# Knowledge Graphs (KG) Assisted Variational Autoencoder (VAE) for Large-Scale Anomaly and Event Detection

Ying Zhao[0000−0001−8350−4033]    yzhao@nps.edu

Naval Postgraduate School, Monterey CA 93943, USA

**Abstract.** This work focuses on general ML/AI assisted analytic processes for monitoring, detection, and classification of anomaly signals from multi-modality sensor data. Specifically, I will show series of variational autoencoders (VAE) including unsupervised AI transformers and related workflow pipelines applied to maritime surveillance in the different time scales of hours, minutes, and seconds. I show these developments using the distributed acoustic sensing (DAS) data set. The DAS data set is from the Sandia National Laboratories. DAS is a special type of fiber optic seafloor communications cables to interrogate the submarine environment at Arctic Alaska. Acoustic heatmaps can be generated to detect waves, ships, marine mammals, and other events. There are 18000 channels and sampled at 1kHZ for the data in 2022. The data set is used to demonstrate the VAE methodology for detecting anomaly, event, and classify objects. The results can enable processing and data analytical capabilities critical to actionable intelligence for mission planning and emerging behavior detection.

**Keywords:** anomaly detection · event detection · machine learning · ML · artificial intelligence · AI · generative AI · AI transformer · variational autoencoder · VAE · distributed acoustic sensing · DAS · heatmap · latent variables · knowledge graph · KG

## 1 Introduction

Current constellations of sensing capabilities including EO imagery, synthetic aperture radar (SAR), radio frequency signals provide an unprecedented ability to persistently observe maritime traffic. Maritime activities of variety vessels of Navy, coast guard, fishing, and civilian/merchant ships can take on many forms. Advanced detection algorithms are needed to help monitor and discriminate the indicators of patterned behaviors. Specifically, one would consider that maritime activities are repetitive behaviors from observed multi-modality data as time series. The effort applies a range of spatio-temporal analytic algorithms for separating patterns and anomalies which can be used to infer further emerging behaviors.

At present, this process is manual, labor intensive, time consuming, and not conducive to rapid decision making, and could benefit from presently available

machine learning (ML) and artificial intelligence (AI) methods. The objective is to develop an improved and automated methodology that can assist in a more rapid and accurate identification of anomalous behavior. Our contribution of this paper is summarized as the following:

- I develop a systemic ML/AI assisted analytic process for monitoring, detection, and classification of multi-modality data. I show an unsupervised AI transformer, namely variable autoencoder (VAE), and related theory and workflow pipelines in the different time scales of hours, minute, and seconds for anomaly and detection.
- I apply a VAE to a big data set collected from a distributed acoustic sensing (DAS) infrastructure in the Arctic area in Alaska. The DAS data set in this paper is a subset from the Sandia National Laboratories 160TB DAS data. The DAS collected can be used to detect sea ice, landfast ice, and other events [1]. The data set is used to show the feasibility of a VAE to detect events such as when waves, ships, and marine mammals pass by the DAS infrastructure. The methodology can be also extended to real-time surveillance, multi-modality anomaly detection applications for sea, land, and space.

## 2   Approaches

### 2.1   Anomaly Detection from Heatmap Time Series

I first map the need of detecting malicious, deceptive, and adversarial behaviors to a general model of anomaly detection. Anomaly detection is important for monitoring and discovering unknown risks of systems with structured and unstructured big data. Anomaly detection for a dynamic system represented as time series [2] assumes that there is a normal state of a system, and the difference between a previously unknown state and the normal state of the system, would be anomaly and require more attention. Event detection can be considered as anomaly detection from time series.

For example, as shown in Fig. 1, for a sensor data set is collected in an interval for a week, it can be pre-processed into a time series of images in seconds. Each second is mapped to a power density spectral density map or heatmap which represents the strength of the observed signals in terms of locations and frequencies of events in that second.

### 2.2   Variational Autoencoder (VAE)

The state-of-the-art anomaly detection techniques are associated with deep learning and analytics [3] for big data surveillance ranging for cyber, social media, and innovation discovery [4]. I consider VAE-based AI transformers, which belong to a type of neural network architecture, capable of modeling cognitive self-attention [5] and translate meaning between different modalities for anomaly detection. VAEs [6, 7] belong to a class of self-supervised and generative AI.
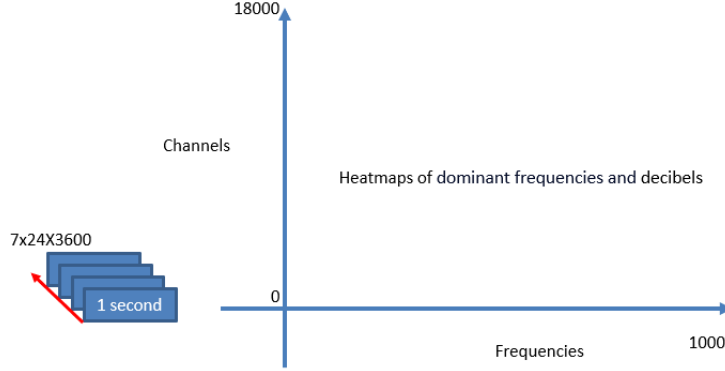
**Fig. 1.** A big data is pre-processed into a time series of images in hours, minutes, or seconds. Each time point is mapped to a power density Spectral density maps or heatmap which represents the strength of the observed signals in terms of locations (channels) and frequencies of events.

VAEs can flexibly incorporate traditional deep learning such as convolutional neural networks (CNNs), commonly used in anomaly detection tools of video like real-time streams such as YoLo [8], and other AI transformers such as BERT Pre-Training of Image Transformers (BEiT) [9].

Fig. 2 shows the workflow and pipeline of using VAE to perform anomaly, event, and object detection from the data.

### 2.3 Variational Autoencoder (VAE) Theory [10, 11]

Unsupervised learning and semi-supervised learning models focus on the characteristics of input data and require less labeled data. Generative AI (GAI) models combine probabilistic models and deep learning [12]. A VAE is an unsupervised learning approach to learn hidden structures and data distributions of the input set. In other words, A VAE learns underlying so called latent variables and their distributions as more fundamental and interpretable representation of an observational data set. VAEs are related to the traditional statistical principal component analysis (PCA), Latent Dirichlet Analysis (LDA) [13, 14], and Hidden Markov Models [15].

In order to reconstruct a corresponding target image with some attribute parameters, let $x, \theta$ denote an input image and the parameters of the underlying true probability distribution $p(x|\theta)$ that generates the image, respectively. A generation model can be realized by maximizing a model's parameters $\theta$ for a log-likelihood function of input data $x$ in Equation (1):

$$\max_{\theta} \mathbf{E}_{x \sim p(x)} log(p(x)|\theta) = \max_{\theta} \int_{x \sim p(x)} log(p(x)|\theta) dx \qquad (1)$$
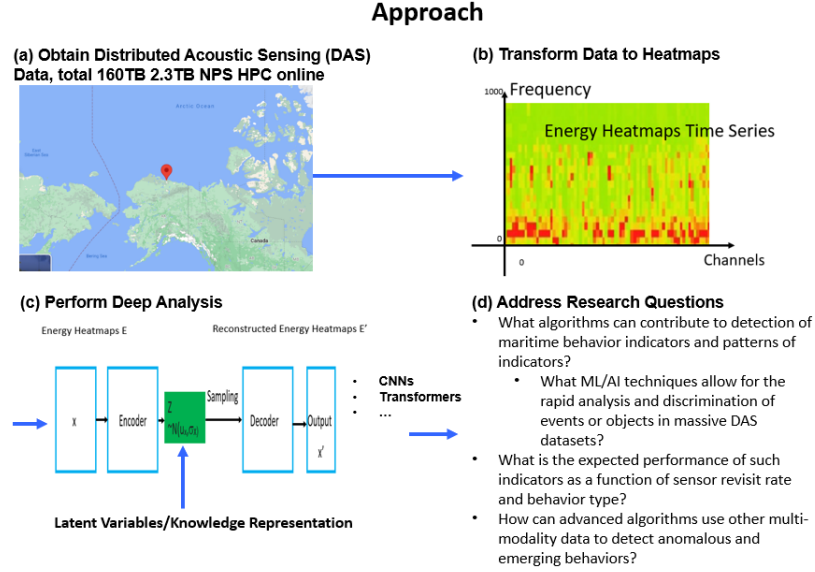
**Approach**



**Fig. 2.** The Workflow and Pipeline of Using VAE to Perform Anomaly, Event, and Object Detection from the Data.

where $\mathbf{E}_{x \sim p(x)}$ represents the expectation of input $x$ generated from the probability distribution $p(x)$. $p(x), p(x|\theta)$ are the prior distribution and posterior distribution of input image $x$. According to the VAE theory, a set of latent variables $z$ is added to $p(x|\theta)$ and $q(z|x)$ is used to represent the posterior distribution of $z$. The log-likelihood function $\log p(x|\theta)$ can be expressed as the sum of an evidence lower bound (ELBO) and the Kullback–Leibler- or KL-divergence between the true distributions $q(z|x)$ and a smooth (i.e., variational) parameterized $q(z|x,\theta)$ as shown in Equation (2):

$$log(p(x|\theta)) = \mathbf{E}_z q(z|x) log \frac{p(x,z|\theta)}{q(z|x)} + \mathbf{KL}(q(z|x), q(z|x,\theta)) \tag{2}$$

The first term in Equation (2) is the ELBO term – because of the non-negativity of KL-divergence, maximizing the ELBO term, which is easier, would result in the maximized total log-likelihood. The KL-divergence term is served as a regularization for learning the latent representation z. The decomposition is proved using the traditional Expectation-Maximization (EM) theory, which is widely used in statistical maximum likelihood estimation, density estimation, HMM, and LDA models. $\theta$ denotes the unknown parameters of of deep neural networks that can be learned from data or with traditional approaches like the EM algorithm.

The ELBO term in Equation (2) can be further decomposed into

$$\mathbf{E}_z q(z|x) log \frac{p(x,z|\theta)}{q(z|x)} = \mathbf{E}_z log(p(x|z)) - \mathbf{KL}(q(z|x), p(z|\theta), \tag{3}$$

where $\mathbf{E}_z$ represents the expectation of $z$ generated from probability density function $q(z|x)$. The first term in Equation (3) maximizes the probability of generating an input image $x$ from the latent variables $z$, i.e., $p(x|z)$, which represents the reconstruction accuracy. The second term in Equation (3) minimizes the difference between the distribution $q(z|x)$ and $p(z|\theta)$, which represents the distribution loss. $q(z|x)$ and $p(z|\theta)$ are usually modeled as normal distribution with different means $\mu$ and standard deviations $\sigma$, who are both modeled as neural networks as shown in Fig. 3.



**Fig. 3.** A VAE Pipeline

### 2.4   VAE for Multi-modality Data

In many anomaly detectipn problems, sensor data sets come in multi-modality such as audio, text, and imagery. They need to be fused and combined for downstream anomaly detection and other event/object classification/decision making tasks. As shown in Figure 3(a). Figure 3(b) shows an exemplar code to implement an image reconstruction VAE.

In Fig. 4 (a), when an input data $\mathbf{x}$ is an image, it is first fed into an encoder of a deep convolutional neural network (CNN) to generate 256 latent variables. A resampling process is used to generate reconstructed $\mathbf{x}'$ from a decoder. Fig. 4 (b), when an input data $\mathbf{x}$ is a text, it is first fed into an encoder of a BERT model to generate 768 latent variables. A resampling process is used to generate reconstructed $\mathbf{x}'$ from a decoder. Fig. 4(c) one can combine the latent variables from (a) and (b) of multi-modality to train a new decoder for downstream event or object classification task. The decoder for classification can be trained with a few labeled data (i.e., a few shots) because the latent variables from multi-modality data has already been trained with a large amount of unlabeled data.

### 2.5   VAE for Cross-modality Data Generation

Further more, when dealing multi-modality data, e.g., image fused with text, one can train separate latent variables (z) in each individual domain and then train a translator between the latent variables between two multi-modalities using semantic links (e.g., captions and images). Such multi-modality VAEs are used
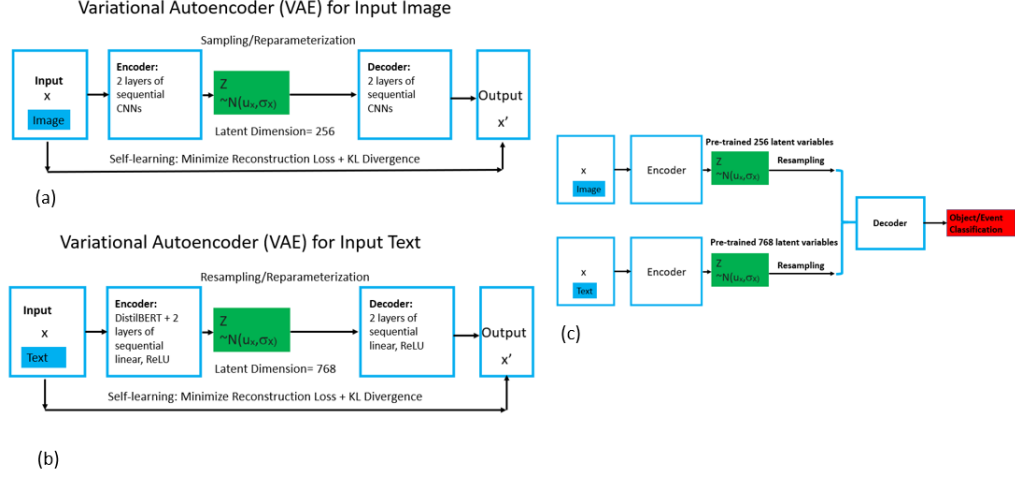
**Fig. 4.** A Multi-modality Pipeline

in the SAR image analysis [16–21] to translate synthetic images to more realistic images.

For example, [13, 7] applied modified VAEs for cross-domain data generation. The task of using synthetic SAR images to generate real SAR images are considered as a cross-domain data generation [19]. A latent representation has the four overall benefits using VAEs [13]: i) implicit latent decomposition into shared and private subspaces, ii) coherent joint generation over all modalities, iii) coherent cross-generation across individual modalities, and iv) improved model learning for individual modalities through multi-modal integration. Generative adversarial neural networks (GANs) [18, 20, 21] need labeled training data and tend to generate sharper images, however, interpretability is not enough.

### 2.6   Knowledge Graph (KG)

My KG is based on on setting of a graph neural network (GNN). In GNNs [22–25], the neural message passing refers to that a hidden embedding such as $z(i)$ is updated for each node $i$ based on the information gathered from its graph neighborhood $G(i)$ in Equation (4):

$$z(t+1)_i = Update(z(t)_i, Message(t)(z(t)_j, \forall j \in G(i)) \tag{4}$$

Update and Message are typically modeled as neural networks. Message generates a message based on the aggregated information from the graph neighborhood of node $i$. In this work, I use the KG for the time dimension in the decoder of a VAE. I regulate the hidden embedding $z(t)$ which represent the normal situations and patterns in the time scales of $z(t+1)$ and reconstruct z to preserve the smoothness between time points. Assume $\mu$ is the mean neural network for

the latent variable $z$ and $x$ is an input data in Equation (5):

$$z(\mu, t+1) = Encoder(x(t)) + KG(z(\mu, t)) \tag{5}$$

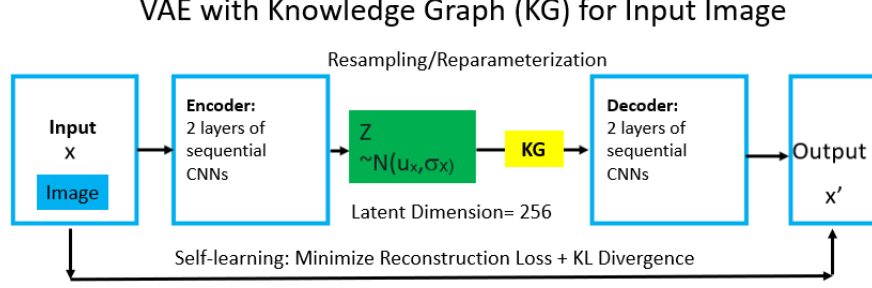## VAE with Knowledge Graph (KG) for Input Image



**Fig. 5.** A VAE-KG Pipeline: Total number of parameters is  500M,

Fig. 5 shows a KG is inserted into a VAE next to the latent variables **z** as the overall architecture.

## 3   Data Set

The data is an open source data set from the Sandia National Laboratories and Distributed Acoustic Sensing (DAS) [1]. DAS is a special fiber optic cable of total 36 kilometers in the Alaska arctic area with 18,000 channels generate measurements. Each channel has a sea location (latitude, longitude, bearing, depth). Each channel covers two meters. Measurements are sampled in 1kHZ (1000 samples per second). Total data include eight weeks of summer and winter data of 160TB. The data set are mainly collected for studying seismic events so data are in Seismic Analysis Code events (SAC) or TDMS formats. This paper focuses on the data set of 2.3TB, 97 hours of 1800 channels online in a high performance computing center (HPC). DAS has the potential to sense higher resolutions for example, detecting waves. Current methods for detecting anomaly and event for a DAS data is to using template matching or k-means clustering [26]. A separate hydrophone network in the same area can be used label a DAS for a supervised learning [27].

## 4   Initial Findings

### 4.1   Pre-processing of DAS into Images

Fig. 6 shows the DAS data set is transformed into images in different time scales. Table 1 shows corresponding numbers of images corresponding different
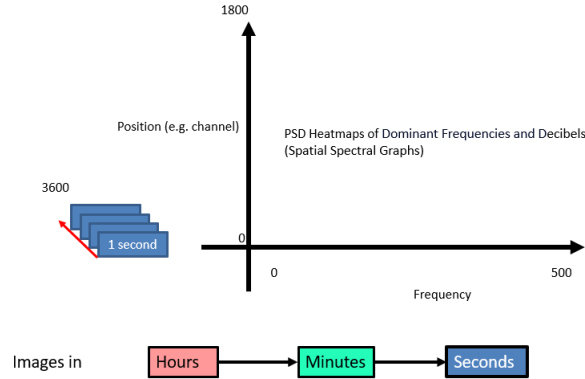
**Fig. 6.** DAS Image Input Heatmaps in Different Time Scales

**Table 1.** DAS Data Set Transformed into Images in Different Time Scales

| Time Scales | Number of Images |
|---|---|
| Hours | 97 |
| Minutes | 5,760 |
| Seconds | 338,184 |

time scales.

Fig. 7 shows the VAE reconstruction error for the 97 Images in the time scale of hours. Higher reconstruction errors indicate possible anomalies and events at the corresponding hours.

### 4.2   VAE Reconstruction Process

Fig. 8 shows an example of high reconstruction error, which indicate events of waves and objects in specific locations and frequencies.

Fig. 9 shows an example of the VAE reconstruction process in iteration 1, 2, 10. The intensity of all the images are scaled between 0 and 55. The axes and labels are kept to illustrate the reconstruction process. The blurring of the input images, compared to the one in high resolution in Fig. 8, is caused by the scale down process in the VAE to keep the number of parameters relatively small. As shown in Fig. 9, initially in Iteration 1 and 2, content of the image are not reconstructed, however in Iteration 10, most content of the image including numbers in axes is reconstructed. Some content, potentially anomalies,for example, in the box Iteration 10 around the bright dot on the top left of the image,is not reconstructed.

Fig. 10 shows the reconstruction errors as the indicators for normal and anomaly hours. Fig. 11 shows the corresponding VAE original images (inputs), reconstructed ones, and high resolution original images. The high resolution
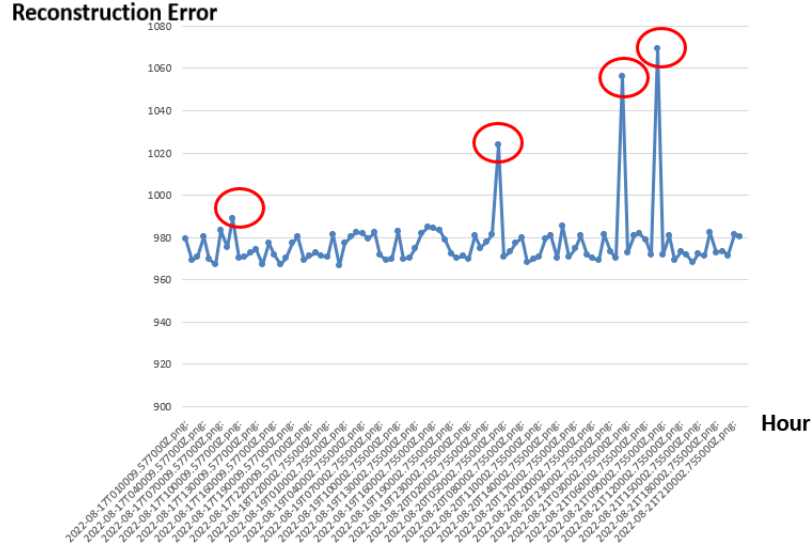
**Fig. 7.** VAE reconstruction error for the 97 Images in the time scale of hours: Higher reconstruction errors show possible anomalies and events at the corresponding hours.
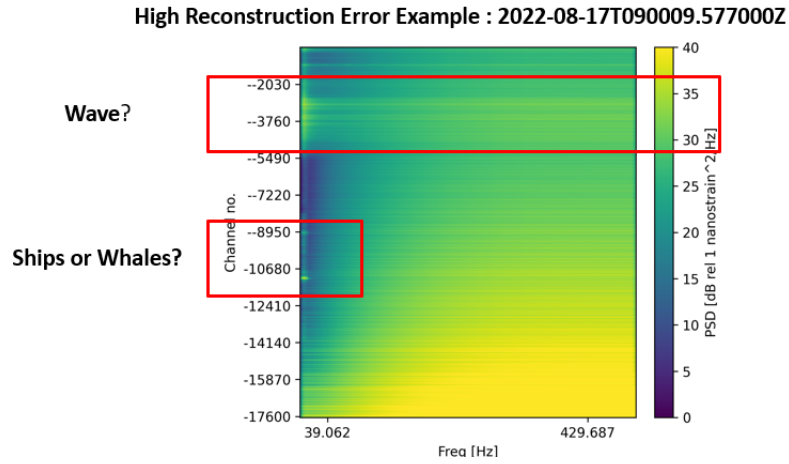


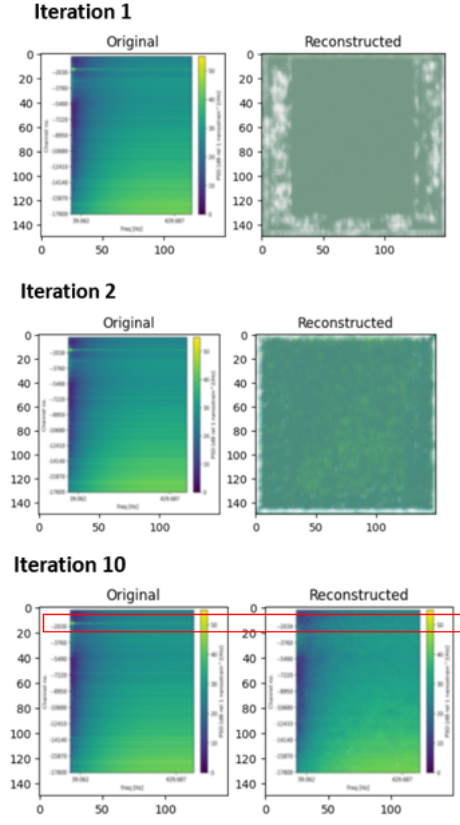**Fig. 8.** An example has a high reconstruction error.

**Fig. 9.** An example illustrate the reconstruction process.

original images are scaled into smaller original images in the VAE algorithm to fit to the number of parameters in the model. Compared to the normal ones, the anomalous original input images (to VAE) have high reconstruction errors in Fig. 10, their reconstructed images are correspondingly brighter than the originals.



**Fig. 10.** High reconstruction errors are indicators of anomalies.

Fig. 12 shows a drill-down to the minutes for the reconstruction errors. Fig. 12 shows corresponding examples of normal and anomaly of original and reconstructed images in Fig. 13. The anomaly image has a bright spot in the original images which is not reconstructed, therefore, an indication of anomaly, for example, an object passing by. The normal one has an accurate reconstruction.

### 4.3   Effects of Adding Knowledge Graphs

Fig. 14 Compares reconstruction errors in the time scale of hours with and without knowledge graphs. Fig. 15 shows the corresponding example of the anomaly of original, reconstructed, and high resolution original images in Fig. 14. The KG model smooths away false alarms.

## 5   Discussion

The VAE-based ML/AI models show the potentials to handle big data. Analyzing such big data for anomaly detection is extremely challenging, where the time series continue through time and space and are tremendous in size. The analysis of such big data can overwhelm and challenge classic analytic methods.
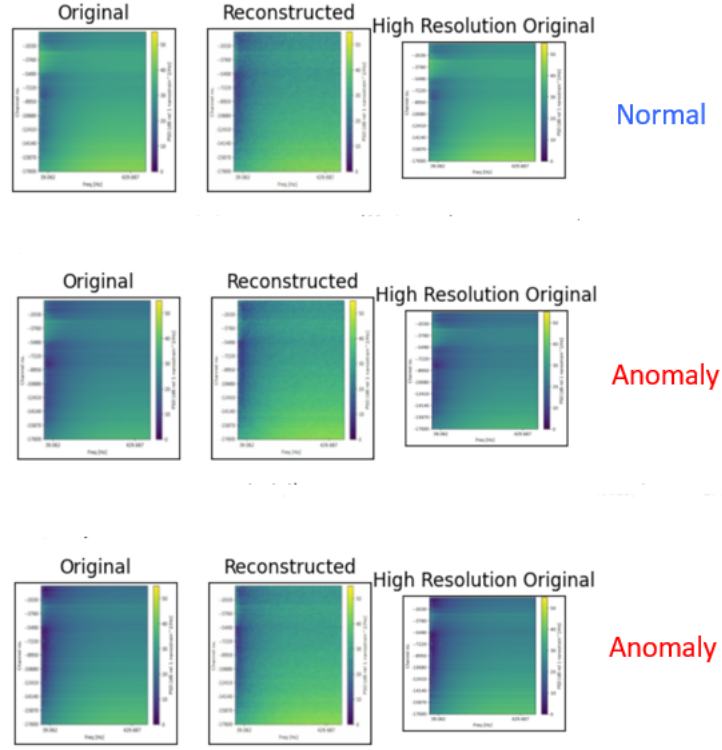
**Fig. 11.** Corresponding examples of original and reconstructed images in Fig. 10: Anomalous original iamges have darker lower right corners.
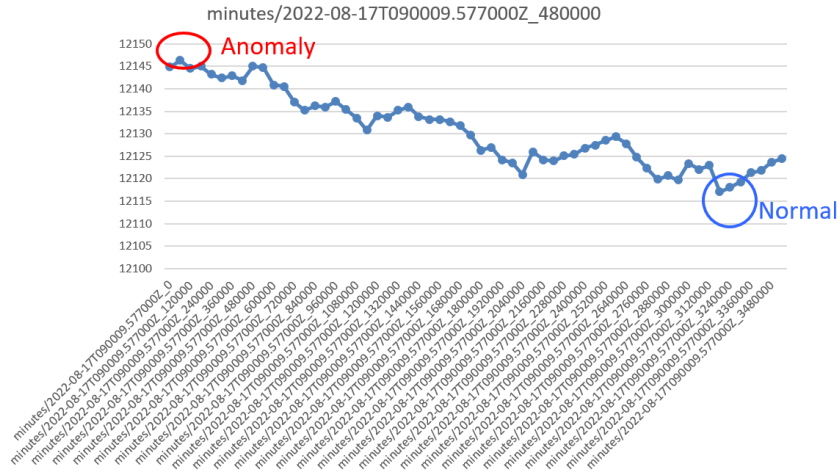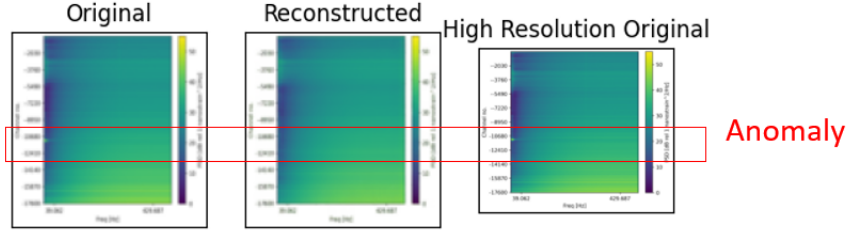


**Fig. 12.** Drill-down to the time scale of minutes.

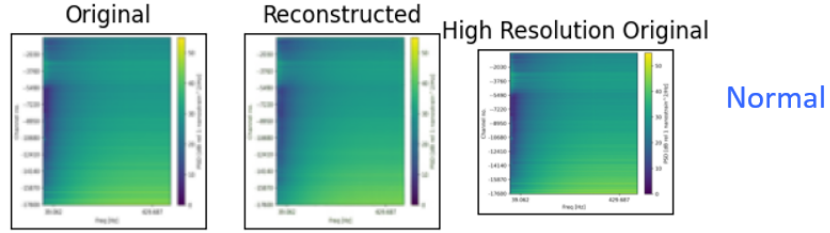## 2022-08-17T090009.577000Z_spatiospectral_60000



**Fig. 13.** Corresponding examples of normal and anomaly of original, reconstructed, and high resolution original images in Fig. 12.
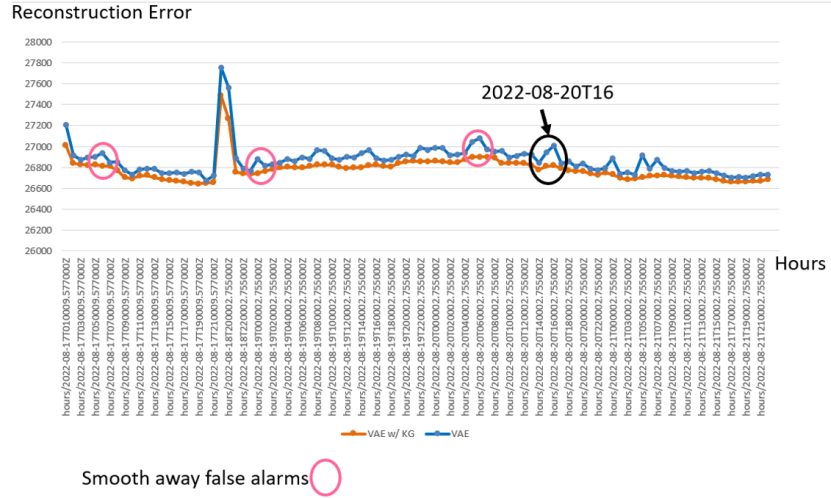


**Fig. 14.** Comparing reconstruction errors in the time scale of hours with and without knowledge graphs: the KG model smooths away false alarms.
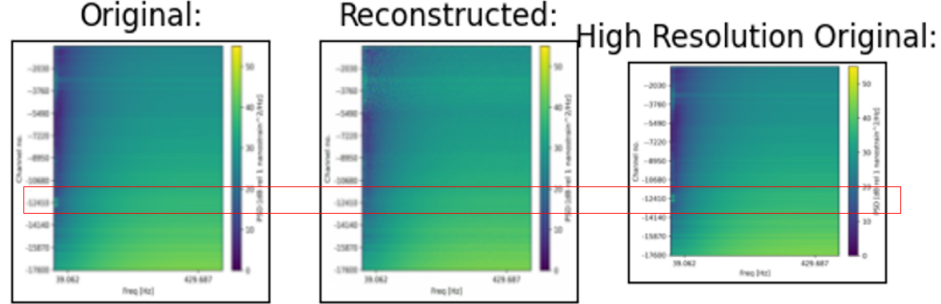
**Fig. 15.** Corresponding example of the true anomaly with its original, reconstructed, and high resolution original images in Fig. 14.

Further heatmap time series images such as DAS are different from the typical video data [8], since objects presented in the data are not similar to typical objects. Therefore pre-trained ML/AI models are rarely available.

The VAE-based ML/AI models show the potentials to handle multi-modality data. Different data types and modalities can also present difficulty for traditional methods. VAEs learn underlying latent variables and distributions which are more fundamental representation of observational data. A VAE architecture is capable of modeling and translating meaning between different modalities, therefore, has the potential to drastically improve the downstream classification and recognition tasks via transfer learning or few-shot learning. For example, different data type of sensor measurement, text, image, and audio data can represent a same entity using an architecture of multiple and mixed deep neural networks and AI transformers.

## 6    Recommendations for Further Research

Further research is needed to validate the results from the ground truth such as matched AIS ship tracks in the same area. Public reports of seismic events and marine mammal activities can be also used as validation sources. Further work is also needed to work on drilling down strategies to more specific areas, frequencies, and resolutions.

## 7    Conclusion

The DAS data set provides a study case for events happened in a different tempo such as passing waves, ships, and marine mammals that could happen in different time scales. We found that VAE models are able to detect anomalies and events in different time scales, when reconstruction error is used as an anomaly indicator. Knowledge graphs can smooth away false alarms.

## Acknowledgment

## References

1. Baker, M. G. and Abbott, R. E. (2022). Rapid refreezing of a marginal ice zone across a seafloor distributed acoustic sensor. In Geophysical Research Letters Volume 49, Issue 24 https://doi.org/10.1029/2022GL099880.
2. Doshi, K. and Yilmaz, Y. (2020). Continual Learning for Anomaly Detection in Surveillance Videos", Computer Vision Foundation,[2004.07941].
3. Welling, M. and Kingma, D. (2019). An Introduction to Variational Autoencoders. Foundations and Trends in Machine Learning. 12 (4): 307392. arXiv:1906.02691.
4. Landauer et al (2022). Deep Learning for Anomaly Detection in Log Data: A Survey. https://arxiv.org/abs/2207.03820v1.
5. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need. In https://arxiv.org/abs/1706.03762.
6. Kingma, D.P. and Welling, M. (2013). Auto-encoding variational bayes. arXiv:1312.6114v10.
7. Rezende, D.J., Mohamed, S. and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. arXiv:1401.4082.
8. Doshi, K. and Yilmaz, Y. (2020). Continual learning for anomaly detection in surveillance videos. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pp. 254-255.
9. Bao, H., Dong, L., Piao, S. and Wei, F. (2023). BEiT: BERT pre-training of image transformers, https://arxiv.org/abs/2106.08254.
10. Shi, Y., Siddharth, N., Paige, B., and Torr, P. (2019). Variational mixture-of-experts autoencoders for multi-modal deep generative models. In https://arxiv.org/abs/1911.03393.
11. Wu, M. and Goodman, N. (2018). Multimodal generative models for scalable weakly-supervised learning. In Advances in Neural Information Processing Systems, pages 5580–5590.
12. Rezende, D.J., Mohamed, S., and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. arXiv:1401.4082.
13. Blei, D. M. and Jordan, M. I. (2006). Variational inference for dirichlet process mixtures. In Bayesian analysis, 1(1):121– 143.
14. Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2017). Variational inference: A review for statisticians. In Journal of the American statistical Association, 112(518):859–877.

15. Zhang, A., Gultekin, S., and Paisley, J. (2016). Stochastic variational inference for the hdp-hmm. In Artificial Intelligence and Statistics, pp. 800–808.

16. Toizumi, T.; Sagi, K.& Senda, Y. (2018). Automatic association between SAR and optical images based on zero-shot learning. In Proceedings of the IGARSS 2018–2018 IEEE International Geoscience and Remote Sensing Symposium, Valencia, Spain, 22–27 July 2018; pp. 17–20.

17. Song, Q., Chen, H., Xu, F., and Cui, T.J. (2020). EM simulation-aided zero-shot learning for SAR automatic target recognition. IEEE Geosci. Remote Sens. Lett. 2020, 17, 1092–1096.

18. Hu, X., Feng, W., Guo, Y., and Wang, Q. (2021). Feature learning for sar target recognition with unknown classes by using cvae-gan. In Remote Sensing September 2021 13(18):3554, DOI:10.3390/rs13183554.

19. Malmgren-Hansen, D., Kusk, A., Dall, J., Nielsen, A., Engholm, R., and Skriver, H. (2017). Improving sar automatic target recognition models with transfer learning from simulated data. In IEEE Geosci. Remote Sens. Lett. 2017, 14, 1484–1488.

20. Larsen, A., Sonderby, S., Larochelle, H., and Winther, O. (2016). Autoencoding beyond pixels using a learned similarity metric. In Int. Conf. Mach. Learn. 2016, 48, 1558–1566.

21. Bao, J., Chen, D., Wen, F., Li, H., and Hua, G. (2017). CVAE-GAN: Fine-grained image generation through asymmetric training. In arXiv 2017, 2745–2754, arXiv:1703.10155.

22. Goodall, R. and Lee, A. (2020). Predicting materials properties without crystal structure: deep representation learning from stoichiometry. In Nat Commun 11, 6280.

23. Barp, A., Costa, L. D., Franca, G., Friston, K., Girolami, M., Jordan, M., and Pavliotis, G. (2022). Geometric methods for sampling, optimization, inference and adaptive agents. In https://arxiv.org/abs/2203.10592

24. Bronstein, M. M., Bruna, J., Cohen, T., and Velickovic, P. (2021). Geometric deep learning, grids, groups, graphs, geodesics, and gauges. In https://arxiv.org/pdf/2104.13478.pdf

25. Ma, Y., Lium, Z., Nan Yang, N., Xu, H., and Zhang, G. (2024). Research on knowledge graph construction and semantic representation of low earth orbit satellite spectrum sensing data. In Electronics 2024, 13(4), 672; https://doi.org/10.3390/electronics13040672.

26. Pena Castro, A. F. F., Schmandt, B., Baker, M., and Abbott, R. (2022). Automated high-resolution tracking of sea ice extent offshore Oliktok Point, Alaska, using distributed acoustic sensing and machine learning. In AGU Fall Meeting, Chicago, IL, 12-16 December 2022.

27. Berry, M. (2024). Tracking sea-ice extent and detecting boats in the Beaufort Sea using Distributed Acoustic Sensing and machine learning. Presentation at the 8th Annual Workshop on Naval Applications of Machine Learning (NAML 2024), 11-14 March, 2024, San Diego, CA.