

# Successful Contrastive Pretraining of Set Transformer

**Divyashree Koti**

Middle Tennessee State University

*dsk2v@mtmail.mtsu.edu*

*Advisor:*

**Dr. Joshua Lee Phillips**

August 20, 2023



Computational Science

# Presentation Overview

## 1 Introduction

Data

## 2 Background

### 3 Methods

Set Transformer

Contrastive Learning using Set Transformer

Fine Tuning

Generalization

Visualization

## 4 Results

Set Transformer

Contrastive Learning using Set Transformer

Fine Tuning

Generalization

Visualization

## 5 Conclusion

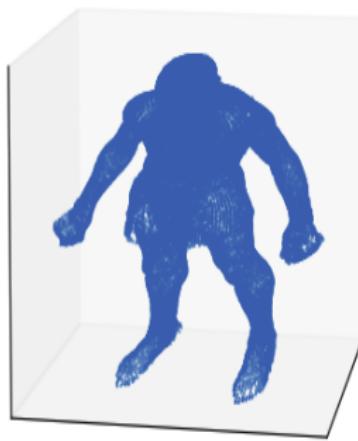
# Introduction

- Transformer - Process sequence of data
- Vision Transformer - Process images
- Set Transformer - Process unordered sets
- Contrastive Learning
- Generalization
  - Weak generalization
  - Strong generalization
- Point cloud 3D Body scan

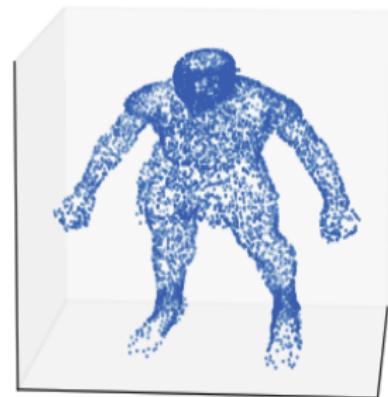


# Data

- 75 Participants
- Collected by Dr. Frederick Steven Cottle
- KX-16 Body Scanner [1]
- Point cloud count - 51000 to 63000



**Figure:** Original 3D Body Scan

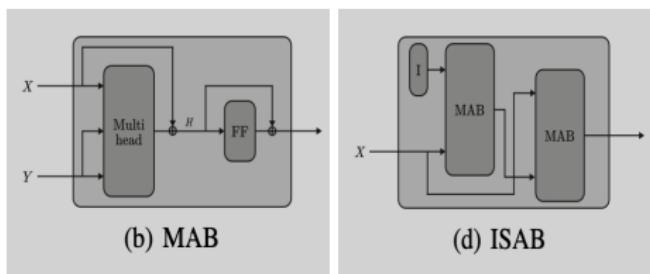


**Figure:** Subsampled 3D Body Scan



# Background

- Set Transformer - Lee, J. et. al. (2018) [2]
  - Induced Set Attention Block(ISAB) - speeds up computation and extracts meaningful features.
  - Pooling layer - parameterized aggregation function, helps to capture the varying contribution of the instances, for better aggregation.



**Figure:** Multi Head Attention Block(MAB)[2]

**Figure:** Induced Set Attention Block(ISAB)[2]



Computational Science

# Background

- Simple framework for Contrastive Learning [3] - Chen, T. et. al. (2020): ResNet-50 on ImageNet dataset, Positive pair/ Negative pair. Using SimCLR self-supervised method shows similar performance to Resnet-50 supervised learning.
- Generalization
  - Kullback-Leibler divergence(KL divergence) [4]

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log \left( \frac{P(x)}{Q(x)} \right) \quad (1)$$



Computational Science

# Methods

- Data Loader
- Batch size - train:4, validation:4
- Sub sample size - train:8000, validation:2048
- Train : Validation = 80:20
- Monte Carlo Simulation - 10
- Optimizer - Adam
- Loss - Sparse categorical cross entropy



# Set Transformer

- Hyperparameters
  - Embedded dimension - 32, 64
  - Number of heads - 4, 16, 32
  - Induce points - 32, 64, 128
  - stack - 2 to 6
  - Dropout - 0.05, 0.2
  - Learning rate - 1e-02 to 1e-04
  - Number of epochs - 250



Computational Science

# Set Transformer

---

## Set Transformer architecture

---

- 1:  $y = \text{Linear}(3, \text{embed\_dim})$
- 2: **for**  $j = 1$  to  $\text{stack}$  **do**
- 3:    $y = \text{InducedSetAttentionBlock}(y)$
- 4: **end for**
- 5:  $y = \text{Dropout}(0.05)(y)$
- 6:  $y = \text{PoolingByMultiHeadAttention}(y)$
- 7:  $y_{\text{embedding}} = \text{Dropout}(0.05)(y)$
- 8:  $y = \text{FinalDense}(\text{numberofclasses})(y_{\text{embedding}})$
- 9: *return*  $y, y_{\text{embedding}}$

---



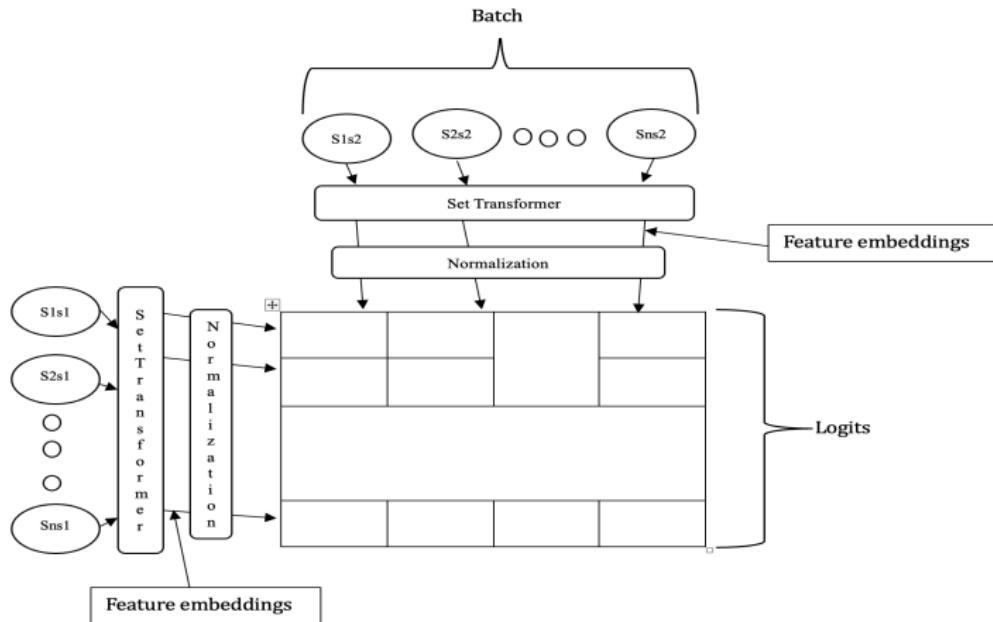
Computational Science

# Contrastive Learning using Set Transformer

- Base model - Set Transformer, to extract the feature embeddings
- For every batch a sample with its augmented self is a positive pair.
- For every batch a sample with the rest of the sample is negative pair.
- Experimented batch size - 4, 16, 32
- Point cloud - 8000, 2048, 1024
- Temperature - 0.5, 1.0, 2.0
- Optimizer - Adam with learning rate 1e-03



# Contrastive Learning using Set Transformer



**Figure:** Contrastive Learning using Set Transformer architecture

# Contrastive Learning using Set Transformer

---

## Contrastive Learning architecture

---

- 1:  $y1, y\_embedding1 = \text{SetTransformer}(\text{batch1})$
  - 2:  $y2, y\_embedding2 = \text{SetTransformer}(\text{batch2})$
  - 3:  $y\_embedding1 = \text{Linear}(\text{embed\_dim}, \text{projection\_dim})(y\_embedding1)$
  - 4:  $y\_embedding2 = \text{Linear}(\text{embed\_dim}, \text{projection\_dim})(y\_embedding2)$
  - 5:  $y\_embedding1 = \text{Norm}(y\_embedding1)$
  - 6:  $y\_embedding2 = \text{Norm}(y\_embedding2)$
  - 7:  $y = \text{Mul}(y\_embedding1, y\_embedding1.T) * \text{Temperature}$
- 



Computational Finance

# Fine Tuning

- Additional Layers:
  - Linear layer - 1024, 256, 128
  - Non-linear activation function - LeakyReLU
  - Dropout - 0.1
- Hyperparameters - same as that of standard Set transformer



# Generalization

- Weak Generalization - Has knowledge about all the category data at the time of training.
- Strong Generalization(leave one out) - Is not aware of a particular class data at a given training period.
- Probability matrix -  $75 \times 75$
- Reassign all diagonal elements in weak generalization 0.
- Normalize both of the matrices
- Apply KL Divergence



# Visualization

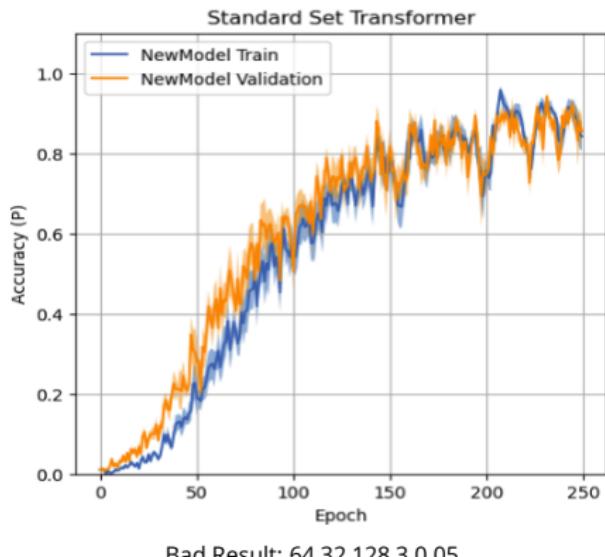
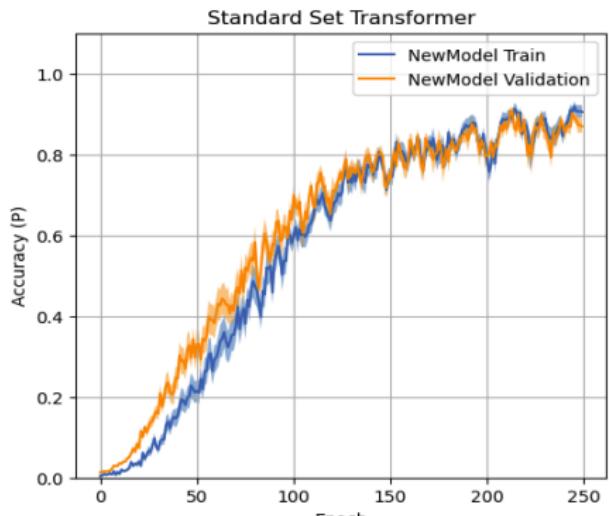
- Get feature embeddings for both standard Set transformer and Contrastive pre-trained model.
- 10 vectors for each sample.
- Applied PCA with 2 principal components to visualize in 2D.
- Applied t-SNE with a 2D component with the perplexity of 20.
- Evaluate along the 2D axis if, distribution is relatable to human distinguish.



# Results

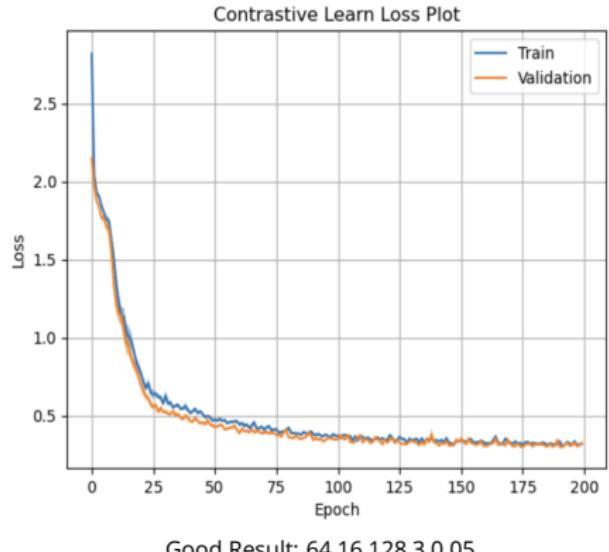


# Standard Set Transformer

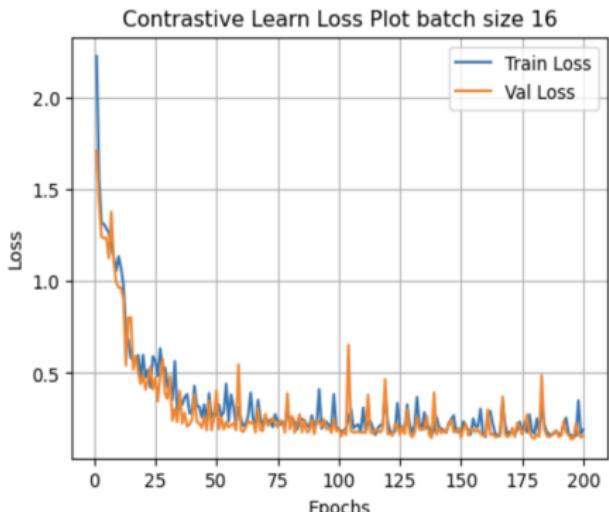


**Figure:** Monte Carlo simulated Standard Set Transformer Accuracy

# Contrastive Learning using Set Transformer



**Figure:** Monte Carlo simulated Contrastive Learning using Set Transformer Loss

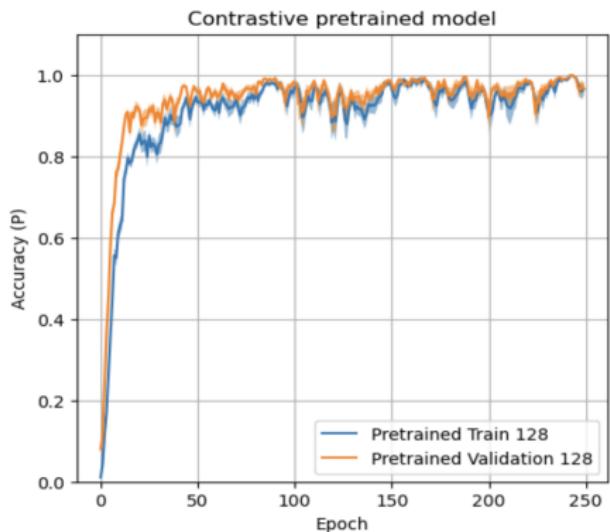


Bad Result: 64,16,128,3,0.05,batch size: 16

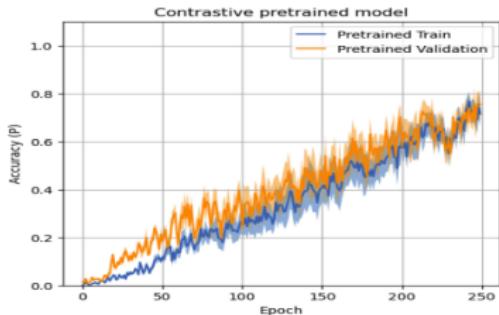


Computational Science

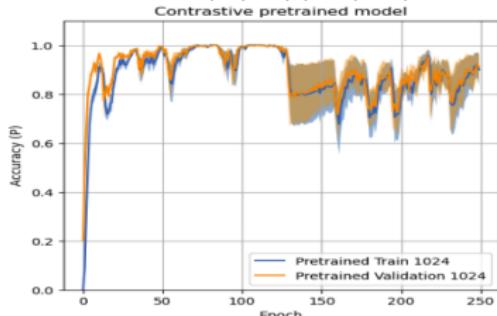
# Fine Tuning



Good Result: 64,16,128,3,0.05, 128, 0.1



Bad Result: 64,32,128,3,0.05, 128, 0.1



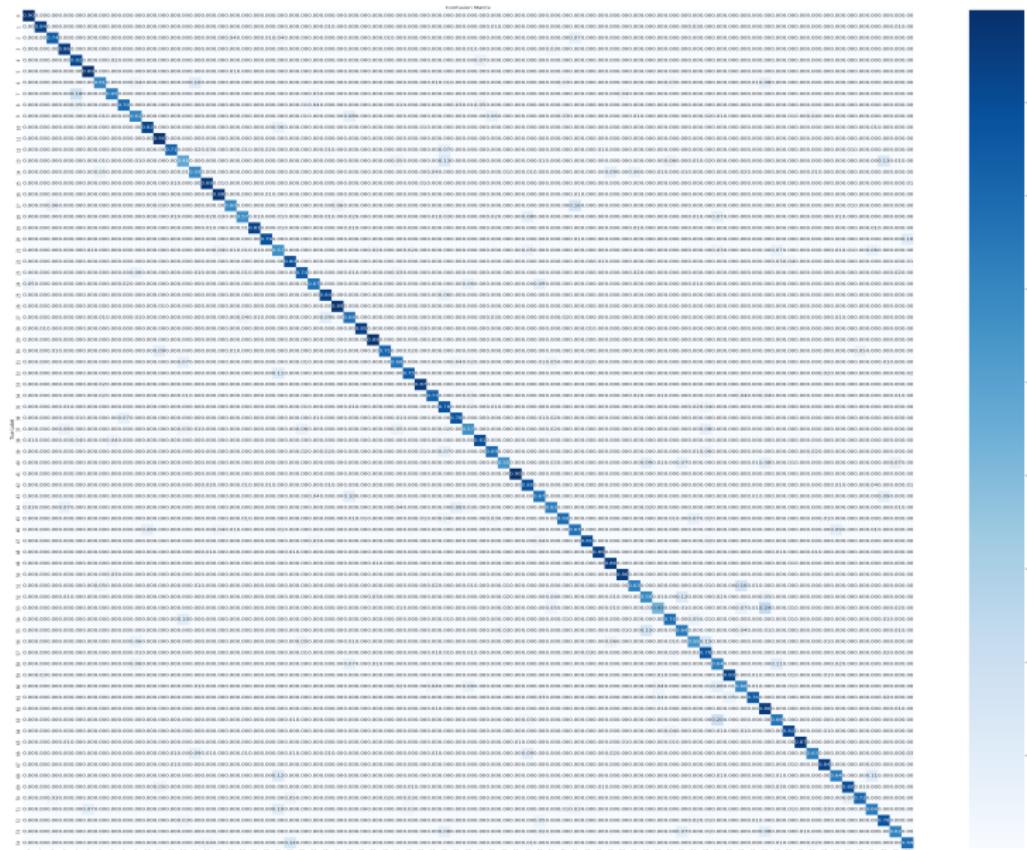
Bad Result: 64,16,128,3,0.05, 1024, 0.1

**Figure:** Monte Carlo simulated Contrastive pre-trained model Accuracy



Computational Science

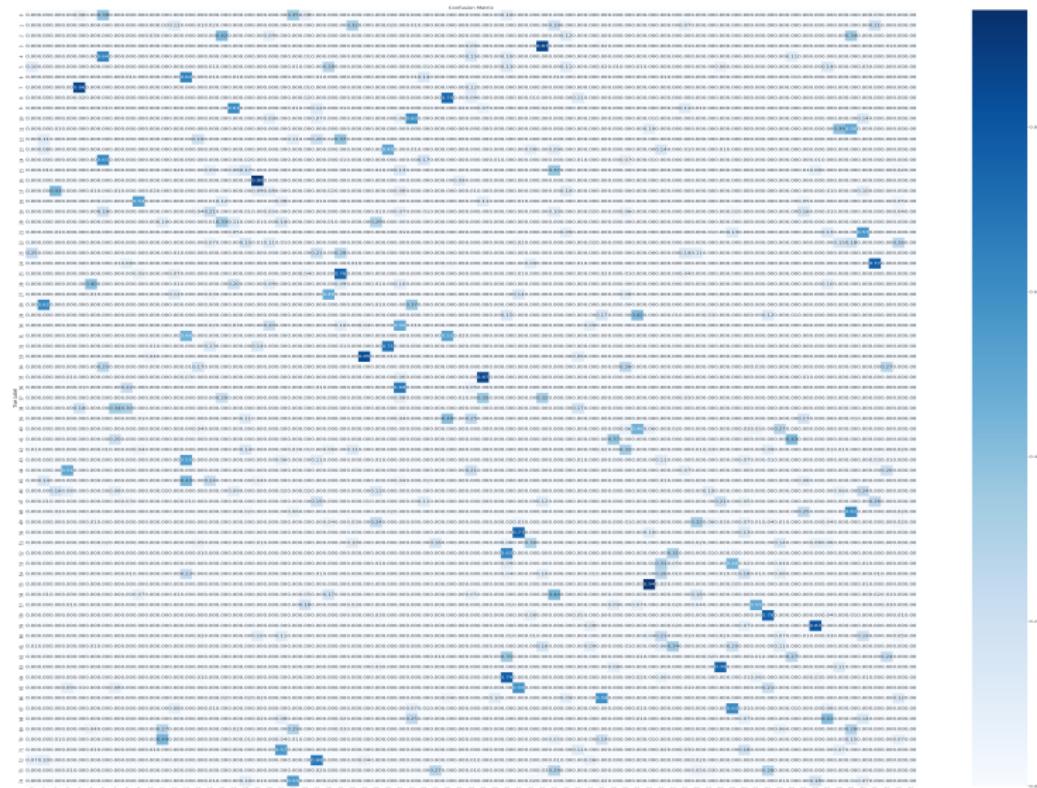
# Weak Generalization: Standard Set Transformer







## Strong Generalization: Contrastive pre-trained model



# Generalization

**Table:** Results for Generalization using KL Divergence

Method	KL Divergence
Standard Set Transformer	6.60
Contrastive Pre-trained model	6.25



Computational Finance

# Visualization

## PCA Visualization for Standard Set Transformer

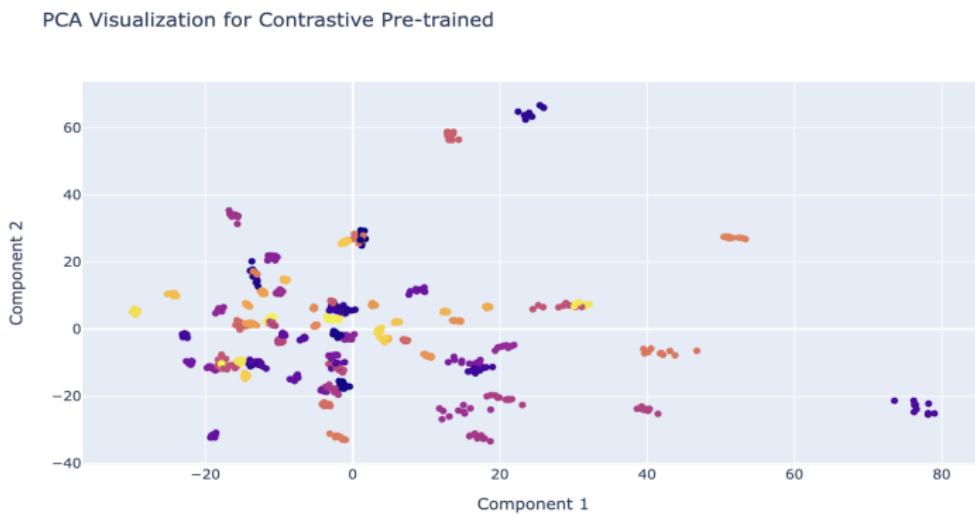


**Figure:** PCA Visualization for Standard Set Transformer

- Principal Component 1: 30.26%
- Principal Component 2: 24.38%
- Total Principal Component: 54.64%

# Visualization

## PCA Visualization for Contrastive Pre-trained Model

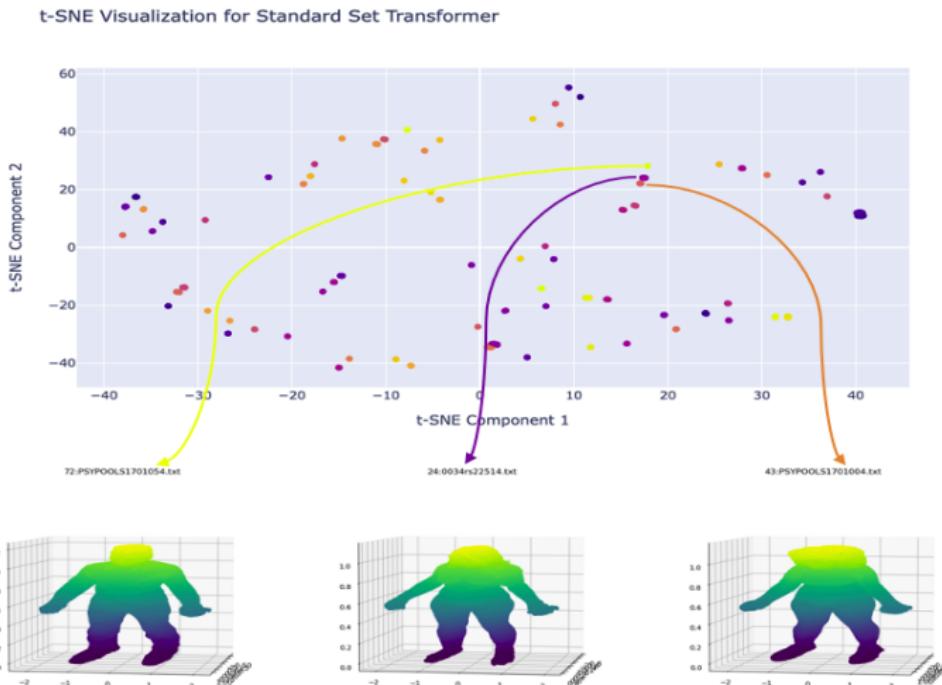


**Figure:** PCA Visualization for Contrastive Pre-trained

- Principal Component 1: 20.50%
- Principal Component 2: 18.35%
- Total Principal Component: 38.85%

# Visualization

## t-SNE Visualization for Standard Set Transformer



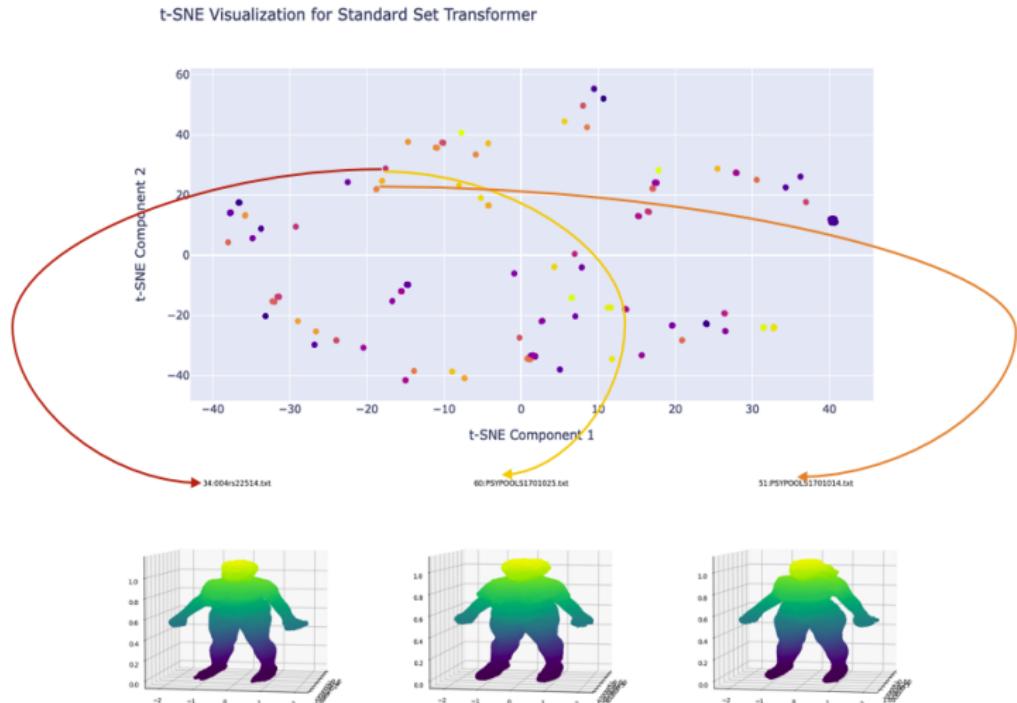
**Figure:** t-SNE Visualization for Standard Set Transformer



Computational Science

# Visualization

## t-SNE Visualization for Standard Set Transformer



**Figure:** t-SNE Visualization for Standard Set Transformer

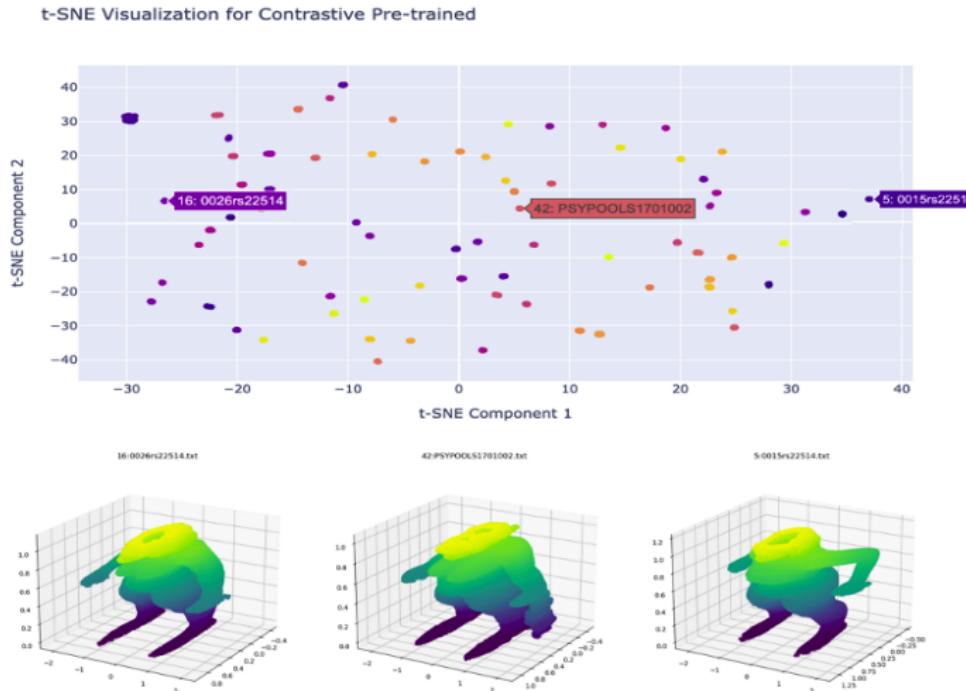


Computational Science



# Visualization

## t-SNE Visualization for Contrastive Pre-trained Model



**Figure:** t-SNE Visualization for Contrastive Pre-trained Model



Computational Science

# Conclusion

- Standard Set transformer and Contrastive Pre-trained model for self-identification task.
- Weak and strong generalization for efficiency of the model.
- Contrastive Pre-trained model excelled Standard Set Transformer performance in terms of accuracy, speed, and stability.
- Contrastive pre-trained model showed improved generalization than the standard Set Transformer.



# Conclusion

## Future work

- Contrastive Pre-trained model can be used as a solid foundation for extended work like Generative Adversarial Set Transformer(GAST).
- GAST approach is to generate point cloud data to form a full or a partial 3D body scan.
- This approach is in the notion to provide support in case of missing scanned parts or lesser point cloud data that helps with increased precision.



# Reference

- [TC]2 Introduces KX-16 Body Scanner. Textile World. (2012, March 20). <https://www.textileworld.com/textile-world/new-products/2012/03/tc2-introduces-kx-16-body-scanner-3/>
- Lee, J., Lee, Y., Kim, J., Kosiorek, A. R., Choi, S., & Teh, Y. W. (2018). Set Transformer: A Framework for Attention-based Permutation-Invariant Neural Networks. ArXiv. /abs/1810.00825
- Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). A Simple Framework for Contrastive Learning of Visual Representations. ArXiv. /abs/2002.05709
- Wikimedia Foundation. (2023, July 16). Kullback–Leibler divergence. Wikipedia. [https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler\\_diverge](https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler_divergence)

# Thank you