

Tractable Probabilistic Multimodal Learning

Jian Vora Pranay Reddy Samala Siddharth Chandak
Stanford University

{jianv, pranayr, chandaks}@stanford.edu

Abstract

Given M modalities or views of a data-point in the form of $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_M$, we aim to learn a joint distribution $p_\theta(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_M)$ that allows for exact inference queries. This is useful as it allows us to sample from a conditional model which allows us to do things like text guided image synthesis and image captioning from the same model itself (with image, text as the modalities). We can also train using missing data, a common issue in problems with multiple modalities, by training marginals rather than the joint.

It is known that probabilistic circuits or sum-product networks [15] as a generative model allow for tractable evaluation of the above inference queries. We train independent auto-encoders for each modality and train a tractable generative model on the joint latent space which is simply the concatenation of the latent spaces of individual modalities. We show that our method generates 1.5 times better results (both quantitative and qualitative) on a variety of datasets such as (MNIST-SVHN, CelebA attributes, CUB-Captions) as compared to prior VAE based methods such as MMVAE [19] and MVAE [25] which cannot compute inference queries and only rely on variational approximations for the same.

1. Introduction

Over the past ten years, deep learning has evolved considerably. From simple tasks such as classification and regression, deep learning is now making strides towards solving more complex and creative endeavors. Developments in generative modeling techniques have enabled this growth and allowed the proliferation of computers for domains that were traditionally considered possible for only humans. Models such as variational auto-encoders (VAEs) [7], generative adversarial networks (GANs) [6], autoregressive models, diffusion models [2] have been tremendously successful, and have widely been used for performing a number of tasks including image generation [18] [16] and text generation [1]. While these methods have been quite powerful, each of these generative models is quite task-specific. For

instance, a generative model trained for generating images from text can only perform this single task, and is unable to go in the backward direction - i.e. generate captions for images. This drawback limits the power of generative models, which we aim to address through our project by designing a generative model for multimodal data.

Multimodal data, i.e., data comprising different sets of features (modalities) belonging to different distributions, is ubiquitous: from collections of heterogeneous unstructured representations of objects (e.g. text, images, audio, categories, to describe a single data point) to sets of homogeneous features providing different views of samples (multi-view learning).

A principled probabilistic treatment of multimodal learning would allow not only to compactly represent multimodal distributions but also to perform inference over them. It would be possible to draw new samples from all or some modalities, or compute the likelihood of some joint assignment. An example application of this could be to use the same model for (a) image captioning, (b) image generation from text, and (c) joint generation of an image with its caption. Although this example constitutes two modalities, the idea is more general and could encompass more than two (for example visual, auditory, and text such as movies).

In this report we investigate how adopting tractable probabilistic models (TPMs), more specifically Sum Product Networks [15], jointly model multimodal data that can enable tractable probabilistic inference over subsets of the modalities at hand as well as scaling multimodal learning. We explore several scenarios where the employed TPMs, in the form of probabilistic circuits (PCs) can deliver implicit or explicit likelihood by aggregating latent representations for different modalities in a “plug&play” fashion.

2. Notation

Upper-case letters X denote random variables (RVs) and lower-case letters their values, i.e., $x \sim X$. Similarly, we denote (ordered) sets of RVs as \mathbf{X} , and their corresponding values as \mathbf{x} . For a general discussion, assume that the set of RVs comprises M modalities—also referred to as views in the literature, i.e., we have a partitioning of the feature space

as $\mathbf{X} = \bigcup_{i=1}^M \mathbf{X}_i$ and $\mathbf{X}_i \cap \mathbf{X}_j = \emptyset$ for any $i \neq j$, where $i, j \in \{1, \dots, M\}$. When there will be the need to refer a particular modality, we will label its corresponding RV set accordingly, e.g., \mathbf{X}_{txt} for some structured representation of text data (like a bag-of-words representation, or some text embedding), \mathbf{X}_{img} for images, and so on.

3. Related Work

Multimodal data has been investigated in a number of fashions for non-probabilistic modeling (deterministic mappings) of low-dimensional spaces for multi-view learning.

Among works dealing with probabilistic mappings, recent research lines involve deep generative models like GANs and VAEs. Both are primarily used as simulators (e.g. to sample) as they do not have an explicit likelihood model (GANs) or if computing the likelihood exactly is hard (VAEs). We briefly list some of them and also cite their main limitations.

3.1. VAEs

All models based on VAEs have issues in modeling a joint evidence lower bound (ELBO): many have to represent explicit inference networks for all subsets of modalities at hand, or resort to heuristics during training to let a single architecture adapt to missing (subsets of) modalities.

1. **Variational methods for Conditional Multimodal Deep Learning** [13]: They introduce **CMMA** which learns one conditional distribution per modality as a conditional VAE.
2. **Deep Variational Canonical Correlation Analysis** [23]: They introduce **BiVCCA**, a deep CCA requiring a network for each subset of modalities.
3. **Joint Multimodal Learning With Deep Generative Models** [20]: **JMVAE** aims to represent each possible subset of modalities by an inference network.
4. **Generative Models of Visually Grounded Imagination** [22]: They use triple ELBO (**TELBO**) but the method does not generalize to more than 3 modalities.
5. **Multimodal Generative Models for Scalable Weakly-Supervised Learning** [25]: They introduce **MVAE** as a joint VAE having a product of experts (PoE) formulation which helps dealing with missing modalities (setting each modality input to 0). This seems the best competitor so far. However, during training they still help the model deal with the missing modalities by generating K masks for random subsets of the M modalities and looking at the D marginals.

They indeed optimize for:

$$\text{ELBO}(\mathbf{X}_1, \dots, \mathbf{X}_M) + \sum_{i=1}^D \text{ELBO}(\mathbf{X}_i) + \sum_{j=1}^K \text{ELBO}(\mathbf{X}_j)$$

6. **Variational Mixture-of-Experts Autoencoders for Multi-Modal Deep Generative Models** [19]: They substitute the PoE in the MVAE with a mixture of univariate experts to have a joint posterior, delivering the **MMVAE**. While not requiring the additional terms in the ELBO as in MVAE, they have to resort to more expensive stratified sampling [17] to average over M modalities.

3.2. GANs

GANs-based models, on the other hand, have the classical issue of pesky adversarial training and we do not use GAN based methods for our comparison as they cannot answer any inference queries even approximately but still mention them here for the sake of completeness.

1. **Adversarially Learned Inference** [3]: The generation network maps samples from stochastic latent variables to the data space while the inference network maps training examples in data space to the space of latent variables.
2. **Triple Generative Adversarial Nets** [9]: Triple-GAN consists of three players—a generator, a discriminator and a classifier. Needs to model all conditional independencies.
3. **Triangle Generative Adversarial Networks** [4]: Δ -GAN is developed for semi-supervised cross-domain joint distribution matching, can be considered as a combination of conditional GAN and ALI.

4. Methods

4.1. Probabilistic Circuits

Probabilistic circuits or sum-product networks [15] are a class of generative models which model the joint distribution as a polynomial over the input leaf distributions. They allow for efficient tractable inference for the following kinds of queries exactly:

Modality sampling (M – SAM). We would like to sample from any subset of modalities given any other subset, that is

$$\mathbf{x}_{M_1}^{\text{new}} \sim p_{\theta}(\mathbf{X}_{M_1} \mid \mathbf{X}_{M_2} = \mathbf{x}_{M_2}) \quad (1)$$

where $\mathbf{X}_{M_1}, \mathbf{X}_{M_2} \subset \mathbf{X}$ and $\mathbf{X}_{M_1} \cap \mathbf{X}_{M_2} = \emptyset$.

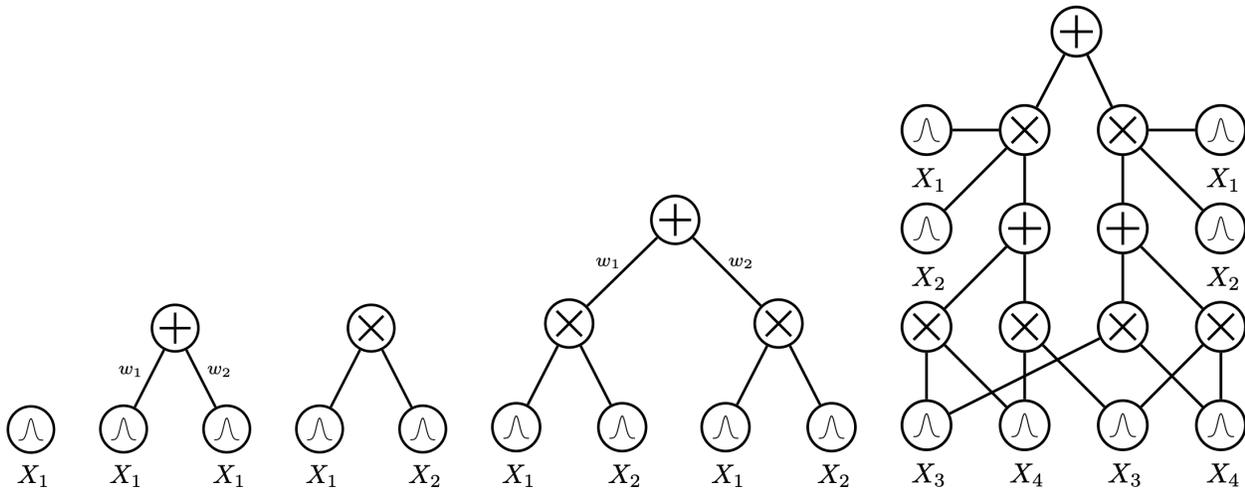


Figure 1. Various components of a probabilistic circuit: Sequentially building a full-fledged PC (rightmost) by composing sum and product nodes as a computational graph. The final output is the joint likelihood of the variables taking the assigned values. The above image was taken from <http://web.cs.ucla.edu/~guyvdb/slides/AAAI20.pdf>.

Modality MAR inference (M – MAR). Here we would like to compute (log-)likelihoods while being able to marginalize over arbitrary sets of modalities (and potentially condition over arbitrary evidence). A prototypical query in the class looks like:

$$p_{\theta}(\mathbf{x}_{M_1} \mid \mathbf{x}_{M_2}) \quad (2)$$

where $\mathbf{X}_{M_1}, \mathbf{X}_{M_2} \subset \mathbf{X}$ and $\mathbf{X}_{M_1} \cap \mathbf{X}_{M_2} = \emptyset$ and we are marginalizing over $\mathbf{X} \setminus \{\mathbf{X}_{M_1} \cup \mathbf{X}_{M_2}\}$.

Modality MAP inference (M – MAP). In this case, we would like to retrieve the mode of the conditional distribution obtained after conditioning on some subsets of modalities

$$\mathbf{x}_{M_1}^* = \operatorname{argmax}_{\mathbf{x}_{M_1}} p_{\theta}(\mathbf{x}_{M_1}, \mathbf{x}_{M_2}) \quad (3)$$

where $\mathbf{X}_{M_1} \cup \mathbf{X}_{M_2} = \mathbf{X}$ and $\mathbf{X}_{M_1} \cap \mathbf{X}_{M_2} = \emptyset$.

Modality marginal MAP inference (M – MMAP). This case is similar to M – MAP but here we marginalize over some modalities which we don’t care about (e.g., modalities missing).

Probabilistic Circuits are implemented as computational graphs consisting of three types of nodes:

1. **Leaf Nodes:** Given a random variable/s at the input of the PC, these nodes parameterise a distribution over that random variable/s. For example, given an input \mathcal{X} , these distributions will output $p_{\mathcal{X}}$ for some parametric distribution p such as a Gaussian.

2. **Sum Nodes:** These nodes take in multiple distributions as inputs and simply output a mixture of these distributions where the weights of the mixture are learnable.
3. **Product Nodes:** These nodes take in multiple distributions as inputs and output the product of these distributions which effectively models the variables in the scope of the product node to be independent.

So eventually, the only learnable parameters of the PC include the mixture weights of all the sum nodes and the parameters of the input distribution. An illustration of a simple PC can be found in Fig 1. In our work, we use the Probabilistic Circuits to model the joint distribution of multiple modalities over their latent space instead of high-dimensional pixel spaces. PCs are trained using EM (expectation-maximization) by maximising the likelihood over the training data.

4.2. Regularized Autoencoders

Informative latent space encoding of the multimodal subspaces is crucial for training probabilistic circuits on multimodal data. Vanilla deterministic autoencoders have a spiky distribution of the latents making maximum likelihood training hard for probabilistic circuits. To learn a Probabilistic Circuit on the fused latent space, we attempt to enforce smoothness in the learnt latent space. To achieve this, instead of training autoencoders with MSE loss between the input and the reconstruction, we add two additional terms to the loss as shown in [5]: l_2 norm on the latent space and l_2 norm on the decoder gradients while training.

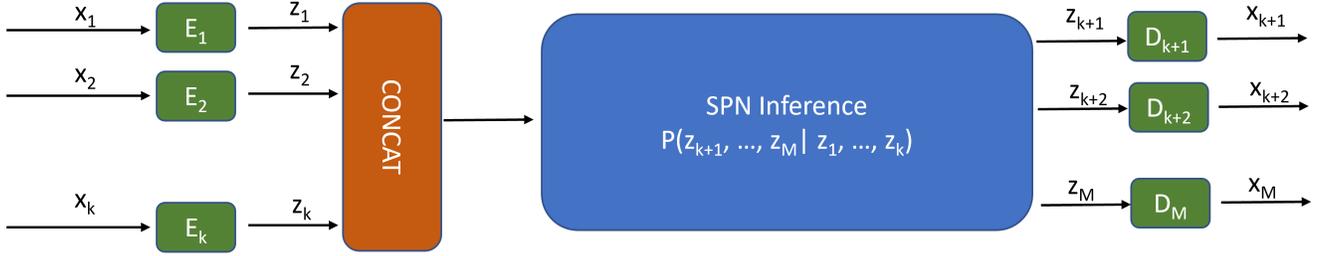


Figure 2. Using the model to sample $\mathbf{x}_{k+1}, \dots, \mathbf{x}_M$ given $\mathbf{x}_1, \dots, \mathbf{x}_k$

4.3. PPPC: Plug&Play Probabilistic Circuits

4.3.1 PCs To “Glue” Modalities

We discuss the simplest generative model for multimodal data: it consists of a set of independent mechanisms, one for each modality $i = 1, \dots, M$. Each mechanism can be modeled as an autoencoder mapping a certain (lower-dimensional) latent space \mathbf{Z}_i to \mathbf{X}_i , the corresponding set of observed RVs. A global joint generative mechanism for $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_M$ is recovered by modeling the joint distribution over latent space $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_M$ via some tractable model p_S , e.g. a PC \mathcal{S} . A PC would give the flexibility to perform inference over subsets of the modalities flawlessly: it would *orchestrate* encoding-decoding over different autoencoders. We codename this architecture PPPC.

This is the most basic idea that was initially proposed for mixed sum-product networks (MSPNs) [11]. The idea of using VAEs instead of PCs has been explored by in [21] in the context of uni-modal data. While it might seem appealing to have a joint ELBO, we lose all the advantages of PPPC as the tractable properties discussed above cannot be achieved and we go back to a harder optimization problem (In [21] they were not able to scale besides MNIST).

4.3.2 Training and Inference Procedure

Let M be the number of modalities for the multimodal data used for training. Our generative model consists of two major components: M encoder-decoder networks (which are deterministic) trained on each modality $[(\mathbf{E}_1, \mathbf{D}_1), (\mathbf{E}_2, \mathbf{D}_2), \dots, (\mathbf{E}_M, \mathbf{D}_M)]$ and one SPN network \mathcal{S} that fits a generative model on the latent spaces of these encoders.

More formally, let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$ represent a *single* datapoint consisting of the M modalities. Let $\mathbf{z}_i = \mathbf{E}_i(\mathbf{x}_i)$ denote the latent representations of the input. Denoting $\mathbf{z} = \text{CONCAT}(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_M)$, the SPN network predicts the likelihood of $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_M$ as $\mathcal{S}(\mathbf{z})$. Thus, our complete likelihood model is simply $\mathcal{S}(\text{CONCAT}(\mathbf{E}_1(\mathbf{x}_1), \mathbf{E}_2(\mathbf{x}_2), \dots, \mathbf{E}_M(\mathbf{x}_M)))$.

As an example (figure 2), suppose we want to perform

inference queries such as predicting $\mathbf{x}_{k+1}, \dots, \mathbf{x}_M$ given $\mathbf{x}_1, \dots, \mathbf{x}_k$. Such a query is useful in the case of text to image generation, where image and text are two of the modalities ($M = 2$ and $k = 1$ in this case). To perform this query, first we encode the given attributes, i.e., compute $\mathbf{z}_1, \dots, \mathbf{z}_k$ in the latent space. Next, we utilize tractable inference of SPNs to compute $\mathcal{S}(\mathbf{z}_{k+1}, \dots, \mathbf{z}_M | \mathbf{z}_1, \dots, \mathbf{z}_k)$. Finally, we sample from this distribution to obtain $\hat{\mathbf{z}}_{k+1}, \dots, \hat{\mathbf{z}}_M$ and then utilize the decoders $\mathbf{D}_{k+1}, \dots, \mathbf{D}_M$ to compute $\hat{\mathbf{x}}_i = \mathbf{D}_i(\hat{\mathbf{z}}_i), i \in \{k+1, \dots, M\}$. Thus, we are able to perform efficient sampling in multimodal data.

4.3.3 Advantages of Plug&Play Learning

The main advantage of PPPC is that each (R)AE could be trained independently from others. This i) greatly simplifies a potentially tough joint optimization problem and ii) provides a single inference machine for all possible inference “directions” (e.g., while conditioning, overcoming the need of other GAN- and VAE-based competitors that have to either train different architectures for each directions or sampling subsets of modalities); iii) allows to leverage sota AE architectures, *out-of-the-box* each tailored for each modalities, while iv) plugging them in (without retraining), out or swapping them more easily. Moreover, v) PCs can flawlessly deal with heterogeneous embeddings. Indeed, we can devise an online learning scheme where we train the PC over \mathbf{Z} by adding some modalities at the time, reusing the partial distributions previously learned.

5. Datasets

We work with the following datasets for this project:

1. We combined digit pairs of the same kind from MNIST [8] and SVHN [12] to create a multimodal dataset for initial experimentation. We pair each image from the MNIST dataset with 20 images from the SVHN dataset. The MNIST dataset is divided into 50k training images and 10k validation and test images each.
2. We also use the CelebA dataset [10] which consists of images and binary attributes. The dataset consists of

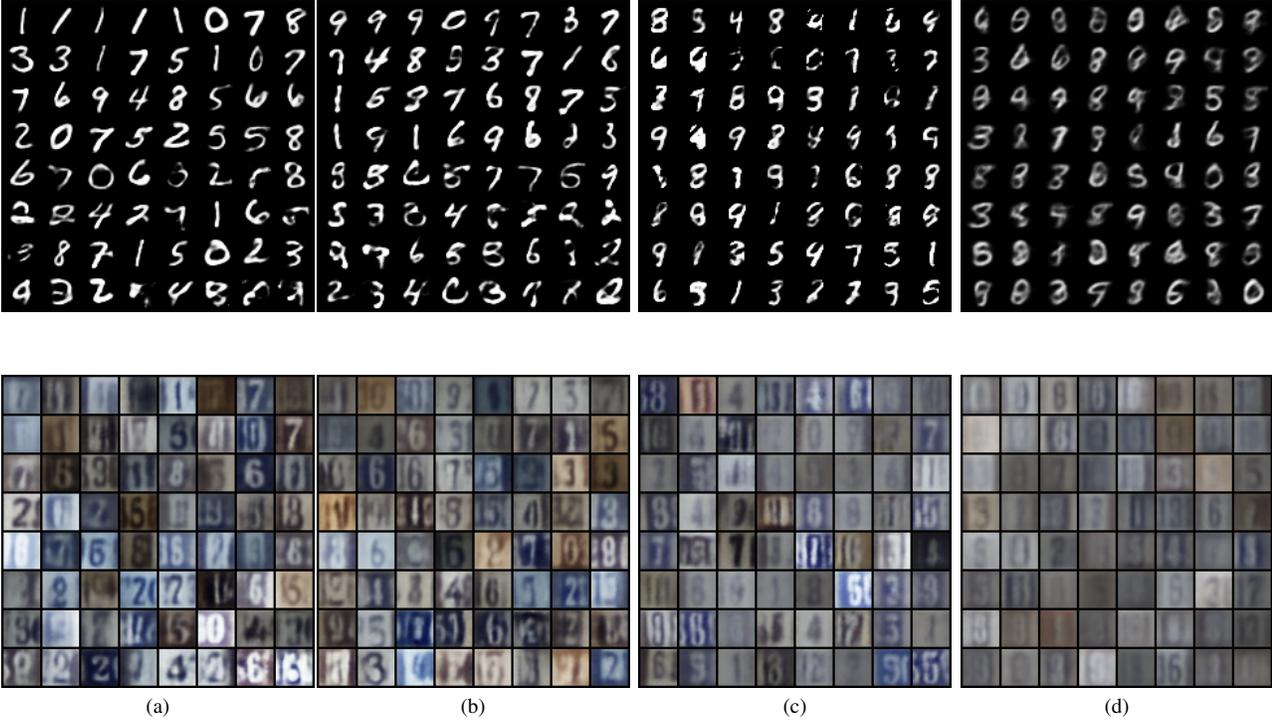


Figure 3. The above figure shows unconditional samples for (a): GMM, (b): SPN, (c): MMVAE, (d): MVAE arranged as a column drawn from the joint distribution $p(x_{\text{mnist}}, x_{\text{svhn}})$. These samples are arranged in order of their likelihoods in a decreasing order

around 200k celebrity faces, each annotated with 40 attributes. We only work with 4 selected attributes - glasses, hair color, gender, open mouth.

3. We also perform experiments on image and caption data from the CUB dataset [24], which consists of images of birds and multiple single-sentence captions for each image. We divide the dataset into around 9,600 images for training and into 1,200 images for validation and test each. The results of this dataset are in the appendix.

6. Results and Discussion

6.1. Label Paired MNIST-SVHN

In this section, we shall show the results on the MNIST-SVHN multiview dataset which involves pairing of MNIST and SVHN images having the same class label.

6.1.1 Method Details

For this particular experiment, the latent space size for both the modalities was chosen to be 32 which led to a dimensionality of 64 over which the probabilistic circuit was to be learnt. The autoencoder architecture was a simple conv arch with 6 layers and Swish activations which lead to smoother image reconstructions as shown in [5]. The latent space

was restricted to lie in a hypercube of $[-1, 1]^{64}$ which shall make it easier to learn the PC and avoid any numerical instability issues. For the PC, we use a randomized GPU implementable architecture as shown in [14] along with Gaussian leaves. For qualitative samples, we also compare with learning a Gaussian Mixture Model (GMM) on the latent space which is a very simple PC with only 1 layer and one sum node.

6.1.2 Unconditional Inference Evaluation

Once we learn the joint distribution, we can sample from the joint $p_{\theta}(x_{\text{mnist}}, x_{\text{svhn}})$ to generate a pair of same class images. We can evaluate the joint query by the following two metrics:

1. Measure the FID scores of the images generated with their corresponding datasets which gives an indication of how well the images look when compared to other samples from the same distribution.
2. Measure coherence, which means how well do the class labels of both the sampled images match with each other. For this, we train two classifiers on MNIST and SVHN separately and the joint predictions of these classifiers are used to measure coherence. The results are in Fig 3.

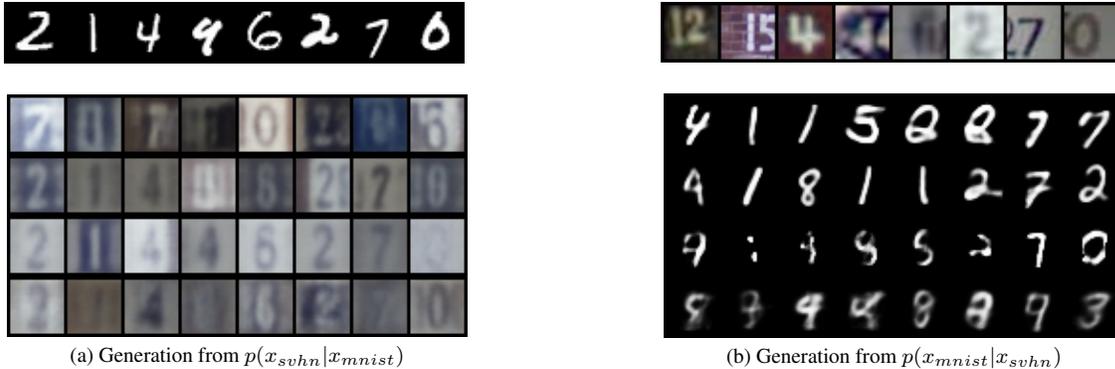


Figure 4. Samples drawn from conditional distribution with each row indicating a model in the order of (1): GMM, (2): SPN, (3): MMVAE, (4): MVAE. Qualitatively samples from SPNs are the most in sync with the digit classes and also show the least blurry artifacts unlike MVAE/MMVAE

6.1.3 Conditional Inference Evaluation

To understand cross-modality generation capabilities of all the models, we condition on one of the modality and sample the other from the conditional distribution $p(x_{mnist}|x_{svhn})$ and $p(x_{svhn}|x_{mnist})$. For PCs, we perform an approximate MAP query to assign the most likely image from the conditional distribution. From the samples in figure 4, it can be clearly seen that the samples from the plug-and-play models offer a better coherence in aligning modalities which is further reinforced by the quantitative scores later.

6.1.4 Quantitative Evaluation

In this subsection, we look at the FID scores and the coherence accuracies of the generated samples and compare those with MVAE and MMVAE. We see PPC outperforms the competitors by a huge margin except of the coherence in the MNIST to SVHN modality case.

Model	Joint		Mod ₁ → Mod ₂		Mod ₂ → Mod ₁	
	Qua(↓)	Coh(↑)	FID(↓)	Acc(↑)	FID(↓)	Acc(↑)
MVAE	220.56	28.15	92.58	54.60	94.27	27.45
MMVAE	112.49	34.75	101.58	72.25	35.98	59.02
PPPC	87.71	38.10	64.60	75.67	21.34	47.35
AE	76.49	78.38	55.52	81.48	17.42	98.15

Table 1. Quantitative Evaluation of generative capacities of various models using FID scores and classification accuracies

The final row AE stands for the autoencoder and can be thought of as the groundtruth/best result that we can obtain for that particular column. Those numbers are indicative of the quality of the learnt autoencoders and the per modality classifiers and hence we cannot expect any generative model to outperform them.

6.2. CelebA-Attributes

In this section, we shall show the results on the CelebA-Attributes dataset which involves a CelebA image (64×64) and associated 4 binary attributes (blonde hair, gender, mouth open or not, glasses on or not) and hence we have a total of 5 modalities.

6.2.1 Method Details

For this particular experiment, the latent space size for images was chosen to be 64 and all the binary attributes did not have any encoder or decoder associated with them and were simply appended in the joint latent space \mathbf{Z} which led to learning a PC over a dimensionality of 68. Here, as the last 4 variables were binary, we modified the PC leaves such that they encoded a Bernoulli distribution for these variables and a Gaussian distribution for the latent codes obtained from the image encoder. The image latent space was also constrained in $[-1, 1]^{64}$ and we also tested with the GMM training on the joint latent space.

6.2.2 Unconditional Inference Evaluation

We sample from the joint distribution and observe the images generated after passing through the decoder in Fig 5. As one can clearly observe, PPC leads to better unconditional samples and are sharper compared to other methods.

6.2.3 Conditional Inference Evaluation

We start off with just conditioning on a single attribute being true and observe the generated image and visually see the coherence between the conditioned attribute and image sampled in Fig 6. Again the coherence of our model is much better compared other counterparts.

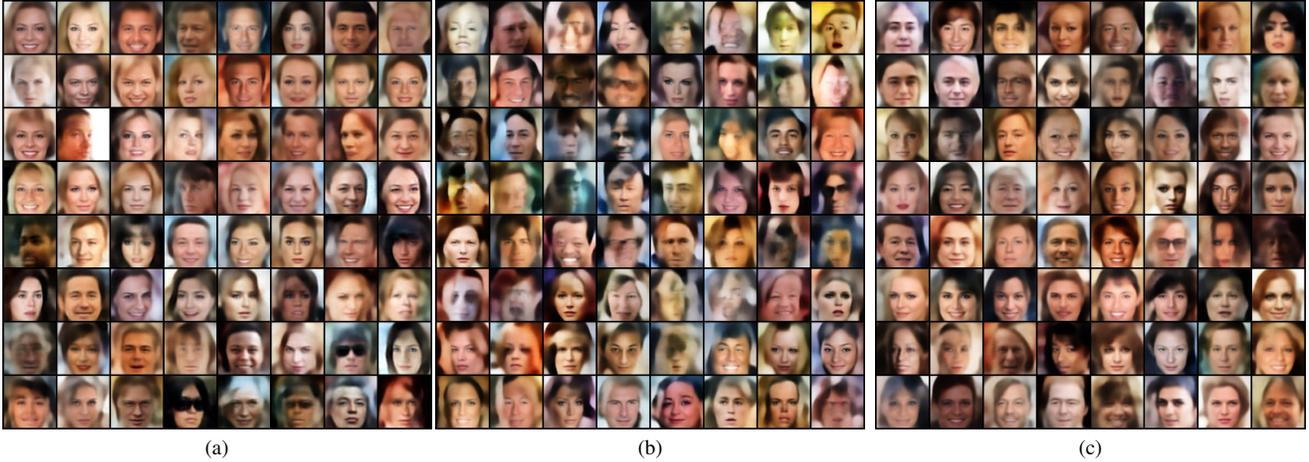


Figure 5. The above figure shows unconditional samples (only the image modality) for (a): PPC, (b): MMVAE, (c): MVAE drawn from the joint distribution arranged in order of decreasing likelihoods for each method

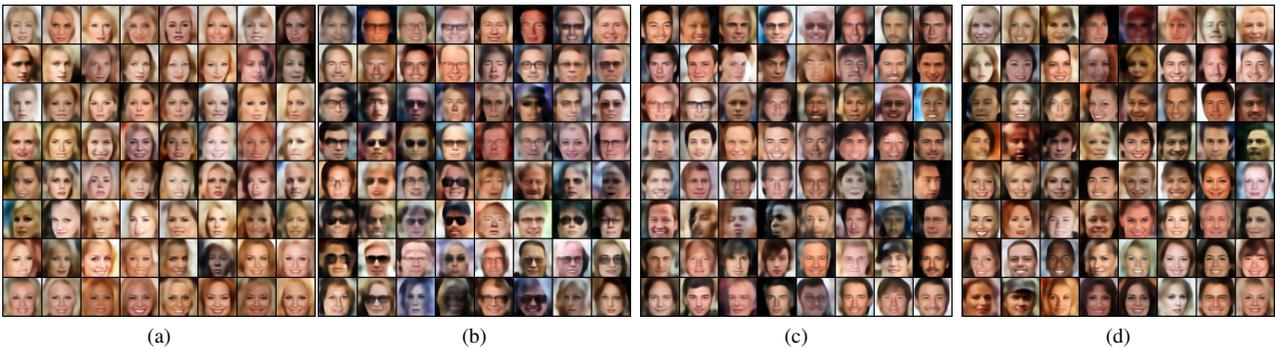


Figure 6. The above figure shows conditional samples for (a): Blond Hair, (b): Glasses, (c): Male, (d): Mouth Open attributes true for PPC. Samples from MVAE and MMVAE are in the Appendix as their performance is significantly worse.

We next condition over multiple modalities which was the main motivating factor of using a tractable probabilistic model. We get the following samples where each row denotes the following features as denoted in figure 7:

1. Blond Hair + Glasses
2. Male + Mouth Open
3. Blond Hair + Woman + Mouth Close
4. Blond Hair + Woman + Mouth Open

Past methods such as MVAE & MMVAE use variational approximations, which leads to strange images.

6.2.4 Quantitative Evaluation

For unconditional queries, we measure the image quality by measuring the FID scores of the generated samples. We also want to measure the coherence between the generated image and the generated attributes for the same. For this, we train 4 separate binary classifiers which take in the image as an input and outputs the class of the image for each of the attribute. Thus, to measure coherence, we pass the

generated image through each of the 4 classifiers and get corresponding labels. As we are sampling from the joint distribution, the model also generates labels and we compare these labels to the labels given by the classifier to measure coherence. Joint coherence is measured by finding the Hamming Loss between these two sets of binary labels (4 in total). The results for unconditional sampling are shown in Table 2.

Model	Joint	
	Qua(↓)	Coh(↑)
MVAE	70.264	0.252
MMVAE	93.031	0.236
PPPC	66.713	0.126
AE	58.871	0.0365

Table 2. Quantitative Evaluation of generative capacities of various models using FID scores and classification accuracies



Figure 7. Samples generated by conditioning on 1) blond hair + glasses, 2) male + mouth open, 3) blond hair + woman + mouth close, 4) blond hair + woman + mouth open, arranged row-wise from top to bottom respectively

Model	Blond		Glasses		Male		Mouth Open	
	FID(↓)	Acc(↑)	FID(↓)	Acc(↑)	FID(↓)	Acc(↑)	FID(↓)	Acc(↑)
MVAE	81.16	20.4	120.74	6.5	93.59	2.9	69.61	19.3
MMVAE	118.02	6.2	137.35	7.1	104.55	41.8	99.72	34.7
PPPC	72.47	81.2	89.27	50.7	74.64	87.4	63.83	57.4

Table 3. Quantitative Evaluation of generative capacities of various models using FID scores and classification accuracies by conditioning on single attributes

For evaluating conditional queries quantitatively, we use the FID scores and classification accuracies (condition on only 1 attribute at a time). We use the same classifiers mentioned above to calculate the classification accuracies. The results are shown in Table 3.

7. Conclusion

In this project, we have developed a general approach to perform multimodal generative modelling that allows for exact inference and sampling queries. Our developed framework is independent of the type of modality, i.e., it works for all kinds of data: fixed length, sequential, image etc. We have demonstrated the efficacy of this framework on two datasets - MNIST-SVHN and CelebA and shown that our method significantly outperforms comparable methods on multimodal learning such as MVAE and MMVAE.

There are several refinements that we can pursue to get better results for our method. Currently, due to constraints on computation and time, we had trained the autoencoder and the Sum Product Networks separately. We envision that joint training would allow for coupled learning of weights between these two modules and improve generalization. Finally, the architecture of the sum-product network that we use was selected randomly - with techniques such as Network Architecture Search we could get better inference/conditioning.

8. Acknowledgments

We broadly divide each experiment into 3 parts – a) training the encoder-decoder and SPN on the fused latent space, b) performing inference to get the qualitative samples and c) getting quantitative metrics. The division of work among team members was as follows (code words: MNIST-SVHN is MS, CelebA is CA and CUB is C):

1. JV: MS a, CA b, C b
2. SC: MS b, CA a, C b
3. PRS: MS c, CA c, C a

The main idea of this work was developed at the UCLA StarAI Lab with Guy Van den Broeck and Antonio Vergari. We are also thankful to Stefano Ermon at Stanford and Isabel Valera at MPI for useful comments. We used the SAIL atlas cluster to run our experiments.

All the implementation was done by the team members only. We used the following open-source implementations:

1. [Einsum Networks](#)
2. [RATSPN](#)
3. [RAE](#)

References

- [1] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- [2] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021.
- [3] Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016.
- [4] Zhe Gan, Liqun Chen, Weiyao Wang, Yuchen Pu, Yizhe Zhang, Hao Liu, Chunyuan Li, and Lawrence Carin. Triangle generative adversarial networks. In *Advances in neural information processing systems*, pages 5247–5256, 2017.
- [5] Partha Ghosh, Mehdi S. M. Sajjadi, Antonio Vergari, Michael Black, and Bernhard Scholkopf. From variational to deterministic autoencoders. In *International Conference on Learning Representations*, 2020.
- [6] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- [7] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2013.
- [8] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, Nov 1998.
- [9] Chongxuan Li, Taufik Xu, Jun Zhu, and Bo Zhang. Triple generative adversarial nets. *Advances in neural information processing systems*, 30, 2017.
- [10] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [11] Alejandro Molina, Antonio Vergari, Nicola Di Mauro, Sri-raam Natarajan, Floriana Esposito, and Kristian Kersting. Mixed sum-product networks: A deep architecture for hybrid domains. In *AAAI*, 2018.
- [12] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [13] Gaurav Pandey and Ambedkar Dukkipati. Variational methods for conditional multimodal deep learning. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 308–315. IEEE, 2017.
- [14] Robert Peharz, Steven Lang, Antonio Vergari, Karl Stelzner, Alejandro Molina, Martin Trapp, Guy Van den Broeck, Kristian Kersting, and Zoubin Ghahramani. Einsum networks: Fast and scalable learning of tractable probabilistic circuits. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, jul 2020.
- [15] Hoifung Poon and Pedro Domingos. Sum-product networks: A new deep architecture, 2012.
- [16] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation, 2021.
- [17] Christian P. Robert and George Casella. *Monte Carlo Statistical Methods (Springer Texts in Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2005.
- [18] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022.
- [19] Yuge Shi, N Siddharth, Brooks Paige, and Philip Torr. Variational mixture-of-experts autoencoders for multi-modal deep generative models. In *Advances in Neural Information Processing Systems*, pages 15692–15703, 2019.
- [20] Masahiro Suzuki, Kotaro Nakayama, and Yutaka Matsuo. Joint multimodal learning with deep generative models. *arXiv preprint arXiv:1611.01891*, 2016.
- [21] Ping Liang Tan and Robert Peharz. Hierarchical decompositional mixtures of variational autoencoders. In *International Conference on Machine Learning*, pages 6115–6124, 2019.
- [22] Ramakrishna Vedantam, Ian Fischer, Jonathan Huang, and Kevin Murphy. Generative models of visually grounded imagination. *arXiv preprint arXiv:1705.10762*, 2017.
- [23] Weiran Wang, Xinchen Yan, Honglak Lee, and Karen Livescu. Deep variational canonical correlation analysis. *arXiv preprint arXiv:1610.03454*, 2016.
- [24] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010.
- [25] Mike Wu and Noah Goodman. Multimodal generative models for scalable weakly-supervised learning. In *Advances in Neural Information Processing Systems*, pages 5575–5585, 2018.

9. Appendix

9.1. MMVAE and MVAE Results for CelebA



Figure 8. The above figure shows conditional samples for (a): Blond Hair, (b): Glasses, (c): Male, (d): Mouth Open attributes being true respectively for MMVAE

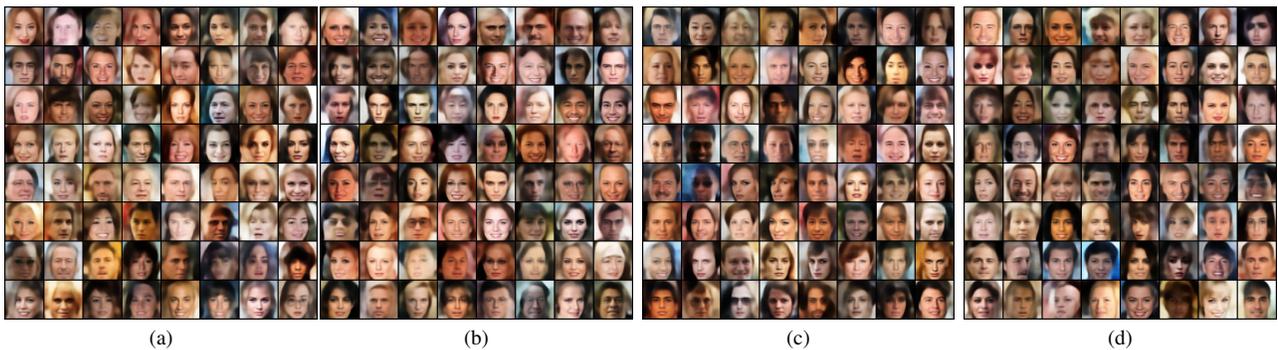


Figure 9. The above figure shows conditional samples for (a): Blond Hair, (b): Glasses, (c): Male, (d): Mouth Open attributes true for MVAE

We can clearly see that coherence is much worse than PPC and even the overall quality of the samples are bad and we can see many artifacts in the images generated as shown in Fig. 8 and Fig. 9.

9.2. CUB Dataset

In this section, we see some results by training SPNs on the CUB dataset. The encoder for the image was a standard CNN encoder and that for the captions was an LSTM and the hidden state was used as a latent representation. While the results for text generation are still good, we observe that the image generation results are not very great. We think this is mainly because of the small size of the dataset due to which encoders and decoders couldn't be trained properly.

9.2.1 Unconditional Sampling



Figure 10. Unconditional Image from PPC

The generated image samples as a result of unconditional sampling are shown in Fig. 10 and the results of unconditional sampling of text are as follows:

this is bird has a brown colored head with long and a yellow stripe on its crown.
this bird has a yellow breast, belly and black stripes on its head and secondaries.
this colorful bird has a grey body and yellow breast, and has black wings on its tail.
this multicolored bird has a white chest and dark brown primaries, with a grey beak and feet.
this bird has a brown back and white belly, also has a dark gray on it.
this bird is has a brown breast with pink feet and a large belly.
this bird has spotted brown wings and tail, with a white and black beak on its crown.
this bird has a white body with dark grey and black wings, a long, pointy bill.

9.2.2 Conditional Sampling

For conditional sampling for generating captions, the test images are shown in Fig. 11 and the generated caption for each of them are given below (in the same order). Conditional sampling in the other direction (text to images) could not be done well, again due to the small size of dataset.



Figure 11. Test images used for text generation

this bird has a dark grey belly and breast, with black superciliaries, and a white tail.
this bird has a white body, grey and dark stripes on its head and a orange beak.
this bird has a white body, grey wings and two long black beak.
this bird has a white body, black and long wings with a red eye ring.
this bird has a large wingspan, white belly and black eyes, with a long and blue beak.
this bird has a black back, white and breast with a orange crown.
this bird has a dark grey belly and breast, with a white eyering and brown crown.
this bird has a black back, wings and long with a yellow tip.