# MultiVae: A Python package for Multimodal Variational Autoencoders on Partial Datasets

**Agathe Senellart**[1,2*¶] **and Stéphanie Allassonnière**[2*]

**1** Université de Paris-Cité **2** Inria **3** Inserm ¶ Corresponding author * These authors contributed equally.

## Summary

In recent years, there has been a major boom in the development of multimodal machine learning models. Among open topics, representation (fusion) and genera- tion of multimodal data are very active fields of research. Recently, Multimodal Variational Autoencoders (VAEs) have been attracting growing interest for both tasks, thanks to their versatility, scalability, and interpretability as probabilistic latent variable models. They are also particularly interesting models in the partially observed setting, as most of them can learn even with missing data. This last point makes them particularly interesting for research fields such as the medical field, where missing data are commonplace.

In this article, we present MultiVae, an open-source Python library for bringing together unified imple- mentations of multimodal VAEs. It has been designed for easy, customizable use of these models on fully or partially observed data. This library also facilitates the development and benchmarking of new algorithms by integrating several popular datasets, variety of evaluation metrics and tools for monitoring and sharing models.

## Multimodal Variational Autoencoders

In Multimodal Machine Learning, two goals are generally targeted: (1) Learn a shared representation from multiple modalities; (2) Learn to generate one missing modality given the ones that are available.

Multimodal Variational Autoencoders aim at solving both issues at the same time. These models learn a latent representation $z$ of all modalities in a lower dimensional common space and learn to *decode* $z$ to generate any modality (Suzuki & Matsuo, 2022). Let $X = (x_1, x_2, ...x_M)$ contain $M$ modalities. In the VAE setting, we suppose that the generative process behind the observed data is the following:

$$z \sim p(z) \qquad\qquad \forall 1 \le i \le M, x_i|z \sim p_\theta(x_i|z) \qquad (1)$$
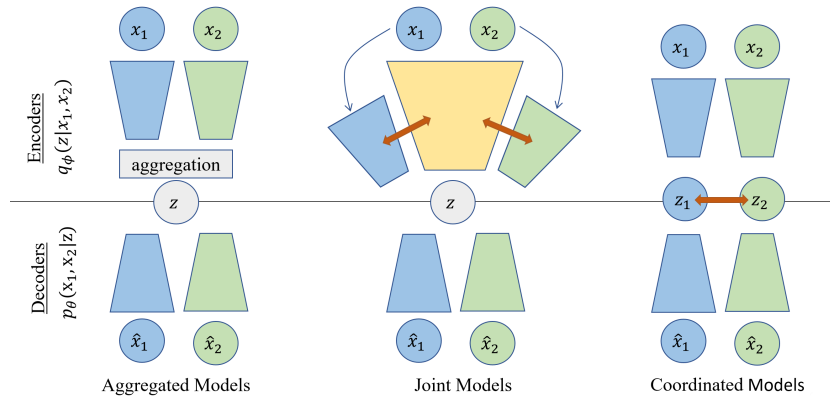
where $p(z)$ is a prior distribution that is often fixed, and $p_\theta(x_i|z)$ are called *decoders* and are parameterized by neural network. Typically, $p_\theta(x_i|z) = \mathcal{N}(x_i, \mu_\theta(z), \sigma_\theta(z))$ where $\mu_\theta, \sigma_\theta$ are neural networks. We aim to learn these *decoders* that translate $z$ into the high dimensional data $x_i$. At the same time, we aim to learn an *encoder* $q_\phi(z|X)$ that map the multimodal observation to the latent space. $q_\phi(z|X)$ is also parameterized by a neural network. Derived from variational inference (Kingma & Welling, 2014), the VAE objective writes:

$$\mathcal{L}(X) = \mathbb{E}_{q_\phi(z|X)}\left(\sum_i \ln(p_\theta(x_i|z))\right) - KL(q_\phi(z|X)|p(z))$$

A simple interpretation of this objective is to see that the first term is a reconstruction loss and the second term is a regularization term that avoids overfitting. A typical training of a

<sub>36</sub> multimodal VAE consists in encoding the data with the encoder, reconstructing each modality
<sub>37</sub> with the decoders and take a gradient step to optimize the loss $L(X)$.

<sub>38</sub> Most multimodal VAEs differ in how they construct the encoder $q_\phi(z|X)$. In Figure **??**, we
<sub>39</sub> summarize several approaches: Aggregated models (Shi et al., 2019; Sutter et al., 2021; Wu &
<sub>40</sub> Goodman, 2018) use a mean or a product operation to aggregate the information coming from
<sub>41</sub> all modalities, where joint models (Senellart et al., 2023; Vedantam et al., 2018; **?**) uses a neural
<sub>42</sub> network taking all modalities as input. Finally coordinated (Tian & Engel, 2019; Wang et al.,
<sub>43</sub> 2017) models uses different latent spaces but add a constraint term in the loss to force them to



Aggregated Models    Joint Models    Coordinated Models

<sub>44</sub> be similar.    ⟷ : Similarity constraint    Recent
<sub>45</sub> extensions of multimodal VAEs include additional terms to the loss, multiple or hierarchical
<sub>46</sub> latent spaces to more comprehensively describe the multimodal data. Aggregated models have
<sub>47</sub> a natural way of learning on incomplete datasets: for an incomplete sample $X$, we use only
<sub>48</sub> the available modalities to encode the data and compute the loss $l(X)$. However, except
<sub>49</sub> in MultiVae, there doesn't exist an implementation of these models that can be used on
<sub>50</sub> incomplete datasets in a straightforward manner.

<sub>51</sub> Another application of these models is data augmentation: from sampling latent codes $z$ and
<sub>52</sub> decoding them, fully synthetic multimodal samples can be generated. Data augmentation has
<sub>53</sub> been proven useful in many deep learning applications. In MultiVae we propose different ways
<sub>54</sub> of sampling latent codes $z$ to further explore the generative abilities of these models.
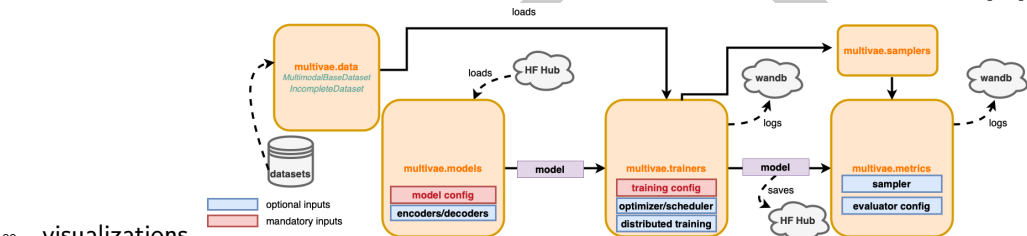
## Statement of need

<sub>56</sub> Although multimodal VAEs have interesting applications in different fields, the lack of easy-to-
<sub>57</sub> use and verified implementations might hinder applicative research. With MultiVae, we offer
<sub>58</sub> unified implementations, designed to be easy to use by non-specialists and even on incomplete
<sub>59</sub> data. To this end, we offer online documentation and tutorials. In order to propose reliable
<sub>60</sub> implementations of each method, we tried to reproduce, whenever possible, a key result from
<sub>61</sub> the original paper. Some works similar to ours have grouped together model implementations:
<sub>62</sub> the Multimodal VAE Comparison Toolkit includes 4 models and the Pixyz library groups 2
<sub>63</sub> multimodal models. The work closest to ours and developed in parallel is the multi-view-ae
<sub>64</sub> library (**?**), which contains a dozen of models. Nevertheless, we are convinced that our library
<sub>65</sub> complements what already exists: our API is quite different, the models implemented are
<sub>66</sub> not all the same, and for those we have in common, our implementation offers additional
<sub>67</sub> parameterization options. Indeed, for each model, we've made sure to offer great flexibility
<sub>68</sub> on parameters and to include all implementation details present in the original codes that
<sub>69</sub> boost results. What's more, our library offers numerous additional features: compatibility
<sub>70</sub> with incomplete data, which we consider essential for real-life applications, and a range of
<sub>71</sub> tools dedicated to the research and development of new algorithms: benchmark datasets,
<sub>72</sub> metrics modules and samplers, for testing and analyze models. Our library also supports
<sub>73</sub> distributed training and straightforward model sharing via HuggingFace Hub. In this way, our

74 work complements existing work and addresses different needs.

# Description of the software

76 Our implementation is based on PyTorch and is inspired by the architecture of (Chadebec et
77 al., 2022). The implementations of the models are collected in the module `multivae.models`.
78 Each model class is accompanied by a configuration dataclass gathering the collection of any
79 relevant hyperparameter which enables them to be saved and loaded straightforwardly. The
80 models are implemented in a unified way, so that they can be easily integrated within the
81 `multivae.trainers`. Trainers are also accompanied by a training configuration dataclass used
82 to specify any training-related hyperparameters (number of epochs, optimizers, schedulers, etc..).
83 Models that have a multistage training [50, 40] benefit from their dedicated trainer that makes
84 them as straightforward to use as other models. Partially observed datasets can be conveniently
85 handled using the `IncompleteDataset` class that contains masks informing on missing or
86 corrupted modalities in each sample. Finally, the MultiVae library also integrates an evaluation
87 pipeline for all models with common metrics such as likelihoods, coherences, FID scores [18] and



88 visualizations.

# List of models and features

| Model or Feature | MultiVae | multi-view-ae |
|---|---|---|
| JMVAE | * | |
| MVAE | * | |
| MMVAE | * | |
| MoPoE | * | |
| DMVAE | | |
| MVTCAE | | |
| MMVAE+ | * | |
| CMVAE | | |
| Nexus | | |
| CVAE | | |
| MHVAE | | |
| TELBO | | |
| JNF | | |
| MCVAE | | |
| mAAE | | |
| DVCCA | | |
| mWAE | | |
| mmJSD | | |
| gPoE | | |
| Support of Incomplete datasets | | |
| GMM Sampler | | |
| MAF Sampler, IAF Sampler | | |

| Model or Feature | MultiVae | multi-view-ae |
|---|---|---|
| Metrics : Likelihoods, Coherences, FIDs, Reconstruction, Clustering | | |
| Inline Datasets | | |
| Model sharing via Hugging Face | | |

# Documentation

The main features are illustrated through tutorials made available either as notebooks or scripts allowing users to get started easily. An online documentation is also made available at https://multivae.readthedocs.io/en/latest.

# Acknowledgements

# References

Chadebec, C., Vincent, L., & Allassonniere, S. (2022). Pythae: Unifying Generative Autoencoders in Python - A Benchmarking Use Case. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, & A. Oh (Eds.), *Advances in Neural Information Processing Systems* (Vol. 35, pp. 21575–21589). Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2022/file/872f0e04ef95be7970d9a9d74b198fdf-Paper-Datasets_and_Benchmarks.pdf

Kingma, D. P., & Welling, M. (2014). *Auto-Encoding Variational Bayes*. arXiv. http://arxiv.org/abs/1312.6114

Senellart, A., Chadebec, C., & Allassonnière, S. (2023). Improving multimodal joint variational autoencoders through normalizing flows and correlation analysis. *arXiv Preprint arXiv:2305.11832*.

Shi, Y., Siddharth, N., Paige, B., & Torr, P. H. S. (2019). Variational Mixture-of-Experts Autoencoders for Multi-Modal Deep Generative Models. *arXiv:1911.03393 [Cs, Stat]*. http://arxiv.org/abs/1911.03393

Sutter, T. M., Daunhawer, I., & Vogt, J. E. (2021). Generalized Multimodal ELBO. *ICLR*.

Suzuki, M., & Matsuo, Y. (2022). A survey of multimodal deep generative models. *Advanced Robotics*, *36*(5-6), 261–278. https://doi.org/10.1080/01691864.2022.2035253

Tian, Y., & Engel, J. (2019). Latent Translation: Crossing Modalities by Bridging Generative Models. *ArXiv*.

Vedantam, R., Fischer, I., Huang, J., & Murphy, K. (2018). Generative Models of Visually Grounded Imagination. *arXiv:1705.10762 [Cs, Stat]*. http://arxiv.org/abs/1705.10762

Wang, W., Yan, X., Lee, H., & Livescu, K. (2017). *Deep Variational Canonical Correlation Analysis*. arXiv. https://doi.org/10.48550/arXiv.1610.03454

Wu, M., & Goodman, N. (2018). Multimodal Generative Models for Scalable Weakly-Supervised Learning. *Advances in Neural Information Processing Systems*, *31*. https://proceedings.neurips.cc/paper/2018/hash/1102a326d5f7c9e04fc3c89d0ede88c9-Abstract.html