

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

Agathe Senellart^{1,2*}¶ and Stéphanie Allasonnière^{2*}

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

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Open Journals](#) ↗

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

In partnership with



AMERICAN
ASTRONOMICAL
SOCIETY

This article and software are linked with research article DOI [10.3847/xxxxx](https://doi.org/10.3847/xxxxx) <- [update this with the DOI from AAS once you know it.](#), published in the *Astrophysical Journal* <- The name of the AAS journal..

Summary

In recent years, there has been a major boom in the development of multimodal machine learning models. Among open topics, representation (fusion) and generation 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 be trained even with missing data. This last point makes them particularly interesting for the medical field, where missing data are commonplace ([Antelmi et al., 2019](#); [Lawry Aguila et al., 2023](#)).

We present MultiVae, an open-source Python library for bringing together unified implementations 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 benchmark datasets, a variety of evaluation metrics and tools for monitoring and sharing models.

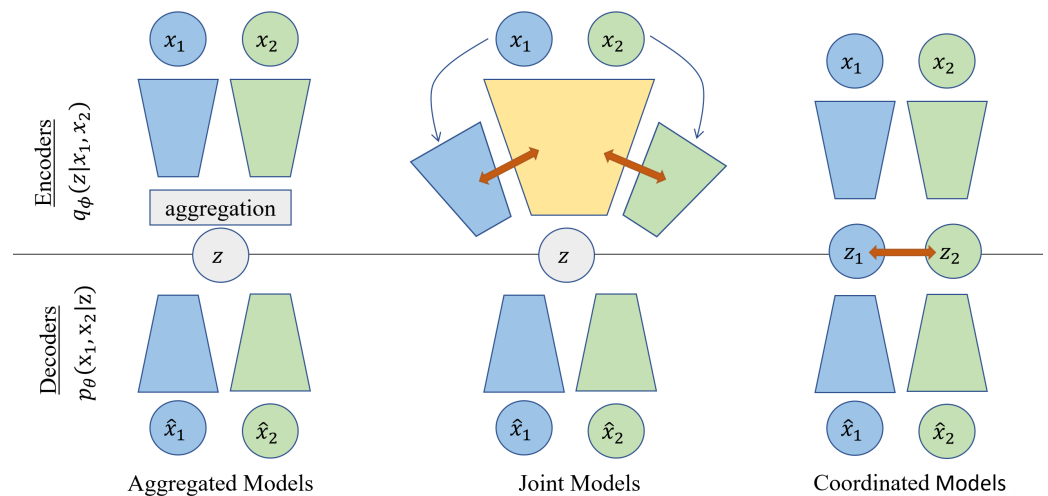
Multimodal Variational Autoencoders

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.

Let $X = (x_1, x_2, \dots, x_M)$ contain M modalities. In the VAE setting, we define an *encoder* distribution $q_\phi(z|X)$ projecting the observations to the latent space, and decoders distributions $(p_\theta(x_i|z))_{1 \leq i \leq M}$ translating the latent code z back to the observations. Those distributions are parameterized by neural networks that we train using variational inference. See ([Kingma & Welling, 2014](#)) to learn more about the VAE framework and ([Suzuki & Matsuo, 2022](#)) for a survey on multimodal VAEs.

Most multimodal VAEs differ in how they construct the encoder $q_\phi(z|X)$. In the figure below, we summarize several approaches: *Aggregated models* ([Shi et al., 2019](#); [Sutter et al., 2021](#); [Wu & Goodman, 2018](#)) use a mean or a product operation to aggregate the information coming from all modalities, where *Joint models* ([Senellart et al., 2023](#); [Suzuki et al., 2016](#); [Vedantam et al., 2018](#)) use a neural network taking all modalities as input. Finally *coordinated models* ([Tian & Engel, 2019](#); [Wang et al., 2017a](#)) use different latent spaces but add a constraint term in the loss to force them to be similar.



39

40 In our library, we implement all these approaches in an unified and modular way.

41 Aggregated models offer a natural way of learning on incomplete datasets: for an incomplete
 42 sample X , we use only the available modalities to encode the data and compute the loss that
 43 we minimize during training. However, except in our library MultiVae, there does not exist an
 44 implementation of these models that can be used on incomplete datasets in a straightforward
 45 manner.

46 Data Augmentation

47 Another application of these models is Data Augmentation (DA): from sampling new latent
 48 codes z and decoding them, *fully synthetic multimodal* samples can be generated to augment
 49 a dataset. DA has been proven useful in many data-intensive deep learning applications
 50 (Chadebec et al., 2023). In a dedicated module `multivae.samplers`, we propose different
 51 ways of sampling latent codes z to further explore the generative abilities of these models.

52 Statement of need

53 Although multimodal VAEs have interesting applications in different fields, the lack of easy-
 54 to-use and verified implementations might hinder applicative research. With MultiVae, we
 55 offer unified implementations, designed to be accessible even for non-specialists. In order to
 56 propose reliable implementations, we reproduced, whenever possible, a key result from the
 57 original paper. Some works similar to ours have grouped together model implementations:
 58 the [Multimodal VAE Comparison Toolkit](#) (Sejnova et al., 2024) includes 4 models and the
 59 [Pixyz](#) (Masahiro Suzuki & Matsuo, 2023) library contains 2 multimodal models. The work
 60 closest to ours and released while we were developping our library is `multi-view-ae` (Aguila
 61 et al., 2023), which contains a dozen of models. We compare in a summarizing table below,
 62 the different features of each work. Our library complements what already exists: our API
 63 is quite different compared to previous work, the models implemented are not all the same,
 64 and for those we have in common, our implementation offers additional parameterization
 65 options. Indeed, for each model, we have made sure to offer great flexibility on parameters
 66 and to include all implementation details present in the original codes. Our library also offers
 67 additional features: **compatibility with incomplete data**, which we consider essential for real-life
 68 applications, and a range of tools dedicated to research and development of new algorithms:
 69 benchmark datasets, **metrics modules** and **samplers**, for testing and analyzing models.

70 **List of models and features**

71 In the Table below, we list available models and features, and compare to previous work. This
72 symbol (✓*) indicates that the implementation include additional options.

Models/ Features	Ours	(Aguila et al., 2023)	(Sejnova et al., 2024)
JMVAE (Suzuki et al., 2016)	✓*	✓	
MVAE(Wu & Goodman, 2018)	✓*	✓	✓
MMVAE(Shi et al., 2019)	✓*	✓	✓
MoPoE(Sutter et al., 2021)	✓*	✓	✓
DMVAE(Lee & Pavlovic, 2021)	✓	✓*	✓
MVTCAE(Hwang et al., 2021)	✓	✓	
MMVAE+(Palumbo et al., 2023)	✓*	✓	
CMVAE(Palumbo et al., 2024)	✓		
Nexus(Vasco et al., 2022)	✓		
CVAE(Kingma & Welling, 2014)	✓		
MHVAE(Dorent et al., 2023)	✓		
TELBO(Vedantam et al., 2018)	✓		
JNF(Senellart et al., 2023)	✓		
CRMVAE(Suzuki & Matsuo, 2023)	✓		
MCVAE(Antelmi et al., 2019)		✓	
mAAE		✓	
DVCCA(Wang et al., 2017b)		✓	
mWAE		✓	
mmJSD(Sutter et al., 2020)		✓	
gPoE(Lawry Aguila et al., 2023)		✓	
Support of Incomplete datasets	✓		
GMM Sampler	✓		
MAF Sampler, IAF Sampler	✓		
Metrics: Likelihood, Coherences, FIDs, Reconstruction, Clustering	✓		
Benchmark Datasets	✓		✓
Model sharing via Hugging Face	✓		

73 An important difference in our user-interface, is that we handle all training and model parameters
74 within python dataclasses while (Aguila et al., 2023; Sejnova et al., 2024) uses independant
75 YAML configuration files.

76 **Code quality and documentation**

77 Our code is available on Github (<https://github.com/AgatheSenellart/MultiVae>) and Pypi
78 and we provide a full online documentation at (<https://multivae.readthedocs.io/>). The main
79 features are illustrated through **tutorials** made available either as notebooks or scripts allowing
80 users to get started easily. To further showcase how to use our library for research applications,
81 we provide detailed case studies in the documentation.

82 **Acknowledgements**

83 We are grateful to the authors of all the initial implementations of the models included in
84 MultiVae.

References

- Aguila, A. L., Jayme, A., Montaña-Brown, N., Heuveline, V., & Altmann, A. (2023). Multi-view-AE: A python package for multi-view autoencoder models. *Journal of Open Source Software*, 8(85), 5093. <https://doi.org/10.21105/joss.05093>
- Antelmi, L., Ayache, N., Robert, P., & Lorenzi, M. (2019). *Sparse multi-channel variational autoencoder for the joint analysis of heterogeneous data*. 97, 302–311. <https://proceedings.mlr.press/v97/antelmi19a.html>
- Chadebec, C., Thibeau-Sutre, E., Burgos, N., & Allasonnière, S. (2023). Data augmentation in high dimensional low sample size setting using a geometry-based variational autoencoder. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(3), 2879–2896. <https://doi.org/10.1109/TPAMI.2022.3185773>
- Dorent, R., Haouchine, N., Kogl, F., Joutard, S., Juvekar, P., Torio, E., Golby, A. J., Ourselin, S., Frisken, S., Vercauteren, T., Kapur, T., & Wells, W. M. (2023). Unified brain MR-ultrasound synthesis using multi-modal hierarchical representations. In *Medical image computing and computer assisted intervention – MICCAI 2023* (pp. 448–458). Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-43999-5_43
- Hwang, H., Kim, G.-H., Hong, S., & Kim, K.-E. (2021). Multi-view representation learning via total correlation objective. *Advances in Neural Information Processing Systems*, 34, 12194–12207.
- Kingma, D. P., & Welling, M. (2014). *Auto-Encoding Variational Bayes*. arXiv. <http://arxiv.org/abs/1312.6114>
- Lawry Aguila, A., Chapman, J., & Altmann, A. (2023). *Multi-modal variational autoencoders for normative modelling across multiple imaging modalities*. 425–434. https://doi.org/10.1007/978-3-031-43907-0_41
- Lee, M., & Pavlovic, V. (2021). *Private-shared disentangled multimodal VAE for learning of latent representations*. 1692–1700. <https://doi.org/10.1109/CVPRW53098.2021.00185>
- Masahiro Suzuki, T. K., & Matsuo, Y. (2023). Pixyz: A python library for developing deep generative models. In *Advanced Robotics* (No. 0; Vol. 0, pp. 1–16). Taylor & Francis. <https://doi.org/10.1080/01691864.2023.2244568>
- Palumbo, E., Daunhawer, I., & Vogt, J. E. (2023). *MMVAE+: ENHANCING THE GENERATIVE QUALITY OF MULTIMODAL VAES WITHOUT COMPROMISES*.
- Palumbo, E., Manduchi, L., Laguna, S., Chopard, D., & Vogt, J. E. (2024). *Deep generative clustering with multimodal diffusion variational autoencoders*. <https://openreview.net/forum?id=k5THrhXDV3>
- Sejnova, G., Vavrecka, M., Stepanova, K., & Taniguchi, T. (2024). *Benchmarking multimodal variational autoencoders: CdSprites+ dataset and toolkit*. <https://arxiv.org/abs/2209.03048>
- Senellart, A., Chadebec, C., & Allasonnière, 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. (2020). Multimodal generative learning utilizing jensen-shannon-divergence. *CoRR*, abs/2006.08242. <https://arxiv.org/abs/2006.08242>
- Sutter, T. M., Daunhawer, I., & Vogt, J. E. (2021). Generalized Multimodal ELBO. *ICLR*.

- 131 Suzuki, M., & Matsuo, Y. (2022). A survey of multimodal deep generative models. *Advanced*
132 *Robotics*, 36(5-6), 261–278. <https://doi.org/10.1080/01691864.2022.2035253>
- 133 Suzuki, M., & Matsuo, Y. (2023). *Mitigating the limitations of multimodal VAEs with*
134 *coordination-based approach*. <https://openreview.net/forum?id=Rn8u4MYgeNJ>
- 135 Suzuki, M., Nakayama, K., & Matsuo, Y. (2016). Joint Multimodal Learning with Deep
136 Generative Models. *arXiv:1611.01891 [Cs, Stat]*. <http://arxiv.org/abs/1611.01891>
- 137 Tian, Y., & Engel, J. (2019). Latent Translation: Crossing Modalities by Bridging Generative
138 Models. *ArXiv*.
- 139 Vasco, M., Yin, H., Melo, F. S., & Paiva, A. (2022). Leveraging hierarchy in multimodal
140 generative models for effective cross-modality inference. *Neural Networks*, 146, 238–255.
- 141 Vedantam, R., Fischer, I., Huang, J., & Murphy, K. (2018). Generative Models of Visually
142 Grounded Imagination. *arXiv:1705.10762 [Cs, Stat]*. <http://arxiv.org/abs/1705.10762>
- 143 Wang, W., Yan, X., Lee, H., & Livescu, K. (2017b). *Deep Variational Canonical Correlation*
144 *Analysis*. *arXiv*. <https://doi.org/10.48550/arXiv.1610.03454>
- 145 Wang, W., Yan, X., Lee, H., & Livescu, K. (2017a). *Deep Variational Canonical Correlation*
146 *Analysis*. *arXiv*. <https://doi.org/10.48550/arXiv.1610.03454>
- 147 Wu, M., & Goodman, N. (2018). Multimodal Generative Models for Scalable Weakly-Supervised
148 Learning. *Advances in Neural Information Processing Systems*, 31. [https://proceedings.](https://proceedings.neurips.cc/paper/2018/hash/1102a326d5f7c9e04fc3c89d0ede88c9-Abstract.html)
149 [neurips.cc/paper/2018/hash/1102a326d5f7c9e04fc3c89d0ede88c9-Abstract.html](https://proceedings.neurips.cc/paper/2018/hash/1102a326d5f7c9e04fc3c89d0ede88c9-Abstract.html)