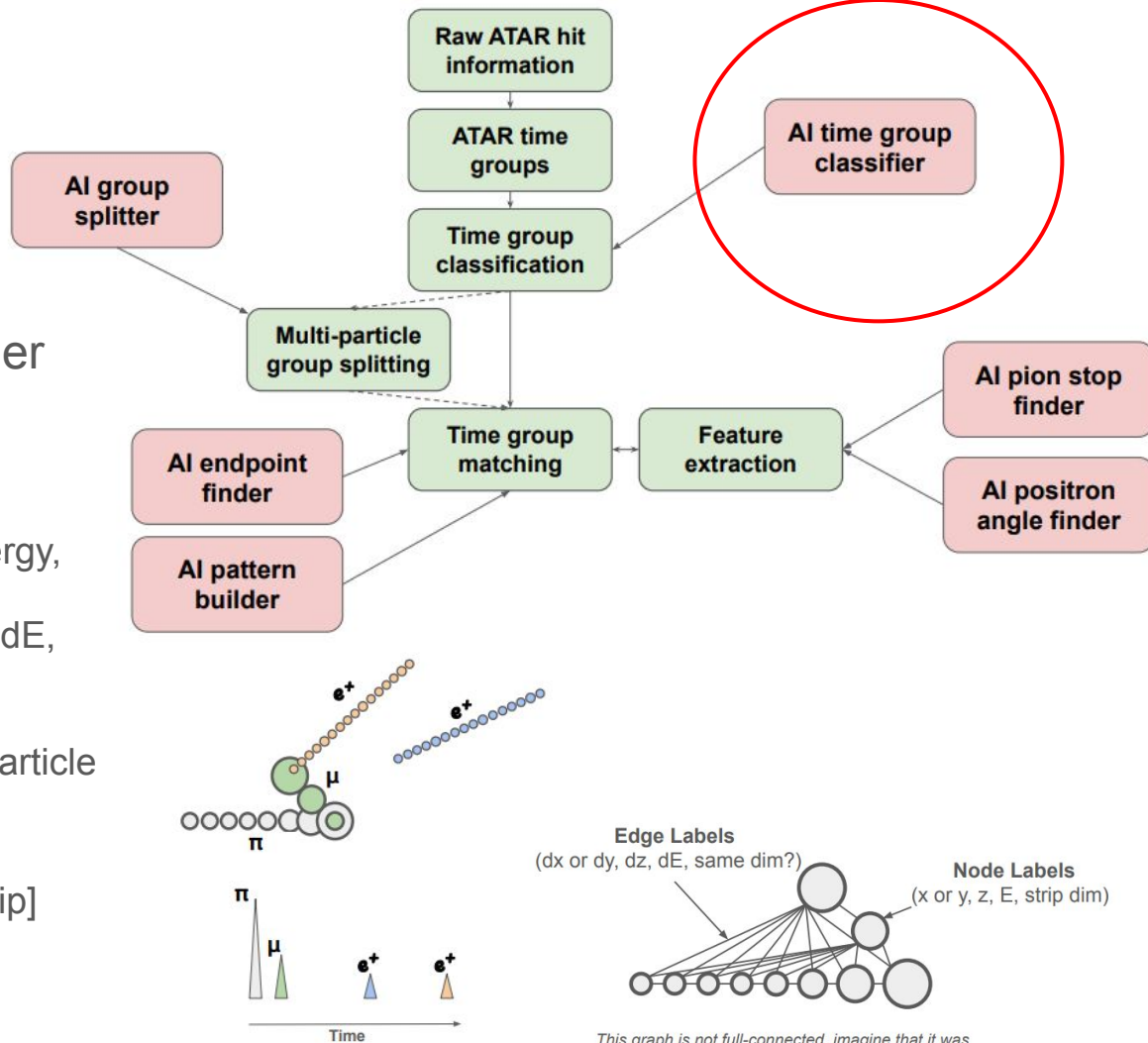# Training and Validation Details For Classification Models for PIONEER Reconstruction

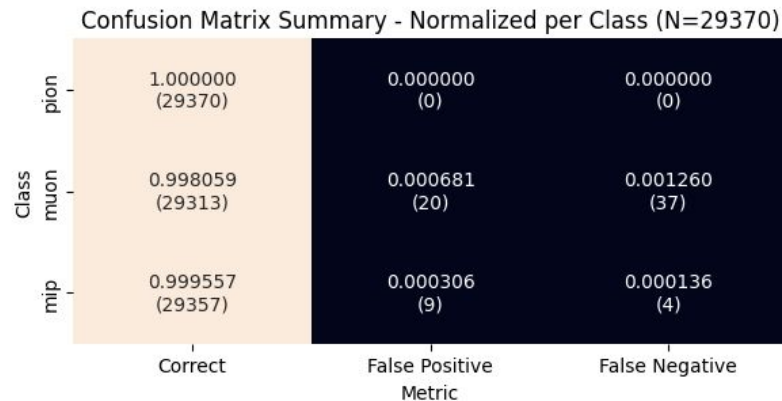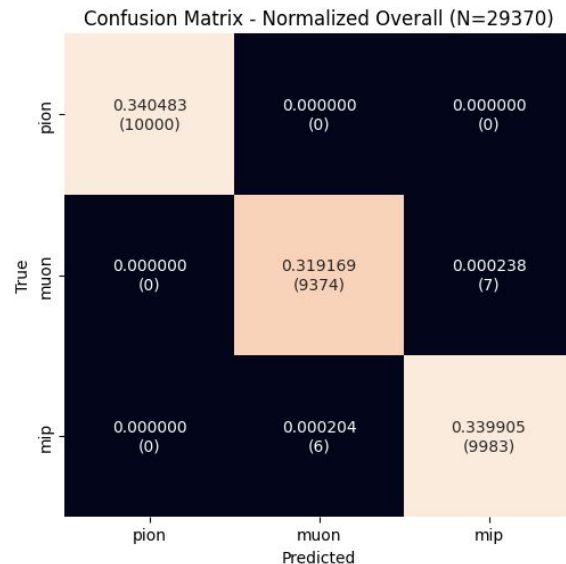Jack Carlton
University of Kentucky

# Time Group Classifier

- I'll mostly be talking about Omar's "time group classifier model"
- Input:
  - Time grouped hit graph
    - Nodes: [x (or y), z, energy, view, group_energy]
    - Edges: [dx (or dy), dz, dE, same_view]
  - Groups split by time (could potentially contain multiple particle types)
- Output
  - Class labels: [muon, pion, mip]
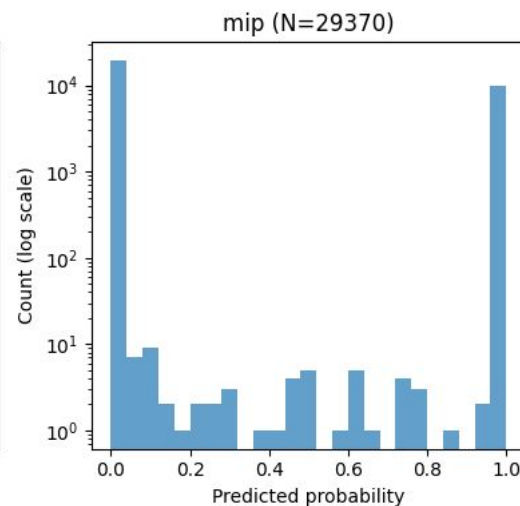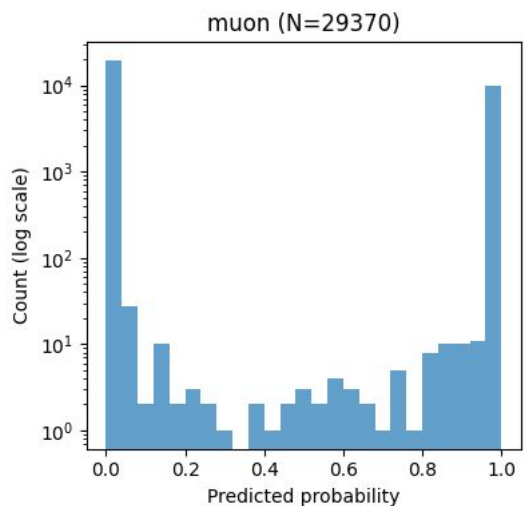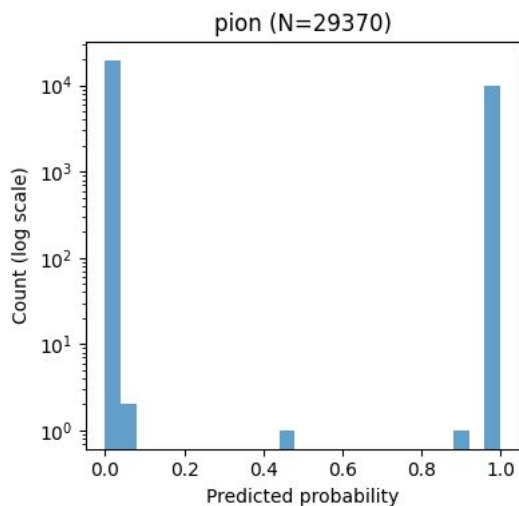- Much better explained in Omar's presentation

# Performance of Model

- Closely matches performance of [Omar's work](#)
- Differences from Omar's work
  - Used my machine to train
  - Did a hyperparameter search using Optuna
  - Incorporated into ZenML framework for creating pipelines in python
- Unsure how this compares to the traditional reco values(?)
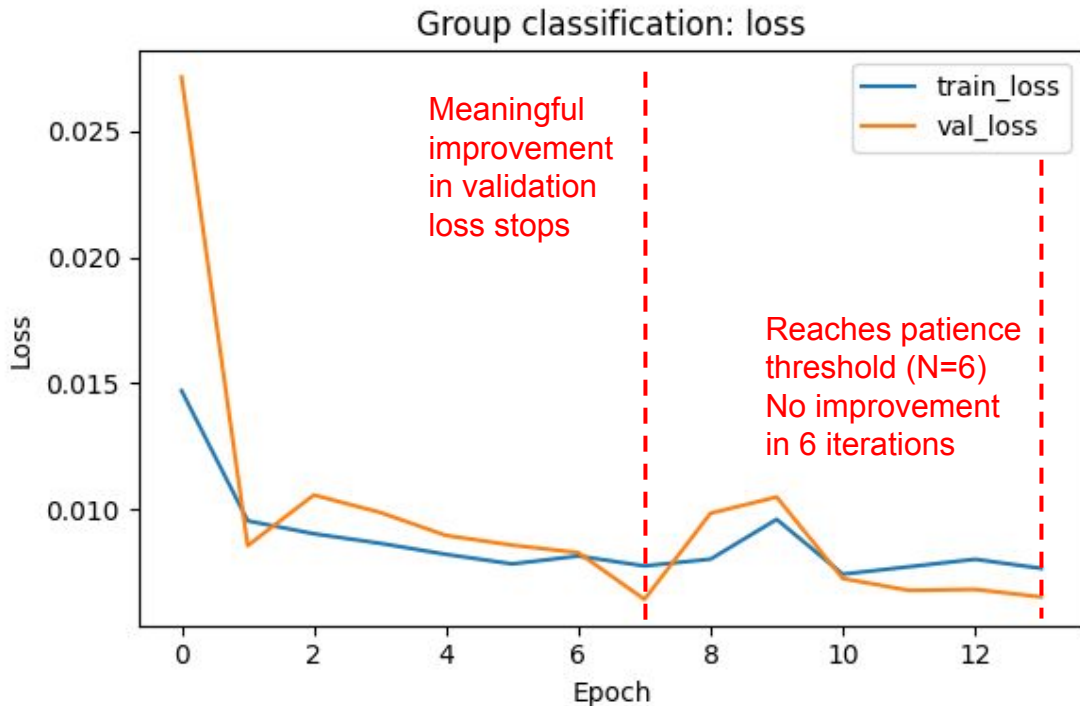- Unsure exact parameters in Omar's data set



Confusion Matrix - Normalized Overall (N=29370)

| True \ Predicted | pion | muon | mip |
|---|---|---|---|
| pion | 0.340483 (10000) | 0.000000 (0) | 0.000000 (0) |
| muon | 0.000000 (0) | 0.319169 (9374) | 0.000238 (7) |
| mip | 0.000000 (0) | 0.000204 (6) | 0.339905 (9983) |



Confusion Matrix Summary - Normalized per Class (N=29370)

| Class \ Metric | Correct | False Positive | False Negative |
|---|---|---|---|
| pion | 1.000000 (29370) | 0.000000 (0) | 0.000000 (0) |
| muon | 0.998059 (29313) | 0.000681 (20) | 0.001260 (37) |
| mip | 0.999557 (29357) | 0.000306 (9) | 0.000136 (4) |

# Performance of Model

- Histograms below shows the models "sureness" of each class
- Want to see large peaks at 0 and 1
  - 0 → this is definitely not this class
  - 1 → this is definitely this class
- The muon groups are the biggest struggle
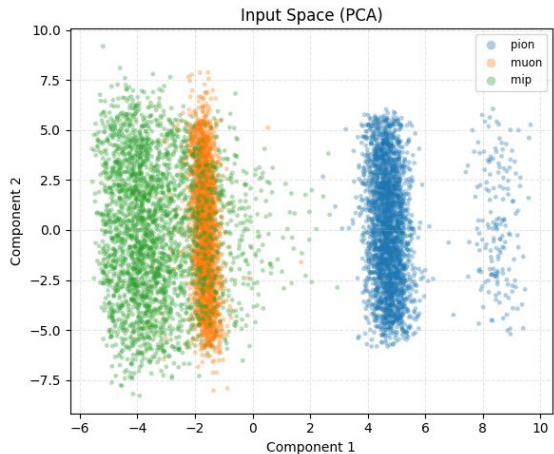
# Preventing Overtraining (Early Stopping)

- Can set early stopping when training loss curves
  - If no meaningful (> δ) improvement within N epochs, stop training
  - N := "patience" usually set to ~5% of expected training epochs
- Potential improvements:
  - Loss curve smoothing
  - "Noise" estimations, only continue training if smoothed improvement > noise of previous iterations
- See documentation

*NOTE: This particular model is overtrained because N set to 6 (too high)



Group classification: loss

Meaningful improvement in validation loss stops

Reaches patience threshold (N=6)
No improvement in 6 iterations

# Preventing Overtraining (PCA)

- Need a way to visualize many dimensions
- [Principal Component Analysis](#) (PCA) good candidate for this
- Compare clusters in [embedding space](#) vs. input space
  - Input space shows data driven groupings
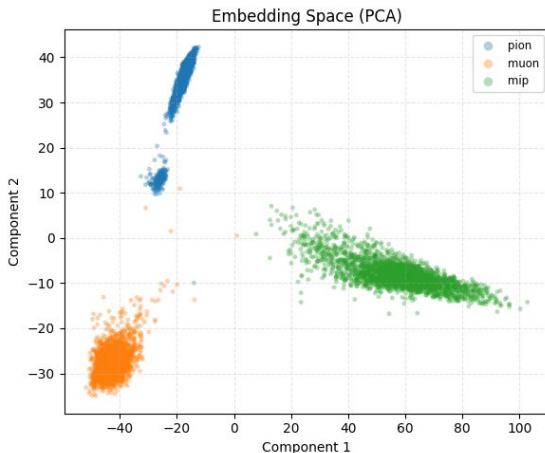  - Embedding space shows learned groupings
  - # of groupings should match!



Average vector of all nodes in graph:
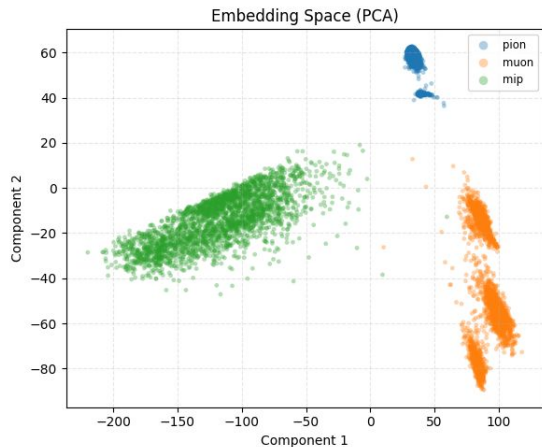[coord,
 z_pos,
energy,
view,
group_energy]

5D → 2D
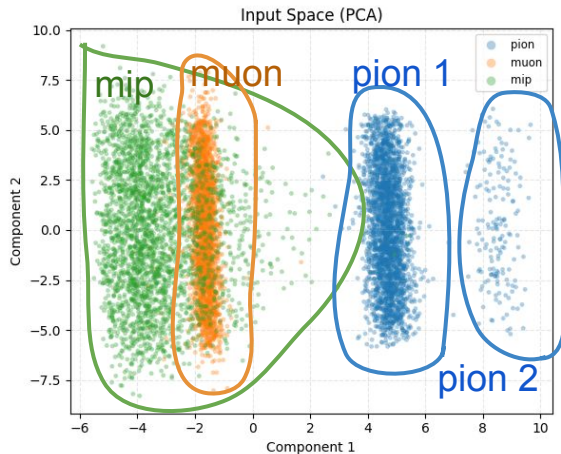
Well-trained          Over-trained



512D → 2D



512D → 2D

# Preventing Overtraining (PCA)

- Need a way to visualize many dimensions
- [Principal Component Analysis](#) (PCA) good candidate for this
  - Choose to project down to d = 2 dimensions this way
- Compare clusters in [embedding space](#) vs. input space
  - Input space shows data driven groupings
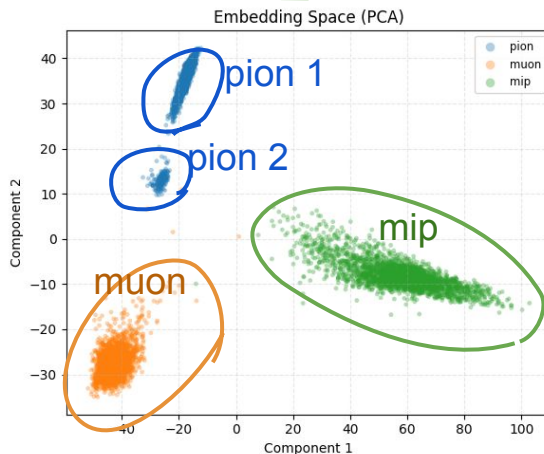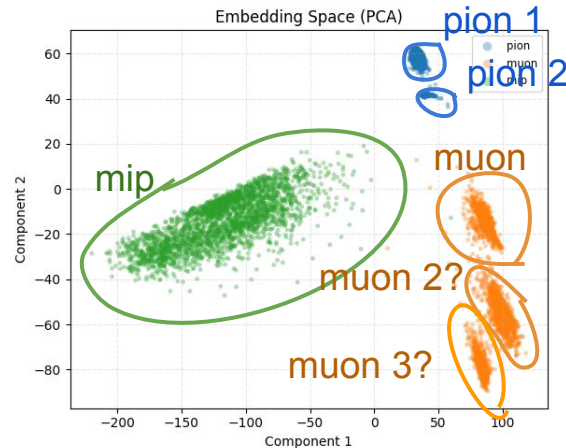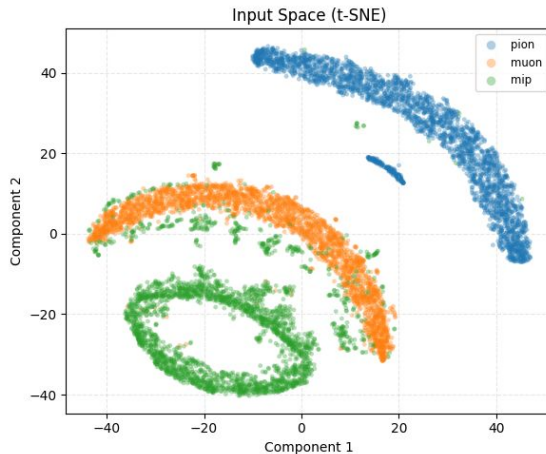  - Embedding space shows learned groupings
  - # of groupings should match!



Input Space (PCA)

mip   muon   pion 1   pion 2

Average vector of all nodes in graph:
[coord,
 z_pos,
 energy,
 view,
 group_energy]

5D → 2D

Well-trained   Over-trained

Embedding Space (PCA)

pion 1   pion 2   mip   muon

512D → 2D

Embedding Space (PCA)

pion 1   pion 2   mip   muon   muon 2?   muon 3?

512D → 2D

# Preventing Overtraining (t-SNE)

- Problem with PCA
  - Global linear may not preserve local neighborhoods
    - May not show groups!
- t-distributed stochastic neighbor embedding (t-SNE) is designed to preserve local clusters
  - Maps N dims → d dims
    - We choose d = 2 for visualization ease
- Same ideas as PCA
  - compare input space and embedding space



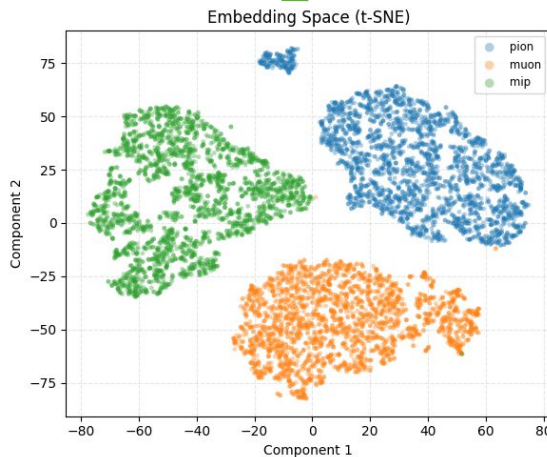Average vector of all nodes in graph:
[coord,
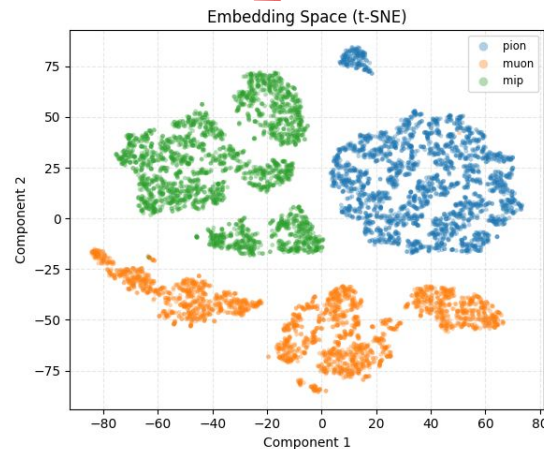 z_pos,
energy,
view,
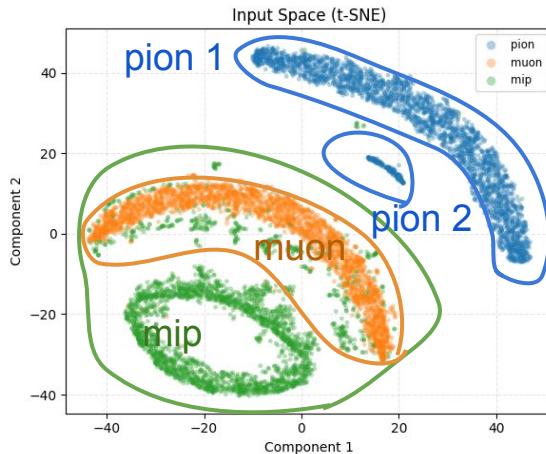group_energy]

5D → 2D

Well-trained          Over-trained



512D → 2D



512D → 2D

# Preventing Overtraining (t-SNE)

- Problem with PCA
  - Global linear may not preserve local neighborhoods
    - May not show groups!
- t-distributed stochastic neighbor embedding (t-SNE) is designed to preserve local clusters
  - Maps N dims → d dims
    - We choose d = 2 for visualization ease
- Same ideas as PCA
  - compare input space and embedding space



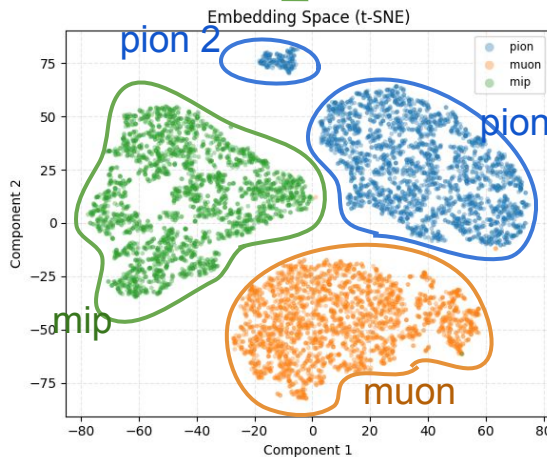Average vector of all nodes in graph:
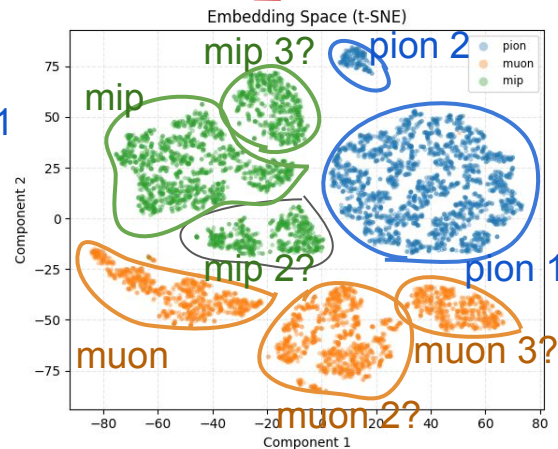[coord, z_pos, energy, view, group_energy]
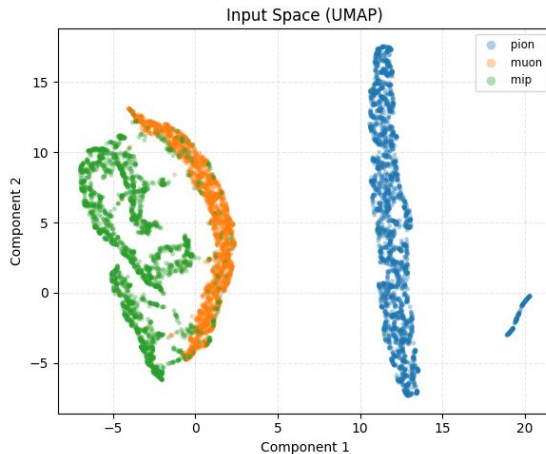
5D → 2D

Well-trained                    Over-trained

512D → 2D                    512D → 2D

# Preventing Overtraining (UMAP)

- Problem with t-SNE
  - Depends on a "Perplexity"
    - Parameter, ~how many neighbors a point can have
  - May artificially split or group clusters
- [Uniform Manifold Approximation and Projection](Uniform Manifold Approximation and Projection) (UMAP) is designed to preserve the underlying manifold structure
  - Tries to preserve both local neighborhoods *and* their global relationships
  - Maps N dims → d dims
    - We choose d = 2 for visualization ease
- Computational expensive
  - Use as a "tie breaker" if PCA and t-SNE disagree



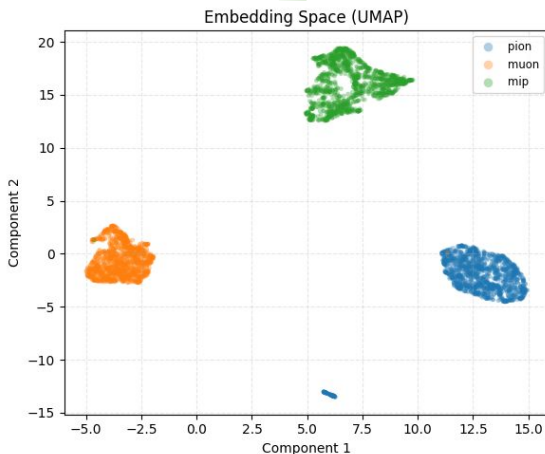Average vector of all nodes in graph:
[coord,
 z_pos,
 energy,
 view,
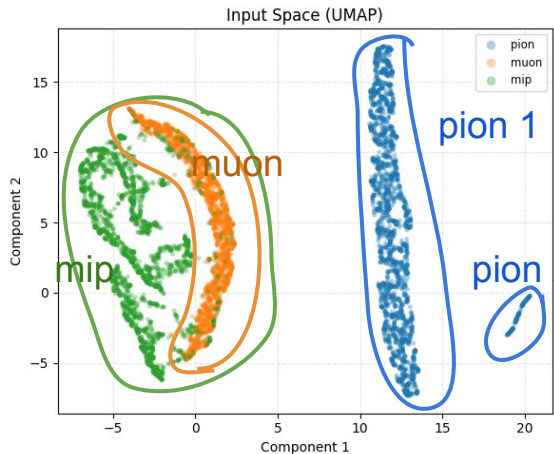 group_energy]

5D → 2D

Well-trained          Over-trained



Sorry!
I don't have an example for this!
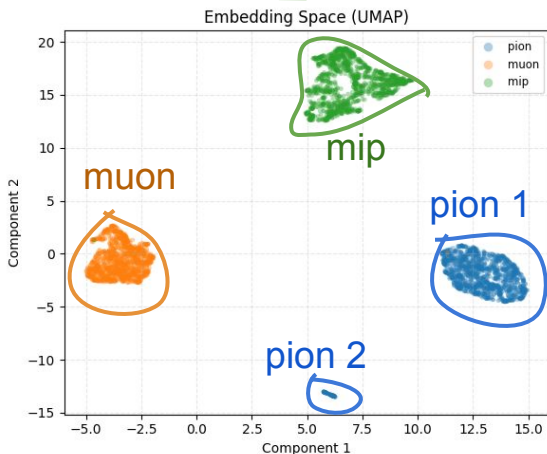
512D → 2D          512D → 2D

# Preventing Overtraining (UMAP)

- Problem with t-SNE
  - Depends on a "Perplexity"
    - Parameter, ~how many neighbors a point can have
  - May artificially split or group clusters
- [Uniform Manifold Approximation and Projection](#) (UMAP) is designed to preserve the underlying manifold structure
  - Tries to preserve both local neighborhoods *and* their global relationships
  - Maps N dims → d dims
    - We choose d = 2 for visualization ease
- Computational expensive
  - Use as a "tie breaker" if PCA and t-SNE disagree



Input Space (UMAP)

Average vector of all nodes in graph: [coord, z_pos, energy, view, group_energy]

5D → 2D

Well-trained          Over-trained

Embedding Space (UMAP)

Sorry!
I don't have an example for this!

512D → 2D          512D → 2D

# Auxiliary Slides