STRUCTURED JOINT ALEATORIC AND EPISTEMIC UN CERTAINTY FOR HIGH DIMENSIONAL OUTPUT SPACES

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023 024

025

Paper under double-blind review

Abstract

Uncertainty estimation plays a vital role in enhancing the reliability of deep learning model predictions, especially in scenarios with high-dimensional output spaces. This paper addresses the dual nature of uncertainty — aleatoric and epistemic — focusing on their joint integration in high-dimensional regression tasks. We introduce an approach to approximate joint uncertainty using a low-rank plus diagonal covariance matrix, which preserves essential output correlations while mitigating the computational complexity associated with full covariance matrices. Specifically, our method reduces memory usage and enhances sampling efficiency and log-likelihood calculations. Simultaneously, our representation matches the true posterior better than factorized joint distributions, offering a clear advancement in reliability and explainability for deep learning model predictions. Furthermore, we empirically show that our method can efficiently enhance out of distribution detection in specific applications.

1 INTRODUCTION

In the realm of deep learning, uncertainty estimation plays a pivotal role in enhancing the reliability of model predictions. This paper delves into the domain of uncertainty estimation, specifically focusing on regression tasks in high-dimensional output spaces.

In these scenarios, model predictions exhibit two types of uncertainty: aleatoric and epistemic
 Kendall & Gal (2017). Heteroscedastic aleatoric or data uncertainty can be modeled as an inherent
 component of the model output. There, the model output consists of the parametrization of an assumed
 distribution, such as a Gaussian distribution in the case of regression or a categorical distribution in
 the case of classification. This distribution is learned through minimizing its negative log-likelihood.
 In contrast, due to its complexity in deep neural networks, epistemic or model uncertainty is typically
 approximated by sampling from an proxy distribution of models Hüllermeier & Waegeman (2021).

Combining both in a single model usually results in a so-called second-order distribution Bengs et al.
 (2023). On the one hand, it consists of a distribution over model weights capturing epistemic uncertainty. On the other hand, it models a distribution over plausible predictions representing aleatoric uncertainty. Sampling from the model weights and performing a transformation (forward pass) of the input data results in another distribution representing the aleatoric uncertainty. The shape of this second-order distribution limits further analysis, as it is difficult to visualize, and it does not allow for calculation of the marginal likelihood of a sample. Therefore, the second-order distribution is typically marginalized and approximated by a single distribution, representing the joint uncertainty.

044 Traditionally, these uncertainties have been jointly modeled without considering correlations between outputs (e.g. pixels), assuming independent factorized univariate Gaussian distributions Kendall 046 & Gal (2017). However, neglecting correlations can limit the comprehensive understanding of 047 uncertainties, especially in scenarios where dependencies between model outputs exist, such as in 048 pixel-wise semantic segmentation Monteiro et al. (2020), pixel-wise regression tasks, such as optical flow estimation or image in-painting, or graph node regression. Figure 1 illustrates the increased representational power of full covariance matrices (right) compared to merely diagonal ones (left). 051 In both cases, samples from the weight space lead to multiple predictions consisting of mean and covariance each. The expected covariance is used for calculating the aleatoric component Σ^a , whereas 052 the covariance of the means is used to calculate the epistemic component Σ^e . The sum of both results in the joint covariance matrix $\Sigma^a + \Sigma^e = \Sigma$.



Figure 1: Visualization of covariance matrices for the 2D case. Three samples with corresponding
means and covariances are depicted (light blue). The columns show the inferred aleatoric (green),
epistemic (blue) and joint uncertainty (red), respectively. On the left, the covariance matrices are
purely diagonal, limiting their representational power. To the right, the same matrices are depicted
with non-diagonal values kept, allowing them to capture the overall uncertainty in greater detail.

Yet, incorporating correlations in high-dimensional output spaces poses a significant challenge, given that the number of correlations between output dimensions scales quadratically in terms of memory complexity $\mathcal{O}(S^2)$ with the total number of outputs S. This leads to large covariance matrices, requiring considerable storage space and making calculations computationally infeasible. This renders many downstream operations like sampling from a normal distribution parameterized by these covariance matrices, which involves Cholesky decomposition $\mathcal{O}(S^3)$, computing the loglikelihood of samples with matrix inversion $\mathcal{O}(S^3)$ and determinant computation (e.g. $\mathcal{O}(S^3)$ with lower-upper (LU) Decomposition) practically impossible.

In summary, an efficient representation of the joint uncertainty containing aleatoric and epistemic uncertainty without neglecting covariances for high-dimensional output spaces has not yet been explored.

082 **Related Work** To estimate epistemic uncertainty, various Bayesian frameworks have been devel-083 oped, including methods like stochastic variational inference Blundell et al. (2015), Monte Carlo 084 dropout Gal & Ghahramani (2016), deep ensembles Lakshminarayanan et al. (2017), stochastic 085 weight averaging Maddox et al. (2019), or Laplace approximation Daxberger et al. (2021). The modeling of heteroscedastic aleatoric uncertainty has been well-established for some time Nix & Weigend 087 (1994); Skafte et al. (2019); Stirn & Knowles (2020); Seitzer et al. (2022). Building upon these works, 880 others have unified epistemic and aleatoric uncertainty in a single model Kendall & Gal (2017); Depeweg et al. (2018); Stirn et al. (2023); Immer et al. (2024). However, all aforementioned methods 089 either evaluate their method only for prediction tasks with a single output value or approximate the 090 marginalized likelihood as a factorized Gaussian, disregarding inter-pixel correlations. 091

092 Covariances for uncertainty estimation have been modeled in various applications, including localization Russell & Reale (2021), human pose estimation Gundavarapu et al. (2019), pixel regression Dorta et al. (2018a;b); Duff et al. (2023), multi class predictions Willette et al. (2021), and segmentation 094 Monteiro et al. (2020). Some approaches that predict full covariance matrices are limited to low 095 dimensional model output spaces Russell & Reale (2021); Gundavarapu et al. (2019). Approaches 096 for handling high-dimensional output spaces typically sparsify the covariance matrix. However, some of these approaches can only model uncertainty in the local neighborhood using a band Cholesky 098 parametrization Dorta et al. (2018a;b); Duff et al. (2023). Some works Salinas et al. (2019); Monteiro 099 et al. (2020); Willette et al. (2021); Nussbaum et al. (2022) use a low-rank plus diagonal (LR+D) 100 parametrization, which is capable of capturing global correlations. Nehme et al. (2024); Yair et al. 101 (2024) learn the low-rank factors of aleatoric uncertainty directly without adding a diagonal and 102 create a rank-deficient semi-definite covariance matrix. This may be sufficient for both sampling and 103 analysis, but it does not provide the positive definiteness required for calculating the log-likelihood. 104 Importantly, all these sparse solutions merely focus on aleatoric uncertainty. Zepf et al. (2023) 105 combine aleatoric and epistemic uncertainty with a LR+D representation. However, by partially using the Maximum a posteriori (MAP) solution as a further approximation, they do not account for the 106 influence of the model uncertainty on the estimation of the aleatoric uncertainty, leading overall to a 107 worse uncertainty estimate. Furthermore, they do not resolve the second-order distribution to provide



Figure 2: Construction of our LR+D matrix. A network predicts values μ_w^P , P_w^P , and D_w for two exemplarily sampled weights w_i , respectively (green and blue). By averaging and stacking these values in a specific manner, we build the diagonal D and the low-rank matrix P as parts of our LR+D representation of Σ . See Section 2 for an in-depth explanation.

a joint representation suitable for further analysis, such as log likelihood calculation, and its usage is
 limited to consecutive sampling.

In conclusion, while significant advancements have been made in modeling covariances for uncertainty estimation, the existing approaches suffer from limitations such as local sparsification, inadequate joint representations, and neglect of weight space uncertainty, indicating a need for further research to develop more comprehensive and globally accurate uncertainty estimation methods.

127 **Contribution** This work is the first to efficiently combine both aleatoric and epistemic uncertainties 128 within any high-dimensional sparse joint representation, leveraging the LR+D framework. Unlike 129 existing approaches that approximate the second-order distribution with a factorized normal distri-130 bution, neglecting correlations between outputs, our method maintains crucial correlations between 131 outputs while avoiding the heavy space and time requirements of a full covariance matrix even for high-dimensional output spaces. We showcase the superior representational power of our approach 132 on multiple high-dimensional regression tasks, e.g. inpainting on the MNIST dataset, colorization of 133 grey-scale versions of the CelebA celebrity faces dataset Liu et al. (2015), and optical flow estimation 134 on the Flying Chairs dataset Dosovitskiy et al. (2015). 135

136 137

138

2 Method

In this paper, we consider supervised learning tasks with high dimensional output spaces $y \in \mathbb{R}^S$, with S denoting the number of output units, e.g. pixels times the number of output channels. We aim to approximate the posterior distribution p(W|X, Y) over model weights W given input-output data pairs (X, Y). As computing p(W|X, Y) directly is generally infeasible for neural networks, we approximate it using Bayesian methods like Monte Carlo Dropout (MCD), Stochastic Variational Inference (SVI), Deep Ensemble (DE) to provide a proxy distribution $q_{\theta}^*(W)$, parametrized by θ .

To represent the joint uncertainty, we decompose the posterior predictive distribution p(y|x, X, Y)145 of an unseen input-output pair (x, y) into two terms: p(y|x, W), representing the likelihood of the 146 output given the input x and the network weights W, and the posterior distribution of weights given 147 the data p(W|X, Y). We model p(y|x, W) as a multivariate Gaussian distribution p(y|x, W) =148 $\mathcal{N}(\mu_W(x), \Sigma_W(x))$, where we keep the spatial complexity of the covariance matrix $\Sigma_W(x)$ low 149 by constructing it in LR+D form. That is, we formulate it as a sum of small matrices, $\Sigma_W(x) =$ 150 $D_W(x) + P_W(x)P_W^{\mathsf{T}}(x)$, with D_W denoting a diagonal matrix of shape $S \times S$ and P_W a tall matrix 151 of shape $S \times R^W$. We choose a rank R^W much lower than the number of outputs $R^W \ll S$, such that 152 only the most important directions of the aleatoric covariance are covered. We further enforce D_W 153 to contain strictly positive diagonal entries and since $P_W P_W^T$ is always symmetric, Σ_W is always 154 symmetric positive definite by construction and thus a valid covariance matrix. The ultimate goal 155 of this work is to calculate an efficient vet representative representation of the posterior predictive 156 distribution p(y|x, X, Y).

157

158 2.1 MODELING THE JOINT UNCERTAINTY

159

We start modeling the parameters of the posterior predictive distribution consisting of mean and covariance by using Monte Carlo integration to approximate the expected model output $\mathbb{E}[y|x, \mathbf{X}, \mathbf{Y}] \approx \mu(x)$. The empirical mean is given as $\mu(x) = \frac{1}{T} \sum_{i}^{T} \mu_{w_i}(x)$, where T repreioint uncertainty

sents the number of weight samples drawn from $w_i \sim q_{\theta}^{\star}(W)$. The joint covariance matrix can be split into epistemic and aleatoric uncertainty using the law of total variance as

$$\underbrace{\operatorname{Cov}\left[y|x, \boldsymbol{X}, \boldsymbol{Y}\right]}_{\sum_{i=1}^{\widetilde{e}} (x)} \approx \underbrace{\operatorname{Cov}_{q_{\theta}^{*}(W)}\left[\mu_{W}(x)\right]}_{\sum_{i=1}^{\widetilde{e}} (x)} + \underbrace{\operatorname{E}_{q_{\theta}^{*}(W)}\left[\Sigma_{W}(x)\right]}_{\sum_{i=1}^{\widetilde{e}} (x)}.$$
(1)

166 167 168

165

This suggests that the mean of covariance matrices across forward pass samples captures aleatoric uncertainty, whereas the covariance of the means represents epistemic uncertainty. We provide a complete derivation of equation 1 in the supplement D.4.

epistemic uncertainty

¹⁷² Our objective is to represent the joint uncertainty $\Sigma(x)$ in LR+D form as the sum of aleatoric and ¹⁷³ epistemic uncertainties,

183

184 185

186

191 192

196 197

$$D + PP^{\mathsf{T}} = (D^{e} + P^{e}P^{e\mathsf{T}}) + (D^{a} + P^{a}P^{a\mathsf{T}}), \tag{2}$$

aleatoric uncertainty

where D^a , D^e , and D are diagonal matrices and P^a , P^e , and P low-rank matrices representing aleatoric, epistemic, and joint uncertainties, respectively. Then, $D = D^e + D^a$ and $P = [P^a P^e]$, where [] denotes columnwise block concatenation. This expression allows us to conveniently represent both aleatoric and epistemic uncertainties in LR+D form, simplifying further analysis and computation. Figure 2 provides an intuitive illustration about the construction of our LR+D matrix components. Starting with Σ^e , we describe in detail the individual components of our LR+D representations in the following sections.

2.2 EPISTEMIC UNCERTAINTY

The epistemic uncertainty is estimated through the distribution over weights. To derive its covariance, we employ empirical sampling from the proxy distribution over model weights as follows:

$$\Sigma^{e}(x) = \frac{1}{T-1} \sum_{i}^{T} \left(\mu_{w_{i}}(x) - \mu(x) \right) \left(\mu_{w_{i}}(x) - \mu(x) \right)^{\mathsf{T}} \quad w_{i} \sim q_{\theta}^{*}(W)$$
(3)

Our objective is to avoid the full covariance matrix and instead seek a representation in LR+D form.

Naive Representation To bring the approximated epistemic covariance matrix into LR+D form, we set the diagonal $D^e(x)$ to zero and rewrite the covariance matrix as $\Sigma^e(x) = P^e(x)P^e(x)^{\mathsf{T}}$, where $P^e(x) \in \mathbb{R}^{S \times R^e}$ has $R^e = T$ columns and is defined as

$$P^{e}(x) = \frac{1}{\sqrt{T-1}} \left[\mu_{w_{1}}(x) - \mu(x) \quad \dots \quad \mu_{w_{T}}(x) - \mu(x) \right].$$
(4)

In high dimensional scenarios, the number of samples will often be much lower than the number of outputs $T \ll S$ rendering Σ^e singular and therefore non-invertible. Computing enough samples for a full rank estimate is usually prohibitive with regard to time and space complexity. In general, one can expect more accurate results from increased sample sizes. However, in this naive representation, larger sample sizes also result in quadratic scaling of computational complexity. Hence, we suggest further approximations to cope with moderately high sample sizes.

205 Truncated Singular Value Decomposition Approximation Assuming that samples are often 206 correlated and exhibit dominant directions of variance, we propose to reduce the dimensionality of 207 $P^{e}(x)$ with truncated Singular Value Decomposition (SVD). Keeping only the most informative 208 columns of $P^{e}(x)$ will improve the efficiency of further computations without losing much informa-209 tion. However, the calculation of SVD comes with its own computational complexity that has to be taken into account. As we use the same nomenclature for further parts of the model, we omit the 210 superscript ()^e and the dependency on (x) for this section. Specifically, we decompose the matrix 211 P as $P^{\intercal} = U\Psi V^{\intercal}$, where U and V are orthogonal matrices, and Ψ is a diagonal matrix containing 212 the singular values in non-decreasing order $\Psi_{1,1} \leq ... \leq \Psi_{S,S}$. Subsequently, we define the matrix 213 $\tilde{P} = V\Psi$ and rewrite the matrix product as $\Sigma = PP^{\mathsf{T}} = \tilde{P}\tilde{P}^{\mathsf{T}}$. To reduce dimensionality, we discard 214 the smallest singular values and their associated columns in V. However, we keep the univariate 215 variance parts of these dropped columns by transferring them to a new diagonal matrix \hat{D} . Hence, the approximated matrix $\hat{\Sigma} = \hat{D} + \hat{P}\hat{P}^{\mathsf{T}}$ keeps all independent variance and the most important covariances of Σ . If we keep the \hat{R} largest singular values, the components of $\hat{\Sigma}$ are

$$\hat{P} = \begin{bmatrix} V_{R-\hat{R}} \cdot \Psi_{R-\hat{R},R-\hat{R}} & \dots & V_R \cdot \Psi_{R,R} \end{bmatrix}$$
(5)

and

$$\hat{D}_{ii} = \sum_{j=1}^{R-\bar{R}-1} V_{ij}^2 \cdot \Psi_{j,j}^2.$$
(6)

The number of columns to keep has to be chosen empirically. The aforementioned approach enables us to effectively represent epistemic uncertainty in the LR+D form.

2.3 Aleatoric Uncertainty

Similar to epistemic uncertainty, the covariance matrix capturing aleatoric uncertainty $\Sigma^{a}(x)$ can be approximated through empirical sampling. We calculate the empirical mean of covariance matrix estimations over all sampled model weights via

$$\Sigma^{a}(x) = \frac{1}{T} \sum_{i}^{T} \Sigma_{w_{i}}(x) \quad w_{i} \sim q_{\theta}^{*}(W).$$

$$\tag{7}$$

We here again intend to represent $\Sigma^{a}(x)$ in LR+D form.

Naive Representation To rewrite the covariance matrix containing the aleatoric uncertainty in LR+D representation, we reformulate $\Sigma^a(x) = D^a(x) + P^a(x)P^a(x)^{\mathsf{T}}$ using

$$D^{a}(x) = \frac{1}{T} \sum_{i}^{T} D_{w_{i}}(x)$$
(8)

$$P^{a}(x) = \frac{1}{\sqrt{T}} \begin{bmatrix} P_{w_{1}}(x) & \dots & P_{w_{T}}(x) \end{bmatrix}.$$
(9)

This yields a $P^a \in \mathbb{R}^{S \times (T \cdot R^W)}$ with $T \cdot R^W$ columns. Although $T \cdot R^W$ generally remains far below S, P^a can still become fairly large as the number of drawn samples increases. Thus, we further reduce the number of columns of Σ^a as for the epistemic case.

Truncated SVD Approximation Like in the previous Subsection 2.2, we reduce the dimensionality of P^a via SVD. This leads for \hat{P}^a to the same equation as 5. For the diagonal matrix, we need to incorporate both, the average calculated in 8 and the removed variance by the SVD given by 6, so the resulting diagonal is given by:

$$\hat{D}^{a}{}_{ii} = D^{a}_{ii} + \sum_{j=1}^{R-\hat{R}-1} V^{2}_{ij} \cdot \Psi^{2}_{j,j}.$$
(10)

EXPERIMENTS

260 3.1 EXPERIMENTAL SETUP

Proposed Method We empirically evaluate our method of joint aleatoric and epistemic uncertainty modeling using our LR+D representation in several experiments. In all experiments, we use variants of the U-Net Ronneberger et al. (2015) architecture. We adapt the U-Net for Bayesian inference by adding dropout, which we use for MCD Gal & Ghahramani (2016) and DE Lakshminarayanan et al. (2017) or by using variational convolutional layer for SVI Blundell et al. (2015) to estimate a distribution over model weights which estimates epistemic uncertainty. However, we note that our approach is compatible with any Bayesian method suitable for large model outputs. We use a single model with multiple outputs for the mean and the covariance prediction parts. To train this model, we gradually we train the mean separately and gradually change increase the weight on the log likelihood loss for the covariance parameters. For architectural details, please refer to our code in the repository. Datasets and Tasks We evaluate our method in different settings on the MNIST, CelebA, and
 Flying Chairs datasets for the tasks of inpainting, colorization, and optical flow.

We train a reconstruction model to inpaint distorted handwritten digits from the MNIST dataset. For the inpainting task, we mask out 5/7 of the image area. We use the official test set and split the training set into 50,000 train and 10,000 validation images. Figure 3.2 (bottom) shows reconstruction results for a single digit.

To evaluate performance on optical flow estimation, we use the Flying Chairs Dosovitskiy et al. (2015)
dataset. This dataset is resized to 192 x 256 and split into 18,297/2,287/2,288 training/validation/test
images. We provide visualizations of the predictions as part of the supplement.

To evaluate our method on facial images, we use the CelebA CelebA-HQ dataset, keeping the original splits from CelebA Liu et al. (2015). The original split contains 24,183 images for training, 2,993
for validation, and 2,824 for testing (image size 256 × 256). We study two tasks on this dataset: colorization and inpainting.

Baselines We compare the performance of non-Bayesian models with different Bayesian approaches. Additionally, we incorporate a diagonal (D) covariance matrix following the method of Kendall & Gal (2017). We extend this by using a low-rank plus diagonal (LR+D) parameterized distribution.

As a further baseline, we also provide results by following the approach by Zepf et al. (2023). Here, we further approximate the aleatoric uncertainty term of Equation 1 to prevent sampling aleatoric Pmatrices. This reduces the number of resulting columns from $T \times (R + 1)$ to T + R. It is achieved by approximating the expectation of the aleatoric uncertainty Σ^a term with the aleatoric covariance prediction of the model with the expected weights:

 $\Sigma^{a} = \mathbb{E}_{q_{a}^{*}(W)}\left[\Sigma_{W}(x)\right] \approx \Sigma_{\mathbb{E}_{a^{*}}[W]}(x)$

To compute this term, we require the expected weights of the Bayesian models to be well-defined. For 295 MCD, this is done by turning dropout off and rescaling the activations accordingly. For SVI, where 296 the weights follow Gaussian distributions, the expected weights are simply the means of the Gaussian 297 distributions. For DE, we are unable to define expected weights, hence this approximation is not 298 evaluated in this case. Note that Zepf et al. (2023) refer to this approach as MAP solution, which 299 coincides with the expected weights solution if the weight uncertainty is modeled with symmetrical 300 unimodal distributions like Gaussians as commonly used by SVI and Laplace Approximation (LA). 301 Furthermore, Zepf et al. (2023) do not provide a joint representation, and log likelihood calculation is 302 only possible using a combination of our methods.

303

284

304 **Hyperparameter** Finally, we evaluate our joint the LR+D parametrization in combination with 305 all three Bayesian methods. For this case, we let the model predict a matrix $P_W \in \mathbb{R}^{S \times R^W}$ of rank 306 $R^W = 8$ and for epistemic models, we draw T = 64 samples. The predictions are multivariate 307 Normal distributions, represented by their LR+D parametrization. Those predictions are joined to a 308 single, LR+D parametrized distribution. For the full joint uncertainty LR+D model, this yields a joint P matrix with $R = T \times (R^W + 1) = 576$ columns, which we jointly compress down to R = 64 with 309 truncated SVD while keeping the diagonal variance of the dropped columns as described in Section 2. 310 For the expected weights baseline, we perform an additional forward pass using the expected weights 311 and concatenate the aleatoric and epistemic columns, which leads to R = 72 in total. All models are 312 trained for the same amount of steps. 313

314 315

3.2 MAIN RESULTS - COMPARISON OF THE FIT OF PREDICTIVE DISTRIBUTIONS

316 **Quantitative Results** To evaluate performance, we use the negative log-likelihood, which measures 317 how well a model predicts the observed data. Lower values imply that the model's predictions are 318 closer to the actual outcomes. Quantitatively, we find that modeling epistemic uncertainty improves 319 the likelihood of unseen test sample predictions, as shown in Table 1. This finding holds for both 320 modeling the diagonal and modeling the LR+D across all experiments. Additionally, including co-321 variances through our LR+D further improves the likelihood of unseen test sample predictions across all experiments. The usefulness of expected weights $\mathbb{E}[W]$ approximation for the aleatoric uncertainty 322 Σ^a is inconsistent. Our approach, to combine both into a joint multivariate uncertainty representation 323 and using SVD, is superior in all performed experiments with all tested Bayesian methods.

324							
325				MNIST	Cele	ebA	Flying Chairs
020	Epistemic	Parameters		Inpainting	Colorization	Inpainting	Optical Flow
320	-			×1	$\times 1000$	×100	×100
321	×	D	-	2665± 7794	-146± 111	153 ± 175	1059± 496
328	X	LR+D	-	2610 ± 24982	$-216\pm$ 78	$134\pm$ 70	$882\pm$ 554
329	MCD	D	-	$-292\pm$ 345	$-152\pm$ 77	$125\pm~146$	1012± 390
330	SVI	D	-	-268 ± 297	-157 ± 35	$125\pm$ 79	$1043\pm$ 381
331	DE	D	-	-308 ± 236	$-158\pm$ 58	$93\pm~101$	$965\pm~291$
332	MCD	LR+D	$\mathbb{E}[W]$	-155± 4257	$-235\pm$ 41	$73\pm$ 62	853± 342
333	SVI	LR+D	$\mathbb{E}[W]$	-327 ± 3232	-225 ± 38	$116\pm$ 35	760 ± 524
334	MCD	LR+D	SVD	-409 ± 302	-241 ± 28	66± 52	844± 308
335	SVI	LR+D	SVD	-383± 2591	-229 ± 31	109 ± 31	748 ± 467
336	DE	LR+D	SVD	$-394\pm$ 118	$-240\pm$ 25	61 ± 44	837 ± 272

Table 1: Quantitative Results. We evaluate the negative log-likelihoods (base 10) of predictions across various dataset-task combinations and its test set variability using standard deviations. Lower values indicate higher likelihood and, therefore, better predictions. Our method is assessed in four experiments: inpainting of removed image parts, colorization, and optical flow estimation. The log-likelihoods scale linearly with the dimensionality of the prediction and, in the case of masking, are evaluated only in the masked area. Results are generated using both Bayesian (MCD,SVD,DE) and non-Bayesian (X) networks, each with diagonal (D) and low-rank plus diagonal (LR+D) covariance parametrization. For the combination of LR+D and Bayesian approaches, we depict both the results by expected weights $\mathbb{E}[W]$ approximation and by our approach using truncated SVD. Our findings demonstrate that employing the LR+D representation and incorporating epistemic uncertainty enhance the posterior predictive distribution and increase the likelihood of the predictions.



Figure 3: Qualitative Results. Random samples from the test sets depict the input, prediction, ground truth, and parameters of the predictive distribution. The top rows show colorization on CelebA images, while the bottom rows display inpainting of MNIST digits. Our model predicts a mean (Pred), the parameter D (Diag), and a low-rank matrix P. In both cases, the predicted joint low-rank matrix P is reduced to the 64 most significant directions (columns) based on their respective eigenvalues. The 10 images in columns 3-8 visualize the 10 most important directions with a random orientation in descending order of the associated eigenvalues. We observe that these columns focus on uncertainty in specific image areas or colors. Additionally, the singular values Ψ measure the importance of the associated direction. For more qualitative results of all datasets and Bayesian methods, please see the Supplement Figures 6, 7, 8, 9, and 10.

378 **Qualitative Results** Figure 3 and Supplementary Figures 6, 7, 8, 9, and 10 provide qualitative 379 results, where we exemplarily visualize how the 10 most important columns in our joint low-rank 380 matrix P describe the areas of correlated uncertainty. E.g. within the CelebA inpanting's (top) 381 visualized P Matrix, the first two eigenvectors focus on full image color shifts, whereas the first 382 one focuses on an axis between orange and blue (inverse color), the second one on the axis between purple to green (inverse color). The third eigenvector focus on the color contrast between foreground and background. The fourth one focuses on hair and eye color. Further, the singular values Ψ provide 384 an interpretation of the importance of these correlations. Visualization of the eigenvectors is only 385 possible with our method, which includes the covariance terms; hence, allowing the identification of 386 image regions with correlated uncertainty. The Parameter D (Diag) captures additional uncertainty, 387 which could not be captured by the Low Rank Covariance Matrix created by PP^{\intercal} . 388

In summary, these qualitative results can help to intuitively describe the underlying relations of uncertainty on an image level.

391 392

3.3 Additional Result - Out of Distribution Detection

We evaluate our method on the task of out-of-distribution detection on the MNIST dataset, where 394 we omit the digit "2" from the training data. Similar to Subsection 3.2, we train four types of 395 models, with and without epistemic uncertainty, using both diagonal and LR+D representations. Our 396 aim is to study whether predictive distributions are more spread out for unseen out-of-distribution 397 (OOD) test samples compared to in-distribution (ID) test samples. A common metric to evaluate the 398 spread of continuous distributions is differential entropy Thomas & Joy (2006). Figure 4 presents 399 histograms of differential entropies of the predictive distributions of two Bayesian models. The 400 left one's output is parametrized using a diagonal normal distribution, whereas the right one uses 401 the LR+D parametrization. We observe that the LR+D parametrized network better separates the ID entropies from the OOD entropies than the diagonally parameterized covariance model. For 402 quantitative analysis, we compute the average entropy for ID and OOD samples in the respective 403 columns. Furthermore, we calculate the Area Under the Curve (AUC) as a metric for the separation 404 of OOD and ID samples using differential entropy. Finally, we employ a Kolmogorov–Smirnov 405 (KS) test to measure a distance between the distribution of differential entropies of ID and OOD 406 samples. Higher values are preferred for both metrics. Our results demonstrate that both 1) modeling 407 epistemic uncertainty and 2) covariances help to improve the differentiation between ID and OOD 408 test samples for the tested dataset. Finally, we present a combination of both, using the expected 409 weights approximation $\mathbb{E}[W]$ and the SVD approximation. Our method, which combine both and 410 compresses using SVD, is superior to all baselines.

411 412

413 414

3.4 SPATIAL AND COMPUTATIONAL COMPLEXITY

Figure 5 presents the memory (left) and time (right) requirements for computing the log-likelihood of different covariance parameterizations: sparse options like diagonal (D) and low-rank plus diagonal

Epistemic	Parameters		ID	OOD	AUC \uparrow	KS↑
X	D		-1134	-988	0.782	0.441
1	D		-751	-561	0.851	0.564
X	LR+D		-1407	-1026	0.842	0.558
1	LR+D	$\mathbb{E}[W]$	-1291	-989	0.859	0.583
1	LR+D	SVD	-844	-593	0.882	0.625

424 Table 2: Results for out-of-distribution (OOD) detection of MNIST digits using the LR+D repre-425 sentation. The digit "2" is excluded from the training set and, therefore, OOD in the test set. The 426 table presents the results for a Bayesian and non-Bayesian network, both combined with D and with LR+D covariance parametrization. The columns in-distribution (ID) and OOD indicate the average 427 differential entropy for the respective sample types. The column AUC measures the Area Under the 428 Curve (AUC) for using differential entropy as a separation criterion. Finally, Kolmogorov-Smirnov 429 (KS) depicts a distance between the distributions of ID and OOD differential entropies. Larger values 430 are preferable for both metrics. The AUC as well as KS demonstrate that both epistemic uncertainty 431 and the LR+D representation are beneficial for the detection of OOD samples.



Figure 4: Distribution of differential entropies predicted by non-Bayesian and Bayesian models parametrized with either a D or an LR+D covariance matrix. The plots depict cumulative kerneldensity plots of the differential entropies of all ID (solid line) and OOD (dashed line) test set samples. Using the LR+D representation, the ID and OOD samples show a better separation based on their entropy compared to using only the variances. The vertical line visualizes the result of the KS test.



Figure 5: Empirical measurement of memory requirement of log likelihood calculation (left) and
Average empirical time measurement of 100 log likelihood calculations (right). Memory requirements
increase linear with the number of variables for LR+D representation and quadratic for full covariance
representations. Full covariance matrix measurements are done until exceeding memory of the GPU.
The LR+D representation enables handling larger covariance matrices than full rank parametrizations.

(LR+D), as well as full covariance Σ using both naive and lower-triangular parameterizations. The complexity is shown as a function of the number of variables in the covariance matrix, with specific points marking the number of variables for each dataset-task combination.

For LR+D, we evaluate various numbers of columns *R* in the low-rank matrix *P*. We limit our
analysis to sizes that fit within a single 48GB GPU. As a result, the full covariance plot stops early
when memory capacity is exceeded. As seen in the figure, the LR+D parameterization (with 64
columns) is significantly more efficient than the naive full covariance, both with respect to memory
and time, for all datasets. In larger datasets like CelebA and Flying Chairs, the full covariance matrix
approaches the GPU memory limit, even without batching or storing the model and its gradients.
Theoretical details on the computational complexity can be found in the Supplementary Material B.

468 469

438

439

440

441

442 443

444

445

446

447

448

449

450 451

452

470 3.5 ABLATION STUDIES

471 For a comprehensive evaluation of our uncertainty framework, we carry out multiple ablations to 472 study which factors are important for optimal model performance. One aspect is the number of 473 samples drawn from the weight distribution and the number of columns R retained in the final 474 representation after performing SVD. We vary the number of samples T, the application of SVD to 475 the low-rank matrix P, and the number of columns retained post-SVD, resulting in R total columns. 476 Table 3 presents the mean negative log-likelihood (NLL) and its test set variability using standard 477 deviations. The results indicate that, generally, a higher number of samples and retained columns improve predictions. However, retaining a high number of columns resulting from sampling without 478 dimensionality reduction can lead to numerical instabilities $(\star, \star\star)$, preventing NLL calculation. In 479 contrast, a low number of samples causes high variability within the test set. We find that optimal 480 results are achieved through dimensionality reduction using SVD, balancing efficient representation 481 with fewer columns against better predictive performance in terms of NLL. 482

We further ablate various design choices. In a first ablation we show that including the update of the diagonal given by the Equations 6 and 10 in Table 6 makes the predictions more robust. Further, we show better performance for the choice to perform SVD on the joint P matrix instead of separated in P^a and P^e in Table 7; next we ablate on the number of columns of P_w directly predicted by the

				MNIST	CelebA		Flying Chairs
				Inpainting	Colorization	Inpainting	Optical Flow
486	T	R	P	$\times 1$	$\times 1000$	$\times 100$	$\times 100$
-100	64	576	-	-455± 994	**	**	*793±274
487	32	288	-	-440 ± 1401	\star -252 \pm 27	53±45	*813±289
488	16	144	-	-411 ± 2035	-245 ± 34	65 ± 51	838±311
100	8	72	-	-352 ± 3096	-237 ± 45	81±59	868 ± 341
489	64	64	SVD(64)	-409 ± 302	-242 ± 35	68 ± 51	844 ± 308
490	64	32	SVD(32)	-372 ± 204	-235 ± 41	82 ± 57	866 ± 322
	64	16	SVD(16)	-345 ± 170	-228 ± 45	96±63	888±336
491	64	8	SVD(8)	-319± 165	-217 ± 48	113±71	915±349

493 Table 3: Comparison of different approximations of the LR+D-parametrized covariance matrix using varying numbers of samples and degrees of dimensionality reduction via SVD. T denotes the number 494 of samples drawn from $q^{\star}(w)$, and R indicates the number of columns in the resulting representation. 495 The column p specifies to what extent the dimensionality is reduced after sampling using SVD. 496 Numbers in brackets represent the retained singular vectors, resulting in columns R. The columns to 497 the right contain the mean negative log-likelihoods (NLL) and their test set variability using standard 498 deviations. We indicate scaling to ease visualization. Some NLLs could not be evaluated due to 499 numerical issues. Cells with more than 10% missing test results are replaced by $\star\star$, those with at 500 least one missing test result (<10%) are marked with \star . The results clearly show that an increasing 501 number of samples improves performance (lower negative log-likelihood (NLL)) and consistency 502 (lower standard deviation). However, it also increases the risk of numerical issues. To harness most of 503 the performance benefits while avoiding numerical instabilities, we can apply SVD for dimensionality 504 reduction. Additionally, retaining more columns after performing SVD results in better outcomes.

model weights in Table 8 and find a moderate improvement at higher computational cost motivating our choice of choosing 8 columns in further experiments. Finally, we provide an extensive list of various design choice combinations in the Ablation Table 9.

- 4
- 510 511

505

506

507

508 509

DISCUSSION

512 **Conclusion** In this work, we have explored the dual nature of uncertainties — aleatoric and 513 epistemic — and their integration in high-dimensional regression tasks. We proposed a novel 514 method that employs a low-rank plus diagonal covariance matrix to approximate joint uncertainty, 515 effectively preserving vital output correlations and significantly reducing the computational demands that are inherent to full covariance matrix representation. Our approach lowers memory usage 516 and improves the efficiency of both sampling and log-likelihood calculations. Empirically, our 517 approach outperforms the commonly used factorized Gaussian representation. It exhibits a lower 518 negative log-likelihood, indicating superior performance in uncertainty estimation, particularly in 519 high-dimensional regression tasks. Furthermore, it excelled in out-of-distribution (OOD) detection 520 on the tested dataset, leveraging the criterion of differential entropy. This success underscores the 521 method's effectiveness in capturing and quantifying uncertainty.

522 523

Limitations Our method conceptually extends to any Bayesian framework; however, for simplicity 524 and computational reasons, we restrict our evaluation to using Monte Carlo Dropout, Stochastic 525 Variational Inference and Deep Ensemble. Further investigations into other Bayesian inference 526 techniques should determine their empirical applicability. We expect that more advanced concepts 527 will lead to better overall uncertainty estimation. The method is flexible with regard to the choice in 528 number of columns utilized in the low-rank plus diagonal parameterization of the covariance matrix. 529 A higher number of columns provides better overall uncertainty estimation but can lead to numerical 530 instabilities and increased computational complexity. Essentially, our method exhibits a trade-off 531 between the quality of uncertainty estimation and these factors, see Table 3. We hypothesize that the numerical instabilities are related to the linear dependencies of the columns in the low-rank and, 532 hence, propose additional investigations to understand and mitigate their impact. Finally, our method 533 builds upon the assumption that uncertainties in output can be modeled by a single multivariate 534 Gaussian, even though this approximation is often used in the literature Kendall & Gal (2017); 535 Monteiro et al. (2020); Duff et al. (2023). However, multivariate Gaussians may not be a suitable 536 approximation for every task, for example, for uncertainties in translation or rotation in images. 537 Exploring epistemic uncertainty under different distributions is a highly promising research question. 538

By more accurately approximating the true posterior than traditional joint distributions, our method enhances both the reliability and explainability of predictions from deep learning models.

540 5 REPRODUCIBILITY

The source code, released under an open-source license, is available via an anonymous public GitHub repository https://anonymous.4open.science/r/structured_joint_
uncertainty/. Checkpoints for all models trained with the proposed method, as well as all
baselines, are available upon request. The datasets used in the experiments are publicly accessible
and links as well as preprocessing scripts are included in the repository. An extensive schematic and
intuitive description of the method, along with proofs, is also included in the Supplementary Material.
Additionally, qualitative examples are provided to enhance understanding of the method.

References

549 550

567

568

- Viktor Bengs, Eyke Hüllermeier, and Willem Waegeman. On second-order scoring rules for epistemic uncertainty quantification. In *International Conference on Machine Learning*, pp. 2078–2091.
 PMLR, 2023.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in
 neural network. In *International conference on machine learning*, pp. 1613–1622. PMLR, 2015.
- Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and
 Philipp Hennig. Laplace redux-effortless bayesian deep learning. Advances in Neural Information
 Processing Systems, 34:20089–20103, 2021.
- Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Finale Doshi-Velez, and Steffen Udluft. Decomposition of uncertainty in bayesian deep learning for efficient and risk-sensitive learning. In *International conference on machine learning*, pp. 1184–1193. PMLR, 2018.
- Garoe Dorta, Sara Vicente, Lourdes Agapito, Neill DF Campbell, and Ivor Simpson. Structured
 uncertainty prediction networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5477–5485, 2018a.
 - Garoe Dorta, Sara Vicente, Lourdes Agapito, Neill DF Campbell, and Ivor Simpson. Training vaes under structured residuals. *arXiv preprint arXiv:1804.01050*, 2018b.
- A. Dosovitskiy, P. Fischer, E. Ilg, P. Häusser, C. Hazırbaş, V. Golkov, P. v.d. Smagt, D. Cremers, and T. Brox. Flownet: Learning optical flow with convolutional networks. In *IEEE International Conference on Computer Vision (ICCV)*, 2015.
- Margaret AG Duff, Ivor JA Simpson, Matthias Joachim Ehrhardt, and Neill DF Campbell. Vaes with structured image covariance applied to compressed sensing mri. *Physics in Medicine & Biology*, 68(16):165008, 2023.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059.
 PMLR, 2016.
- Nitesh B Gundavarapu, Divyansh Srivastava, Rahul Mitra, Abhishek Sharma, and Arjun Jain. Structured aleatoric uncertainty in human pose estimation. In *CVPR Workshops*, volume 2, pp. 2, 2019.
- Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning:
 An introduction to concepts and methods. *Machine learning*, 110(3):457–506, 2021.
- Alexander Immer, Emanuele Palumbo, Alexander Marx, and Julia Vogt. Effective bayesian het eroscedastic regression with deep neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems*, 30, 2017.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
 uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.

594 595 596	Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In <i>Proceedings of the IEEE international conference on computer vision</i> , pp. 3730–3738, 2015.
597 598 599	Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon Wilson. A simple baseline for bayesian uncertainty in deep learning. <i>Advances in neural information processing systems</i> , 32, 2019.
600 601 602 603	Miguel Monteiro, Loïc Le Folgoc, Daniel Coelho de Castro, Nick Pawlowski, Bernardo Marques, Konstantinos Kamnitsas, Mark van der Wilk, and Ben Glocker. Stochastic segmentation networks: Modelling spatially correlated aleatoric uncertainty. <i>Advances in neural information processing systems</i> , 33:12756–12767, 2020.
604 605 606	Elias Nehme, Omer Yair, and Tomer Michaeli. Uncertainty quantification via neural posterior principal components. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
607 608 609	David A Nix and Andreas S Weigend. Estimating the mean and variance of the target probability distribution. In <i>Proceedings of 1994 ieee international conference on neural networks (ICNN'94)</i> , volume 1, pp. 55–60. IEEE, 1994.
610 611 612 613	Frank Nussbaum, Jakob Gawlikowski, and Julia Niebling. Structuring uncertainty for fine-grained sampling in stochastic segmentation networks. <i>Advances in Neural Information Processing Systems</i> , 35:27678–27691, 2022.
614 615 616 617	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In <i>Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18</i> , pp. 234–241. Springer, 2015.
618 619 620	Rebecca L Russell and Christopher Reale. Multivariate uncertainty in deep learning. <i>IEEE Transac-</i> <i>tions on Neural Networks and Learning Systems</i> , 33(12):7937–7943, 2021.
621 622 623	David Salinas, Michael Bohlke-Schneider, Laurent Callot, Roberto Medico, and Jan Gasthaus. High- dimensional multivariate forecasting with low-rank gaussian copula processes. <i>Advances in neural</i> <i>information processing systems</i> , 32, 2019.
624 625 626	Maximilian Seitzer, Arash Tavakoli, Dimitrije Antic, and Georg Martius. On the pitfalls of heteroscedastic uncertainty estimation with probabilistic neural networks. <i>arXiv preprint arXiv:2203.09168</i> , 2022.
627 628 629	Nicki Skafte, Martin Jørgensen, and Søren Hauberg. Reliable training and estimation of variance networks. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
630 631	Andrew Stirn and David A Knowles. Variational variance: Simple, reliable, calibrated heteroscedastic noise variance parameterization. <i>arXiv preprint arXiv:2006.04910</i> , 2020.
632 633 634 635	Andrew Stirn, Harm Wessels, Megan Schertzer, Laura Pereira, Neville Sanjana, and David Knowles. Faithful heteroscedastic regression with neural networks. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 5593–5613. PMLR, 2023.
636	MTCAJ Thomas and A Thomas Joy. <i>Elements of information theory</i> . Wiley-Interscience, 2006.
638 639	Jeffrey Willette, Hae Beom Lee, Juho Lee, and Sung Ju Hwang. Meta learning low rank covariance factors for energy-based deterministic uncertainty. <i>arXiv preprint arXiv:2110.06381</i> , 2021.
640 641 642	Omer Yair, Elias Nehme, and Tomer Michaeli. Uncertainty visualization via low-dimensional posterior projections. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 11041–11051, 2024.
643 644 645 646	Kilian Zepf, Selma Wanna, Marco Miani, Juston Moore, Jes Frellsen, Søren Hauberg, Aasa Feragen, and Frederik Warburg. Laplacian segmentation networks: Improved epistemic uncertainty from spatial aleatoric uncertainty. <i>arXiv preprint arXiv:2303.13123</i> , 2023.

648 A SYMBOLS AND ACRONYMS 649

650 A.1 LIST OF SYMBOLS 651

652	Symbol	Remark
653	\overline{D}	Diagonal matrix used for LR+D
654	P	Tall matrix P used for LR+D
655	$R_{\tilde{\alpha}}$	Number of columns of tall matrix P
656	S	Number of outputs P
657	$\frac{T}{U}$	Number of epistemic samples
658	U JU	Singular vectors
659	V^{Ψ}	Figenvectors of <i>PP</i> ^T
660	Ŵ	Model parameters
661	θ	Parameters of the proxy distribution
662	Х	Training Data
663	Y	Training Labels
664	x	Test data sample
665	y	Test label sample
666	p	Probability distribution
667	q	proxy distribution
668	\mathcal{L}	LOSS Normal distribution
669	v. ⊡	Expectation
670	Cov	Covariance Matrix
671	[]	Column wise Block Concatenation
672	[.]	Stop Gradient Function
674	(.) ^T	Transposed Matrix
675	$(.)^a$	Symbol representing aleatoric uncertainty only
676	$(.)^{c}$	Symbol representing epistemic uncertainty only
677	$(\cdot)_W$	Symbol is a function of the weights
678	$(.)_{w_i}$	Symbol uses the set of weights w_i
679	$(.)_{i_{-}}$	ith as here a financial
680	$(.)_{i}$	
681		
682	A 2 LIST	L OF ACRONYMS
683	11.2 115	
684	ACRONY	MS
685		
686	AUC Area	a Under the Curve
687	CelebA C	elebFaces Attributes
688	D diagona	1
689	DF Deen	Ensemble
690	Electron Ch	
691	Flying Ch	airs
692	GT groun	d truth
693	ID in-dist	ribution
694	KS Kolmo	ogorov–Smirnov
695	LA Lapla	ce Approximation
696		
697	LK+D lov	v-rank plus diagonal
698	LU lower-	upper
699	MAP Max	kimum a posteriori
700	MC Mont	e Carlo
701	MCD Mo	nte Carlo Dropout
		······································

702 MNIST Modified National Institute of Standards and Technology database 703

- 704 NLL negative log-likelihood
- **OOD** out-of-distribution
- 706 SVD Singular Value Decomposition
- 707 SVI Stochastic Variational Inference
 708

B COMPUTATION, TIME AND SPACE COMPLEXITY

		Captured	Parametr.	Memory		Time	
Туре	Parametrization	Corr.	Represent.		$x\Sigma^{-1}x^{T}$	$ \Sigma $	Sampling
full Russell & Reale (2021)	correlation	all	Σ	S^2	S^3	S^3	S^3
full Gundavarapu et al. (2019)	Cholesky	all	Σ	S^2	S^2	S	S^2
sparse Dorta et al. (2018a;b)	inv. band Cholesky	local	Σ^{-1}	SR	SR	S	SR^2
sparse Monteiro et al. (2020)	LR+D	global	Σ	SR	SR^2	SR^2	SR
factorized Kendall & Gal (2017)	diagonal	none	Σ	S	S	S	S

Table 4: This table depicts the computational complexity for calculations using different parametriza-tions for covariance matrices. We use the sparse LR+D parametrization as the basis for our method. This reduces time and spacial complexity in comparison to the naive or Cholesky decomposition and allows for global correlation in comparison to the sparse inverse band Cholesky parametrization. The type and amount of correlations of different parametrization is different (Captured Corr). Furthermore, the used representation enables for efficient calculation of Σ or Σ^{-1} (Parametr. Representation) and needs different amount of memory. The time complexity is given for calculation of the mahalanobis distance $x\Sigma^{-1}x^{\intercal}$, determinant $|\Sigma|$ as well as sampling.

Table 4 give the theoretical time and memory complexities of various covariance parametrizations and calculations. The sparse representations are more efficient in terms of memory and computational complexity. However, they do not provide all degrees of freedom of a covariance matrix and are limited to either local or the most important global correlations.

C ADDITIONAL RESULTS

733 C.1 QUALITATIVE RESULTS

We provide additional qualitative results for every performed tasks. Figures 6 presents optical flow on
Flying Chairs, 7 depicts CelebA inpainting, 8 shows CelebA colorization, and 9 illustrates MNIST
inpainting. Figure 10 compares both, eigenvectors in both random orientations as well as the different
used Bayesian methods.



Figure 6: Additional Qualitative Results, Visualizing Flying Chairs' Optical Flow. Random samples 792 from the test sets showing input, prediction, ground truth and parameters of the predictive distribution. 793 The task here is optical flow estimation in the Flying Chairs dataset. The model predicts a mean 794 (Pred), and the parameter D (Diag), as well as a low-rank matrix P. In all cases, the predicted 795 joint low-rank matrix P is reduced to the 64 most significant directions (columns) and displayed 796 using the 10 most significant ones in descending order of associated eigenvalues. We can clearly see 797 that the columns focus on uncertainty in certain images areas or colors. Furthermore, the singular 798 values Ψ give a measure of importance of the associated direction. Note that the orientation of the 799 singular vectors is arbitrarily chosen and can be inverted, which results in opposite colors (left) and 800 brightness (right). These eigenvectors are only possible to visualize when modeling covariances and show the direction of maximum variability of the data and helps to understand the underlying factors. 801 Furthermore, we show the upper bound of the angles between the directions of the eigenvectors of 802 PP^{\intercal} and the eigenvectors of Σ . 803

- 804
- 805
- 806
- 807
- 808



Figure 7: Additional Qualitative Results, Visualizing Inpainting of Eyes of CelebA Faces. Random samples from the test sets showing input, prediction, ground truth and parameters of the predictive distribution. The task here is inpainting of the eyes region of the CelebA faces dataset. The model predicts a mean (Pred), and the parameter D (Diag), as well as a low-rank matrix P. In all cases, the predicted joint low-rank matrix P is reduced to the 64 most significant directions (columns) and displayed using the 10 most significant ones in descending order of associated eigenvalues.



Figure 8: Additional Qualitative Results, Visualizing the Colorization of CelebA Faces. Random samples from the test sets showing input, prediction, ground truth and parameters of the predictive distribution. The task here is colorization of the CelebA faces dataset. The model predicts a mean (Pred), and the parameter D (Diag), as well as a low-rank matrix P. In all cases, the predicted joint low-rank matrix P is reduced to the 64 most significant directions (columns) and displayed using the 10 most significant ones in descending order of associated eigenvalues.



Figure 9: Additional Qualitative Results, Visualizing Inpainting of MNIST Digits. Random samples from the test sets showing input, prediction, ground truth and parameters of the predictive distribution. The task is inpainting MNIST digits. The model predicts a mean (Pred), and the parameter D (Diag), as well as a low-rank matrix P. In all cases, the predicted joint low-rank matrix P is reduced to the 64 most significant directions (columns) and displayed using the 10 most significant ones in descending order of associated eigenvalues.



Figure 10: Additional Qualitative Results, comparing Bayesian Methods. A Random sample from the test sets showing input, prediction, ground truth and parameters of the predictive distribution with various Bayesian Methods. For the first method, we also show the eigenvectors with inverse sign. Both signs are mathematically equivalent and one of them is randomly chosen for the visualizations. The task here is colorization of the CelebA faces dataset. The model predicts a mean (Pred), and the parameter D (Diag), as well as a low-rank matrix P. In all cases, the predicted joint low-rank matrix P is reduced to the 64 most significant directions (columns) and displayed using the 10 most significant ones in descending order of associated eigenvalues.

1026 C.2 QUANTITATIVE PREDICTION ERRORS

Table 5 lists the predicted errors for all bayesian mehtods. We aim for similar predictive errors for all
 models to get mainly evaluate the quality of the uncertainty using negative log-likelihood (NLL).

1030											
1031				MN	IST		Celeb	A		Flying Chairs	
1020				Inpai	inting	Colori	zation	Inpai	nting	Opt. Flow	
1002	Epistemic	Param.	R_W	$L_1 \downarrow$	$L_2 \downarrow$	$L_1\downarrow$	$L_2 \downarrow$	$L_1\downarrow$	$L_2 \downarrow$	$L_1\downarrow$	$L_2\downarrow$
1033	X	D	0	0.0331	0.0475	1.9002	2.4168	2.64	5.55	0.150	0.346
1034	MCD	D	0	0.0335	0.0477	1.8976	2.4159	2.58	5.62	0.148	0.325
1035	SVI	D	0	0.0366	0.0501	1.9074	2.4176	2.71	5.71	0.161	0.330
1036	DE	D	0	0.0339	0.0483	1.8977	2.4159	2.57	5.63	0.147	0.325
1037	X	LR+D	8	0.0352	0.0494	1.9330	2.4244	2.59	5.52	0.176	0.345
1007	MCD	LR+D	4	0.0379	0.0530	1.9116	2.4183	2.51	5.52	0.178	0.338
1038	MCD	LR+D	8	0.0340	0.0481	1.9191	2.4206	2.56	5.52	0.157	0.328
1039	MCD	LR+D	16	0.0338	0.0480	1.9182	2.4202	2.61	5.57	0.155	0.328
1040	SVI	LR+D	8	0.0384	0.0529	1.9494	2.4297	2.67	5.65	0.162	0.333
1041	DE	LR+D	8	0.0348	0.0489	1.9219	2.4213	2.60	5.60	0.158	0.329

1042 Table 5: Comparison of reconstruction or prediction errors of all methods. We use the same loss 1043 for the prediction between those methods. The last convolutional layer of models with LR+D 1044 parametrization has more channels in comparison to models with diagonal parametrization. Further-1045 more, the uncertainty channels receive gradients from different negative log likelihood functions. 1046 Bayesian models (Epistemic) include additional Dropout layer or variational convolutional layer and 1047 are evaluated using 64 weight samples. Essentially, the presented study shows our robust, better 1048 uncertainty quantification towards the quality of the prediction. This is important to evaluate because 1049 the negative-log likelihood is affected by both prediction and uncertainty estimation.

1050 1051

1053

1052 C.3 ADDITIONAL ABLATION STUDY

One component of our proposed method is to retain the variance of the removed columns after dimensionality reduction using SVD, see 6. In Table 6 we ablate this design choice. The column \hat{D} indicates whether the diagonal D is updated (\checkmark) after performing SVD according to Equations 6 and 10, or if the original D is retained (\checkmark) as per Equation 8. Our ablation shows that updating the diagonal D appears to slightly improve the average, NLL while also enhancing prediction consistency and reducing test set variability. We show that this is consistent for three different configurations of dimensionality reductions, see Table 6.

Table 7 compares the performance of performing SVD independently on the epistemic and aleatoric low-rank matrices P^e and P^a versus on the combined low-rank matrix P. The findings show that

				MNIST Celeb		bA	Flying Chairs
R	P^a	P^e	\hat{D}	Inpainting	Colorization	Inpainting	Optical Flow
				×1	$\times 1000$	×100	×100
64	SVD(32)	SVD(32)	1	-379± 202	-239±35	71±52	849±308
32	SVD(16)	SVD(16)	1	-350± 164	-232 ± 39	$84{\pm}57$	870 ± 322
64	joint S	VD(64)	1	-409 ± 302	-242 ± 35	68 ± 51	844 ± 308
64	SVD(32)	SVD(32)	X	-408 ± 2087	-241 ± 41	$70{\pm}55$	845 ± 324
32	SVD(16)	SVD(16)	X	-368 ± 2985	-233 ± 52	85 ± 62	869 ± 348
64	joint SVD(64)		X	-410 ± 1980	-242 ± 39	68±53	841±321

1071 1072

1063 1064 1065

Table 6: Comparison between adapting the diagonal D after performing the SVD according to Equations 6 and 10 or not. Here R denotes the number of columns in the resulting representation. The columns P^a and P^e denote if and to what degree the dimensionality is reduced after sampling using SVD. The numbers in brackets denote the kept singular vectors, which result in columns R. The column \hat{D} indicates whether the diagonal D is updated (\checkmark) after performing SVD according to Equations 6 and 10, or if the original D is retained as per Equation 8, despite the dimensionality reduction of P. We show in the first three rows that updating the diagonal D appears to slightly improve the average NLL while also enhancing prediction consistency and reducing test set variability.

1080 1081	R	P^{a}	P^e	MNIS Inpain ×1	ST ting	Coloriza	Cele tion	bA Inpainting $\times 100$		Flying Chairs Optical Flow ×100
1082	64	SVD(32)	SVD(32)	-379±	202	-239±	35	71±	52	849±308
1084	32 16	SVD(16) SVD(8)	SVD(16) SVD(-8)	$-350\pm$ $-322\pm$	164 158	$-232\pm$ $-222\pm$	39 41	$^{84\pm}_{99\pm}$	57 64	870 ± 322 892 ± 333
1085	64	joint S	VD(64)	-409±	302	-242±	35	$68\pm$	51	844 ± 308
1086	32	joint S	$-372\pm$	204	$-235\pm$	41	$82\pm$ 06+	57 63	866 ± 322	
1087 1088	8	joint S	-345⊥ -319±	165	-228± -217±	48	$113\pm$	03 71	915 ± 349	

Table 7: Comparison of different approximations of the LR+D-parametrized covariance matrix 1090 using varying degrees of dimensionality reduction via SVD. T denotes the number of samples drawn 1091 from $q^{\star}(w)$, and R indicates the number of columns in the resulting representation. The columns 1092 P^{e} and P^{a} denote on which matrix SVD is performed and the degree of dimensionality reduction. 1093 The numbers in brackets represent the retained singular vectors, resulting in columns R. If the SVD 1094 is performed on the joint matrix P instead of the epistemic and aleatoric submatrices P^e/P^a , it is 1095 marked as joint. We observe that with the same number of columns R, performing SVD on the joint 1096 matrix P yields better performance. Conversely, when retaining the same number of columns per 1097 SVD, the separate version performs slightly better. In both cases, a higher number of columns is preferable. 1098

1100				MNIS	ST		Cele	bA		Flying Chairs
1101	$R^W \mid R$		P	Inpain	ting	Colorization		Inpainting		Optical Flow
1102				×1		$\times 1000$		$\times 100$		×100
1102	4	16	joint SVD(16)	-193±	63	-198±	39	$147\pm$	62	892±412
1103	8	16	joint SVD(16)	-345±	170	$-228\pm$	45	$96\pm$	63	888±336
1104	16	16	joint SVD(16)	-363±	156	$-206\pm$	37	$224\pm$	136	895±322
1105	4	32	joint SVD(32)	-241±	73	$-207\pm$	32	$113\pm$	52	870±387
1106	8	32	joint SVD(32)	-372±	204	$-235\pm$	41	$82\pm$	57	866±322
1107	16	32	joint SVD(32)	-393±	196	$-217\pm$	35	$162\pm$	122	863±310
1100	8	64	joint SVD(64)	-409±	302	$-242\pm$	35	$68\pm$	51	844 ± 308
1100	16	64	joint SVD(64)	-423±	254	$-228\pm$	31	$114\pm$	111	833 ± 302

1109

1099

1110Table 8: Comparison between different number of columns used during training the model. While1111 R^W denotes the number of columns produced by the model without sampling., R denotes the number1112of columns in the resulting representation. The column P show how SVD is used dimensionality is1113reduced. The number in the brackets denote the kept singular vectors, which result as columns R.1114The results tend to be better with a higher number of learned columns R^W . In genreal, increasing1115 R^W can cause increasing training time. However, we also experienced instabilities during training1116for 3 of those tasks when using $R^W = 32$.

1117

1118

combined SVD yields better results for the same number of columns. However, for a fixed number of singular vectors retained per SVD, the independent approach performs slightly better. Note that combined SVD complicates the post-hoc separation of aleatoric and epistemic contributions.

1122Furthermore, Table 8 evaluates the impact of the number of columns predicted by the model's1123output layer. Increasing the number of columns generally benefits the NLL. However, this increases1124computational complexity and may lead to numerical instabilities during training, as observed in11253 out of 4 tasks failing with $R^W = 32$ columns. Balancing these trade-offs, we chose a rank of 8,1126which aligns with the choice of Monteiro et al. (2020).

1127 Finally, Table 9 presents an extended ablation of various parameters, non-Bayesian with various 1128 Bayesan Models (epis) and both kinds of distribution parametrizations (Param). Therefore, it 1129 compares a purely diagonal (D) uncertainty with our LR+D parametrization. The representation took 1130 *T* Bayesian samples, and results in *R* columns of the low-rank matrix. The number of samples can be 1131 reduced for the aleatoric matrix (P^a) the epistemic matrix (P^e) or both together (in the center of both 1132 columns). The columns are either not reduced (-) reduced using Singular Value Decomposition (SVD) 1133 with the remaining number of columns in the brackets, or for the aleatoric covariance matrix using 1136 the expected weights $\mathbb{E}[W]$. The column \hat{D} indicates whether the diagonal *D* is updated (\checkmark) after

1134								MNIST	Cel	ebA	Flying Chairs
1135	Epis	Param	R^W	T	R	$P^a P^e$	\hat{D}	Inpainting	Colorization	Inpainting	Optical Flow
1155	-							×1	$\times 1000$	×100	×100
1136	×	D	0	0 + 1	0	- 0	-	2665 ± 7794	-146± 111	153 ± 175	1059 ± 496
1137	× MCD	LR+D D	8	0 + 1 64 + 1	8	- 0 F[W] -	-	2610 ± 24982 -253 ± 589	-210 ± 78 -151 ± 80	134 ± 70 128 ± 152	882 ± 554 1014 ± 402
1100	MCD	D		64 + 0	Ő		-	-292 ± 345	-152 ± 77	120 ± 132 125 ± 146	1014 ± 402 1012 ± 390
1138	MCD	LR+D	8	64 + 1	72	$\mathbb{E}[W]$ -	-	-155± 4257	$-235\pm$ 41	$73\pm$ 62	853 ± 342
1139	MCD	LR+D	8	64 + 0	576		-	-455± 994	**	**	**
11/0	MCD	LR+D	8	32 + 0	288		-	-440 ± 1401	$*-250\pm 24$	52 ± 46	813 ± 289
1140	MCD	LR+D LR+D	8	10 + 0 8 + 0	72		-	-411 ± 2033 -352 ± 3095	-244 ± 29 -237 ± 37	80 ± 51	858 ± 311 868 + 341
1141	MCD	LR+D	8	64 + 0	64	SVD(32) SVD(32) 🗸	$-379\pm$ 202	-238 ± 29	$69\pm$ 52	849 ± 308
1142	MCD	LR+D	8	64 + 0	32	SVD(16) SVD(16) 🗸	-350± 164	-231± 34	83± 58	$870\pm~322$
1144	MCD	LR+D	8	64 + 0	16	SVD(8) SVD(8) 🗸	$-322\pm$ 158	-222 ± 36	97± 67	892 ± 333
1143	MCD MCD	LR+D	8	64 + 0 64 + 0	64 64	joint SVD(64)	<i>×</i>	-409 ± 302 410 ± 1070	-241 ± 28 242 ± 30	66 ± 52	844 ± 308 841 ± 321
1144	MCD	LR+D	8	64 + 0 64 + 0	32	ioint SVD(32)	2	-372 ± 204	-242 ± 30 -234 ± 32	80 ± 53 80 ± 57	866 ± 322
11/5	MCD	LR+D	8	64 + 0	32	joint SVD(32)	x	-372± 2773	-236± 36	80± 61	864± 345
1140	MCD	LR+D	8	64 + 0	16	joint SVD(16)	1	-345± 170	-227± 37	95± 64	888 ± 336
1146	MCD	LR+D	8	64 + 0	16	joint SVD(16)	X	-326 ± 3914	-229 ± 47	96 ± 71	890 ± 371
1147	MCD MCD	LR+D LR+D	8	64 ± 0 64 ± 0	8	joint SVD(-8)	×	-319 ± 165 -273 ± 5523	-218 ± 40 -219 ± 59	111 ± 74 116 ± 86	915 ± 349 930 ± 417
	MCD	LR+D	4	64 + 1	68	$\mathbb{E}[W]$ -	<u>_</u>	-376 ± 1677	-219 ± 31	$77\pm$ 55	859 ± 422
1148	MCD	LR+D	4	64 + 0	320		-	**	*-230± 21	51± 38	807± 299
1149	MCD	LR+D	4	32 + 0	160		-	*-396± 420	-227 ± 24	$62\pm$ 42	831 ± 325
1150	MCD	LR+D	4	16 + 0	80		-	-396 ± 848	-222 ± 27	78 ± 50	862 ± 369
1150	MCD	LR+D LR+D	4	64 ± 0	40 64	SVD(32) SVD(32	· -	-382 ± 1390 -245 ± 73	-210 ± 31 -222 ± 26	98 ± 00 71 + 47	894 ± 412 849 ± 351
1151	MCD	LR+D	4	64 + 0	32	SVD(32) SVD(32 SVD(16) SVD(16		-198 ± 62	-217 ± 28	87± 53	872 ± 379
1152	MCD	LR+D	4	64 + 0	16	SVD(8) SVD(8		-161± 55	$-211\pm$ 30	104± 59	895 ± 405
1152	MCD	LR+D	4	64 + 0	32	joint SVD(32)	1	$-241\pm$ 73	-218 ± 28	84± 52	870 ± 387
1153	MCD	LR+D	4	64 + 0	32	joint SVD(32)	X	-400 ± 1092	-219 ± 29	84 ± 54	870 ± 406
1154	MCD	LR+D LR+D	4	64 + 0 64 + 0	16	joint SVD(16)	×	-195 ± 05 -399 ± 1372	-213 ± 30 -214 ± 32	101 ± 59 102 ± 62	892 ± 412 896 ± 454
	MCD	LR+D	4	64 + 0	8	joint SVD(8)	1	$-157\pm$ 56	$-206\pm$ 33	102 ± 02 117± 65	914 ± 432
1155	MCD	LR+D	4	64 + 0	8	joint SVD(8)	x	-395± 1635	-208± 37	$120\pm$ 72	$924\pm~492$
1156	MCD	LR+D	16	64 + 1	80	$\mathbb{E}[W]$ -	-	-208 ± 3830	-240 ± 80	56 ± 73	820 ± 365
1157	MCD	LR+D LR+D	16	64 ± 0 32 ± 0	1088 544		-	-483 ± 725 -473 ± 956	**	* 27+ 47	* 760+ 283
1157	MCD	LR+D	16	16 + 0	272		-	-446 ± 1417	*-259± 33	$36\pm$ 52	791 ± 303
1158	MCD	LR+D	16	8 + 0	136		-	-403± 2126	-251± 45	50± 60	819± 330
1159	MCD	LR+D	16	64 + 0	64	SVD(32) SVD(32		-399± 190	-244 ± 37	57± 60	843 ± 299
1100	MCD	LR+D	16	64 + 0	32	SVD(16) SVD(16		-368 ± 155	-234 ± 42	71 ± 67	872 ± 306
1160	MCD	LR+D LR+D	16	64 + 0 64 + 0	64	joint SVD(64)	×	-423 ± 234 -424 ± 1930	-247 ± 33 -248 ± 45	54 ± 60 54 ± 66	833 ± 302 817 ± 332
1161	MCD	LR+D	16	64 + 0	32	joint SVD(32)	1	-393± 196	-239± 40	$71\pm$ 68	863 ± 310
1160	MCD	LR+D	16	64 + 0	32	joint SVD(32)	×	-375± 2892	-238± 61	72± 81	843 ± 360
1102	MCD	LR+D	16	64 + 0	16	joint SVD(16)	1	-363 ± 156	-229 ± 46	91± 74	895 ± 322
1163	SVI	LK+D D	16	64 ± 0 64 ± 1	16	$\mathbb{F}[W]$	<i>.</i>	$-339 \pm 354/$ -263 ± 311	-225 ± 101 -157 ± 35	$9/\pm 90$ 126 ± 81	$\frac{87}{\pm} \frac{410}{1045\pm}$
1164	SVI	D		64 + 0	0		-	-268 ± 297	-157 ± 35	120 ± 01 125 ± 79	1043 ± 380 1043 ± 381
	SVI	LR+D	8	64 + 1	72	$\mathbb{E}[W]$ -	-	-327± 3232	$-225\pm$ 38	116± 35	760 ± 524
1165	SVI	LR+D	8	64 + 0	576		-	-381 ± 2638	**	100 ± 29	**
1166	SVI	LR+D LR+D	8	32 ± 0 16 ± 0	288		-	$-3/1\pm 2810$ -358 ± 3056	*-233-+ 28	104 ± 30 109 ± 32	$* 723 \pm 399$ 733 ± 463
1167	SVI	LR+D	8	10 ± 0 8 ± 0	72		-	-342 + 3340	-226 ± 35	109 ± 32 118 ± 34	733 ± 403 770 ± 508
1107	SVI	LR+D	8	64 + 0	64	SVD(32) SVD(32) 🗸	-396± 2374	$-226\pm$ 32	111 ± 32	754± 479
1168	SVI	LR+D	8	64 + 0	32	SVD(16) SVD(16) 🗸	-404± 2121	-217 ± 33	122 ± 36	787 ± 501
1169	SVI	LR+D	8	64 + 0	16	SVD(8) SVD(8		-399 ± 1603	-194 ± 23	135 ± 38	826 ± 509
1170	SVI	LR+D LR+D	8	64 + 0 64 + 0	64	joint SVD(64)	×	-363 ± 2391 -348 ± 3251	-229 ± 31 -229 ± 33	109 ± 31 108 ± 31	748 ± 407 744 ± 505
1170	SVI	LR+D	8	64 + 0	32	joint SVD(32)	1	-396 ± 2355	-222± 35	118± 34	780 ± 506
1171	SVI	LR+D	8	64 + 0	32	joint SVD(32)	×	-333± 3514	-223± 39	118 ± 34	778 ± 567
1172	SVI	LR+D	8	64 + 0	16	joint SVD(16)	<i>✓</i>	-404 ± 2051	-212 ± 37	132 ± 39	816± 529
1112	SVI	LR+D LR+D	8	64 ± 0 64 ± 0	16	joint SVD(16)	×	$-319\pm 3/7/$	-214 ± 48 -101 ± 25	132 ± 40 147 ± 43	$81/\pm 609$ 860 ± 537
1173	SVI	LR+D	8	64 + 0	8	ioint SVD(8)	x	-303 ± 4061	-205 ± 57	147 ± 45 146 ± 46	863 ± 650
1174	DE	D		64 + 0	0		-	-308± 236	-158± 58	93± 101	965± 291
	DE	LR+D	8	64 + 0	576		-	-483± 337	**	* 43± 37	**
11/5	DE	LR+D	8	32 + 0	288		-	-487 ± 604	*-252± 23	47 ± 41	805 ± 274
1176	DE	LR+D LR+D	8	10 ± 0 8 ± 0	72		-	-405 ± 1155 -339 ± 3213	-245 ± 27 -237 ± 36	60 ± 48 80 ± 60	820 ± 297 868 ± 344
1177	DE	LR+D	8	64 + 0	64	SVD(32) SVD(32) 🗸	-358± 93	-233± 25	64± 43	845± 270
	DE	LR+D	8	64 + 0	32	SVD(16) SVD(16) 🗸	-321± 82	-212± 24	86± 39	$873\pm~269$
1178	DE	LR+D	8	64 + 0	16	SVD(8) SVD(8) 🗸	-280 ± 72	-184 ± 22	112 ± 34	911 ± 264
1179	DE	LK+D IR±D	8	64 ± 0	64 64	joint SVD(64)	✓ ×	-394 ± 118 -490 ±702	-240 ± 25 -243 ± 28	61 ± 44 60 ± 45	$83/\pm 2/2$ 827 ± 281
4400	DE	LR+D	8	64 + 0	32	joint SVD(32)	2	-351 ± 93	-229 ± 27	79 ± 48	862 ± 276
0811	DE	LR+D	8	64 + 0	32	joint SVD(32)	×	-491± 908	-235± 34	76± 54	846± 295
1181	DE	LR+D	8	64 + 0	16	joint SVD(16)	1	-313± 82	-210± 25	100 ± 42	892± 275
1182	DE	LR+D	8	64 + 0	16	joint SVD(16)	×	-487 ± 1093	-230 ± 39	90 ± 61	862 ± 309
1102	DE	LK+D IR±D	8	64 ± 0	8	joint SVD(-8)	✓ ¥	$-2/3 \pm /3$ -477 ± 1220	-183 ± 22 -217 ± 45	123 ± 35 103 ± 67	$92/\pm 209$ 878 ± 326
1183	DE	LITU	0	04 ± 0	0		^	1 -7// 1229	-217 1 45	1051 0/	070± 520

Table 9: Extended Ablation. We compare non-Bayesian networks with aleatoric uncertainty only and various Bayesian networks with both kind of uncertainties and various hyperparameters.

1188
1189performing SVD according to Equations 6 and 10, or if the original D is retained as per Equation 8,
despite the dimensionality reduction of P. Drawing many samples without any rank compression can
make the approach numerically instable. Here, \star marks ($\star\star$ replaces) values, where less (more) than
2% of the test set results run into numerical errors. To alleviate this, we reduce the dimensionality
of the representation and remove the eigenvectors associated with smaller singular values from the
low-rank matrix.

1194 1195

1196

1204 1205 1206

1213 1214 1215

D DERIVATIONS IN DETAIL

1197 D.1 EXPLOITING LR+D FOR EFFICIENT COMPUTATION OF MATRIX DETERMINANT AND INVERSE

Both the likelihood function $p(y|x, W) = \mathcal{N}(\mu_W^a(x), \Sigma_W^a(x))$ as well as the approximate posterior predictive distribution $p(y|x, X, Y) \approx \mathcal{N}(\mu(x), \Sigma(x))$ are multivariate normal distributions parametrized by covariance matrices Σ_W^a and Σ , respectively, where in the following, we only consider Σ for clarity. Denoting by *S* the output dimension, the normal distribution is then defined as

$$\mathcal{N}(\mu(x), \Sigma(x)) = \frac{1}{\sqrt{|\Sigma(x)|(2\pi)^S}} \exp\left(-\frac{1}{2}(\mu(x) - y)^{\mathsf{T}}\Sigma^{-1}(x)(\mu(x) - y)\right)$$
(11)

1207 which requires computation of the covariance matrix' determinant $|\Sigma|$ and inverse Σ^{-1} for sampling 1208 and evaluation of the log likelihood. For full covariance matrices $\Sigma \in \mathbb{R}^{S \times S}$ with large *S*, these are 1209 very expensive, if not impossible, to compute directly. Instead, we exploit our LR+D representation 1210 for efficient computation of the matrix determinant and inverse.

We compute the determinant as

$$|\Sigma| = |D + PP^{\mathsf{T}}| \tag{12}$$

$$= |I_{R} + P^{\mathsf{T}} D^{-1} P||D| \tag{13}$$

$$= |C||D| \tag{14}$$

1216 1217 where we first substituted Σ with its LR+D representation and subsequently applied the matrix 1218 determinant lemma. With $D \in \mathbb{R}^{S \times S}$, $P \in \mathbb{R}^{S \times R}$ and $I_R \in \mathbb{R}^R$, the so-called capacitance 1219 $C = I_R + P^{\mathsf{T}} D^{-1} P$ is an $R \times R$ matrix. Since $R \ll S$, the determinant of the capacitance matrix is 1220 very cheap to compute.

To compute the inverse, we use the Woodbury matrix identity, again by exploiting the LR+D representation.

$$\Sigma^{-1} = (D + PP^{\mathsf{T}})^{-1} \tag{15}$$

$$= D^{-1} - D^{-1}P(I_R + P^{\mathsf{T}}D^{-1}P)^{-1}P^{\mathsf{T}}D^{-1}$$
(16)

$$= D^{-1} - D^{-1} P C^{-1} P^{\mathsf{T}} D^{-1} \tag{17}$$

1226 1227 1228

1229

1223 1224 1225

As before, the capacitance matrix $C \in \mathbb{R}^{R \times R}$ is very small and thus its inverse easy to compute.

1230 D.2 FULL DERIVATION OF SVD

1231 We apply dimensionality reduction using SVD on our tall P matrices. This involves decomposing 1232 into three separate matrices: U, Ψ , and V^{\intercal} . The U matrix represents an arbitrary not further used 1233 rotation, Ψ is a diagonal matrix containing the singular values, and V^{\intercal} contains the columns of the 1234 transformed matrix.

By selecting the top R singular values and corresponding vectors, we can approximate the original matrix. This approximation is achieved by truncating the matrices U and V^{T} to retain only the top Rsingular values and vectors. This reduces the dimensionality of the data while preserving its essential structure.

The reduced dimensionality representation, denoted as \hat{P} , is computed by taking the product of the truncated matrices V and Ψ . Additionally, a diagonal matrix \hat{D} captures the by the dimensionality reduction removed variance of PP^{\intercal} , with each element representing the contribution of the omitted singular values to the overall uncertainty. We use \hat{D} to update our diagonal for the final LR+D representation. 19//

$$P^{\mathsf{T}} = U\Psi V^{\mathsf{T}} \tag{18}$$

$$PP^{\mathsf{T}} = (U\Psi V^{\mathsf{T}})^{\mathsf{T}} (U\Psi V^{\mathsf{T}}) \tag{19}$$

$$= V\Psi U^{\mathsf{T}} U\Psi V^{\mathsf{T}} \tag{20}$$

(21)

$$= V\Psi\Psi V^{\intercal}$$

$$PP^{\mathsf{T}} = \hat{D} + \hat{P}\hat{P}^{\mathsf{T}} \tag{22}$$

$$\hat{P} = \begin{bmatrix} V_{R-\hat{R}} \cdot \Psi_{R-\hat{R},R-\hat{R}} & \dots & V_R \cdot \Psi_{R,R} \end{bmatrix}$$
(23)

$$\hat{D}_{ii} = \sum_{j=1}^{R-R-1} V_{ij}^2 \cdot \Psi_{j,j}^2$$
(24)

D.3 LOSS DEFINITION

For regression problems we intend to maximize the data likelihood $p(\mathbf{Y}|\mathbf{X}, w) = \prod_{i} p(y_i|x_i, w)$, where we assumed all dataset samples to be i.i.d. Equivalently, we can minimize the negative log likelihood $p(\mathbf{Y}|\mathbf{X}, w) = \sum_{i} -\log p(y_i|x_i, w)$. We further assume the network predictions to be distributed around the true value y following a Gaussian distribution with mean $\mu_w(x)$ and covariance $\Sigma_w(x).$

The loss function for a single training sample can then be defined as

$$\mathcal{L} = -\frac{1}{S} \log p(y|x, w)$$

$$= -\frac{1}{S} \log \left(\frac{1}{\sqrt{|\Sigma_w|(2\pi)^S}} \exp\left(-\frac{1}{2}(\mu_w - y)^{\mathsf{T}} \Sigma_w^{-1}(\mu_w - y)\right)\right)$$

$$= \frac{1}{S} \left(\log \sqrt{|\Sigma_w|(2\pi)^S} + \frac{1}{2}(\mu_w - y)^{\mathsf{T}} \Sigma_w^{-1}(\mu_w - y)\right)$$

$$= \frac{1}{S} \left(\frac{1}{2} \log |\Sigma_w| + \frac{S}{2} \log(2\pi) + \frac{1}{2}(\mu_w - y)^{\mathsf{T}} \Sigma_w^{-1}(\mu_w - y)\right)$$

where we normalized by the output dimensionality S.

Dropping constant terms, we are left with:

$$\mathcal{L} = \frac{1}{2S} \log |\Sigma_w| + \frac{1}{2S} (\mu_w - y)^{\mathsf{T}} \Sigma_w^{-1} (\mu_w - y)$$

We can see that evaluating \mathcal{L} involves computing the determinant and inverse of the covariance matrix. To achieve this, we exploit our LR+D representation as described in the previous section D.1

D.4 FULL DERIVATION OF MEAN VECTOR AND COVARIANCE MATRIX

The expectation of the posterior predictive distribution is given by:

$$\mathbb{E}[y|x, \boldsymbol{X}, \boldsymbol{Y}] = \mathbb{E}_{p(W|Y, X)}\left[\mathbb{E}\left[y|x, W\right]\right]$$
(25)

$\approx \mathbb{E}_{a^*(W)} \left[\mathbb{E} \left[y | x, W \right] \right]$ (26)

1291
1292
$$= \mathbb{E}_{q_{\theta}^{*}(W)} \left[\mu_{W}^{a}(x) \right]$$
(27)

1293
1294
$$\approx \frac{1}{T} \sum_{i}^{T} \mu_{w_i}^a(x) \quad w_i \sim q_{\theta}^*(W)$$
 (28)

$$= \mu(x)$$
(29)

Loss weighting λ_M λ_D, λ_{LR} λ_{LRD} Step $\cdot 10^4$ Figure 11: Loss Factors Variation during Training.

The covariance of the posterior predictive distribution is given by:

$$\operatorname{Cov}\left[y|x, \boldsymbol{X}, \boldsymbol{Y}\right] = \operatorname{Cov}_{p(W|Y,X)}\left[\mathbb{E}_{p(W|Y,X)}\left[y|x, W\right]\right] + \mathbb{E}_{p(W|Y,X)}\left[\operatorname{Cov}\left[y|x, W\right]\right]$$
(30)

$$\approx \operatorname{Cov}_{q_{\theta}^{*}(W)} \left[\mathbb{E}_{q_{\theta}^{*}(W)} \left[y | x, W \right] \right] + \mathbb{E}_{q_{\theta}^{*}(W)} \left[\operatorname{Cov} \left[y | x, W \right] \right]$$
(31)

$$=\underbrace{\operatorname{Cov}_{q_{\theta}^{*}(W)}\left[\mu_{W}^{a}(x)\right]}_{epistemic} +\underbrace{\mathbb{E}_{q_{\theta}^{*}(W)}\left[\Sigma_{W}^{a}(x)\right]}_{aleatoric}$$
(32)

$$\approx \Sigma^{e}(x) + \Sigma^{a}(x)$$
(33)
= $\Sigma(x)$ (34)

$$=\Sigma(x) \tag{3}$$

In above transformations of expectation and variance, we applied the law of total expectation or variance, respectively, and subsequently approximated them using the proxy distribution $q_{\alpha}^{*}(W)$. The expectation over y, denoted by $\mathbb{E}[y|x, W]$, is given by the mean of the predicted normal distribution $\mu_W^a(x)$, whereas the covariance over y, denoted as $\operatorname{Cov}[y|x,W]$, is given by the covariance matrix of the predicted normal distribution. Finally, the expectation – and in some suggested solutions also the covariance – over the proxy distribution $q_{\theta}^{*}(W)$ is approximated using Monte Carlo integration.

Ε **TRAINING SETUP**

To stabilize, training, we train our model outputs separately with different loss terms and change the weight of the loss terms over time. We reparametrize the diagonal $D(x) = \text{diag}\left(\exp(s) + 10^{-4}\right)$ to ensure that the diagonal has positive entries and at least a standard deviation of 0.01 in the normalized image domains. Therefore, we combine four losses using the stop gradient operator $| \cdot |$ with factors which change over time. Figure 11 shows the gradual increase and decrease losses factors over time to finally train the joint distribution. The whole training lasts 60000 steps where the number of steps per epoch is 1511 for CelebA, 1143 for Flying Chairs, and 195 for MNIST.

$$\mathcal{L} = \lambda_M \, \mathcal{L}_M + \lambda_D \, \mathcal{L}_D + \lambda_{LR} \, \mathcal{L}_{LR} + \lambda_{LRD} \, \mathcal{L}_{LRD} \tag{35}$$

$$\mathcal{L}_{M} = \log \mathcal{N}\left(\mu\left(x\right), I\right) \tag{36}$$

$$\mathcal{L}_{D} = \log \mathcal{N}\left(\lfloor \mu\left(x\right) \rfloor, D\left(x\right)\right) \tag{37}$$

$$\mathcal{L}_{L} = \log \mathcal{N}\left(\left\lfloor \mu\left(x\right) \right\rfloor, I + P\left(x\right) P^{\mathsf{T}}\left(x\right)\right)$$
(38)

$$\mathcal{L}_{LRD} = \log \mathcal{N} \left(\lfloor \mu \left(x \right) \rfloor, D \left(x \right) + P \left(x \right) P^{\mathsf{T}} \left(x \right) \right)$$
(39)

We run all experiments on a single Quattro RTX8000 NVIDIA GPU with 48GB RAM.