

Generalized Product-of-Experts for Learning Multimodal Representations in Noisy Environments

Abhinav Joshi
ajoshi@cse.iitk.ac.in
IIT Kanpur
Kanpur, India

Naman Gupta*
namang@cse.iitk.ac.in
IIT Kanpur
Kanpur, India

Jinang Shah*
jinang.iitk@gmail.com
IIT Kanpur
Kanpur, India

Binod Bhattarai†
b.bhattarai@ucl.ac.uk
University College London
London, U.K.

Ashutosh Modi‡
ashutoshm@cse.iitk.ac.in
IIT Kanpur
Kanpur, India

Danail Stoyanov‡
danail.stoyanov@ucl.ac.uk
University College London
London, U.K.

ABSTRACT

A real-world application or setting involves interaction between different modalities (e.g., video, speech, text). In order to process the multimodal information automatically and use it for an end application, Multimodal Representation Learning (MRL) has emerged as an active area of research in recent times. MRL involves learning reliable and robust representations of information from heterogeneous sources and fusing them. However, in practice, the data acquired from different sources are typically noisy. In some extreme cases, a noise of large magnitude can completely alter the semantics of the data leading to inconsistencies in the parallel multimodal data. In this paper, we propose a novel method for multimodal representation learning in a noisy environment via the generalized product of experts technique. In the proposed method, we train a separate network for each modality to assess the credibility of information coming from that modality, and subsequently, the contribution from each modality is dynamically varied while estimating the joint distribution. We evaluate our method on two challenging benchmarks from two diverse domains: multimodal 3D hand-pose estimation and multimodal surgical video segmentation. We attain state-of-the-art performance on both benchmarks. Our extensive quantitative and qualitative evaluations show the advantages of our method compared to previous approaches.

KEYWORDS

Multimodal Representations; Multimodal Fusion; Cross-modal Processing; Deep Learning Architectures; Machine Learning

1 INTRODUCTION

Humans interact with the real world by conveying and perceiving information using multiple modalities. For example, when two people talk to each other, besides the primary modality of verbal communication via context (text), they also use additional modalities like the tone of speech (audio) and facial/hand gestures (video). The information from these varied modalities may either overlap or complement each other. Moreover, the signals captured in the real world are often noisy. Humans in real-world interaction may not be able to capture signals from all the modalities efficiently,

often leading to noisy signals. However, humans tend to fuse the noisy information efficiently and learn about one modality from another modality. For example, the emotions of a person speaking in an unknown language can be predicted using their voice tone and physical gestures without any context information. Based on this intuition, researchers have focused on developing methods for exploiting shared information between the different modalities for self-supervision and combining the complementary information to improve a machine learning model’s generalization capability. The area of *Multimodal Representation Learning* (MRL) involves fusing the information coming from varied sources to learn representations that are robust and generalizable in different settings (for example, the case of a missing modality). However, the existing methods in MRL often tend to assume the training dataset to be noise-free. In this work, we propose a model that efficiently handles the noise present in different modalities in the dataset and compare it with the existing modality fusing mechanisms.

In recent years, generative models like Variational Auto-Encoders (VAE)s have attracted colossal interest in modality fusing mechanisms [11, 20, 23, 30]. The ability of VAEs to create information bottlenecks for the available modalities makes it easier to fuse information coming from different sources. The learned unimodal posteriors are combined using a fusing mechanism, for example, some of the widely popular fusing mechanisms include Product of Experts (PoE) [37], and Mixture of Experts (MoE) [31]. PoE multiplies unimodal posteriors, while MoE sums them up. Both of these approaches have their own merits and demerits, and we refer the reader to [35] for details. In particular, compared to MoE, PoE can aggregate any subset of modalities, which provides an efficient way of dealing with the missing modalities. This property has attracted its usage in various tasks such as multimodal 3D hand-pose estimation [3], and multimodal-GAN to generate naturally realistic images [15]. In this paper, we, too, use a variant of a PoE method for learning multimodal representations in a noisy setting.

Another noteworthy component in mixing information from multiple modalities is dynamically deciding the importance of a modality sample. Humans tend to figure out the noisy signals from a specific modality easily and fuse the information by giving less weightage to the modalities with the noisy signals. For example, a person parking a vehicle with a faulty rear-view camera (noisy

*Both authors contributed equally to this research.

†Corresponding author

‡Both authors were senior supervisors

image) will rely more on the rear-view mirrors for making parking decisions. Or a person communicating via video conferencing where the audio signals are noisy will shift their attention towards lip reading from the visuals to understand the speaker efficiently. Dynamically deciding the credibility of a sample from a modality makes the fusion of multimodal information more efficient and robust in such noisy cases. Considering the importance of dynamic weightage in modality fusing mechanisms, we propose to estimate the weightage contribution for each of the modalities in our architecture.

PoE multiplies unimodal posteriors assuming a uniform contribution from every modality. It results in a peaky joint distribution if all marginal posteriors have high density and vice-versa, as shown in Figure 1, where, the green and the red curves show the distribution of two different modalities, and a PoE combination gives a dashed black curve. This is an optimal combination if the information from each modality is equally credible. However, in a real-world scenario, this may not hold. Various types of equipment/sensors, such as depth sensors, RGB cameras, event cameras, LIDAR, etc., are used to capture different modalities. Equipment encounters different levels of noise asynchronously. The varying noise in each source corrupts it differently, and sometimes it can destroy the semantics of that modality and thus affecting the combination of modalities. Consequently, simply combining the posterior of each modality can be detrimental and eventually deteriorate the performance. Thus, there is a need to assess the credibility of information in every constituent modality, reweigh the parameters based on the credibility, and fuse them. For example, the effects of reweighing of unimodal posteriors on their joint posterior are shown in Figure 1 via the dashed blue curves. The demo¹ shows how introducing the alpha parameter to PoE helps reweigh the contribution of two modalities, allowing the freedom to appoint different weightage to different modalities dynamically.

Inspired by recent work on re-scaling Gaussian processes [4], we propose training different networks for each modality to check the credibility of the data coming from each source. In particular, we propose a cross-modal VAE based architecture (§3) for learning latent representation for each modality and then using a generalized product of experts for taking a weighted combination of the latent representations. We evaluate the proposed method on two tasks: multimodal hand pose estimation and multimodal semantic segmentation in surgical videos. Since both datasets are synthetic, to emulate a real-world setting, we add different noise levels to different modalities and perform an extensive set of experiments on the corrupted input. Our results show improvements over the existing PoE-based methods and attain state-of-the-art performance on the tasks.

2 RELATED WORKS

Multimodal Generative Learning: For the case of two modalities, several variants of variational autoencoders (VAEs) [17, 19] have been proposed to learn generative models for a uni-directional conditional inference. However, we are more interested in approaches that learn joint latent space for better modeling of data distribution, which can be used for conditional inference interchangeably across

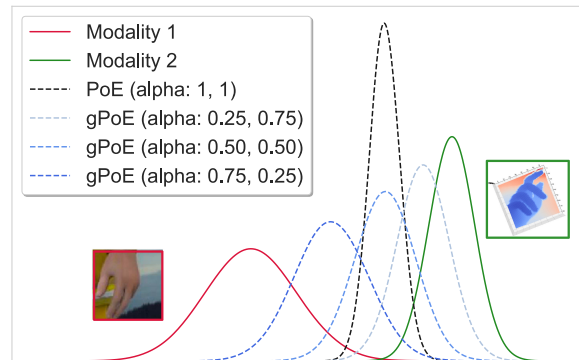


Figure 1: The figure shows a comparison between the joint distribution obtained by Product-of-Experts (PoE) and Generalized Product-of-Experts (gPoE) while fusing two modalities. The alpha introduced in gPoE helps scale the freedom to appoint different weightage to different modalities dynamically. The blue curves show the multiple distributions obtained with varying weightage (alpha) values. The darker blue shade (alpha: 0.75, 0.25) gives more weightage to Modality 1 (red) whereas the lighter blue shade (alpha: 0.25, 0.75) gives more weightage to Modality 2 (green).

all modalities. Joint multimodal variational autoencoder (JMVAE) proposed in Suzuki et al. [36] attempts to learn a joint distribution explicitly but by training separate inference networks for each possible subset of present modalities, which becomes intractable as more modalities get involved. Wu and Goodman [37] introduces multimodal variational autoencoder (MVAE), which leverages Product-of-Experts (PoE), whereas Multimodal Mixture-of-Experts VAE (MMVAE), as proposed by Shi et al. [31], uses a Mixture of Experts (MoE) to learn joint distribution efficiently. These approaches only measure their robustness in terms of their ability to handle missing modalities during inference, which is necessary but, in turn, do not attempt to address their limitations against noisy multimodal input.

Dealing with Noise: To address the issue of noisy inputs in multimodal settings, numerous approaches have been proposed [6, 22, 34] for specific tasks, but they largely remain restricted to discriminative models. Thus, leaving a gap in research to address limitations of current approaches focused on learning a joint distribution in the presence of noisy multimodal inputs. Since PoE has been widely used and documented for multimodal generative networks across various tasks, we use PoE to investigate further how noisy inputs limit the performance for specific tasks of 3D Hand Pose Estimation and Semantic Segmentation.

3D Hand Pose Estimation: Generative methods attempt to learn the distribution of the hand model from the given observations and remain highly susceptible to initialization. Although using depth or 3D data (especially 3D point clouds [10, 21]) provides the most accurate results [25, 26], their availability usually remains uncertain during training or inference. Thus, recent works Cai et al. [3], Spurr

¹<https://www.desmos.com/calculator/l4y75hedez>

et al. [33], Yang et al. [38] leverage depth information for training while restricting testing exclusively with RGB images. While Cai et al. [3] utilizes rendered depth maps from poses to regularize RGB image-based training, Spurr et al. [33] proposes a VAE-based method to learn a shared latent space, which suffers from its alternating training strategy for different modalities and a relatively slower convergence. Yang et al. [38] proposes a multimodal VAE-based method, where, authors rather attempt at aligning latent spaces of diverse modalities, including 3D poses, point clouds, and heat maps using PoE. Rather than forcing learning of a joint latent space, aligning latent spaces from different modalities allows a much faster convergence and also a better handling with non-corresponding data from other available modalities.

Semantic Segmentation: As semantic segmentation plays a vital role in medical diagnosis and treatment, the particular domain has been well studied. In recent years, there have been numerous deep learning-based approaches suggested to utilize the multimodal data for medical image segmentation as detailed in Zhou et al. [42]. Yet, semantic segmentation using multimodal generative learning and its analysis with noisy inputs remains an understudied area for the medical domain. We, thus, consider a recently published Surgical Video-Sim2Real dataset [29] for semantic segmentation with multiple modalities, which includes RGB images and depth maps. More details about the dataset are mentioned in (§4.1).

3 METHODOLOGY

For our setting, we consider a data acquisition pipeline with multiple sensors for capturing information tuples (e.g., audio, video, text) emanating from multiple sources. Since information from every modality can be noisy, it can lead to corruption of the sample. Consider the acquisition of a sample instance i having M different modalities: $\{m_1^{(i)}, m_2^{(i)}, \dots, m_M^{(i)}\}$. Consider the probability of noise in respective sensors to be $\{p_{m_1}, p_{m_2}, \dots, p_{m_M}\}$, with increasing number of sensors, the probability of multiple noisy modalities occurring simultaneously decreases exponentially $\prod_{k=1}^M p_{m_k} \ll 1$. In contrast, the number of corrupted data samples increases with the number of sensors. The primary assumption of clean, noise-free data samples from all the modalities limits the current deep learning architectures to use corrupted data samples when even one modality is noisy. Recently, a wide variety of approaches (using MVAE [37]) have tried to efficiently handle the missing modalities by aligning the learned latent spaces. However, the major limitation of such an approach is the underlying assumption that the available modalities are always noise-free.

Moreover, PoE giving equal weightage to all modalities results in a distribution dominated by highly confident modality-specific experts. For noisy examples, the experts produced by the encoders might not represent the required task-specific underlying data distribution. For such noisy cases, it is better to rely on noise-free experts. To control the credibility of each of the modalities, we introduce the use of *Generalized Product-of-Experts* [4], where the information present in different modalities controls the contribution of an expert dynamically. In the next section, we describe the architecture details. We start with the formulation of standard-VAE and extend it to crossmodal domain. Further, we introduce the use

of PoE for mixing information from multiple modalities. In the end, we formulate the generalized PoE that dynamically captures every modality’s contribution to help learn multimodal representations in a noisy environment effectively.

3.1 Architecture

Standard VAEs are encoder-decoder based generative architectures that maximise the evidence lower bound (ELBO) of the data log-likelihood. The amortized variational inference scheme with reparameterization trick [19, 28] helps formulating the approximate posterior $q(\phi(z|x))$ and the likelihood $p_\theta(x|z)$ distributions using deep neural networks with parameters ϕ and θ respectively.

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathcal{D}_{KL}[q_\phi(z|x) \| p(z)] \quad (1)$$

where \mathcal{D}_{KL} represents the Kullback-Leibler divergence and $p(z)$ represents a prior distribution which can vary from a standard Gaussian [19] to more expressive priors like normalizing flows [2, 12, 18].

Crossmodal-VAE: A standard VAE can be extended for modelling crossmodal inference, i.e. translating information from one modality to another. Consider two modalities, $m_{input} = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ and $m_{target} = \{y^{(1)}, y^{(2)}, \dots, y^{(n)}\}$, the ELBO in crossmodal VAE is minimized to generate the target modality ($y^{(i)}$) from its corresponding pair in input modality ($x^{(i)}$).

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(y|z)] - \mathcal{D}_{KL}[q_\phi(z|x) \| p(z)] \quad (2)$$

Further, for our case we assume the presence of other auxiliary modalities $M_{aux} = \{m_1, m_2, \dots, m_N\}$ during training and extend the minimization objective for crossmodal VAE. Thus, we additionally have $M_{aux} = \{m_1, m_2, \dots, m_N\}$, and the total number of available modalities becomes $N + 2$ including the primary input and output modalities.

$$M = \{m_1, m_2, \dots, m_N, m_{input}, m_{target}\} \quad (3)$$

Generalized Product-of-Experts: For mixing the features from multiple modalities, existing approaches [9, 37] use the product of experts (PoE) for computing the joint latent representation, which is proportional to the individual modality distributions.

$$P(z) = \frac{1}{Z} \prod_i p_i(z) \quad (4)$$

$$\mu_{PoE}(z) = \left(\sum_i \mu_i(z) T_i(z) \right) \left(\sum_i T_i(z) \right)^{-1} \quad (5)$$

$$\Sigma_{PoE}(z) = \left(\sum_i T_i(z) \right)^{-1} \quad (6)$$

where $T_i(z) = \Sigma_i^{-1}(z)$ is the precision of the i^{th} Gaussian expert in z , and Z is the normalization constant. PoE generates distribution dominated by highly confident experts compared to the less confident ones. The presence of noise in the dataset can cause an expert to produce erroneously low predicted variance along a latent dimension which further develops a strong bias in the joint predictions. To overcome this issue, we propose using a generalized formulation of the product of experts (gPoE) [4] which introduces a

weighing mechanism for scaling down such overconfident experts as formulated below.

$$P(z) = \frac{1}{Z} \prod_i p_i^{\alpha_i(z)}(z) \quad (7)$$

$$\mu_{gPoE}(z) = \left(\sum_i \mu_i(z) \alpha_i(z) T_i(z) \right) \left(\sum_i \alpha_i(z) T_i(z) \right)^{-1} \quad (8)$$

$$\Sigma_{gPoE}(z) = \left(\sum_i \alpha_i(z) T_i(z) \right)^{-1} \quad (9)$$

For handling the noisy information present in the dataset, we formulate the modality-specific scaling factors α_i as a direct function of input modalities ($\mathcal{F}(m_i)$). We estimate the scaling parameters α for each dimension in the latent space using independent modality encoders. An overview of our architecture is shown in Figure 2, where features from different modalities are concatenated and passed through standard feed-forward layers. For a normalized scaling of predicted gaussian distributions, we use a softmax function to distribute the importance across available modalities such that for each latent dimension $\sum_{i=1}^{N+2} \alpha_i = 1$. Further, using gPoE, we define the joint representation of available modalities as,

$$z_{joint} \sim \mathcal{N}(\mu_{gPoE}, \Sigma_{gPoE}) \quad (10)$$

$$\mathcal{L}(\Theta, \Phi) = \mathcal{L} \left(\{ \theta^{(y)}, \theta^{(m_1)}, \dots, \theta^{(m_N)} \}, \{ \phi^{(x)}, \phi^{(m_1)}, \dots, \phi^{(m_N)} \} \right) \quad (11)$$

For generating the target modality y , the reconstruction loss is defined as follows:

$$\begin{aligned} \mathcal{L}_{target} &= \mathbb{E}_{z_{joint}} [\log p_{\theta}(y | z_{joint})] + \mathbb{E}_{z_x} [\log p_{\theta}(y | z_x)] \\ &+ \sum_{i=1}^N \mathbb{E}_{z_{m_i}} [\log p_{\theta}(y | z_{m_i})] \end{aligned} \quad (12)$$

Similarly, for generating other auxiliary modalities, the reconstruction loss is,

$$\begin{aligned} \mathcal{L}_{aux} &= \sum_{k=1}^N \left(\mathbb{E}_{z_{joint}} [\log p_{\theta}(m_k | z_{joint})] \right. \\ &+ \mathbb{E}_{z_x} [\log p_{\theta}(m_k | z_x)] \\ &+ \left. \sum_{i=1}^N \mathbb{E}_{z_{m_i}} [\log p_{\theta}(m_k | z_{m_i})] \right) \end{aligned} \quad (13)$$

Furthermore, the Kullback-Leibler divergence for all the latent spaces is,

$$\mathcal{L}_{KL} = \mathcal{D}_{KL}[q_{\phi_x}(z_x | x) \| p(z)] + \sum_{i=1}^N \mathcal{D}_{KL}[q_{\phi_{m_i}}(z_{m_i} | m_i) \| p(z)] \quad (14)$$

Adding the losses in 12, 13 14, we obtain the minimising objective of our formulation

$$\mathcal{L}(\Theta, \Phi) = \mathcal{L}_{target} + \mathcal{L}_{aux} + \beta \mathcal{L}_{KL} \quad (15)$$

where β is a hyperparameter [14] used to balance the loss between reconstruction and Kullback-Leibler divergence. Note that the objective of our work is to leverage the information present in corresponding modality pairs during training that helps make unimodal inference robust towards noisy data samples.

4 EXPERIMENTS

4.1 Datasets

To evaluate our method in the noisy environment, we choose two synthesized datasets from different domains, which provide multiple modalities that help formulate a unimodal prediction task. The two publicly available datasets include the Rendered Hand Pose Dataset (RHD) [43] for predicting 3D hand pose from RGB images and Surgical Video-Sim2Real Dataset [29] for predicting segmentation masks from RGB images.

RHD is a synthesized dataset of rendered hand images with 320×320 resolution. The dataset was created from 20 subjects, performing 39 actions that were rendered using Blender [8]. It consists of 41258 training and 2728 testing samples. The dataset provides RGB images, along with the corresponding depth maps, segmentation masks, and 3D keypoints for hand pose prediction.

Surgical Video-Sim2Real dataset consists of 21000 randomly sampled views for 7 simulated, surgical 3D scenes, which were rendered using the liver meshes obtained from the 3D-IRCADb dataset [32] composed of 3D CT scans of liver. The surgical dataset provides realistic translation views and the corresponding depth maps, camera poses, and segmentation masks. The provided segmentation mask consists of 5 classes: liver, fat/stomach, abdominal wall, gallbladder, and ligament.

4.2 Preprocessing

For hand pose prediction, we follow a preprocessing scheme similar to Yang et al. [38] and crop the hand region of the image. Further, using the available hand segmentation mask, we extract the corresponding hand region from the depth image and convert it to point clouds using the provided camera intrinsic parameters. For segmentation prediction in the Surgical Sim2Real dataset, we center-crop the RGB image for visual modality and resize it to 256×256 . The available depth images are converted to point clouds using the provided camera intrinsic parameters.

4.3 Noise Simulation

Modeling imaging sensor noise is a fundamental problem in image processing. The practical application pipeline of camera sensors is highly complex and usually has different modules that cause various types of noise inductions in an acquired image. Many existing approaches have proposed statistical noise models to simulate real-world noise in the acquired images. Recently, a wide variety of deep-learning methods have been introduced to simulate the real-world noise [1, 5, 7, 16, 24, 41]. In contrast, various methods use simple diffusion models to introduce Gaussian noise in images which are further used to make the deep learning architectures more robust towards a noisy environment. In general, noise in a visual modality can be interpreted as information loss in terms

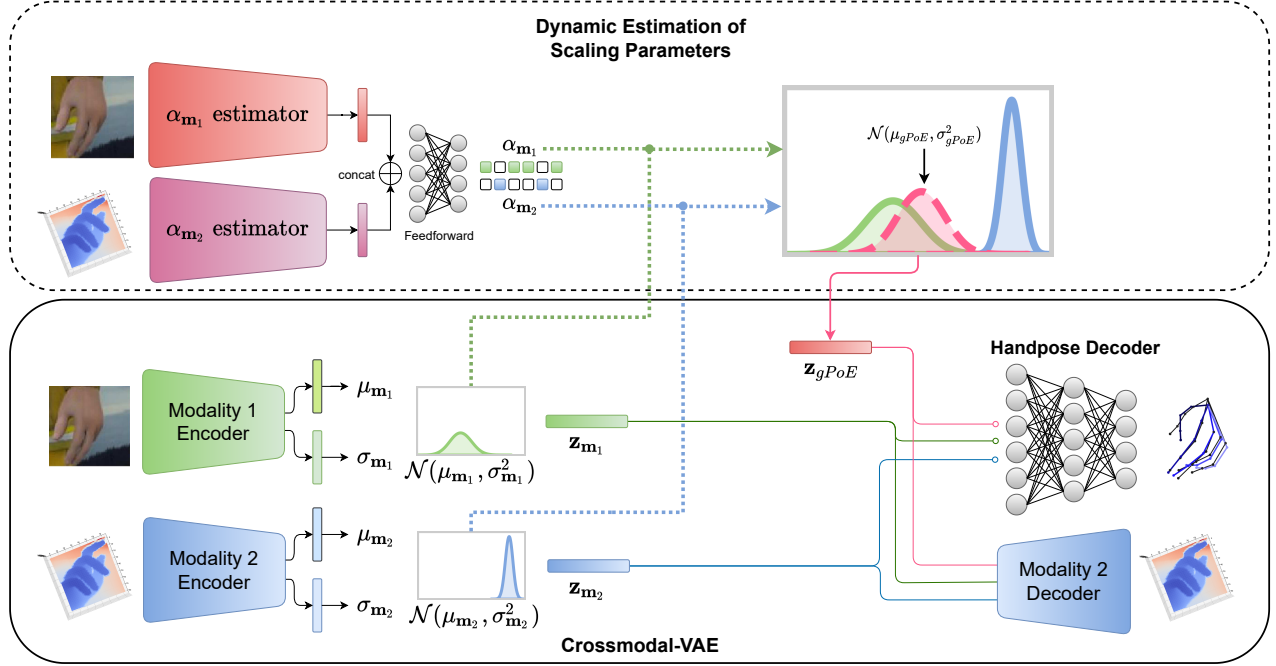


Figure 2: Crossmodal-VAE architecture with gPoE used for 3D hand pose estimation. The (top) estimated scaling parameters $\{\alpha_1, \alpha_2\}$ are used with their corresponding learned modal distributions to identify the aligned distribution using gPoE as shown in Eq.(8) and Eq.(9). The proposed framework although shown for two modalities here can be easily extended to any N number of modalities present.

of pixels. To simulate noise for our synthetic environment, we choose a simple scheme of corrupting the pixels by replacing the pixel values with random values. Figure 3 shows an example of an image with various percentages of pixel corruption added as noise. In real-world applications, sensors like LiDAR, RADAR, etc., have shown poor performance in adverse weather conditions resulting in deviated computation of point clouds. For generalizing our methods to multiple modalities, we consider the auxiliary modality to be noisy and simulate noise using standard Gaussian noise applied to the observed values. Figure 3 shows an example of Gaussian noise added to a depth point cloud sample.

4.4 Evaluation Metrics

For evaluation of the predicted hand poses in the RHD dataset, we use the standard metrics [38]: mean end-point-error (EPE) and area under the curve (AUC) on the percentage of correct key points (PCK) curve (higher is better). EPE is computed using the euclidean distance between the predicted and ground-truth 3D keypoints of hand joints (lower is better). PCK represents the percentage of predicted key points that fall within certain error thresholds of EPE. To evaluate the predicted segmentation masks in the Surgical Sim2Real dataset, we use standard IoU and F1 metrics. IoU or Jaccard similarity coefficient computes the size of the intersection divided by the size of the union of two label sets, and F1 score is interpreted as a harmonic mean of the precision and recall, where an F1 score reaches its best value at 1 and worst score at 0.

4.5 Implementation Details

For a fair comparison and showing the effect of gPoE over PoE over the noisy environment, we follow the encoder and decoder architectures inspired from Yang et al. [38]. We consider RGB as our primary modality and use a pre-trained ResNet-18 [13] with fully connected layers to generate μ_{RGB} and σ_{RGB} . For encoding point clouds, we use the ResPCL network [21], for decoding point clouds, we use the Folding-Net decoder [40]. We use a set of linear layers in the RHD dataset to decode hand keypoints and standard DC-GAN architecture [27] to predict segmentation maps in the surgical dataset, keeping the rest of the architecture the same in both the datasets.

5 RESULTS-ANALYSIS

5.1 Quantitative Results

To validate the effectiveness of our proposed architecture in a noisy setting, we compare the performance of unimodal encoder models trained in similar noisy settings over various ranges of data corruption. The increasing deviation in performance shown in Figure 4 (first from the left) represents the robustness of generalized PoE towards noisy data samples. The difference in performance increases almost linearly with the amount of data corruption depicting the performance gain when the larger proportion of the test dataset is noisy. Moreover, to judge the efficacy of the proposed architecture, we also test the unimodal encoders on a varying range of pixel corruption. For this setting, we corrupt all the test set examples

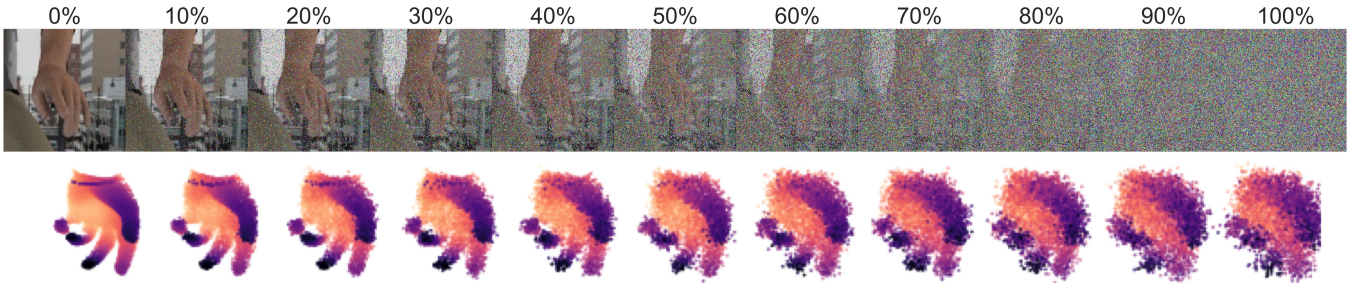


Figure 3: (top) Random pixel noise added to image modality. (bottom) Gaussian noise (increasing left to right) added to a point cloud.

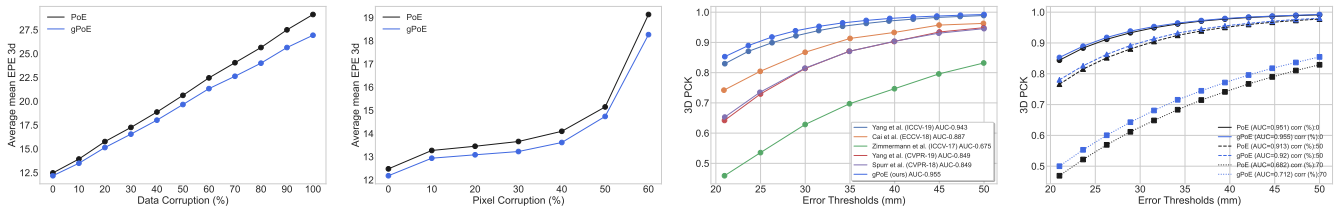


Figure 4: The first two figures from the left show the Average mean EPE 3D (mm) comparison between PoE and gPoE for varying levels of data corruption (first) and pixel corruption (second), respectively, on the RHD dataset. The third figure shows the 3D PCK vs. Error Thresholds (mm) plot comparison against prior works, and the fourth figure shows the performance across different noise levels on the RHD dataset (Zoom in for a better view).

Pixel Corruption	Method	IOU	F1
0%	PoE	0.6272	0.7692
	gPoE	0.7138	0.8323
50%	PoE	0.6221	0.7653
	gPoE	0.7078	0.8281
70%	PoE	0.6128	0.7576
	gPoE	0.6866	0.8131

Table 1: Comparison of gPoE and PoE for semantic segmentation on the surgical dataset with varying degree of pixel corruption.

with a fixed percentage of pixel corruption and observe the performance for both architectures. Figure 4 (second from the left) shows the Average Mean EPE varying across the pixel corruption percentage. As observed, gPoE outperforms PoE over the entire range of pixel corruption, highlighting the generalizability of gPoE in a noisy natural environment.

Furthermore, to quantitatively judge the effect of pixel corruption, we compare the performance of PoE and gPoE using the PCK curve for a considerable amount of pixel corruption. To show this effect, we choose 3 settings, 0%, 50%, and 70% pixel corruption. Figure 4 (fourth from the left) shows the observed PCK curves for these settings. The effectiveness of our method is observable when the noise in the dataset increases. The significant performance gap

between the 50% and 70% pixel corruption highlights the adverse effects of noise in the test environment. The increasing difference in AUC values between gPoE and PoE (3% difference in 70% pixel corruption) shows the robustness of our system in a noisy environment. The encoder-decoder models trained with gPoE are less prone to noise for unimodal inference. In contrast, the encoder-decoder trained with PoE performs inadequately when pixel corruption increases in the test set.

We also test our approach by experimenting with the Surgical Video-Sim2Real dataset to predict pixel-level segmentation masks using RGB images and corresponding depth information. In Table 1, we provide a comparison between PoE and gPoE for three different levels of pixel corruption on the surgical dataset for semantic segmentation. In this setting, gPoE outperforms PoE by a significant margin showing the robustness towards noisy samples present during training. The introduction of dynamic scaling in gPoE not only helps in training the unimodal encoders simultaneously but also results in unimodal encoders, which are more robust against real-world noisy data samples.

5.2 State-of-the-Art Comparison:

Though the main focus of our work is tackling a realistic noisy setting, for fair analysis, we briefly highlight the comparison with the other existing works on hand pose prediction from RGB images. To validate our results on the RHD dataset, we compare our method with the existing 3D hand pose prediction approaches. Spurr et al. [33] and Yang et al. [38] are related to our method as they use VAE based architectures for hand pose prediction. Out of the two

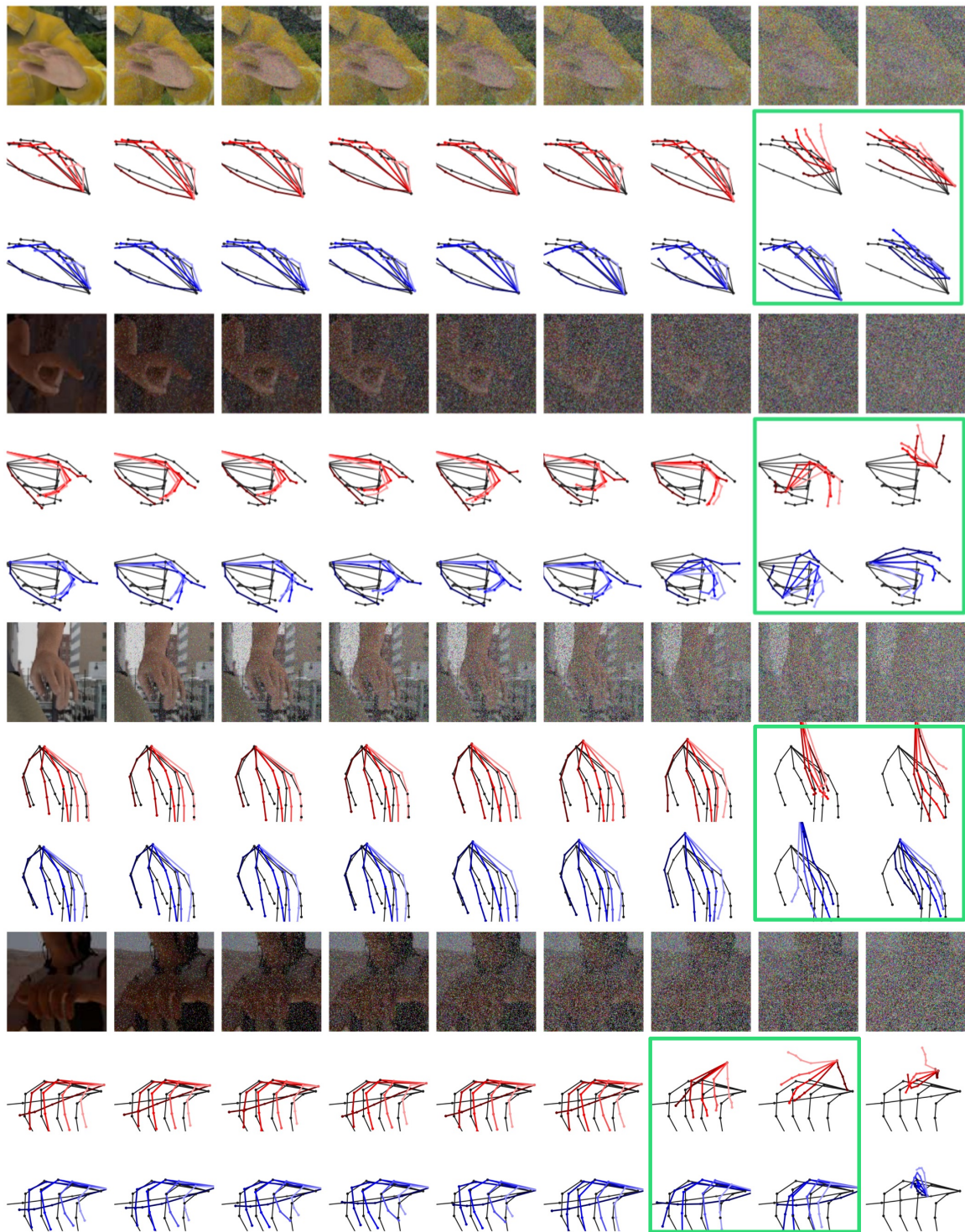


Figure 5: Hand Pose Estimation using PoE (red) and gPoE (blue) on the RHD dataset. The first row in each sample shows the input RGB image, and the ground truth keypoints are shown in (black). The second and third row shows the predicted masks from PoE and gPoE, respectively. Noise intensity increases going from left to right. Green boxes highlight significant improvement in predicted hand poses for extremely noisy inputs.

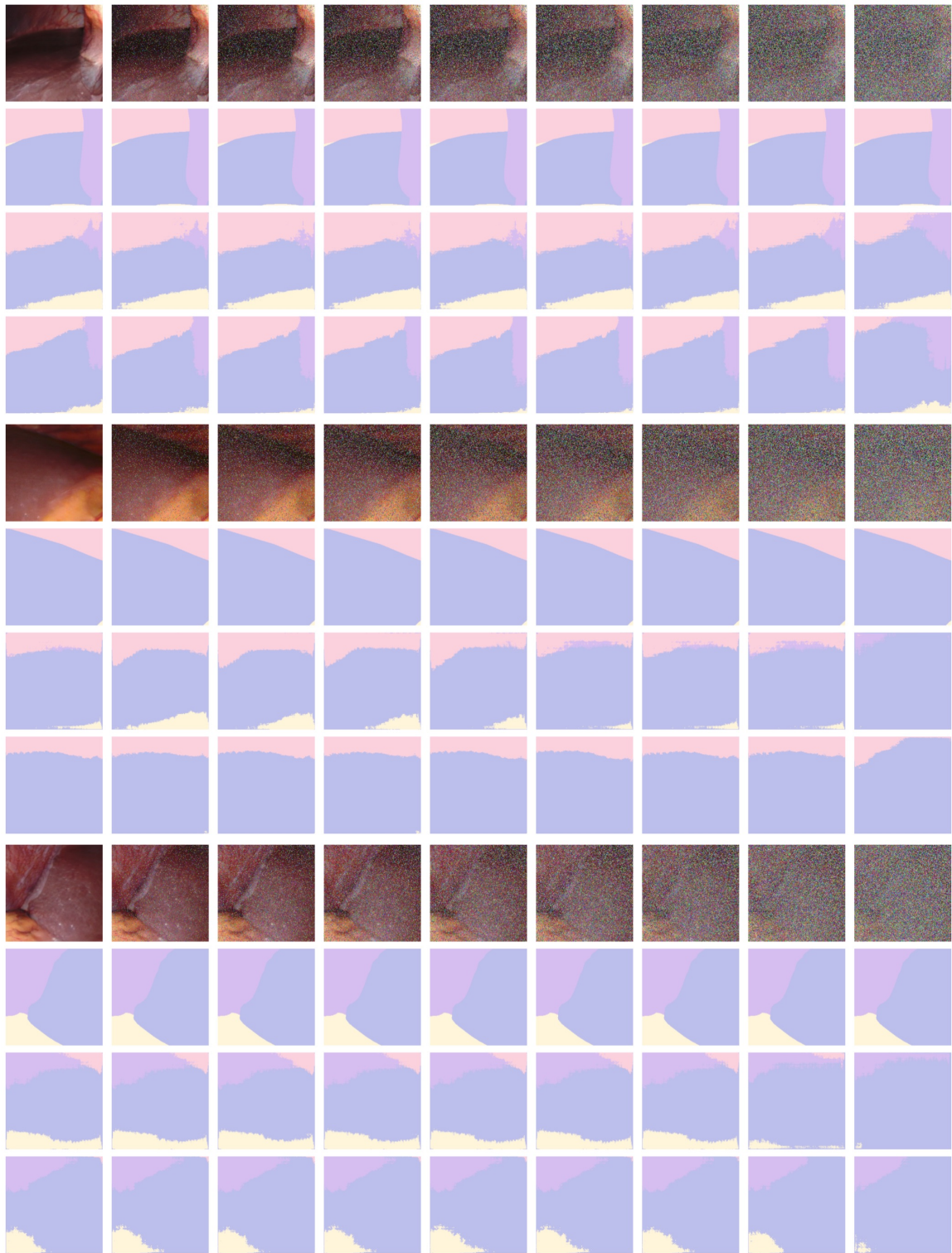


Figure 6: Surgical Semantic Segmentation comparison on architectures trained with PoE and gPoE in a noisy setting (noise increasing from left to right). The first row in a sample represents the input image. The second row represents the ground truth segmentation map, third and fourth row shows the predicted masks from PoE and gPoE, respectively.

Method	mean EPE (mm)AUC	
Zimmerman et al. [44] (ICCV-17)	-	0.675
Spurr et al. [33] (CVPR-18)	19.73	0.849
Cai et al. [3] (ECCV-18)	-	0.887
Yang et al. [39] (CVPR-19)	19.95	0.849
Yang et al. [38] (ICCV-19)	13.14	0.943
gPoE (ours)	12.18	0.955

Table 2: Comparing gPoE with prior works for 3D Hand Pose Estimation on RHD dataset. Lower mean end-point-error (EPE) is better, higher AUC is better.

approaches, Yang et al. [38] come closer to our approach and act as our baseline for highlighting the drawbacks of the product of experts in a noisy environment. Table 2 compares the related works on 3D hand pose prediction in the RHD dataset. We also compare the PCK curve obtained by our method with the existing works. Figure 4 (third from the left) shows a comparison with the existing works on the RHD dataset. The encoder-decoder models trained using gPoE show comparable performance on unimodal prediction when compared to the baseline Yang et al. [38].

5.3 Modality Fusing Mechanism

In order to analyze if the proposed modality fusing mechanism plays a significant role in learning unimodal encoders robust to the noise present in the collected dataset, we compare the proposed training mechanism with other existing fusing mechanisms. Since our method builds upon cross-modal VAE architecture, we consider two widely favored methods for fusing information in the latent space, Product-of-Experts, and Mixture-of-Experts. For a fair comparison between different fusing mechanisms, we only change the fusing mechanism module and keep the modality encoders and decoders same in all the settings.

Product-of-Experts (PoE) proposes a fusing mechanism where the product of the present unimodal posteriors helps formulate the joint posterior. The detailed formulation for PoE can be found in the equations 4, 5 and 6. For an unbiased comparison, we use the learning objective similar to our architecture (equation 15), the only difference being the method of computing joint posterior.

Mixture-of-Experts (MoE) as proposed in mixture-of-experts multimodal variational autoencoder (MMVAE) [31], factorizes joint posterior by weighted averaging of individual posteriors. Rather than learning individual weights for the modalities, MMVAE suggests giving equal weightage to each modality present to avoid a dominant-modality issue as in PoE. More details can be found in MMVAE Shi et al. [31]. Keeping the training strategy same, the learning objective (ELBO) for MoE as proposed in MMVAE is formulated as follows:

$$\begin{aligned}
 \text{ELBO} &= \frac{1}{M} \sum_i^M \mathbb{E}_{z_{m_i} \sim q_{\phi_{m_i}}(z|m_i)} \left[\log \frac{p_{\Theta}(z_{m_i}, m_{1:M})}{q_{\Phi}(z_{m_i} | m_{1:M})} \right] \\
 &= \frac{1}{M} \sum_i^M \left(\mathbb{E}_{z_{m_i} \sim q_{\phi_{m_i}}(z|m_i)} \left[\log p_{\Theta}(m_{1:M} | z_{m_i}) \right] \right. \\
 &\quad \left. - \text{KL} [q_{\Phi}(z_{m_i} | m_{1:M}) \| p(z)] \right] \quad (16)
 \end{aligned}$$

Moreover, to capture the effect of noise on various modality fusing mechanisms, we compare them on different degrees of pixel corruption. We choose no corruption as our baseline and compare them against two additional settings, with 25% pixel corruption and 50% pixel corruption. Table 3 shows the comparison results on the RHD dataset. The lower mean EPE with a high AUC score highlights the performance boost obtained using the proposed generalized-Product-of-Experts as a modality fusing mechanism. We perform the same set of inferences on the segmentation mask prediction task in the Surgical Video-Sim2Real dataset. Table 4 shows the results of the fusing mechanisms on varying degrees of noise. The proposed modality fusing mechanism shows higher IOU and F1 scores in all the settings when compared to PoE and MoE in multiple degrees of noise present in the test set. Both the tables clearly illustrate the significance of the proposed fusing mechanism in learning robust unimodal encoders and decoders.

A noteworthy advantage of the proposed fusing mechanism is its ability to train unimodal predictors in a noisy dataset with extra computation cost only during training. During inference, the same architecture results in parameters more robust against noisy samples without any computational overhead.

5.4 Qualitative Results

For assessing the quality of predictions made by gPoE, we compare it with an architecture using standard Product of Experts for mixing the modalities. We train the architectures in the same noisy setting using multiple modalities and evaluate them using only the primary input modality (RGB images). Further, the learned latent spaces are tested using the increasing amount of noise in the RGB modality. Using the learned encoders for both the methods, z_{rgb} is generated and used as input to the learned hand pose decoder. Appendix Figure 5 shows the predicted hand poses using the learned latent space z_{rgb} in both the architectures. The predicted 3D keypoints are similar for both the architectures when pixel corruption is small. However, as the pixel corruption increases, gPoE outperforms PoE depicting the importance of dynamic scaling in aligning the latent spaces. For Surgical Video-Sim2Real dataset z_{rgb} is used to predict the segmentation using the learned segmentation decoder. Appendix Figure 6 shows the comparison between the architectures trained with PoE and gPoE. The predicted segmentations of gPoE are less noisy when compared to PoE, conveying the robustness for training in a noisy setting.

Modality Fusing Mechanism	No Pixel Corruption		25% Pixel Corruption		50% Pixel Corruption	
	mean EPE (mm)	AUC	mean EPE (mm)	AUC	mean EPE (mm)	AUC
Mixture-of-Experts (MoE)	12.48	0.950	14.12	0.927	15.01	0.914
Product-of-Experts (PoE)	12.47	0.951	13.50	0.938	15.14	0.913
Generalized-Product-of-Experts (gPoE)	12.18	0.955	13.14	0.942	14.73	0.918

Table 3: Comparing gPoE with other modality fusing mechanisms for 3d Handpose Estimation from RGB images in RHD [43] dataset.

Modality Fusing Mechanism	No Pixel Corruption		25% Pixel Corruption		50% Pixel Corruption	
	IOU	F1	IOU	F1	IOU	F1
Mixture-of-Experts (MoE)	0.4477	0.6143	0.4598	0.6260	0.4545	0.6208
Product-of-Experts (PoE)	0.6272	0.7692	0.6198	0.7636	0.6221	0.7653
Generalized-Product-of-Experts (gPoE)	0.7138	0.8323	0.7030	0.8249	0.7078	0.8281

Table 4: Comparing gPoE with other modality fusing mechanisms for predicting segmentation masks from RGB images in Surgical Video-Sim2Real Dataset [29] dataset.

6 CONCLUSION

This paper addresses the problem of mixing modalities in a real-world setting where training data captured from one or more modalities are noisy. We propose a novel method for multimodal representation learning in a noisy environment via the generalized product of experts. We test our architecture on two publicly available datasets, Rendered Hand Pose Dataset (RHD) and Surgical Video-Sim2Real, for unimodal prediction tasks in different domains. We show our method’s effectiveness in a noisy setting by qualitative and quantitative analysis on various noise levels. We observe that the architecture trained by our method leads to unimodal encoders that are more robust toward noisy data samples. Though the proposed method can be generalized to fuse information from N different modalities, we have only tested the architecture considering two input modalities. In future work, we plan to test our architecture on prediction tasks where more than two modalities are available.

7 ACKNOWLEDGMENTS

We would like to thank reviewers for their insightful comments. Ashutosh Modi is supported in part by SERB India (Science and Engineering Board) (SRG/2021/000768).

Binod Bhattarai and Danail Stoyanov are funded by in whole, or in part, by the Wellcome/EPSRC Centre for Interventional and Surgical Sciences (WEISS) (203145/Z/16/Z), Engineering and Physical Sciences Research Council (EPSRC) (EP/P012841/1), the Royal Academy of Engineering Chair in Emerging Technologies scheme, and EndoMapper project by Horizon 2020 FET (GA863146).

REFERENCES

- [1] A. Abdelhamed, Marcus A. Brubaker, and M. S. Brown. 2019. Noise Flow: Noise Modeling With Conditional Normalizing Flows. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* (2019), 3165–3173.
- [2] Rianne van den Berg, Leonard Hasenclever, Jakub M Tomczak, and Max Welling. 2018. Sylvester normalizing flows for variational inference. *arXiv preprint arXiv:1803.05649* (2018).
- [3] Yujun Cai, Lihao Ge, Jianfei Cai, and Junsong Yuan. 2018. Weakly-supervised 3D Hand Pose Estimation from Monocular RGB Images. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- [4] Yanshui Cao and David J. Fleet. 2014. Generalized Product of Experts for Automatic and Principled Fusion of Gaussian Process Predictions. *CoRR* abs/1410.7827 (2014). arXiv:1410.7827 <http://arxiv.org/abs/1410.7827>
- [5] Ke-Chi Chang, Ren Wang, Hung-Jin Lin, Yu-Lun Liu, Chia-Ping Chen, Yu-Lin Chang, and Hwann-Tzong Chen. 2020. Learning Camera-Aware Noise Models. In *Proceedings of European Conference on Computer Vision (ECCV)*.
- [6] Cheng Chen, Qi Dou, Yueming Jin, Hao Chen, Jing Qin, and Pheng-Ann Heng. 2020. Robust Multimodal Brain Tumor Segmentation via Feature Disentanglement and Gated Fusion. *CoRR* (2020). arXiv:2002.09708 <https://arxiv.org/abs/2002.09708>
- [7] Jingwen Chen, Jiawei Chen, Hongyang Chao, and Ming Yang. 2018. Image Blind Denoising with Generative Adversarial Network Based Noise Modeling. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2018), 3155–3164.
- [8] Blender Online Community. 2018. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam. <http://www.blender.org>
- [9] Imant Daunhawer, Thomas M. Sutter, Ricards Marcinkevics, and Julia E. Vogt. 2020. Self-supervised Disentanglement of Modality-Specific and Shared Factors Improves Multimodal Generative Models. In *GCPR*. 459–473. https://doi.org/10.1007/978-3-030-71278-5_33
- [10] Lihao Ge, Yujun Cai, Junwu Weng, and Junsong Yuan. 2018. Hand PointNet: 3D Hand Pose Estimation Using Point Sets. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8417–8426. <https://doi.org/10.1109/CVPR.2018.00878>
- [11] Eloy Geenjaer, Noah Lewis, Zening Fu, Rohan Venkatdas, Sergey Plis, and Vince Calhoun. 2021. Fusing multimodal neuroimaging data with a variational autoencoder. In *2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE, 3630–3633.
- [12] Will Grathwohl, Ricky TQ Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. 2018. Fjord: Free-form continuous dynamics for scalable reversible generative models. *arXiv preprint arXiv:1810.01367* (2018).
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [14] Irina Higgins, Loic Matthey, Arka Pal, Christopher P. Burgess, Xavier Glorot, Matthew M. Botvinick, Shakir Mohamed, and Alexander Lerchner. 2017. beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. In *ICLR*.
- [15] Xun Huang, Arun Mallya, Ting-Chun Wang, and Ming-Yu Liu. 2021. Multimodal Conditional Image Synthesis with Product-of-Experts GANs. *arXiv preprint arXiv:2112.05130* (2021).

- [16] Dong-Wook Kim, Jae-Ryun Chung, and Seung-Won Jung. 2019. GRDN: Grouped Residual Dense Network for Real Image Denoising and GAN-Based Real-World Noise Modeling. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)* (2019), 2086–2094.
- [17] Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. 2014. Semi-supervised Learning with Deep Generative Models. In *Advances in Neural Information Processing Systems*, Vol. 27.
- [18] Diederik P. Kingma, Tim Salimans, Rafal Józefowicz, Xi Chen, Ilya Sutskever, and Max Welling. 2016. Improving Variational Autoencoders with Inverse Autoregressive Flow. In *NIPS*. 4736–4744. <http://papers.nips.cc/paper/6581-improving-variational-autoencoders-with-inverse-autoregressive-flow>
- [19] Diederik P Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. arXiv:1312.6114 [stat.ML]
- [20] Mihee Lee and Vladimir Pavlovic. 2021. Private-shared disentangled multimodal vae for learning of latent representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1692–1700.
- [21] Shile Li and Dongheui Lee. 2019. Point-To-Pose Voting Based Hand Pose Estimation Using Residual Permutation Equivariant Layer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [22] Y. Liang, F. Ma, Y. Li, and S. Huang. 2021. Person Recognition with HGR Maximal Correlation on Multimodal Data. In *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE Computer Society, 1–8. <https://doi.org/10.1109/ICPR48806.2021.9413164>
- [23] Navonil Majumder, Soujanya Poria, Gangeshwar Krishnamurthy, Niyati Chhaya, Rada Mihalcea, and Alexander Gelbukh. 2019. Variational fusion for multimodal sentiment analysis. *arXiv preprint arXiv:1908.06008* (2019).
- [24] Ioannis Marras, Grigorios G. Chrysos, Ioannis Alexiou, Gregory G. Slabaugh, and Stefanos Zafeiriou. 2020. Reconstructing the Noise Variance Manifold for Image Denoising. In *ECCV*.
- [25] Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. 2017. V2V-PoseNet: Voxel-to-Voxel Prediction Network for Accurate 3D Hand and Human Pose Estimation from a Single Depth Map. *CoRR* (2017). arXiv:1711.07399
- [26] Markus Oberweger, Paul Wohlhart, and Vincent Lepetit. 2015. Hands Deep in Deep Learning for Hand Pose Estimation. *CoRR* (2015). arXiv:1502.06807 <http://arxiv.org/abs/1502.06807>
- [27] Alec Radford, Luke Metz, and Soumith Chintala. 2016. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. arXiv:1511.06434 [cs.LG]
- [28] Danilo Jimenez Rezende and Shakir Mohamed. 2015. Variational Inference with Normalizing Flows. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37 (Lille, France) (ICML '15)*. JMLR.org, 1530–1538.
- [29] Dominik Rivoir, Micha Pfeiffer, Reuben Docea, Fiona Kolbinger, Carina Riediger, Jurgen Weitz, and Stefanie Speidel. 2021. Long-Term Temporally Consistent Unpaired Video Translation From Simulated Surgical 3D Data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 3343–3353.
- [30] Yuge Shi, Brooks Paige, Philip HS Torr, and N Siddharth. 2020. Relating by contrasting: A data-efficient framework for multimodal generative models. *arXiv preprint arXiv:2007.01179* (2020).
- [31] Yuge Shi, N. Siddharth, Brooks Paige, and Philip H. S. Torr. 2019. Variational Mixture-of-Experts Autoencoders for Multi-Modal Deep Generative Models. arXiv:1911.03393 [stat.ML]
- [32] L. Soler, A. Hosteller, V. Agnus, A. Charnoz, I. Fasquel, I. Moreau, A. Osswald, M. Bouhadjar, and I. Marescaux. 2010. 3D image reconstruction for comparison of algorithm database: A patient specific anatomical and medical image database. <https://www.ircad.fr/research/data-sets/liver-segmentation-3d-ircadb-01/>
- [33] Adrian Spurr, Jie Song, Seonwook Park, and Otmar Hilliges. 2018. Cross-modal Deep Variational Hand Pose Estimation. arXiv:1803.11404 [cs.CV]
- [34] Viswanath P. Sudarshan, Uddeshya Upadhyay, Gary F. Egan, Zhaolin Chen, and Suyash P. Awate. 2021. Towards lower-dose PET using physics-based uncertainty-aware multimodal learning with robustness to out-of-distribution data. *Medical Image Anal.* 73 (2021), 102187. <https://doi.org/10.1016/j.media.2021.102187>
- [35] Thomas M Sutter, Imant Daunhawer, and Julia E Vogt. 2020. Generalized Multimodal ELBO. In *International Conference on Learning Representations*.
- [36] Masahiro Suzuki, Kotaro Nakayama, and Yutaka Matsuo. 2016. Joint Multimodal Learning with Deep Generative Models. *arXiv preprint arXiv:1611.01891* (2016).
- [37] Mike Wu and Noah Goodman. 2018. Multimodal Generative Models for Scalable Weakly-Supervised Learning. arXiv:1802.05335 [cs.LG]
- [38] Linlin Yang, Shile Li, Dongheui Lee, and Angela Yao. 2019. Aligning Latent Spaces for 3D Hand Pose Estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [39] Linlin Yang and Angela Yao. 2019. Disentangling Latent Hands for Image Synthesis and Pose Estimation. arXiv:1812.01002 [cs.CV]
- [40] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. 2018. FoldingNet: Point Cloud Auto-encoder via Deep Grid Deformation. arXiv:1712.07262 [cs.CV]
- [41] Zongsheng Yue, Qian Zhao, Lei Zhang, and Deyu Meng. 2020. Dual Adversarial Network: Toward Real-world Noise Removal and Noise Generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- [42] Tongxue Zhou, Su Ruan, and Stéphane Canu. 2019. A review: Deep learning for medical image segmentation using multi-modality fusion. *Array* 3-4 (2019), 100004. <https://doi.org/10.1016/j.array.2019.100004>
- [43] Christian Zimmermann and Thomas Brox. 2017. *Learning to Estimate 3D Hand Pose from Single RGB Images*. Technical Report. arXiv:1705.01389. <https://lmb.informatik.uni-freiburg.de/projects/hand3d/> <https://arxiv.org/abs/1705.01389>.
- [44] Christian Zimmermann and Thomas Brox. 2017. Learning to Estimate 3D Hand Pose from Single RGB Images. arXiv:1705.01389 [cs.CV]