymbol{a}}$
$$\geq \int \log \mathcal{N}(\boldsymbol{a} \mid \hat{\boldsymbol{a}}, \sigma_a^2 \mathbf{I}) \mathcal{N}(\hat{\boldsymbol{a}} \mid \boldsymbol{\mu_a}, \boldsymbol{\Sigma_a}) d\hat{\boldsymbol{a}}$$

= $\log \mathcal{N}(\boldsymbol{a} \mid \boldsymbol{\mu_a}, \sigma_a^2 \mathbf{I}) - \frac{1}{2\sigma_a^2} \operatorname{tr}(\boldsymbol{\Sigma_a}).$ (12)

We apply this bound to the expectation to get

$$\mathbb{E}_{q(\boldsymbol{u})}[\log p(\boldsymbol{a} | \boldsymbol{u})] \ge \mathbb{E}_{q(\boldsymbol{u})}[\log \mathcal{N}(\boldsymbol{a} | \boldsymbol{\mu}_{\boldsymbol{a}}, \sigma_a^2 \mathbf{I})] - \frac{1}{2\sigma_a^2} \operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{a}}), \text{ with}$$
(13)

$$\mathbb{E}_{q(\boldsymbol{u})}[\log \mathcal{N}(\boldsymbol{a} \mid \boldsymbol{\mu}_{\boldsymbol{a}}, \sigma_{a}^{2}\mathbf{I})] = \log \mathcal{N}(\boldsymbol{a} \mid \boldsymbol{K}_{\boldsymbol{au}}\boldsymbol{K}_{\boldsymbol{uu}}^{-1}\boldsymbol{m}, \sigma_{a}^{2}\mathbf{I}) + \frac{1}{2\sigma_{a}^{2}}\operatorname{tr}(\boldsymbol{K}_{\boldsymbol{au}}\boldsymbol{K}_{\boldsymbol{uu}}^{-1}\boldsymbol{S}\boldsymbol{K}_{\boldsymbol{uu}}^{-1}\boldsymbol{K}_{\boldsymbol{ua}}).$$
(14)

Resubstituting this result into (11) yields the final bound

$$\log p(\boldsymbol{a} \mid \boldsymbol{X}) \geq \log \mathcal{N}(\boldsymbol{a} \mid \boldsymbol{K_{au}} \boldsymbol{K_{uu}}^{-1} \boldsymbol{m}, \sigma_a^2 \mathbf{I}) - \operatorname{KL}(q(\boldsymbol{u}) \parallel p(\boldsymbol{u})) - \frac{1}{2\sigma_a^2} \operatorname{tr}(\boldsymbol{\Sigma_a}) - \frac{1}{2\sigma_a^2} \operatorname{tr}(\boldsymbol{K_{au}} \boldsymbol{K_{uu}}^{-1} \boldsymbol{S} \boldsymbol{K_{uu}}^{-1} \boldsymbol{K_{ua}}).$$
(15)

This bound, which depends on the hyper parameters of the kernel and likelihood $\{\theta, \sigma_a\}$ and the variational parameters $\{Z, m, S\}$, can be calculated in $\mathcal{O}(NM^2)$ time. It factorizes along the data points which enables stochastic optimization.

In order to obtain a bound on the full model, we apply the same techniques to the other processes. Since the alignment processes a_d are assumed to be independent, we have $\log p(a_1, \ldots, a_D | X) = \sum_{d=1}^{D} \log p(a_d | X)$, where every term can be approximated using the bound in (15). However, for all subsequent layers, the bound is not directly applicable, since the inputs are no longer known but instead are given by the outputs of the previous process. It is therefore necessary to propagate their uncertainty and also handle the interdependencies between the layers introduced by the latent function values a, f and g. **The Second and Third Layer** Our next goal is to derive a bound on the outputs of the second layer

$$\log p(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}) = \log \int p(\boldsymbol{f}, \boldsymbol{a}, \boldsymbol{u}_{\boldsymbol{a}} | \boldsymbol{u}_{\boldsymbol{f}}) \, \mathrm{d}\boldsymbol{a} \, \mathrm{d}\boldsymbol{u}_{\boldsymbol{a}}, \tag{16}$$

that is, an expression in which the uncertainty about the different a_d and the cross-layer dependencies on the $u_{a,d}$ are both marginalized. While on the first layer, the different a_d are conditionally independent, the second layer explicitly models the cross-covariances between the different outputs via convolutions over the shared latent processes w_r . We will therefore need to handle all of the different f_d , together denoted as f, at the same time.

We start by considering the relevant terms from (9) and apply (12) to marginalize a in

$$\log p(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}, \boldsymbol{u}_{\boldsymbol{a}}) = \log \int p(\boldsymbol{f}, \boldsymbol{a} | \boldsymbol{u}_{\boldsymbol{f}}, \boldsymbol{u}_{\boldsymbol{a}}) \, \mathrm{d}\boldsymbol{a}$$

$$\geq \log \int \tilde{p}(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}, \boldsymbol{a}) \tilde{p}(\boldsymbol{a} | \boldsymbol{u}_{\boldsymbol{a}}) \cdot \exp \left(-\frac{1}{2\sigma_{a}^{2}} \operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{a}}) - \frac{1}{2\sigma_{f}^{2}} \operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{f}})\right) \, \mathrm{d}\boldsymbol{a} \quad (17)$$

$$\geq \mathbb{E}_{\tilde{p}(\boldsymbol{a} | \boldsymbol{u}_{\boldsymbol{a}})} [\log \tilde{p}(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}, \boldsymbol{a})] - \mathbb{E}_{\tilde{p}(\boldsymbol{a} | \boldsymbol{u}_{\boldsymbol{a}})} \left[\frac{1}{2\sigma_{f}^{2}} \operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{f}})\right] - \frac{1}{2\sigma_{a}^{2}} \operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{a}}),$$

where we write $\tilde{p}(\boldsymbol{a} | \boldsymbol{u}_{\boldsymbol{a}}) = \mathcal{N}(\boldsymbol{a} | \boldsymbol{\mu}_{\boldsymbol{a}}, \sigma_a^2 \mathbf{I})$ to incorporate the Gaussian noise in the latent space. Due to our assumption that $\boldsymbol{u}_{\boldsymbol{a}}$ is a sufficient statistic for \boldsymbol{a} we choose

$$q(\boldsymbol{a} \mid \boldsymbol{u}_{\boldsymbol{a}}) = \tilde{p}(\boldsymbol{a} \mid \boldsymbol{u}_{\boldsymbol{a}}), \text{ and}$$
$$q(\boldsymbol{a}) = \int \tilde{p}(\boldsymbol{a} \mid \boldsymbol{u}_{\boldsymbol{a}}) q(\boldsymbol{u}_{\boldsymbol{a}}) d\boldsymbol{u}_{\boldsymbol{a}},$$
(18)

and use another variational approximation to marginalize u_a . This yields

$$\log p(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}) = \log \int p(\boldsymbol{f}, \boldsymbol{u}_{\boldsymbol{a}} | \boldsymbol{u}_{\boldsymbol{f}}) d\boldsymbol{u}_{\boldsymbol{a}}$$

$$= \log \int p(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}, \boldsymbol{u}_{\boldsymbol{a}}) p(\boldsymbol{u}_{\boldsymbol{a}}) d\boldsymbol{u}_{\boldsymbol{a}}$$

$$\geq \int q(\boldsymbol{u}_{\boldsymbol{a}}) \log \frac{p(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}, \boldsymbol{u}_{\boldsymbol{a}}) p(\boldsymbol{u}_{\boldsymbol{a}})}{q(\boldsymbol{u}_{\boldsymbol{a}})} d\boldsymbol{u}_{\boldsymbol{a}}$$

$$= \mathbb{E}_{q(\boldsymbol{u}_{\boldsymbol{a}})}[\log p(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{a}}, \boldsymbol{u}_{\boldsymbol{f}})] - \mathrm{KL}(q(\boldsymbol{u}_{\boldsymbol{a}}) \parallel p(\boldsymbol{u}_{\boldsymbol{a}}))$$

$$\geq \mathbb{E}_{q(\boldsymbol{u}_{\boldsymbol{a}})}[\mathbb{E}_{\tilde{p}(\boldsymbol{a} | \boldsymbol{u}_{\boldsymbol{a}})}[\log \tilde{p}(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}, \boldsymbol{a})]] - \mathrm{KL}(q(\boldsymbol{u}_{\boldsymbol{a}}) \parallel p(\boldsymbol{u}_{\boldsymbol{a}}))$$

$$- \frac{1}{2\sigma_{a}^{2}} \operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{a}}) - \mathbb{E}_{q(\boldsymbol{u}_{\boldsymbol{a}})}\left[\mathbb{E}_{\tilde{p}(\boldsymbol{a} | \boldsymbol{u}_{\boldsymbol{a}})\left[\frac{1}{2\sigma_{f}^{2}} \operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{f}})\right]\right]$$

$$\geq \mathbb{E}_{q(\boldsymbol{a})}[\log \tilde{p}(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}, \boldsymbol{a})], - \mathrm{KL}(q(\boldsymbol{u}_{\boldsymbol{a}}) \parallel p(\boldsymbol{u}_{\boldsymbol{a}}))$$

$$- \frac{1}{2\sigma_{a}^{2}} \operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{a}}) - \frac{1}{2\sigma_{f}^{2}} \mathbb{E}_{q(\boldsymbol{a})}[\operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{f}})],$$
(19)

where we apply Fubini's theorem to exchange the order of integration in the expected values. The expectations with respect to q(a) involve expectations of kernel matrices, also called Ψ -statistics, in the same way as in [8] and are given by

$$\psi_{f} = \mathbb{E}_{q(a)}[tr(K_{ff})],$$

$$\Psi_{f} = \mathbb{E}_{q(a)}[K_{fu}],$$

$$\Phi_{f} = \mathbb{E}_{q(a)}[K_{uf}K_{fu}].$$
(20)

These Ψ -statistics can be computed analytically for multiple kernels, including the squared exponential kernel. In Appendix A.3 we show closed-form solutions for these Ψ -statistics for the implicit kernel defined in the CP layer. To obtain the final formulation of the desired bound for $\log p(f | u_f)$ we substitute (20) into (19) and get the analytically tractable bound

$$\log p(\boldsymbol{f} | \boldsymbol{u}_{\boldsymbol{f}}) \geq \log \mathcal{N} \left(\boldsymbol{f} \left| \boldsymbol{\Psi}_{\boldsymbol{f}} \boldsymbol{K}_{\boldsymbol{u}_{\boldsymbol{f}}\boldsymbol{u}_{\boldsymbol{f}}}^{-1} \boldsymbol{m}_{\boldsymbol{f}}, \sigma_{\boldsymbol{f}}^{2} \mathbf{I} \right) - \mathrm{KL}(\mathbf{q}(\boldsymbol{u}_{\boldsymbol{a}}) \parallel \mathbf{p}(\boldsymbol{u}_{\boldsymbol{a}})) - \frac{1}{2\sigma_{a}^{2}} \operatorname{tr}(\boldsymbol{\Sigma}_{\boldsymbol{a}}) - \frac{1}{2\sigma_{f}^{2}} \left(\psi_{\boldsymbol{f}} - \operatorname{tr} \left(\boldsymbol{\Psi}_{\boldsymbol{f}} \boldsymbol{K}_{\boldsymbol{u}_{\boldsymbol{f}}\boldsymbol{u}_{\boldsymbol{f}}}^{-1} \right) \right) - \frac{1}{2\sigma_{f}^{2}} \operatorname{tr} \left(\left(\boldsymbol{\Phi}_{\boldsymbol{f}} - \boldsymbol{\Psi}_{\boldsymbol{f}}^{\mathsf{T}} \boldsymbol{\Psi}_{\boldsymbol{f}} \right) \boldsymbol{K}_{\boldsymbol{u}_{\boldsymbol{f}}\boldsymbol{u}_{\boldsymbol{f}}}^{-1} \left(\boldsymbol{m}_{\boldsymbol{f}} \boldsymbol{m}_{\boldsymbol{f}}^{\mathsf{T}} + \boldsymbol{S}_{\boldsymbol{f}} \right) \boldsymbol{K}_{\boldsymbol{u}_{\boldsymbol{f}}\boldsymbol{u}_{\boldsymbol{f}}}^{-1} \right) \right)$$

$$(21)$$

The uncertainties in the first layer have been propagated variationally to the second layer. Besides the regularization terms, $f \mid u_f$ is a Gaussian distribution. Because of their cross dependencies, the different outputs f_d are considered in a common bound and do not factorize along dimensions. The third layer warpings g_d however are conditionally independent given f and can therefore be considered separately. In order to derive a bound for $\log p(y \mid u_g)$ we apply the same steps as described above, resulting in the final bound, which factorizes along the data, allowing for stochastic optimization methods:

$$\log p(\boldsymbol{y} \mid \boldsymbol{X}) \geq \sum_{d=1}^{D} \log \mathcal{N} \left(\boldsymbol{y}_{d} \mid \boldsymbol{\Psi}_{\boldsymbol{g}, \boldsymbol{d}} \boldsymbol{K}_{\boldsymbol{u}_{\boldsymbol{g}, \boldsymbol{d}} \boldsymbol{u}_{\boldsymbol{g}, \boldsymbol{d}}}^{-1} \boldsymbol{m}_{\boldsymbol{g}, \boldsymbol{d}}, \sigma_{\boldsymbol{y}, \boldsymbol{d}}^{2} \mathbf{I} \right) - \sum_{d=1}^{D} \frac{1}{2\sigma_{\boldsymbol{x}, \boldsymbol{d}}^{2}} \operatorname{tr} \left(\boldsymbol{\Sigma}_{\boldsymbol{a}, \boldsymbol{d}} \right) \\ - \frac{1}{2\sigma_{f}^{2}} \left(\psi_{f} - \operatorname{tr} \left(\boldsymbol{\Phi}_{f} \boldsymbol{K}_{\boldsymbol{u}_{f} \boldsymbol{u}_{f}}^{-1} \right) \right) - \sum_{d=1}^{D} \frac{1}{2\sigma_{\boldsymbol{y}, \boldsymbol{d}}^{2}} \left(\psi_{\boldsymbol{g}, \boldsymbol{d}} - \operatorname{tr} \left(\boldsymbol{\Phi}_{\boldsymbol{g}, \boldsymbol{d}} \boldsymbol{K}_{\boldsymbol{u}_{\boldsymbol{g}, \boldsymbol{d}} \boldsymbol{u}_{\boldsymbol{g}, \boldsymbol{d}}}^{-1} \right) \right) \\ - \sum_{d=1}^{D} \operatorname{KL}(q(\boldsymbol{u}_{\boldsymbol{a}, \boldsymbol{d}}) \parallel p(\boldsymbol{u}_{\boldsymbol{a}, \boldsymbol{d}})) - \operatorname{KL}(q(\boldsymbol{u}_{f}) \parallel p(\boldsymbol{u}_{f})) - \sum_{d=1}^{D} \operatorname{KL}(q(\boldsymbol{u}_{\boldsymbol{y}, \boldsymbol{d}}) \parallel p(\boldsymbol{u}_{\boldsymbol{y}, \boldsymbol{d}})) \quad (22) \\ - \frac{1}{2\sigma_{f}^{2}} \operatorname{tr} \left(\left(\boldsymbol{\Phi}_{f} - \boldsymbol{\Psi}_{f}^{\mathsf{T}} \boldsymbol{\Psi}_{f} \right) \boldsymbol{K}_{\boldsymbol{u}_{f} \boldsymbol{u}_{f}}^{-1} \left(\boldsymbol{m}_{f} \boldsymbol{m}_{f}^{\mathsf{T}} + \boldsymbol{S}_{f} \right) \boldsymbol{K}_{\boldsymbol{u}_{f} \boldsymbol{u}_{f}}^{-1} \right) \\ - \sum_{d=1}^{D} \frac{1}{2\sigma_{\boldsymbol{y}, d}^{2}} \operatorname{tr} \left(\left(\boldsymbol{\Phi}_{\boldsymbol{g}, \boldsymbol{d}} - \boldsymbol{\Psi}_{\boldsymbol{g}, \boldsymbol{d}}^{\mathsf{T}} \boldsymbol{\Psi}_{\boldsymbol{g}, \boldsymbol{d}} \right) \boldsymbol{K}_{\boldsymbol{u}_{\boldsymbol{g}, \boldsymbol{u}, \boldsymbol{g}, \boldsymbol{d}}^{-1} \left(\boldsymbol{m}_{\boldsymbol{g}, \boldsymbol{d}} \boldsymbol{m}_{\boldsymbol{g}, \boldsymbol{d}}^{\mathsf{T}} + \boldsymbol{S}_{\boldsymbol{g}, \boldsymbol{d}} \right) \boldsymbol{K}_{\boldsymbol{u}_{\boldsymbol{g}, \boldsymbol{d}} \boldsymbol{u}_{\boldsymbol{g}, \boldsymbol{d}}}^{-1} \right) \right)$$

A.3 Convolution Kernel Expectations

In Section 2 we assumed the latent processes w_r to be white noise processes and the smoothing kernel functions $T_{d,r}$ to be squared exponential kernels, leading to an explicit closed form formulation for the covariance between outputs shown in (3). In this section, we derive the Ψ -statistics for this generalized squared exponential kernel needed to evaluate (22).

The uncertainty about the first layer is captured by the variational distribution of the latent alignments a given by $q(a) \sim \mathcal{N}(\mu_a, \Sigma_a)$, with $a = (a_1, \ldots, a_d)$. Every aligned point in a corresponds to one output of f and ultimately to one of the y_d . Since the closed form of the multi output kernel depends on the choice of outputs, we will use the notation $\hat{f}(a_n)$ to denote $f_d(a_n)$ such that a_n is associated with output d.

For notational simplicity, we only consider the case of one single latent process w_r . Since the latent processes are independent, the results can easily be generalized to multiple processes. Then, ψ_f is given by

$$\psi_{f} = \mathbb{E}_{q(\boldsymbol{a})}[\operatorname{tr}(\boldsymbol{K}_{ff})]$$

$$= \sum_{n=1}^{N} \mathbb{E}_{q(\boldsymbol{a}_{n})} \left[\operatorname{cov} \left[\hat{f}(\boldsymbol{a}_{n}), \hat{f}(\boldsymbol{a}_{n}) \right] \right]$$

$$= \sum_{n=1}^{N} \int \operatorname{cov} \left[\hat{f}(\boldsymbol{a}_{n}), \hat{f}(\boldsymbol{a}_{n}) \right] q(\boldsymbol{a}_{n}) \, \mathrm{d}\boldsymbol{a}_{n}$$

$$= \sum_{n=1}^{N} \hat{\sigma}_{nn}^{2}.$$
(23)

Similar to the notation $\hat{f}(\cdot)$, we use the notation $\hat{\sigma}_{nn'}$ to mean the variance term associated with the covariance function $\cos[\hat{f}(a_n), \hat{f}(a_{n'})]$. The expectation $\Psi_f = \mathbb{E}_{q(a)}[K_{fu}]$ connecting the alignments and the pseudo inputs is given by

$$\Psi_{f} = \mathbb{E}_{q(\boldsymbol{a})}[\boldsymbol{K}_{f\boldsymbol{u}}], \text{ with}$$

$$(\Psi_{f})_{ni} = \int \operatorname{cov}\left[\hat{f}(\boldsymbol{a}_{n}), \hat{f}(\boldsymbol{Z}_{i})\right] q(\boldsymbol{a}_{n}) \,\mathrm{d}\boldsymbol{a}_{n}$$

$$= \hat{\sigma}_{ni}^{2} \sqrt{\frac{(\boldsymbol{\Sigma}_{\boldsymbol{a}})_{nn}^{-1}}{\hat{\ell}_{ni} + (\boldsymbol{\Sigma}_{\boldsymbol{a}})_{nn}^{-1}}} \cdot \exp\left(-\frac{1}{2} \frac{(\boldsymbol{\Sigma}_{\boldsymbol{a}})_{nn}^{-1} \hat{\ell}_{ni}}{(\boldsymbol{\Sigma}_{\boldsymbol{a}})_{nn}^{-1} + \hat{\ell}_{ni}} \left((\boldsymbol{\mu}_{\boldsymbol{a}})_{n} - \boldsymbol{Z}_{i}\right)^{2}\right)$$
(24)

where $\hat{\ell}_{ni}$ is the combined length scale corresponding to the same kernel as $\hat{\sigma}_{ni}$. Lastly, $\Phi_f = \mathbb{E}_{q(a)}[K_{uf}K_{fu}]$ connects alignments and pairs of pseudo inputs with the closed form

$$\begin{split} \boldsymbol{\Phi}_{\boldsymbol{f}} &= \mathbb{E}_{q(\boldsymbol{a})}[\boldsymbol{K}_{\boldsymbol{u}\boldsymbol{f}}\boldsymbol{K}_{\boldsymbol{f}\boldsymbol{u}}], \text{ with} \\ (\boldsymbol{\Phi}_{\boldsymbol{f}})_{ij} &= \sum_{n=1}^{N} \int \operatorname{cov}\left[\hat{f}(\boldsymbol{a}_{\boldsymbol{n}}), \hat{f}(\boldsymbol{Z}_{\boldsymbol{i}})\right] \cdot \operatorname{cov}\left[\hat{f}(\boldsymbol{a}_{\boldsymbol{n}}), \hat{f}(\boldsymbol{Z}_{\boldsymbol{j}})\right] q(\boldsymbol{a}_{\boldsymbol{n}}) \, \mathrm{d}\boldsymbol{a}_{\boldsymbol{n}} \\ &= \sum_{n=1}^{N} \hat{\sigma}_{ni}^{2} \hat{\sigma}_{nj}^{2} \sqrt{\frac{(\boldsymbol{\Sigma}_{\boldsymbol{a}})_{nn}^{-1}}{\hat{\ell}_{ni} + \hat{\ell}_{nj} + (\boldsymbol{\Sigma}_{\boldsymbol{a}})_{nn}^{-1}}} \cdot \exp\left(-\frac{1}{2} \frac{\hat{\ell}_{ni} \hat{\ell}_{nj}}{\hat{\ell}_{ni} + \hat{\ell}_{nj}} (\boldsymbol{Z}_{\boldsymbol{i}} - \boldsymbol{Z}_{\boldsymbol{j}})^{2} \right)^{2} \\ &- \frac{1}{2} \frac{(\boldsymbol{\Sigma}_{\boldsymbol{a}})_{nn}^{-1} + \hat{\ell}_{ni} + \hat{\ell}_{nj}}{(\boldsymbol{\Sigma}_{\boldsymbol{a}})_{nn}^{-1} + \hat{\ell}_{ni} + \hat{\ell}_{nj}} \cdot \left((\boldsymbol{\mu}_{\boldsymbol{a}})_{n} - \frac{\hat{\ell}_{ni} \boldsymbol{Z}_{\boldsymbol{i}} + \hat{\ell}_{nj} \boldsymbol{Z}_{\boldsymbol{j}}}{\hat{\ell}_{ni} + \hat{\ell}_{nj}}\right)^{2} \right). \end{split}$$

Note that the Ψ -statistics factorize along the data and we only need to consider the diagonal entries of Σ_a . If all the data belong to the same output, the Ψ -statistics of the standard squared exponential kernel can be recovered as a special case. It is used to propagate the uncertainties through the output-specific warpings g.

A.4 Approximative Predictions

Using the variational lower bound in (5), our model can be fitted to data, resulting in appropriate choices of the kernel hyper parameters and variational parameters. Now assume we want to predict approximate function values $g_{d,\star}$ for previously unseen points $X_{d,\star}$ associated with output d, which are given by $g_{d,\star} = g_d(f_d(a_d(X_{d,\star})))$.

Because of the conditional independence assumptions in the model, other outputs $d' \neq d$ only have to be considered in the shared layer f. In this shared layer, the belief about the different outputs and the shared information and is captured by the variational distribution $q(u_f)$. Given $q(u_f)$, the different outputs are conditionally independent of one another and thus, predictions for a single dimension in our model are equivalent to predictions in a single deep GP with nested variational compression as presented by Hensman and Lawrence [11].

B Joint models for wind experiment

In the following, we show plots with joint predictions for the models discussed in Section 4.2. Similar to Section 4.1, we trained a standard GP in Figure 6, a multi-output GP in Figure 7, a deep GP in Figure 8 and our model in Figure 9. All models were trained until convergence and multiple runs result in very similar models. For all models we used RBF kernels or dependent RBF kernels where applicable.

Each plot shows the data in gray and two mean predictions and uncertainty bands. The first violet uncertainty band is the result of the variational approximation of the respective model. The second green or blue posterior is obtained via sampling. For both the GP and MO-GP, we used the SVGP approximation [12] and since the models are shallow, the approximation is almost exact.

Figure 8 showcases the difficulty of training a deep GP model and the shortcomings of the nested variational compression. The violet variational approximation is used for training and approximates

the data comparatively well. As discussed above, the deep GP cannot share information, so the test sets cannot be predicted. However, as discussed in more detail in [12], the approximation tends to underestimate uncertainties when propagating them through the different layers and because of this, uncertainties obtained via sampling tend to vary considerably more. Because during model selection sample performance does not matter, the true posterior can be (and in this case is) considerably different.

Our approach in principle has the same problem as the deep GP. However, because of the strong interpretability of the different parts of the hierarchy, uncertainties within the model are never placed arbitrarily and because of this, the variational posteriors and true posteriors look much more similar. They tend to disagree in places when there is high uncertainty about the alignment.



Figure 6: GP



Figure 7: MO-GP



Figure 8: DGP



Figure 9: AMO-GP (Ours)