



**HAL**  
open science

## Exploring unseen 3D scenarios of physics variables using machine learning-based synthetic data: An application to wave energy converters

César Quilodrán-Casas, Qian Li, Ningbo Zhang, Sibó Cheng, Shiqiang Yan, Qingwei Ma, Rossella Arcucci

### ► To cite this version:

César Quilodrán-Casas, Qian Li, Ningbo Zhang, Sibó Cheng, Shiqiang Yan, et al.. Exploring unseen 3D scenarios of physics variables using machine learning-based synthetic data: An application to wave energy converters. *Environmental Modelling and Software*, 2024, 177, pp.106051. 10.1016/j.envsoft.2024.106051 . hal-04660519

**HAL Id: hal-04660519**

**<https://hal.science/hal-04660519>**

Submitted on 24 Jul 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Highlights

### **Exploring unseen 3D scenarios of physics variables using machine learning-based synthetic data: an application to wave energy converters**

César Quilodrán-Casas, Qian Li, Ningbo Zhang, Sibó Cheng, Shiqiang Yan, Qingwei Ma, Rossella Arcucci

- We propose a model surrogate for Wave Energy Converters Computational Fluid Dynamics.
- Generative models are used to create unseen scenarios of velocity and dynamic viscosity responses.
- Two different generative models are employed and expand the design space.

# Exploring unseen 3D scenarios of physics variables using machine learning-based synthetic data: an application to wave energy converters

César Quilodrán-Casas<sup>a,b</sup>, Qian Li<sup>c</sup>, Ningbo Zhang<sup>c</sup>, Sibó Cheng<sup>a</sup>, Shiqiang Yan<sup>c</sup>, Qingwei Ma<sup>c</sup> and Rossella Arcucci<sup>a,b</sup>

<sup>a</sup>Data Science Institute, Imperial College London, London, United Kingdom

<sup>b</sup>Department of Earth Science and Engineering, Imperial College London, London, United Kingdom

<sup>c</sup>City University, London, United Kingdom

---

## ARTICLE INFO

### Keywords:

Generative models  
Synthetic data  
Wave energy converters  
Model surrogate

## ABSTRACT

This work aims to use machine learning to produce synthetic data of wave energy converters from time-expensive 3D simulations based on computational fluid dynamics models. The simulations to analyse the response of these systems to incoming waves are lengthy and computationally expensive to obtain. Here, we explore the use of a beta-VAE and a Principal Components-based adversarial autoencoder for generating new synthetic data. The compression plus the generation of synthetic data introduces an exceptionally fast model surrogate of the original simulation and delivers more samples of either dynamic viscosity and velocity fields, enlarging the design space. The newly generated synthetic samples can have a speed up from 5 to 6 orders of magnitude. The new design space can be used to improve the prediction of dynamic viscosity given the velocity fields. The generative model has the potential to capture the transition and the new physical phenomena under extreme initial conditions.

---

## 1. Introduction

1 Due to the increasing demand for clean energy, various renewable energy resources are being explored, among  
2 which wave energy is one of the topics with the greatest potential (Glendenning, 1977). Various forms of oscillating  
3 Wave Energy Converter (WEC) devices have been developed to capture wave energy to generate electricity (Antonio,  
4 2010).

5 WECs are devices that convert the kinetic and potential energy associated with a moving ocean wave into useful  
6 mechanical or electrical energy. A point absorber is a floating structure that absorbs energy from all directions through  
7 its movements at or near the water's surface. It converts the motion of the buoyant top relative to the base into electrical  
8 power. An interesting interaction to predict is how the dynamic viscosity around the point absorber reacts to the stimulus  
9 of an incoming wave, and this can be predicted with the velocity fields and the dynamic viscosity from the previous  
10 time level (Jin, Patton, and Guo, 2018).

11 In the process of studying a complete WEC system, it is essential to obtain a general and applicable hydrodynamic  
12 description of how the device interacts with the incident waves. This mathematical description is important to suggest  
13 the design of the power take-off (PTO), as well as the development of the control system since these WEC subsystems  
14 are influenced by the dynamic interaction that the WEC device has with the movement of the waves (Son and Yeung,  
15 2017). However, Computational Fluid Dynamics (CFD) simulations of how these systems behave due to ocean wave  
16 perturbations can be computationally expensive to run and time-consuming. Often it can produce a small number of  
17 samples to work with, yielding poor predictors for a predictive model of the dynamic viscosity due to the lack of a  
18 high number of samples. Therefore, an attractive solution is to generate new simulations, at a considerable speedup,  
19 that learn from the original CFD simulations.

---

\*Corresponding author

\*\*Principal corresponding author

✉ c.quilodran@imperial.ac.uk (C. Quilodrán-Casas); qian.li.8@city.ac.uk (Q. Li); ningbo.zhang.2@city.ac.uk (N. Zhang); sibó.cheng@imperial.ac.uk (S. Cheng); shiqiang.yan@city.ac.uk (S. Yan); q.ma@city.ac.uk (Q. Ma); r.arcucci@imperial.ac.uk (R. Arcucci)

ORCID(s):

20 A classic approach to this difficult problem is Data Augmentation (DA). In DA, the network is trained using  
21 additional synthetic data. DA was introduced in object recognition in LeCun, Bottou, Bengio, and Haffner (1998).  
22 The advantages of DA are that it increases the size of the training data; eliminates the overfitting problem; and it  
23 makes the network more robust to data variations that may exist in any real-world application. The basic idea behind  
24 DA is to apply transformations so that the semantics of the labels associated with the data does not change. By training  
25 the network with this extra data, one would expect, its performance on unseen data to be enhanced.

26 Quilodrán-Casas, Arcucci, Mottet, Guo, and Pain (2021) used generative models to create stable rollouts for air  
27 pollution. The experimental design space of microfluidics has also been augmented using generative networks (Chagot,  
28 Quilodrán-Casas, Kalli, Kovalchuk, Simmons, Matar, Arcucci, and Angeli, 2022). Other successful implementations  
29 of generative-based augmentation are Karras, Aila, Laine, and Lehtinen (2017); Berthelot, Schumm, and Metz (2017);  
30 Radford, Metz, and Chintala (2015).

31 The generation of high-quality synthetic data allowed the augmentation of small-sample data sets (Forestier,  
32 Petitjean, Dau, Webb, and Keogh, 2017; Hoffmann, Bar-Sinai, Lee, Andrejevic, Mishra, Rubinstein, and Rycroft,  
33 2019). Although the use of the synthetic data needs to be developed and adapted for each case (Chen, Lu, Chen,  
34 Williamson, and Mahmood, 2021), it can be a powerful tool to increase the robustness and adaptability of data-driven  
35 models (Yoon, Jordon, and Schaar, 2018; Quilodrán-Casas et al., 2021). However, Machine Learning (ML) often  
36 requires representative data to be effective (Zhou, Pan, Wang, and Vasilakos, 2017; Li, Zhang, Chen, Shen, and Niu,  
37 2023; Sadeghi, Nguyen, Hsu, and Sorooshian, 2020; Razavi, 2021).

38 Recent advances in ML have shown strong predictive power to determine complex correlations and find patterns  
39 between inputs and outputs (Goodfellow, Bengio, and Courville, 2016). ML has been employed in wave energy studies  
40 previously. Rodriguez-Delgado, Bergillos, and Iglesias (2019) have used neural networks to assess the efficiency of  
41 WECs. Li, Yuan, and Gao (2018) used Deep Learning (DL) for assessing the energy absorption of a WEC. Sarkar,  
42 Contal, Vayatis, and Dias (2016). Sclavounos and Ma (2018) and Mousavi, Ghasemi, Dehghan Manshadi, and Mosavi  
43 (2021) used ML for forecasting the time series response of WEC and wave energy conversion rates.

44 However, generative methods to augment datasets have not been used to augment data and reproduce unseen physics  
45 conditions. Here, we use data from a high-fidelity CFD simulating 3D velocity and dynamic viscosity of WECs in some  
46 scenarios and we use ML models to develop a surrogate model to reproduce unseen physics conditions for WECs. Due  
47 to the problem complexity and the high dimensionality, running this CFD simulation can be computationally expensive  
48 and time-consuming. To tackle this bottleneck, in this paper, a  $\beta$ -Variational Autoencoder (VAE) and a Principal  
49 Components-based Adversarial Autoencoder (PC-AAE) are used to generate synthetic data to enlarge the experimental  
50 dataset, explore unseen scenarios of this Three-dimensional (3D) simulation in a fast manner. The generated data can  
51 also be used to train a ML surrogate model for predicting future wave dynamics. To the best of our knowledge, this is  
52 the first reported attempt to obtain synthetic data of WECs 3D CFD simulations using generative models.

53 The contribution of this paper lies in the use of generative networks such as  $\beta$ -VAE and PC-AAE for the generation  
54 of unseen synthetic data to expand the relationship between velocity and dynamic viscosity in CFD models, bypassing  
55 running a different CFD simulation in a supercomputer to create more samples. These generative models allow  
56 increasing the design space and access to sampled data from a matched distribution of the original data.

57 This paper is organised as follows. Section 2 describes the implementation of the  $\beta$ -VAE and the PC-AAE. Section  
58 3 describes the irregular wave test case of a point absorber WEC. Section 4 shows the results and discussion of  
59 the samples of velocity fields and dynamic viscosity generated by these networks. And finally, Section 5 presents  
60 a summary, conclusion and future work.

## 61 2. Generative models

62 In this paper, two different generative models are used:  $\beta$ -VAE and PC-AAE. These two methods were chosen as  
63 they are well-established and well-documented for producing high-fidelity results.

### 64 2.1. $\beta$ -Variational autoencoder ( $\beta$ -VAE)

65 Autoencoder (AE)s were developed to reconstruct high-dimensional data using a neural network model composed  
66 of an encoder and a decoder. AEs can also reduce the dimensionality of the system with the encoder mapping the  
67 input onto a bottleneck layer. Furthermore, a  $\beta$ -VAE instead of mapping onto a fixed vector, maps the input onto an  
68 arbitrary distribution (Higgins, Matthey, Pal, Burgess, Glorot, Botvinick, Mohamed, and Lerchner, 2016). The  $\beta$ -VAE  
69 is a modification of the VAE with a special emphasis to discover disentangled latent factors. Following the same

70 incentive in VAE, the probability of generating real data is maximised whilst maintaining the distance between the real  
71 and estimated posterior distributions small.

72 Let  $Q$  and  $P$  denote the encoder and decoder, respectively. Moreover, let  $q(\mathbf{z}|\mathbf{x})$  and  $p(\tilde{\mathbf{x}}|\mathbf{z})$  denote the encoding and  
73 decoding distributions, respectively., where  $\mathbf{x}$  is the input vector,  $\tilde{\mathbf{x}}$  is the reconstructed input, and  $\mathbf{z}$  is the latent space.  
74 As suggested by Makhzani, Shlens, Jaitly, Goodfellow, and Frey (2015), a Gaussian posterior can be used assuming  
75 that  $q(\mathbf{z}|\mathbf{x})$  is a Gaussian distribution, where its mean  $\mu$  and variance  $\sigma$  are predicted by the encoder  $Q$  by adding two  
76 dense layers of means  $\mu$  and  $\log \sigma$  to the final layer of the encoder  $Q$ , and return  $\mathbf{z}$  as a vector of samples (Kingma and  
77 Welling, 2013). To ensure that  $\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mu, \sigma^2)$ , the aggregated posterior, the reparameterisation trick described  
78 by Kingma and Welling (2013) was used for backpropagation.

The minimisation of the Kullback-Leibler Divergence Score (KL) loss ( $\mathcal{L}^{KL}$ ) quantifies how much the probability  
distribution  $a$  differs from the probability distribution  $b$  as:

$$\mathcal{L}^{KL}(a, b) = - \sum a \log \left( \frac{b}{a} \right) \quad (1)$$

where, in this case,  $a = q(\mathbf{z}|\mathbf{x})$  and  $b = g(\mathbf{z}) = \mathcal{N}(0, \mathbf{I})$ , the arbitrary prior, and  $\mathbf{I}$  is the identity matrix. In other words,  
we expect the latent distribution  $g(\mathbf{z})$  to approximate a centred and standard Gaussian distribution. Adam is used as  
the optimiser (Kingma and Welling, 2013). The total loss  $\mathcal{L}^\theta$  is then defined as  $\mathcal{L}^\theta = \lambda \mathcal{L}^{KL} + \mathcal{L}^{mse}$ , where  $\lambda = 0.001$   
acts as a regulariser. The reconstruction error  $\mathcal{L}^{mse}$  is the mean squared error defined as:

$$\mathcal{L}^{mse} = \|\mathbf{x} - \tilde{\mathbf{x}}\|^2 \quad (2)$$

79 where  $\tilde{\mathbf{x}}$  is the reconstructed input of experimental data, defined as  $\tilde{\mathbf{x}} = P(Q(\mathbf{x}))$ . The inputs were scaled between 0  
80 and 1.

81 The implementation of the  $\beta$ -VAE is in Python using `pytorch` (Paszke, Gross, Massa, Lerer, Bradbury, Chanan,  
82 Killeen, Lin, Gimelshein, Antiga et al., 2019) and the `pytorch-lightning` wrapper (Falcon et al., 2019). The  
83 algorithm of the generation of synthetic samples using  $\beta$ -VAE is shown in Algorithm 1.

---

#### Algorithm 1: $\beta$ -VAE

---

**Require** :  $\theta^Q$  and  $\theta^P$  trainable parameters for encoder  $Q$ , and decoder  $P$ , respectively; batch sizes  $m^Q$  for  
 $Q$ ; number of epochs  $k$

**for**  $epochs = 0, \dots, k$  **do**

- Match a latent vector of size  $m_Q$  to a normal distribution  $\mathbf{z} \sim \mathbf{q}(\mathbf{z}|\mathbf{x})$
- Calculate the  $\mathcal{L}^{KL}(q(\mathbf{z}|\mathbf{x}), g(\mathbf{z}))$
- Calculate the  $\mathcal{L}^{mse}$
- Reconstruction error  $\mathcal{L}^\theta = \mathcal{L}^{KL} + \mathcal{L}^{mse}$
- Update the  $\beta$ -VAE parameters ( $\theta^Q$  and  $\theta^P$ ) via:  $\theta^{\beta\text{-VAE}} \leftarrow \text{Adam}(\mathcal{L}^{\beta\text{-VAE}})$

until convergence detach Decoder  $P$  and use it to generate samples

---

## 84 2.2. PC-based adversarial AE (PC-AAE)

As described by Lever, Krzywinski, and Altman (2017), Principal Component Analysis (PCA) is an unsupervised  
learning method that simplifies high-dimensional data by transforming it into fewer dimensions. The PCA consists in  
decomposing  $\mathbf{x}$  as  $\mathbf{x} = \mathbf{P}\mathbf{\Pi} + \bar{\mathbf{x}}$  where  $\mathbf{P} \in \mathcal{R}^{n \times n}$  are the Principal Components (PC)s of  $\mathbf{x}$ ;  $\mathbf{\Pi} \in \mathcal{R}^{n \times m}$  are the Empirical  
Orthogonal Functions (EOF)s; and  $\bar{\mathbf{x}}$  is the mean vector of the model. The dimension reduction of the system comes  
from truncating  $\mathbf{P}$  at the first  $\tau$  PCs as  $\mathbf{x}_\tau = \mathbf{P}_\tau \mathbf{\Pi}_\tau + \bar{\mathbf{x}}$ , with  $\mathbf{P}_\tau \in \mathcal{R}^{n \times \tau}$  and  $\mathbf{\Pi}_\tau \in \mathcal{R}^{\tau \times m}$ . To further reduce the  
system dimension, a number of recent researches combined PCA with deep learning AEs (see Cheng, Jin, Harrison,  
Quilodr an-Casas, Prentice, Guo, and Arcucci (2022); Cheng, Chen, Anastasiou, Angeli, Matar, Guo, Pain, and Arcucci  
(2023); Gong, Cheng, Chen, Li, Quilodr an-Casas, Xiao, and Arcucci (2022)). In this study, the principle components  
are used to train an Adversarial Autoencoder (AAE) (Makhzani et al., 2015). The functional of our PC-AAE is defined  
as:

$$f^{PC\text{-AAE}} : \mathbf{P}_{t_k} \rightarrow \tilde{\mathbf{P}}_{t_k} \quad (3)$$

85 where  $\mathbf{P}_{t_k}$  are the scaled PCs time series between -1 and 1 at time-level  $k$ . The AE consists of an encoder  $\mathcal{Q}$  and a  
 86 decoder  $\mathcal{P}$ , both mirrored fully-connected networks, where the scaled reconstructed PCs  $\tilde{\mathbf{P}}_{t_k} = \mathcal{P}(\mathcal{Q}(\mathbf{P}_{t_k}))$ . Let  $q(\mathbf{z}|\mathbf{P})$   
 87 and  $p(\tilde{\mathbf{P}}|\mathbf{z})$  be the encoding and decoding distributions, respectively. As suggested by Makhzani et al. (2015), we use a  
 88 Gaussian posterior and assume that  $q(\mathbf{z}|\mathbf{P})$  is a Gaussian distribution, where its mean and variance are predicted by the  
 89 encoder  $\mathcal{Q}$ . This is achieved by adding two dense layers of means  $\mu$  and  $\log \sigma$  to the final layer of the encoder  $\mathcal{Q}$ , and  
 90 return  $\mathbf{z}$  as a vector of samples. This is achieved similarly to the reparameterisation trick described above for  $\beta$ -VAE.

The adversarial training of PC-AAE includes a discriminator  $\mathcal{D}^A$  to distinguish between the real samples, given  
 by an arbitrary prior  $g(\mathbf{z})$ , and fake samples, given by  $q(\mathbf{z}|\mathbf{P})$ . Therefore, the adversarial autoencoder is regularised by  
 matching  $g(\mathbf{z})$  to  $q(\mathbf{z}|\mathbf{P})$ . The ground truth PCs  $\mathbf{P}$  are fed to the discriminator as real sequences (ground truth). Let,  
 $\mathcal{D}^A(\alpha, \gamma)$  represent the discriminator function with an input  $\alpha$  and a target label  $\gamma$  such that, for  $\alpha = \mathbf{z} \sim q(\mathbf{z}|\mathbf{P})$ ,  $\gamma = 1$   
 and for  $\alpha = \hat{\mathbf{z}} \sim p(\mathbf{z})$ ,  $\gamma = 0$ , where  $\hat{\mathbf{z}}$  is the latent space sampled from  $g(\mathbf{z})$ . The training of  $\mathcal{D}^A$  is based on the  
 minimisation of the binary cross-entropy loss ( $\mathcal{L}^{bce}$ ), using the Nesterov Adam optimizer (Nadam) (Dozat, 2016). The  
 adversarial losses  $\mathcal{L}^{adv}$  for  $\mathcal{D}^A$  and  $f^{PC-AAE}$  are then defined as:

$$\mathcal{L}_{\mathcal{D}^A}^{adv}(\mathbf{P}) = \mathcal{L}_{\hat{\mathbf{z}} \sim g(\mathbf{z})}^{bce}(\mathcal{D}^A(\hat{\mathbf{z}}, 1)) + \mathcal{L}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{P})}^{bce}(\mathcal{D}^A(\mathbf{z}, 0)) \quad (4)$$

$$\mathcal{L}_{f^{PC-AAE}}^{adv}(\mathbf{P}) = \mathcal{L}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{P})}^{bce}(\mathcal{D}^A(\mathbf{z}, 1)) + \mathcal{L}^{mse}(\tilde{\mathbf{P}}, \mathbf{P}) \quad (5)$$

91 where  $\mathcal{L}^{mse}$  is the Mean Squared Error (MSE) between  $\tilde{\mathbf{P}}$  and  $\mathbf{P}$ . The algorithm of the generation of synthetic samples  
 92 using PC-AAE is shown in Algorithm 2.

---

#### Algorithm 2: PC-AAE

---

**Require** :  $\theta^{\mathcal{Q}}$ ,  $\theta^{\mathcal{P}}$ , and  $\theta^{\mathcal{D}^A}$  trainable parameters for encoder  $\mathcal{Q}$ , decoder  $\mathcal{P}$ , and discriminator  $\mathcal{D}^A$ ,  
 respectively; number of discriminator iterations per AE iteration ( $n_{\mathcal{D}^A}$ ); batch sizes  $m_{\mathcal{Q}}$  for  $\mathcal{Q}$ ;  
 number of epochs  $k$

**for**  $epochs = 0, \dots, k$  **do**

- Discriminator training: **for**  $i = 0, \dots, n_{\mathcal{D}^A}$  **do**
  - Sample a latent vector of size  $m_{\mathcal{Q}}$  from a normal distribution  $\hat{\mathbf{z}} \sim g(\mathbf{z})$ , plus Gaussian noise
  - Fake samples  $\mathbf{z} \leftarrow \mathcal{Q}_{\theta^{\mathcal{Q}}}(\hat{\mathbf{z}})$
  - Update the discriminator  $\mathcal{D}^A$  to differentiate between real and fake sample via:  $\theta^{\mathcal{D}^A} \leftarrow \text{NAdam}(\mathcal{L}_{\mathcal{D}^A}^{adv})$
- Update the AE PC-AAE parameters ( $\theta^{\mathcal{Q}}$  and  $\theta^{\mathcal{P}}$ ) via:  $\theta^{PC-AAE} \leftarrow \text{NAdam}(\mathcal{L}_{PC-AAE}^{adv})$

until convergence Detach Decoder  $\mathcal{P}$  and use it to generate samples of  $\mathbf{P}$

---

Then the reconstruction  $\tilde{\mathbf{x}}$  to the physical space is given by:

$$\tilde{\mathbf{x}} = \tilde{\mathbf{P}}\mathbf{\Pi} + \tilde{\mathbf{x}} \quad (6)$$

93 The implementation of the PC-AAE is in Python using tensorflow (Abadi, Agarwal, Barham, Brevdo, Chen,  
 94 Citro, Corrado, Davis, Dean, Devin et al., 2016) and the keras wrapper (Chollet et al., 2015).

### 95 3. Dataset generation methodology

96 The dataset applied in this study for the ML training is a 3D unstructured grid CFD simulation of a point absorber  
 97 WEC. Here  $\mathbf{x} = [U, \nu_i]$  where  $U = [u, v, w]$  is the whole velocity field and  $u, v$  and  $w$  are the velocity components in  
 98 the X-axis, Y-axis, and Z-axis, respectively; and  $\nu_i$  is the dynamic viscosity.

99 These high-fidelity CFD simulating results are used for data generation and also as reference results for validation.  
 100 An in-house domain decomposition hybrid solver, qaLeFOAM, based on the open-source platform OpenFOAM (Jasak,  
 101 Jemcov, Tukovic et al., 2007) is adopted in this study where a two-phase incompressible Navier-Stokes (NS) solver  
 102 couples with the Quasi Lagrangian-Eulerian Finite Element Method (QALE-FEM) solver based on the fully nonlinear  
 103 potential theory (Ma and Yan, 2006, 2009; Yan and Ma, 2007, 2010; Yan, Ma, and Cheng, 2011). The main target  
 104 of this novel model is to boost the computational efficiency while maintaining the efficiency and details about this  
 105 hybrid model can be found in (Li, Wang, Yan, Gong, and Ma, 2018) and (Yan, Wang, Wang, Ma, and Xie, 2020).

106 Besides, several investigations using this method in various working scenarios have demonstrated the effectiveness  
 107 of the proposed model Yan, Li, Wang, Ma, Xie, and Stoesser (2019) and Yan et al. (2020). Even though significant  
 108 improvement has been seen by applying this hybrid method, the simulation of turbulent effect for the real engineering  
 109 issue is still very computationally expensive. This contributes to the main motivation of coupling the CFD code with  
 110 the technology of machine learning.

In the Reynolds-averaged Navier Stokes (RANS) model employed in this study, an ensemble averaging method  
 is applied to the unsteady turbulent flow modelling. This hypothesis introduces the macroscopic representations of  
 the micro-scale fluctuating flow. It offers access to model the overall effects of small vortexes by correlations and  
 meanwhile, resolves the larger eddies through the numerical simulation. where  $\nu$  is the constant molecular viscosity and  
 $\nu_T(d_i, t)$  is the spatial-temporal dependent turbulent/eddy viscosity, and together they compose the effective viscosity  
 $\nu_{eff}(d_i, t)$ :

$$\nu_{eff}(d_i, t) = \nu + \nu_T(d_i, t) \quad (7)$$

111 where  $d_i$  and  $t$  represent a point in space at time-step  $t$ .

As we are more interested in water physics, the kinematic viscosity is supposed to transfer to the dynamic viscosity  
 $\nu_t$  before being used in the ML training of generative methods by the following equation:

$$\nu_t = \nu [\rho_{air}(1 - \alpha) + \rho_{water}\alpha] \quad (8)$$

112 where water is playing a dominant role and  $\alpha$  is the phase fraction in the two-phase flow;  $\nu$  is the kinematic viscosity  
 113 and  $\rho_{water}$  and  $\rho_{air}$  are the water and air density, respectively. Regarding the boundary condition setting, the no-slip  
 114 boundary condition is applied on the bottom boundary with the total pressure specified on the top boundary in the  
 115 NS domain. For the subdomain configuration, in the fully nonlinear potential domain, there are wave generation and  
 116 absorption boundaries employed at the inlet and outlet boundary, respectively. Benefiting from the proposed hybrid  
 117 model, the turbulent viscosity is only considered in a refined relatively smaller zone which is the NS domain.

118 In this CFD simulation, a  $k - \omega$  Shear Stress Transport (SST) turbulent model that belonged to the RANS equation  
 119 catalogue is applied in this study. Here,  $k$  is the turbulence kinetic energy and  $\omega$  is the specific rate of dissipation of  
 120 the turbulence kinetic energy  $k$  into internal thermal energy. The WEC surface wall treatment is always one of the  
 121 biggest challenges raised in turbulent flow simulation, which can be classified into two categories: the Low-Reynolds  
 122 number (LR) models and High-Reynolds number (HR) models. The LR approach accompanied by a wall function is  
 123 targeting at the sublayer where exists a local low turbulent Reynolds number. One alternative to wall functions is to  
 124 adopt a fine-grid configuration that allows the application of a laminar flow boundary condition. To reach the viscous  
 125 sublayer, the normalised distance ( $y^+$ ) from the first mesh cell centre to the body surface is supposed to be around 1,  
 126 where  $y^+ = u_* y_w / \nu_{eff}$ . In numerical practice, the desired  $y^+$  is usually obtained through consistent trials. However, the  
 127 HR model can cope with a much larger  $y^+$  ( $\sim 30$ ) which integrates with a log-law to estimate the gradient approaching  
 128 the body wall. It should be noted that the first computational mesh should be placed either in the log-layer or the viscous  
 129 sublayer but not in-between (Utyuzhnikov, 2005), since none of the categories can deal with the buffer layer where  
 130 both viscous and Reynolds stresses are significant. All these factors further contribute to the complicity in the turbulent  
 131 flow modelling and also can be a source of error since it highly relies on the experiences of the user. Therefore, by  
 132 cooperating with ML, this study is aimed to deal with these difficulties by employing a surrogate prediction for RANS  
 133 turbulence eddy viscosity.

### 134 3.1. CFD Simulation Configurations

135 In the numerical simulation, a rectangular computational domain is adopted. A Two-dimensional (2D) side view  
 136 of the computational domains is given in Figure 1. A series of numerical simulations targeted at the hydrodynamics  
 137 performances of the three principal motions (surge, heave, and pitch) are conducted.

138 As the important role played by the CFD dataset, its accuracy is examined before transferring into the ML model. To  
 139 test the viability of the results generated by the adopted qaleFOAM solver, the numerical results have been compared to  
 140 the experimental measurement carried out by Todalshaug, Ásgeirsson, Hjalmarsson, Maillet, Möller, Pires, Guérinel,  
 141 and Lopes (2016) for the same wave energy devices under the regular wave conditions, in which good agreements have  
 142 been achieved in the structure motion response. Figure 2 demonstrates the wave propagation near the WEC and also the  
 143 pressure profile on the WEC surface. Besides, the mesh configuration of the computational domain is generated by the  
 144 SnappyHexMesh tool with the refined zone around the free surface and structure. The CFD work is highly dependent

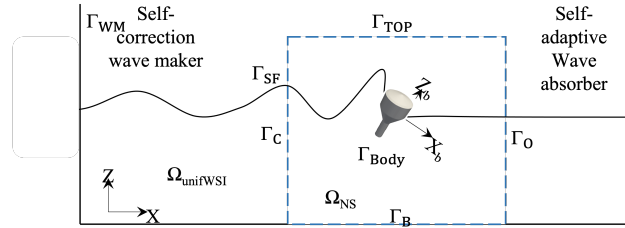


Figure 1: Sketch of the computational domain with boundaries

145 on the mesh resolution. Therefore, for each wave state, the convergence test against mesh resolution is performed to  
 146 identify the optimised mesh configuration with a minimal computational cost. In the turbulent model, the initial values  
 147 of  $k$  and  $\omega$  are set at the inlet boundary with the WEC surface treatment using the low-Reynolds-number approach  
 148 accompanied by a wall function.

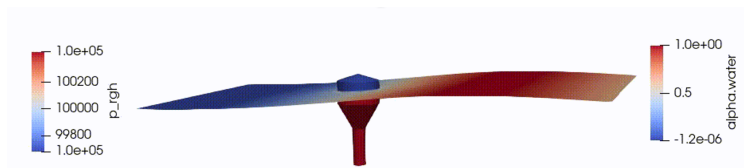
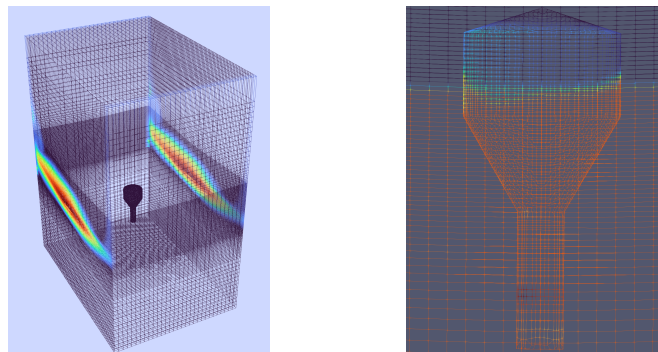


Figure 2: Wave surface profile around the buoy

149 The working condition considered in this simulation is an irregular wave generated by the JONSWAP spectrum  
 150 (Hasselmann, Barnett, Bouws, Carlson, Cartwright, Enke, Ewing, Gienapp, Hasselmann, Kruseman et al., 1973). The  
 151 significant wave height is  $H=8.8$  m and the significant wave period is  $T=1$  s with the numerical wave tank depth  $d=50$   
 152 m. Within certain mesh configuration, the time step size  $\Delta t$  is automatically determined by using the fixed Courant  
 153 number  $C_0$  ( $C_0 = (u_{max} \Delta d) / \Delta t$ , where  $\Delta d$  is the mesh size and  $u_{max}$  is the largest velocity value at the current time  
 154 step). This simulation has  $n = 60$  time-levels, with  $M = 851101$  nodes. Each node contains a scalar feature kinematic  
 155 viscosity  $\nu_i$  and a three-dimensional velocity vector  $U$ . The 4 fields account for each node containing  $m = 3404404$   
 156 features. A snapshot of the dataset is shown in Figure 3 showcasing the regular mesh and the irregular mesh around  
 the buoy.



(a) 3D-view of the WEC CFD unstructured mesh (b) Close-up to the unstructured mesh around the WEC within the global mesh

Figure 3: 3D and 2D view of a snapshot of the WEC CFD simulation

157



## 4. Results

Four different experiments were performed. For the PC-AAE we constructed the PC-space for inputs in three different ways:

- PC-AAE (Unut\_sep): full velocity field  $U$  and  $v_t$  separated
- PC-AAE (uvwnut): 3 components of the velocity field  $u, v, w$  and the dynamic viscosity  $v_t$  separated.
- PC-AAE (Unut): Full velocity field  $U$  and  $v_t$  together

The rationale behind this design is to show how the PC-space can be constructed and its effects on the training and reconstruction of the physical space. If the velocity field  $U$  is considered fully along with the  $v_t$  then the values of  $U$ , being larger than the ones of  $v_t$ , will give more weight towards  $U$ . However, decomposing the PC-space into  $U$  and  $v_t$ ; or  $u, v, w$  and  $v_t$ , the balance the reduced-space. For the  $\beta$ -VAE, no dimension reduction like PCA was applied and each variable was independently standardised:

- VAE: 3 components of the velocity field  $u, v, w$  and the dynamic viscosity  $v_t$  separated.

Therefore, the networks with PC-AAE are small and deal with tens or hundreds of inputs, whilst the  $\beta$ -VAE is large and has  $\sim 3.4M$  inputs.

For each one of these combinations, the data were standardised by its mean and standard deviation of the field. The training data set is 80% of the original simulation data and the test data is 20%, These datasets were shuffled randomly.

For each one of these experiments we tested their accuracy in the ground truth reconstruction (see Section 4.1), then we used them as a model surrogate and we tested the accuracy and efficiency of seen and unseen scenarios (see Section 4.2)

### 4.1. Reconstruction of ground truth using generative networks

To assess the effectiveness of the generative networks  $\beta$ -VAE and PC-AAE, the test data of the ground truth is reconstructed via the AEs. The test data is composed of 20% of the time-levels from the original CFD simulation, i.e. 12 time-levels.

Figure 4 shows the Kernel density estimator (KDE) for ground truth (blue), and the reconstructed test data, predicted by different generative methods, of the average values for each node over all samples. There is a good agreement between the reconstruction given by the generative models and the ground truth. Moreover, Figure 5 shows the excellent overlapping per time-level between the ground truth and the generative methods when the test data is reconstructed.

The mutual information (Mutual Information (MI)) between distributions was calculated using:

$$MI(U; V) = \sum_{i=1}^{|U|} \sum_{j=1}^{|V|} \frac{U_i \cap V_j}{N} \log \frac{U_i \cap V_j}{|U_i| |V_j|} \quad (9)$$

where  $U$  and  $V$  are the distributions of ground truth values of the test dataset and the predicted values by the generative methods, respectively. Here, a Normalised Mutual Information (NMI) (Strehl and Ghosh, 2002) is used where normalisation of the MI score scales the results between 0 (no mutual information) and 1 (perfect correlation). The NMI is then defined as:

$$NMI(U; V) = \frac{MI(U; V)}{H(U)H(V)} \quad (10)$$

where  $H(U)$  and  $H(V)$  are the entropies of  $U$  and  $V$ , respectively.

Another metric to assess the fidelity of the reconstruction of the test dataset is Normalised Root Mean Squared Error (NRMSE). The NRMSE for all points at all time levels is defined by:

$$NRMSE = \frac{\sqrt{\|y_{GT} - y_{pred}\|^2}}{\max(y_{GT}) - \min(y_{GT})} \quad (11)$$

where  $y_{GT}$  and  $y_{pred}$  are the unravelled values of the ground truth and predicted values by the generative methods, respectively.

A summary of the different metrics is presented in Table 1.

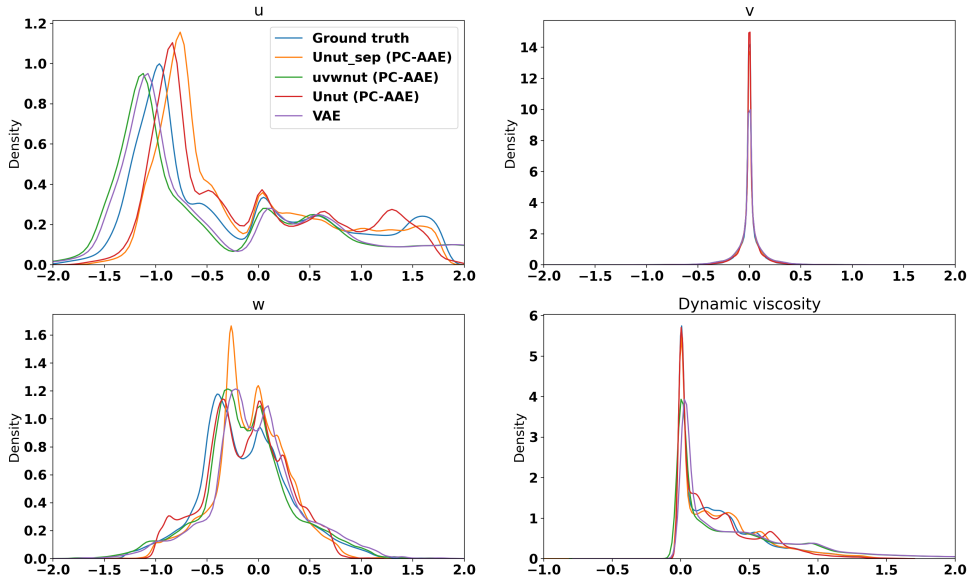


Figure 4: KDE for ground truth (blue), and the reconstructed test data, predicted by different generative methods, of the average values for each node over all samples.

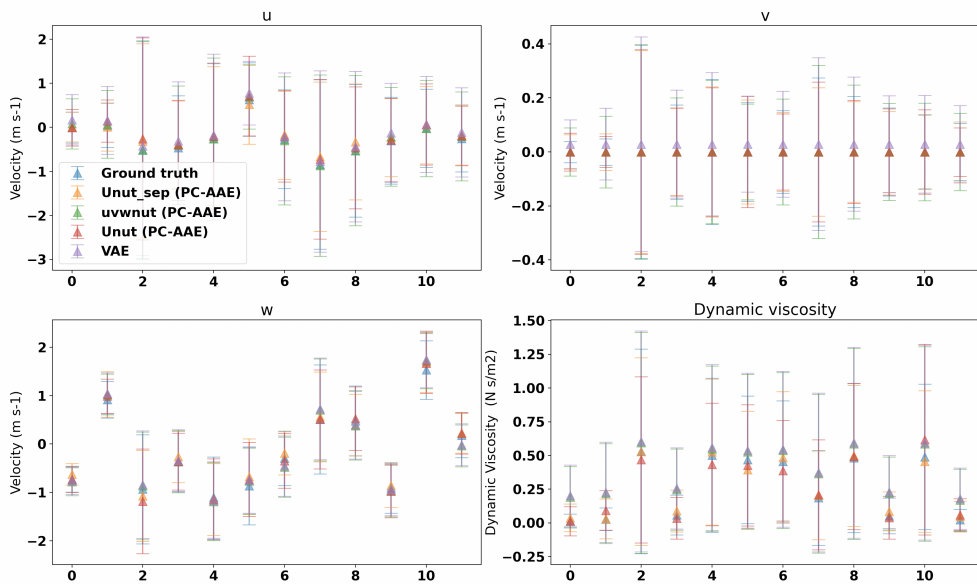


Figure 5: Error bars for reconstructed test data, averaged over all data points.

#### 4.2. Model surrogates of seen and unseen scenarios

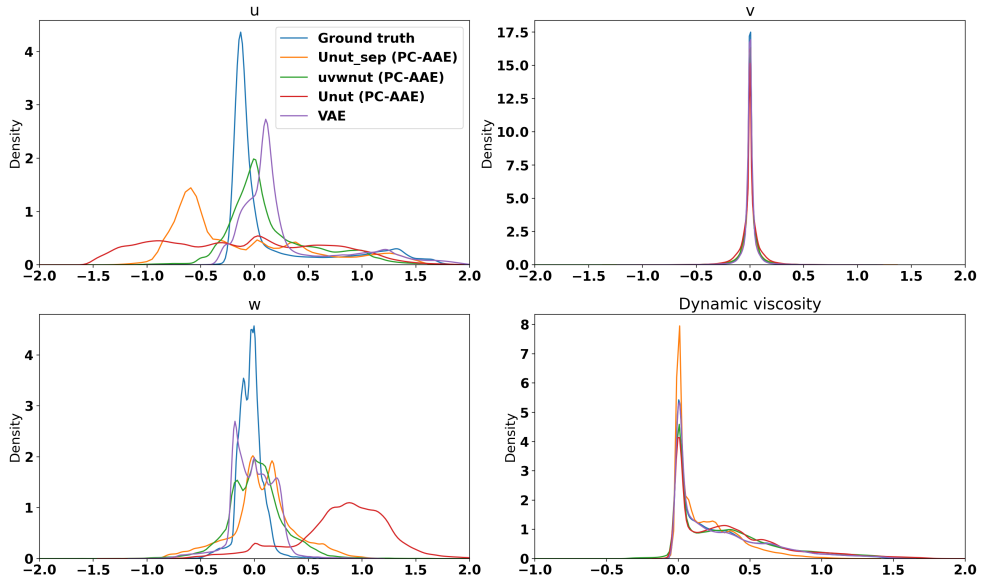
High-fidelity synthetic data were generated using the  $\beta$ -VAE (see Section 2.1) and PC-AAE (see Section 2.2). These techniques permit to augment experimental datasets, which can be costly and time-consuming to acquire.

Figure 6 shows the KDE of  $u$ ,  $v$ ,  $w$  and  $v_t$ , comparing ground truth and synthetic data generated by the aforementioned different methods. This is the mean for each node across all samples. For  $u$ , PC-AAE (uvwnut) and  $\beta$ -VAE show the best agreement with the ground truth. For  $v$ , all experiments exhibit an excellent match with the observed velocity. This can be explained by the fact that  $v$  is almost negligible and tends to 0, as there are minimal wave movement in that direction. For  $w$ , the movement in that axis is more restricted than in the X-axis which yields

**Table 1**

NMI between the ground truth values of the test dataset and the predicted values by the generative methods

Experiments	NRMSE				NMI			
	u	v	w	$v_t$	u	v	w	$v_t$
Unut_sep	0.029	0.017	0.023	0.015	0.942	0.968	0.942	0.945
uvwnut	0.040	0.022	0.028	0.038	0.942	0.968	0.942	0.945
Unut	0.030	0.018	0.021	0.022	0.942	0.968	0.942	0.945
$\beta$ -VAE	0.022	0.028	0.038	0.170	0.942	0.954	0.942	0.945



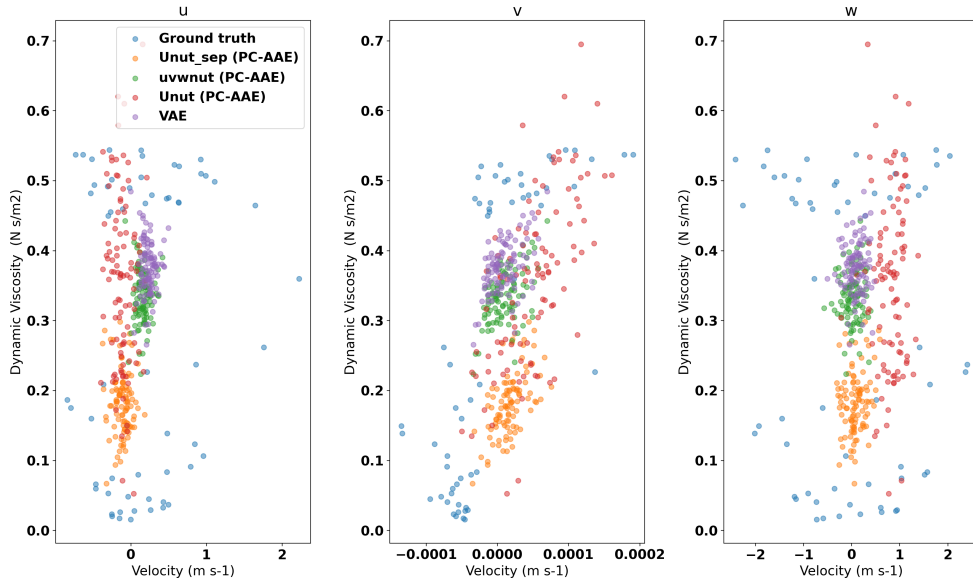
**Figure 6:** KDE for ground truth (blue), and 100 synthetic data, generated by different generative methods, of the average values for each node over all samples.

203 smaller velocities, and here is where the different experiments show greater variations with the ground truth in general.  
 204 Finally, for  $v_t$ , all experiments show an excellent agreement with the ground truth distribution.

205 In general, the experiments that show the highest fidelity to the ground truth data come from  $\beta$ -VAE and  
 206 PC-AAE (uvwnut). This can be explained because each field was standardised by its mean and standard deviation.  
 207 This is not the case for the other experiments where the occasional higher values shown in the velocity fields could  
 208 weight the PC space towards them.

209 In particular, the generated data with PC-AAE (Unut) shows the biggest distribution spread which expands the  
 210 physical scenarios and explores more extreme velocities, due to a covariate shift. This can help understand how the  
 211 point absorber WEC reacts to larger unseen incoming waves. This characteristic is fundamental for predicting unseen  
 212 scenarios (Yang, Zhou, Li, and Liu, 2021) and we will expand the analysis of this aspect in later.

213 Whilst generative networks aim to maximise the probability of generating real data whilst maintaining the distance  
 214 between the real and estimated posterior distributions small, this might not represent any physical meaning. This is  
 215 extremely important in test cases like WEC. In Figure 7, the averaged values of velocities  $u$ ,  $v$  and  $w$  are plotted  
 216 against the  $v_t$  to showcase its relationship. It is portrayed that this relationship is preserved in all experiments where  
 217 the averaged, over the number of data points, samples clusters overlap. Figure 7 also shows the spread over time-level  
 218 samples from the different experiments. The larger spread is given by PC-AAE (Unut) and the smaller spread of samples  
 219 is yielded by  $\beta$ -VAE. The latter can be explained due to the large size of the network required to train the  $\beta$ -VAE which  
 220 has an input and output of 3.4M points making it difficult to train due to memory allocation, rather than the PC-AAE  
 221 experiments which use tens or hundreds of points for input. However, the PC-AAE-related experiments need to store



**Figure 7:** Scatter values of velocities  $u$ ,  $v$ , and  $w$  against  $v_t$  of the average values per time samples over all nodes.

222 the inverse mapping to the physical space which is of a similar storage size to the original simulation used for training  
 223 data.

224 To understand how the distributions of the newly generated samples are related to the distribution of the ground  
 225 truth data, we obtained a t-Student Stochastic Neighbour Embedding (t-SNE) projection (Van der Maaten and Hinton,  
 226 2008). The t-SNE projection is depicted in Figure 8. For velocities, it is clear that the ground truth ( $U_{GT}$ ) clusters  
 227 together with the synthetic samples for all experiments. Similar behaviour can be observed for  $v_t^{GT}$  (nut\_GT) and its  
 228 synthetic samples. Moreover, the relationship between the velocity fields and the dynamic viscosity observed in the  
 229 ground truth samples is observed and preserved in the synthetic data. This is another example of how the physical  
 230 relationship between  $U$  and  $v_t$  behaves and how it is preserved in the different generative methods.

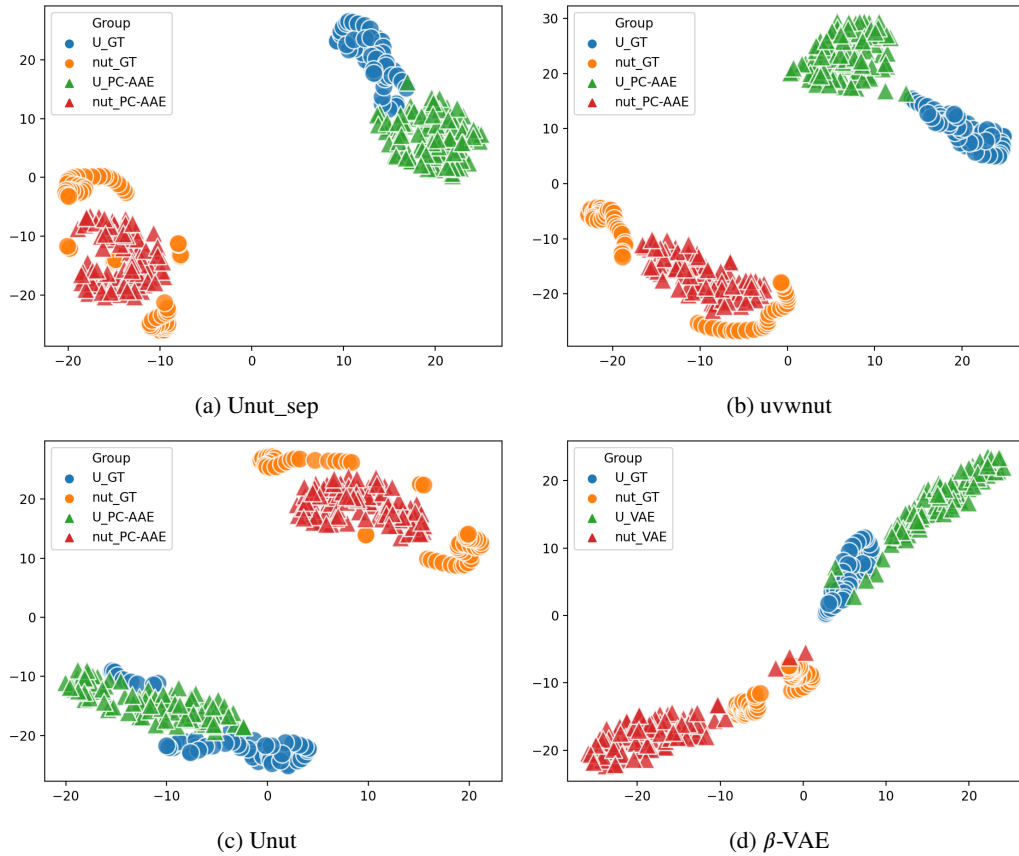
231 The synthetic data are generated by sampling random Gaussian noise, in these cases with a sample size of 16,  
 232 before feeding this into the decoder  $\mathcal{P}$ .

233 As aforementioned, PC-AAE (Unut) is more able to expand the design space, whilst maintaining the physical  
 234 relationship between  $U$  and  $v_t$ . Therefore, for analysing unseen scenarios, only PC-AAE (Unut) will be discussed.

235 To further demonstrate the synthetic data's physical meaning, Figure 9 shows the relationship of real and synthetic  
 236  $u$ ,  $v$ , and  $w$  (X-axis) against  $v_t$  (Y-axis). The scatter plots show  $m = 851101$  points and how the synthetic data  
 237 overlaps with the ground truth data. The PC-AAE overlaps and engulfs the spread shown by the ground truth data.  
 238 This demonstrates how the design space is expanded by PC-AAE. As aforementioned, the generated data with  
 239 PC-AAE (Unut) shows the biggest distribution spread which expands the physical scenarios and explores more extreme  
 240 velocities. The access to this part of the design space comes from how the PC was built. In PC-AAE (Unut) all fields  
 241 were considered together giving  $u$  and  $w$  a larger weight than  $v$  and  $v_t$ . Thus, PC-AAE can expand the design space  
 242 with larger incoming waves whilst preserving the physical relationship among variables.

243 Figures 10a and 10b show 10 newly generated samples of the  $U$  and  $v_t$  in the XZ planes, respectively, using  
 244 PC-AAE. For these Figures, a new set of 100 samples were generated, however, only 10 samples are shown for display  
 245 purposes.

246 As shown, for high-fidelity scenarios it is better to separate all variables and preprocess them individually, i.e  
 247 standardisation, to preserve the statistics without assigning a larger weight to specific fields. However, for expanding  
 248 the design space, a larger weight can be given to the fields with larger values when constructing the PC space on all  
 249 fields, specifically, to the fields mostly affected by the variations of irregular incoming waves.



**Figure 8:** t-SNE projection of the ground truth data (GT) and 100 synthetic data for the Unut\_sep, uvwnut, Unut, and  $\beta$ -VAE experiments each experiment using the PC-AAE and  $\beta$ -VAE.

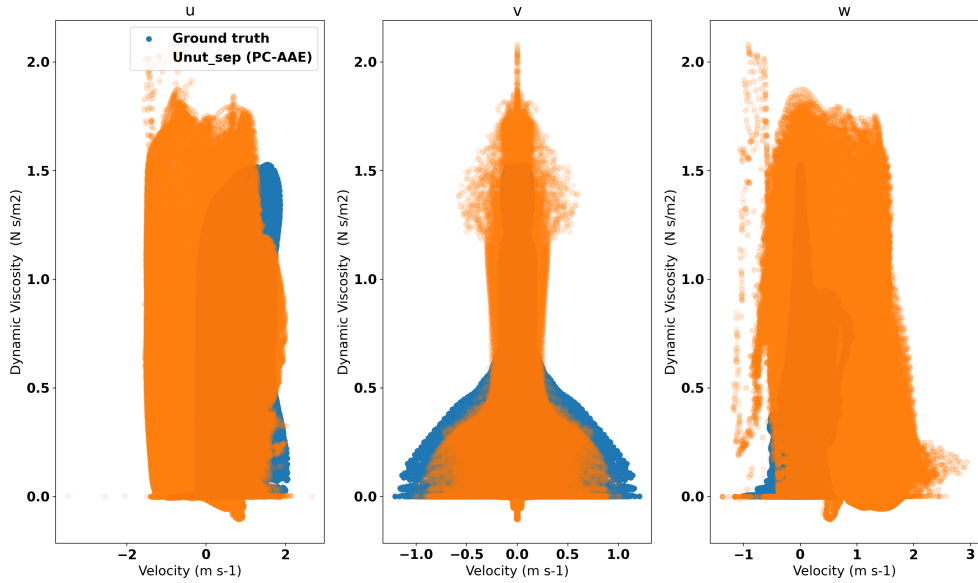
250 This has tremendous implications for expanding the design space and exploring unseen scenarios of the WEC CFD  
 251 simulation. The generative model has the potential to capture the transition and the new physical phenomenons under  
 252 new extreme initial conditions, as shown in the synthetic data generated by PC-AAE (Unut).

### 253 4.3. Model architectures

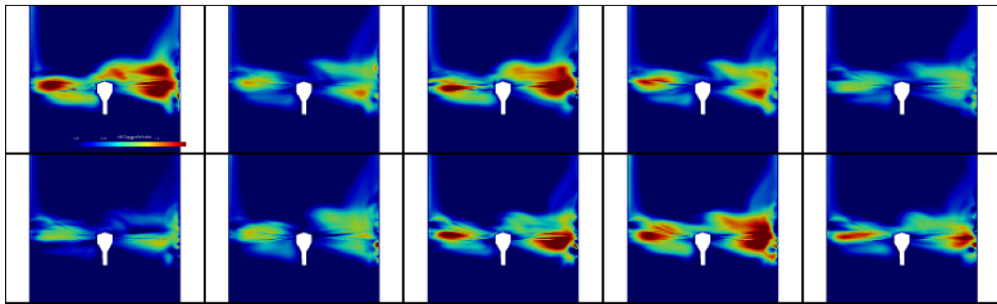
254 The architectures of both generative models are shown in Table 2 where the networks have been trained with batch  
 255 normalisation before and a dropout of 0.5 between layers. For the PC-AAE-related architectures, the discriminator  
 256 updates 10 times before the generator during training, and Nadam parameters with Learning rate  $lr = 10^{-3}$ ,  $\beta_1 =$   
 257  $0.9$ ,  $\beta_2 = 0.999$

### 258 4.4. Performance

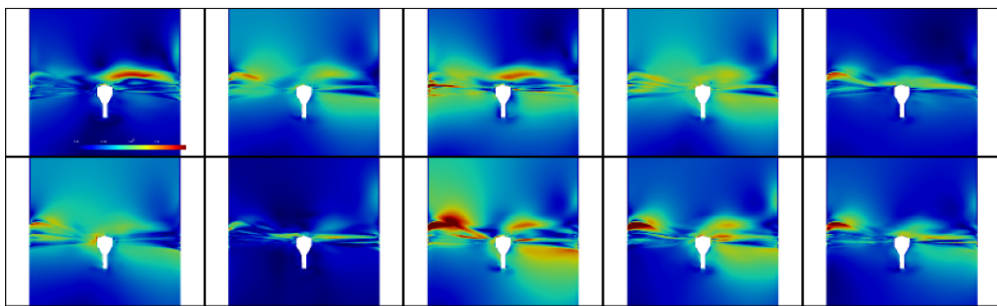
259 The most notable capability of these approaches is the speed of execution and how fast new synthetic samples can  
 260 be obtained. The runtimes, averaged over 10 times, to generate 100 synthetic samples with the different methods are  
 261 10.46 s, 5.23 s, 10.81 s, and 0.64s for PC-AAE: Unut\_sep, uvwnut, Unut, and  $\beta$ -VAE, respectively. For the PC-AAE  
 262 experiments, it only takes  $\sim 0.05$  s to generate 100 samples in the PC-space before projecting these onto the physical  
 263 space. These times were obtained using a 2.3 GHz 8-Core Intel Core i9 processor. The speed of these generative  
 264 models compares to  $\sim 1$  week of simulation using OpenFOAM for the original 3D CFD simulation, which only outputs  
 265 60 time-levels. This is a speed-up from 5 to 6 orders of magnitude.



**Figure 9:** Scatter values of velocities  $u$ ,  $v$ , and  $w$  against  $v_t$  of the average values per point for PC-AAE (Unut)



(a) 10 newly generated samples of  $v_t$  using the PC-AAE. Shown in the  $XZ$  plane.



(b) 10 newly generated samples of  $U$  using the PC-AAE. Shown in the  $XZ$  plane.

**Figure 10:** Example of newly generated samples of  $U$  and  $v_t$  using the PC-AAE. Shown in the  $XZ$  plane.

## 266 5. Conclusion and Future work

267 We have shown how two different ML-based generative models,  $\beta$ -VAE and PC-AAE, can generate synthetic data  
 268 of WECs quickly and at reduced computational cost. These models can generate several samples in a fast manner,  
 269 compared to 1 week time of simulation using OpenFOAM. The t-SNE projection of the ground truth and generated

**Table 2**

Architectures of the different generative networks  $\beta$ -VAE and PC-AAE. The encoder  $Q$ , the decoder  $\mathcal{P}$  and the discriminator  $D^A$  are fully-connected layers.

Experiments AE	Enc $Q$	Dec $\mathcal{P}$	Disc $D^A$
PC-AAE (Unut_sep)	96	16	16
	32	32	
	16	96	
PC-AAE (uvwnut)	192	16	16
	32	32	
	16	192	
PC-AAE (Unut)	48	16	16
	32	32	
	16	48	
	Enc $Q$	Dec $\mathcal{P}$	
$\beta$ -VAE	3404404	16	
	64	32	
	32	64	
	16	3404404	

270 samples are located closely spatially, preserving the relationship between the velocity and dynamic viscosity fields.  
 271 Furthermore, this is also shown by the relationship of  $u$ ,  $v$ , and  $w$  with  $v$ .

272 The reconstruction of the ground truth using these generative models agrees well with the test dataset. As shown,  
 273 for high-fidelity scenarios it is better to separate all variables and preprocess them individually to preserve the statistics  
 274 without assigning a larger weight to specific fields. However, for expanding the design space, a larger weight can be  
 275 given to the fields with larger values when constructing the PC space on all fields.

276 This has tremendous implications for expanding the design space and exploring unseen scenarios of the CFD  
 277 simulation. The generative model has the potential to capture the transition and the new physical phenomenons under  
 278 extreme initial conditions, as shown in the synthetic data generated by PC-AAE.

279 Moreover, these generative models can generate snapshots of  $\sim 3.4$  M features at a fraction of the cost of the original  
 280 simulation. Future work will include using these generated samples to improve the prediction of  $v_t$  in future time steps,  
 281 by augmenting the design space with physically plausible samples. Additionally, other newer generative models like  
 282 latent diffusion models, due to the number of nodes in these simulations, could be used.

## 283 6. Data and Code availability

284 The code used in this paper is available in <http://github.com/c-quilo/WaveSuite> and the data is available  
 285 on request.

## 286 7. Acknowledgements

287 The authors would like to acknowledge support from the UK Engineering and Physical Sciences Research  
 288 Council (EPSRC) Programme Grant PREMIERE (EP/T000414/1), the EPSRC grant EP/T003189/1 Health assessment  
 289 across biological length scales for personal pollution exposure and its mitigation (INHALE), and the EPSRC  
 290 grant EP/V040235/1 New Generation Modelling Suite for the Survivability of Wave Energy Convertors in Marine  
 291 Environments (Wave-Suite).

## 292 CRedit authorship contribution statement

293 **César Quilodrán-Casas:** Conceptualisation of this study, Methodology, Software, Formal analysis, Investigation,  
 294 Writing - Original draft preparation, Writing - Review & Editing. **Qian Li:** Data Curation, Writing - Original Draft,  
 295 Writing - Review & Editing. **Ningbo Zhang:** Writing - Review & Editing. **Sibo Cheng:** Writing - Review & Editing.  
 296 **Shiqiang Yan:** Supervision, Funding acquisition, Writing - Review & Editing. **Qingwei Ma:** Supervision, Funding  
 297 acquisition, Writing - Review & Editing. **Rossella Arcucci:** Supervision, Funding acquisition, Writing - Review &  
 298 Editing.

## 299 References

- 300 I. Glendenning, Ocean wave power, *Applied Energy* 3 (1977) 197–222.
- 301 F. d. O. Antonio, Wave energy utilization: A review of the technologies, *Renewable and sustainable energy reviews* 14 (2010) 899–918.
- 302 S. Jin, R. J. Patton, B. Guo, Viscosity effect on a point absorber wave energy converter hydrodynamics validated by simulation and experiment, *Renewable energy* 129 (2018) 500–512.
- 303 D. Son, R. W. Yeung, Optimizing ocean-wave energy extraction of a dual coaxial-cylinder wec using nonlinear model predictive control, *Applied energy* 187 (2017) 746–757.
- 304 Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to document recognition, *Proceedings of the IEEE* 86 (1998) 2278–2324.
- 305 C. Quilodrán-Casas, R. Arcucci, L. Mottet, Y. Guo, C. Pain, Adversarial autoencoders and adversarial lstm for improved forecasts of urban air pollution simulations, *arXiv preprint arXiv:2104.06297* (2021).
- 306 L. Chagot, C. Quilodrán-Casas, M. Kalli, N. M. Kovalchuk, M. J. Simmons, O. K. Matar, R. Arcucci, P. Angeli, Surfactant-laden droplet size prediction in a flow-focusing microchannel: a data-driven approach, *Lab on a Chip* 22 (2022) 3848–3859.
- 307 T. Karras, T. Aila, S. Laine, J. Lehtinen, Progressive growing of gans for improved quality, stability, and variation, *arXiv preprint arXiv:1710.10196* (2017).
- 308 D. Berthelot, T. Schumm, L. Metz, Began: Boundary equilibrium generative adversarial networks, *arXiv preprint arXiv:1703.10717* (2017).
- 309 A. Radford, L. Metz, S. Chintala, Unsupervised representation learning with deep convolutional generative adversarial networks, *arXiv preprint arXiv:1511.06434* (2015).
- 310 G. Forestier, F. Petitjean, H. A. Dau, G. I. Webb, E. Keogh, Generating synthetic time series to augment sparse datasets, in: 2017 IEEE international conference on data mining (ICDM), IEEE, 2017, pp. 865–870.
- 311 J. Hoffmann, Y. Bar-Sinai, L. M. Lee, J. Andrejevic, S. Mishra, S. M. Rubinstein, C. H. Rycroft, Machine learning in a data-limited regime: Augmenting experiments with synthetic data uncovers order in crumpled sheets, *Science advances* 5 (2019).
- 312 R. J. Chen, M. Y. Lu, T. Y. Chen, D. F. Williamson, F. Mahmood, Synthetic data in machine learning for medicine and healthcare, *Nature Biomedical Engineering* (2021) 1–5.
- 313 J. Yoon, J. Jordon, M. Schaar, Radialgan: Leveraging multiple datasets to improve target-specific predictive models using generative adversarial networks, in: International Conference on Machine Learning, PMLR, 2018, pp. 5699–5707.
- 314 L. Zhou, S. Pan, J. Wang, A. V. Vasilakos, Machine learning on big data: Opportunities and challenges, *Neurocomputing* 237 (2017) 350–361.
- 315 H. Li, C. Zhang, M. Chen, D. Shen, Y. Niu, Data-driven surrogate modeling: Introducing spatial lag to consider spatial autocorrelation of flooding within urban drainage systems, *Environmental Modelling & Software* (2023) 105623.
- 316 M. Sadeghi, P. Nguyen, K. Hsu, S. Sorooshian, Improving near real-time precipitation estimation using a u-net convolutional neural network and geographical information, *Environmental Modelling & Software* 134 (2020) 104856.
- 317 S. Razavi, Deep learning, explained: Fundamentals, explainability, and bridgeability to process-based modelling, *Environmental Modelling & Software* 144 (2021) 105159.
- 318 I. Goodfellow, Y. Bengio, A. Courville, *Deep learning*, MIT press, 2016.
- 319 C. Rodriguez-Delgado, R. J. Bergillos, G. Iglesias, An artificial neural network model of coastal erosion mitigation through wave farms, *Environmental Modelling & Software* 119 (2019) 390–399.
- 320 L. Li, Z. Yuan, Y. Gao, Maximization of energy absorption for a wave energy converter using the deep machine learning, *Energy* 165 (2018) 340–349.
- 321 D. Sarkar, E. Contal, N. Vayatis, F. Dias, Prediction and optimization of wave energy converter arrays using a machine learning approach, *Renewable Energy* 97 (2016) 504–517.
- 322 P. D. Sclavounos, Y. Ma, Wave energy conversion using machine learning forecasts and model predictive control, in: 33rd International Workshop on Water Waves and Floating Bodies, Brest, France, 2018.
- 323 S. M. Mousavi, M. Ghasemi, M. Dehghan Manshadi, A. Mosavi, Deep learning for wave energy converter modeling using long short-term memory, *Mathematics* 9 (2021) 871.
- 324 I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, A. Lerchner, beta-vae: Learning basic visual concepts with a constrained variational framework (2016).
- 325 A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, B. Frey, Adversarial autoencoders, *arXiv preprint arXiv:1511.05644* (2015).
- 326 D. P. Kingma, M. Welling, Auto-encoding variational bayes, *arXiv preprint arXiv:1312.6114* (2013).
- 327 A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learning library, *Advances in neural information processing systems* 32 (2019).
- 328 W. Falcon, et al., Pytorch lightning, GitHub. Note: <https://github.com/PyTorchLightning/pytorch-lightning> 3 (2019).
- 329 J. Lever, M. Krzywinski, N. Altman, Points of significance: Principal component analysis, 2017.
- 330 S. Cheng, Y. Jin, S. P. Harrison, C. Quilodrán-Casas, I. C. Prentice, Y.-K. Guo, R. Arcucci, Parameter flexible wildfire prediction using machine learning techniques: Forward and inverse modelling, *Remote Sensing* 14 (2022) 3228.
- 331 S. Cheng, J. Chen, C. Anastasiou, P. Angeli, O. K. Matar, Y.-K. Guo, C. C. Pain, R. Arcucci, Generalised latent assimilation in heterogeneous reduced spaces with machine learning surrogate models, *Journal of Scientific Computing* 94 (2023) 11.
- 332 H. Gong, S. Cheng, Z. Chen, Q. Li, C. Quilodrán-Casas, D. Xiao, R. Arcucci, An efficient digital twin based on machine learning svd autoencoder and generalised latent assimilation for nuclear reactor physics, *Annals of Nuclear Energy* 179 (2022) 109431.
- 333 T. Dozat, Incorporating Nesterov momentum into adam (2016).
- 334 M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al., Tensorflow: Large-scale machine learning on heterogeneous distributed systems, *arXiv preprint arXiv:1603.04467* (2016).
- 335 F. Chollet, et al., Keras, 2015. URL: <https://github.com/fchollet/keras>.
- 336 H. Jasak, A. Jemcov, Z. Tukovic, et al., Openfoam: A c++ library for complex physics simulations, in: International workshop on coupled methods in numerical dynamics, volume 1000, IUC Dubrovnik Croatia, 2007, pp. 1–20.



- 362 Q. Ma, S. Yan, Quasi ale finite element method for nonlinear water waves, *Journal of computational physics* 212 (2006) 52–72.
- 363 Q. Ma, S. Yan, Qale-fem for numerical modelling of non-linear interaction between 3d moored floating bodies and steep waves, *International*
- 364 *Journal for Numerical Methods in Engineering* 78 (2009) 713–756.
- 365 S. Yan, Q. Ma, Numerical simulation of fully nonlinear interaction between steep waves and 2d floating bodies using the qale-fem method, *Journal*
- 366 *of Computational physics* 221 (2007) 666–692.
- 367 S. Yan, Q. Ma, Qale-fem for modelling 3d overturning waves, *International Journal for Numerical Methods in Fluids* 63 (2010) 743–768.
- 368 S. Yan, Q. Ma, X. Cheng, Fully nonlinear hydrodynamic interaction between two 3d floating structures in close proximity, *International Journal of*
- 369 *Offshore and Polar Engineering* 21 (2011).
- 370 Q. Li, J. Wang, S. Yan, J. Gong, Q. Ma, A zonal hybrid approach coupling fnpt with openfoam for modelling wave-structure interactions with action
- 371 of current, *Ocean Systems Engineering* 8 (2018) 381–407.
- 372 S. Yan, J. Wang, J. Wang, Q. Ma, Z. Xie, Ccp-wsi blind test using qalefoam with an improved passive wave absorber, *International Journal of*
- 373 *Offshore and Polar Engineering* 30 (2020) 43–52.
- 374 S. Yan, Q. Li, J. Wang, Q. Ma, Z. Xie, T. Stoesser, Comparative numerical study on focusing wave interaction with fpso-like structure, *International*
- 375 *Journal of Offshore and Polar Engineering* 29 (2019) 149–157.
- 376 S. Utyuzhnikov, Generalized wall functions and their application for simulation of turbulent flows, *International journal for numerical methods in*
- 377 *fluids* 47 (2005) 1323–1328.
- 378 J. H. Todalshaug, G. S. Ásgeirsson, E. Hjalmarsson, J. Maillet, P. Möller, P. Pires, M. Guérinel, M. Lopes, Tank testing of an inherently phase-
- 379 controlled wave energy converter, *International Journal of Marine Energy* 15 (2016) 68–84.
- 380 K. Hasselmann, T. P. Barnett, E. Bouws, H. Carlson, D. E. Cartwright, K. Enke, J. Ewing, A. Gienapp, D. Hasselmann, P. Kruseman, et al.,
- 381 Measurements of wind-wave growth and swell decay during the joint north sea wave project (jonswap)., *Ergaenzungsheft zur Deutschen*
- 382 *Hydrographischen Zeitschrift, Reihe A* (1973).
- 383 A. Strehl, J. Ghosh, Cluster ensembles—a knowledge reuse framework for combining multiple partitions, *Journal of machine learning research* 3
- 384 (2002) 583–617.
- 385 J. Yang, K. Zhou, Y. Li, Z. Liu, Generalized out-of-distribution detection: A survey, *arXiv preprint arXiv:2110.11334* (2021).
- 386 L. Van der Maaten, G. Hinton, Visualizing data using t-sne., *Journal of machine learning research* 9 (2008).

## 387 **Acronyms**

388 **2D** Two-dimensional

389

390 **3D** Three-dimensional

391

392 **AE** Autoencoder

393

394 **CFD** Computational Fluid Dynamics

395

396 **VAE** Variational Autoencoder

397

398 **PC** Principal Components

399

400 **PCA** Principal Component Analysis

401

402 **WEC** Wave Energy Converter

403

404 **AAE** Adversarial Autoencoder

405

406 **t-SNE** t-Student Stochastic Neighbour Embedding

407

408 **NS** Navier-Stokes

409

410 **RANS** Reynolds-averaged Navier Stokes  
411

412 **KL** Kullback-Leibler Divergence Score  
413

414 **KDE** Kernel density estimator  
415

416 **MSE** Mean Squared Error  
417

418 **PC-AAE** Principal Components-based Adversarial Autoencoder  
419

420 **MI** Mutual Information  
421

422 **NMI** Normalised Mutual Information  
423

424 **LR** Low-Reynolds number  
425

426 **HR** High-Reynolds number  
427

428 **EOF** Empirical Orthogonal Functions  
429

430 **QALE-FEM** Quasi Lagrangian-Eulerian Finite Element Method  
431

432 **NRMSE** Normalised Root Mean Squared Error  
433

434 **SST** Shear Stress Transport  
435

436 **DA** Data Augmentation  
437

438 **DL** Deep Learning  
439

440 **ML** Machine Learning  
441