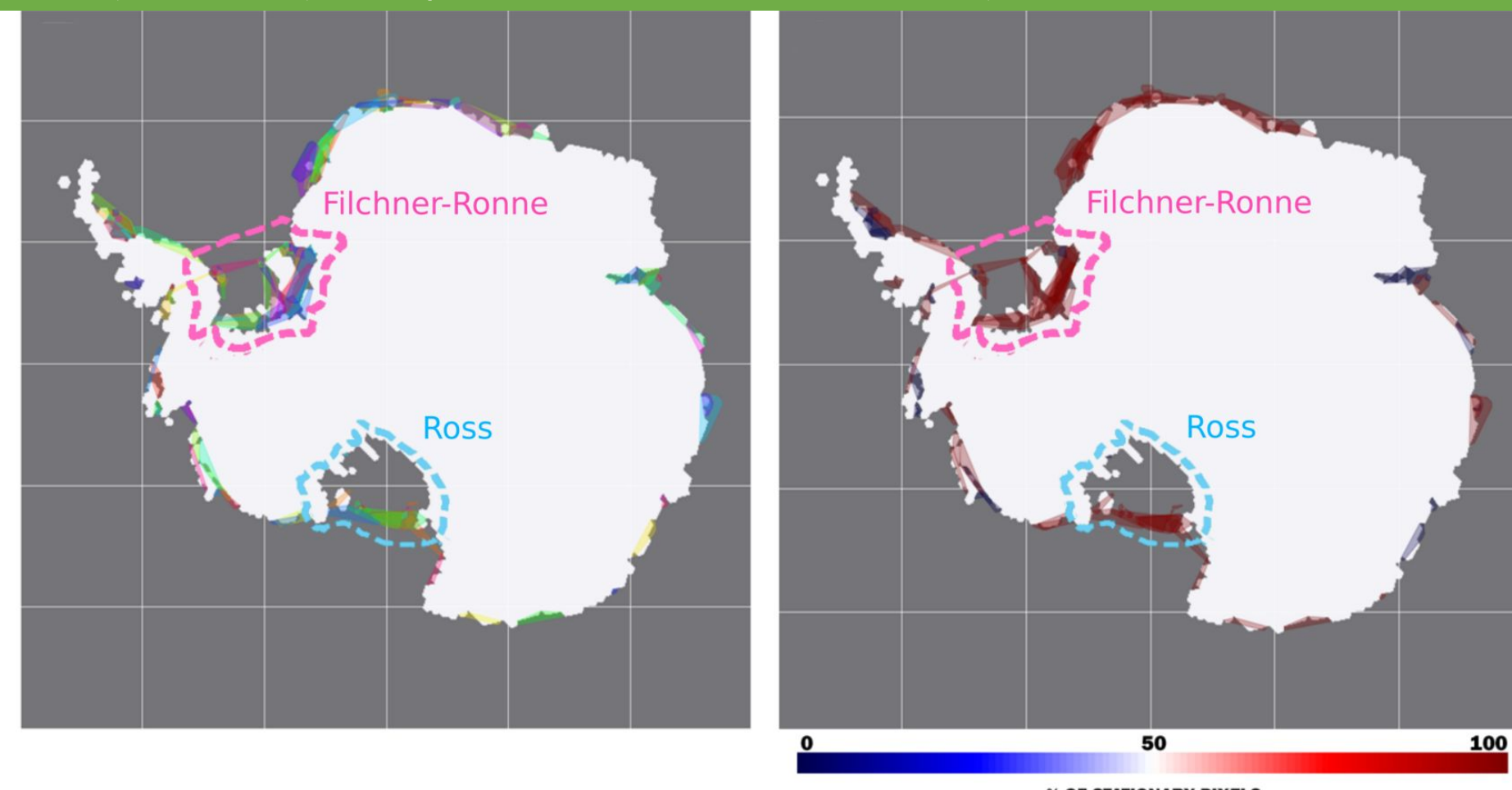


Generate Antarctic sub-shelf melt using recurrent neural network-based Generative Adversarial Models on pixel clusters

Jacquelyn A. Shelton¹, Alexander Robel², Matthew Hoffman³, and Stephen Price³
¹Hong Kong Polytechnic University ²Georgia Institute of Technology ³Los Alamos National Laboratory

- Antarctic Ice Sheet **ice loss** accelerated by surrounding ocean's extreme warming over last 30 years The DOE E3SM v1.2 Cryosphere Configuration: Description and Simulated Antarctic Ice-Shelf Basal Melting → dominant contributor to **global sea level rise**
- Unknown:** How much ice loss due to **anthropogenic changes** and to **internal variability** [1]?
- Goal:** develop/apply machine learning (ML) method(s) for **data generation** using **limited model output** (single realization) from **expensive** state-of-the-art Earth System model (E3SM) [2]
- Previous work:** identified **stationary subspaces** (via customized hierarchical agglomerative clustering) of the input data that are both **realistic (physically consistent)** and **representative** of its **complex spatiotemporal dynamics** [3]
- Current work:** results show **TimeGAN can generate additional realizations** of variable Antarctic sub-ice shelf melt that **preserves the temporal dynamics** and **stationarity** → **three metrics** (PCA [5], t-SNE [6], KPSS [7]) employed for evaluation of original model data and synthetic generated data

Approach: Identify Stationary Subspaces for Data Generation via Dynamic Agglomerative Clustering [3]



Prerequisites prior to data generation: identify individual subspaces in data that are:

- representative** of its spatiotemporal dynamics,
- realistic** in terms of consistency with physically observed dynamics, and
- stationary** over the entire time-series, regardless of the behavior of the data within each subspace relative to any other's (may vary independently over arbitrary time-scales)

Idea: construct dynamic hierarchical clustering pipeline to adaptively learn stationary subspaces, the number of which can grow or shrink in a data-driven fashion according to the data dynamics while simultaneously incorporating relevant prior domain knowledge (e.g. physical observations, problem setting).

Notation: Let $p_1 = (x_1, y_1)$ and $p_2 = (x_2, y_2)$ be 2 pixel locations given by their 2D coordinates on the Antarctic Ice Sheet. Note that each pixel spans 10 km. Ice melt flux time-series at these locations: $F_1 = (f_1^1, \dots, f_1^M)$ and $F_2 = (f_2^1, \dots, f_2^M)$, where M denotes the number of time-steps in the simulation, which for this data is a monthly resolution over 150 years, for $M = 1800$ time-steps. Normalized versions denoted: $\hat{F}_k = \frac{F_k}{\max F_k}$ for $k \in \{1, 2\}$, and the spatial-distance threshold denoted: d_{thr}^s .

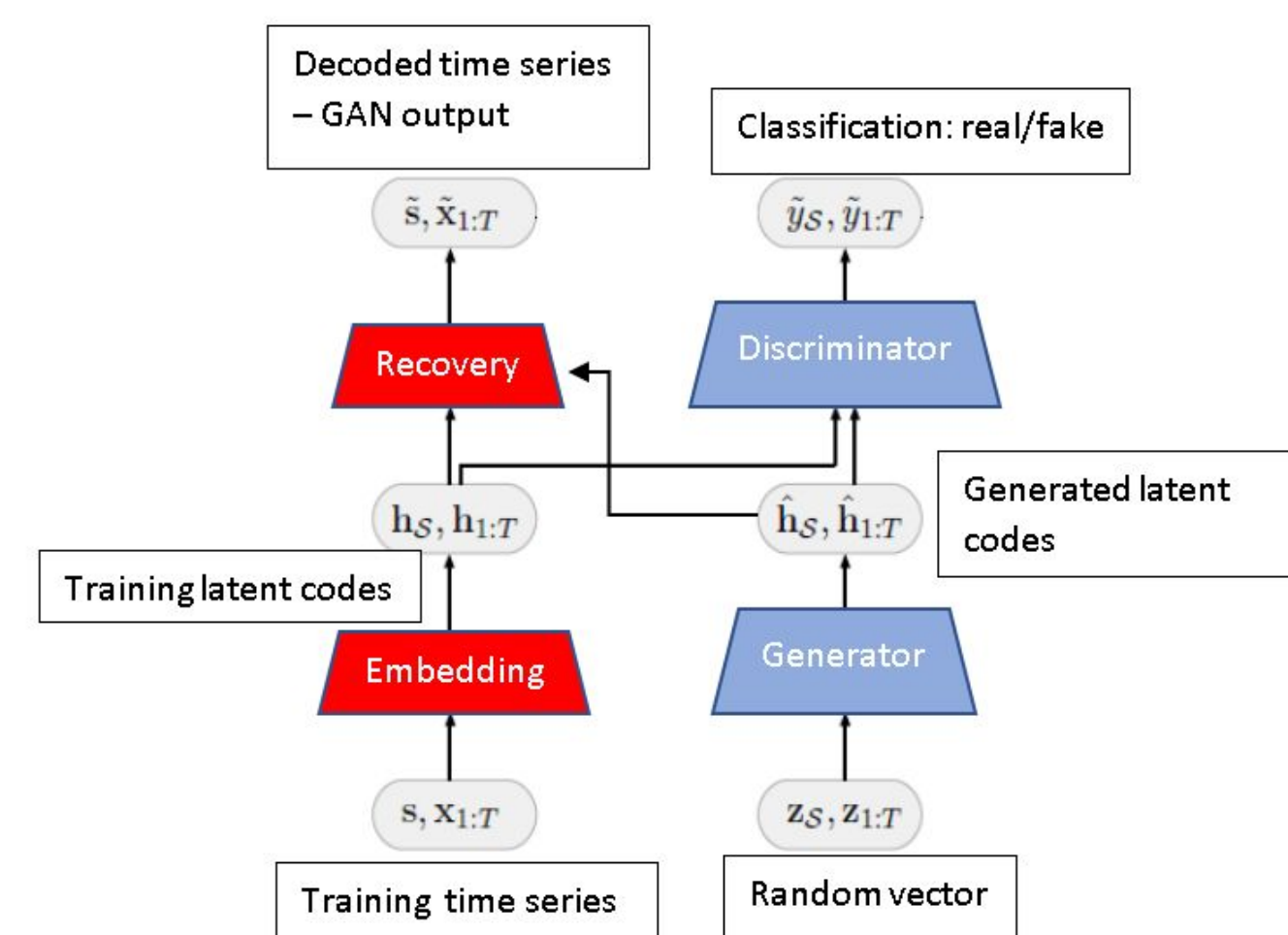
Aggregate Spatiotemporal Distance Criterion

$$d_{s,t}(p_1, p_2) = \begin{cases} \sum_{k=1}^M |\hat{F}_1(k) - \hat{F}_2(k)| & \text{if } \|p_1 - p_2\| \leq d_{thr}^s \\ +\infty & \text{else} \end{cases}$$

Stationarity: Evaluate the **stationarity of each identified cluster** using the Kwiatkowski-Phillips-Schmidt-Shin (KPSS) hypothesis test with the null hypothesis that the time-series is stationary around the mean. A cluster is considered to be stationary if the majority of its pixels pass the KPSS test.

Method: Time-series Generative Adversarial Network (TimeGAN) [4]

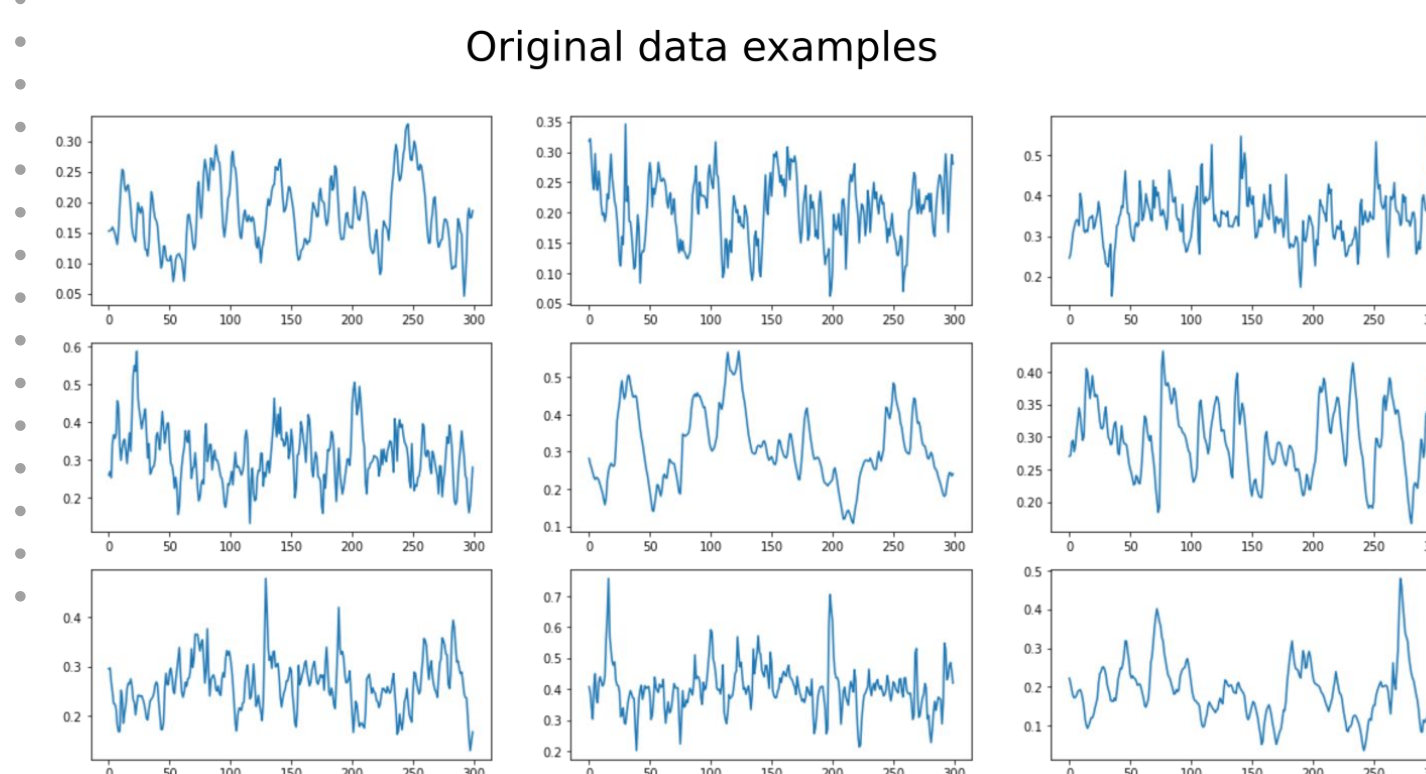
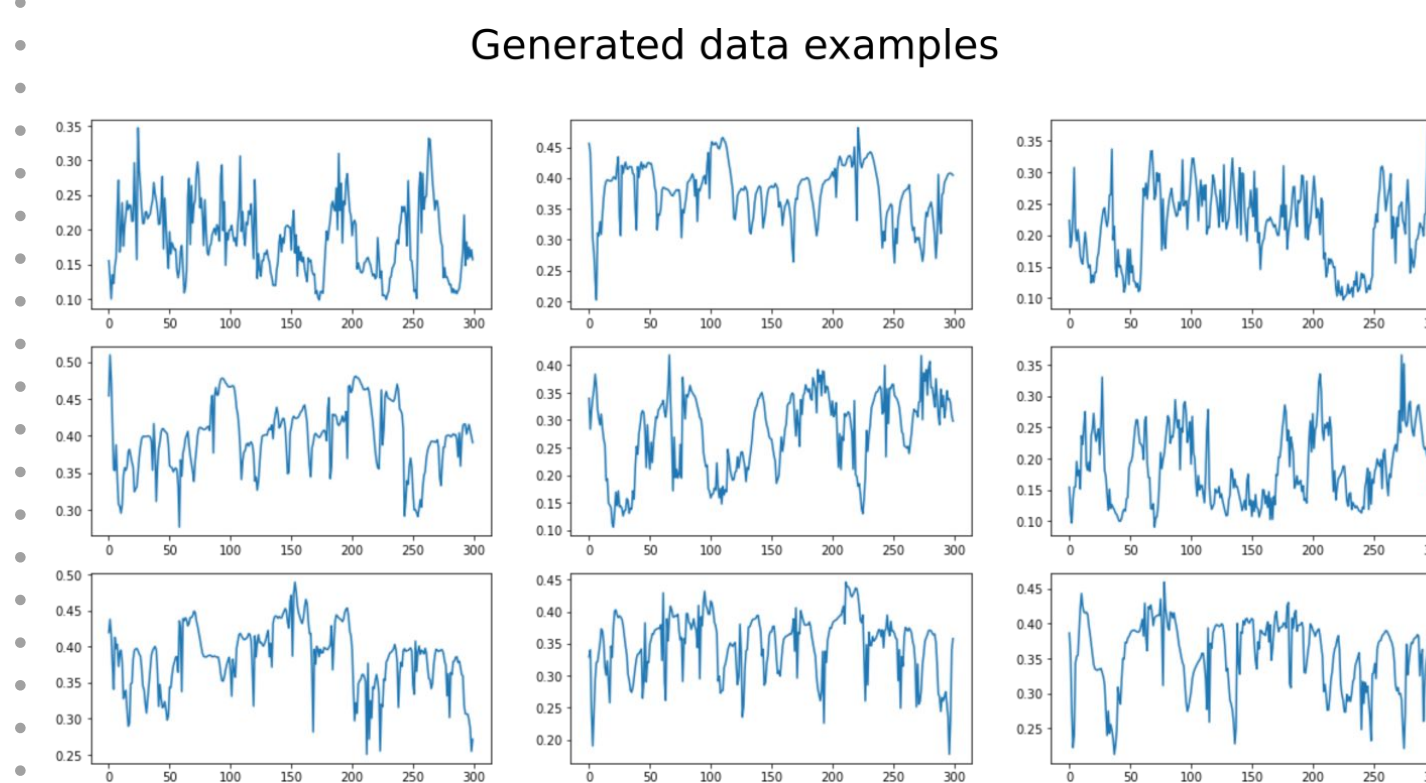
TimeGAN Network Architecture



- TimeGAN learns:** good **generative model** for time-series data that **preserves temporal dynamics** → new sequences respect original relationships between variables across time [4]
- Consists of **four unique Recurrent Convolutional Neural Networks (RCNNs;** e.g. LSTM, GRU): **embedder, generator, discriminator** and **recovery RCNNs**

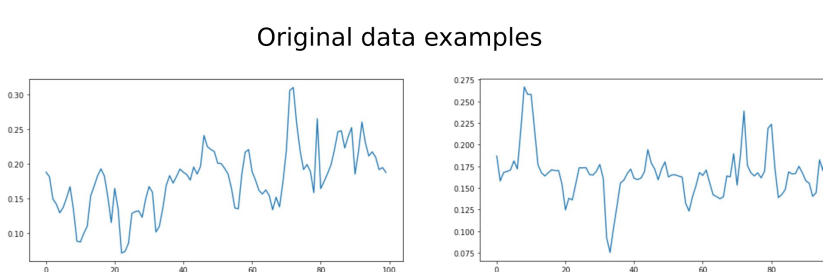
Data Generation: Examples of Original E3SM and TimeGAN Generated Data

Cluster 1 – gen seq. length 300 time-steps:



- Generate data** from identified **stationary clusters** → **train TimeGAN** on individual cluster's respective time-series/pixels
- Each cluster** ⇒ treated as its **own independent data distribution** for GAN to form generative model of

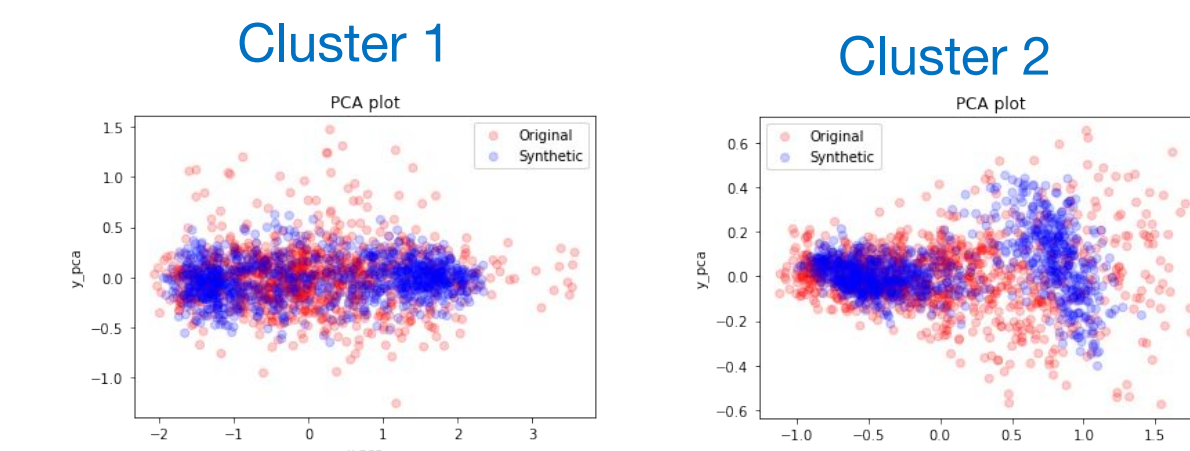
Cluster 2 – gen. seq. length 100:



Evaluation: Qualitative and Quantitative Metrics of TimeGAN Generated Data

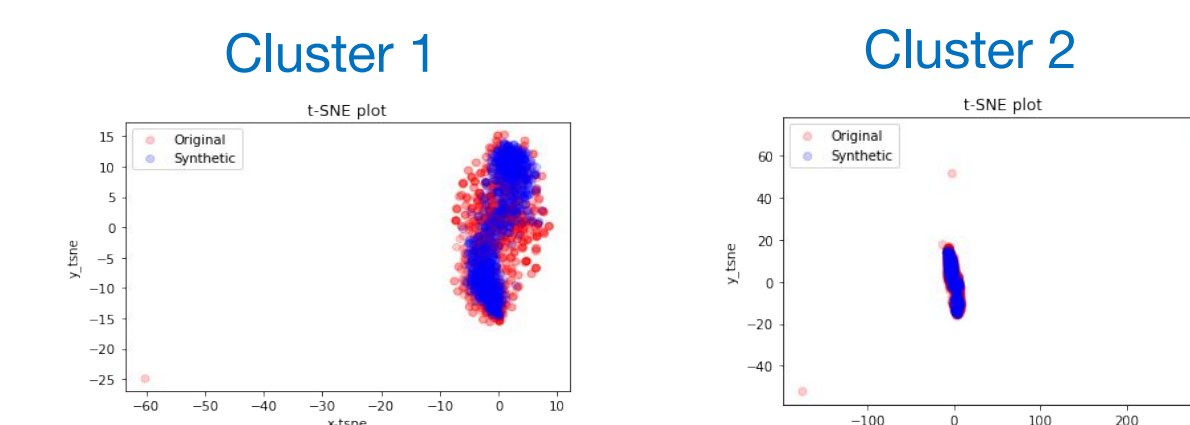
Assess **consistency/similarity** of the distributions of **generated data** vs. **input data:**

1. Principal Components Analysis (PCA) [5]



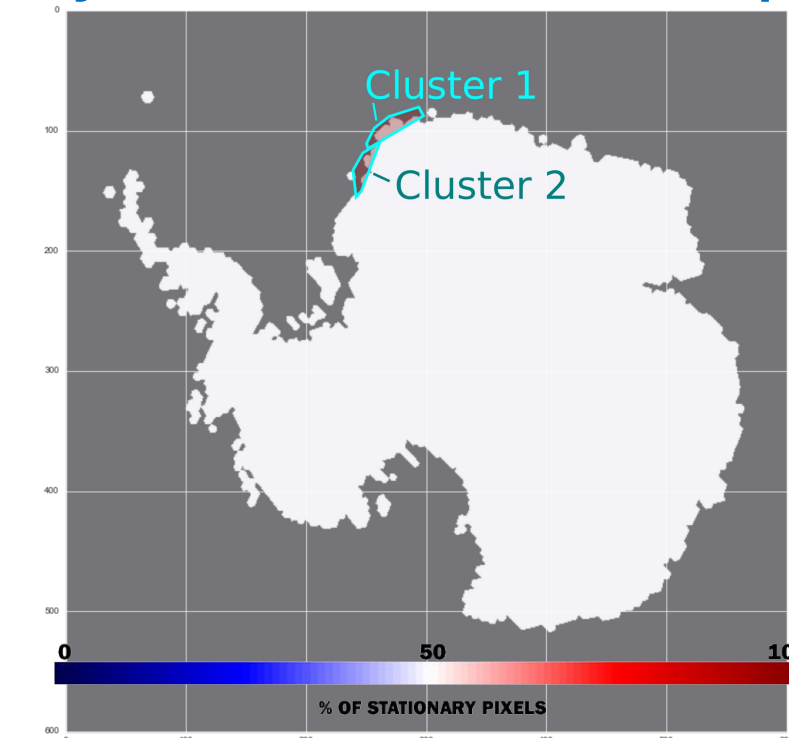
→ Flatten temporal dimension s.t. we can visualize data in 2D space: project all data onto the first 2 PCs of original data
→ **Shows:** significant overlap of and spread of input/generated data pts of **both clusters 1 and 2**

2. t-Distributed Stochastic Neighbor Embedding (t-SNE) [6]



→ Method to quantify/visualize similarity of data – **capable of retaining local structure** of (high dimensional) data and revealing **important global structure**
→ **Shows:** significant overlap of clusters

3. Stationarity: KPSS Kwiatkowski-Phillips-Schmidt-Shin (KPSS) hypothesis test [7]



→ **Shows:** Cluster 1 and 2 **both stationary** according to KPSS with the null hypothesis that the data is stationary around the mean

Conclusions, Impact, and Outlook

- Preliminary results** show TimeGAN can generate additional realizations of variable Antarctic sub-ice shelf melt that **preserves the temporal dynamics** and **stationarity**
- Evaluation summary:** all 3 metrics show promise that the TimeGAN can produce data similar to the input data, both in terms of **preserving temporal dynamics** and **stationarity** → **PCA** and **t-SNE** show the input/generated data have similar temporal dynamics in a lower dimensional space, **KPSS** shows that the generated data retains the input data's stationarity
- This **work addresses the general/pervasive problem of data scarcity** in the climate sciences → far more **computationally affordable** than running climate model
- Current ongoing work** investigating **other time-series GANs** [e.g. 8] and **advanced discriminator functions** for built-in non-parametric high-dimensional distribution comparisons [e.g. 9]

- References**
- [1] Robel, A., Seroussi, H., Roe, G. Marine ice sheet instability amplifies and skews uncertainty in projections of future sea-level rise. In PNAS, 116(30), 2019.
 - [2] Hoffman, M., Price, S. The DOE E3SM v1.2 Cryosphere Configuration: Description and Simulated Antarctic Ice-Shelf Basal Melting. In J. of Advances in Modeling Earth Systems, 2022.
 - [3] Shelton, J.A., Robel, A.A., Hoffman, M., and Price, S.: Towards generating stationary realizations of simulated Antarctic ice shelf melt rates from limited model output. Climate Informatics 2022. (additionally, expanded preprint available soon)
 - [4] Yoon, J., Jarrett, D., van der Schaar, M. Time-series Generative Adversarial Networks. In Proceedings of Advances in Neural Information Processing Systems (NeurIPS), 2019.
 - [5] H. Hotelling. Analysis of a complex of statistical variables into principal components. Journal of Educational Psychology, 24:417–441, 1933.
 - [6] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of Machine Learning Research, 9(Nov):2579–2605, 2008.
 - [7] Kwiatkowski, D., Phillips, P. C., Schmidt, P., Shin, Y. Testing the null hypothesis of stationarity against the alternative of a unit root: How sure are we that economic time series have a unit root? In Journal of Econometrics, 54(1–3), 159–178, 1992.
 - [8] Cristóbal Esteban, Stephanie L. Hyland, C.: RatschReal-valued (Medical) Time Series Generation with Recurrent Conditional GANs. In arXiv:1706.02633v2, 2017.
 - [9] Binkowski, M., Sutherland, D., Arbel, M., Gretton, A. Demystifying MMD GANs. In Proceedings of the ICLR 2018 Conference Blind Submission, 2018.