# Single-cell transcriptome profiling (scRNA-seq) is an important technique to study cellular heterogeneity

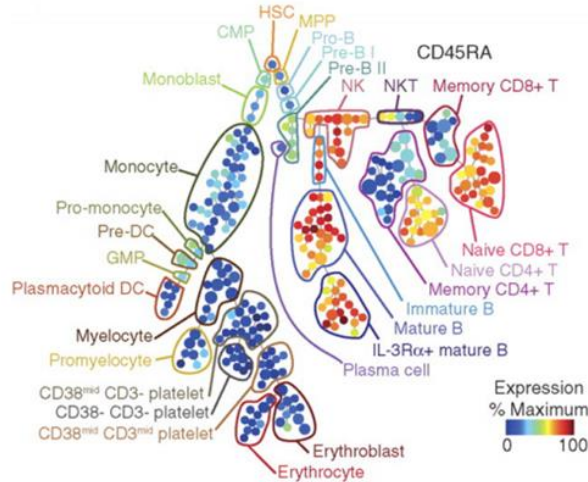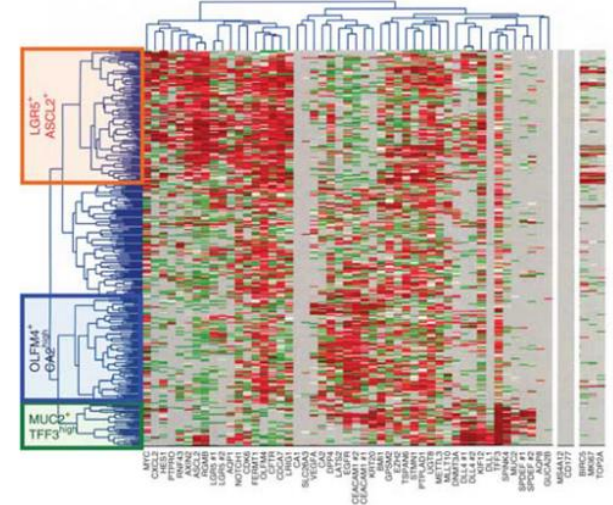

**Cellular heterogeneity**

**Differentiation trajectories**

Bendall et al. (2011), Science

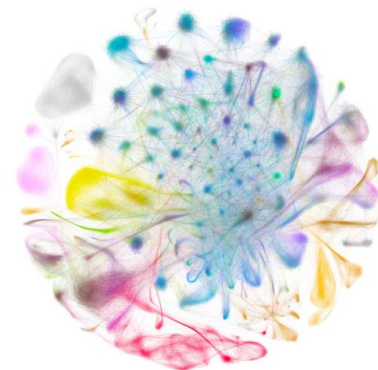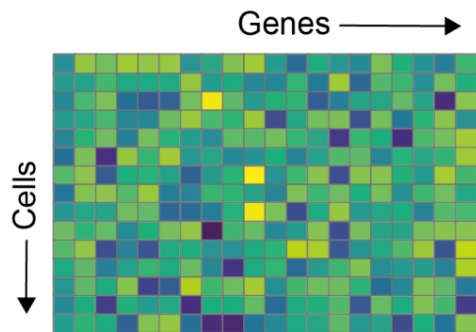**Within-cell-type differences**

Dalerba et al. (2011), Nature Biotech

# Analysis of scRNA-seq data often involves manifold embedding

**Individual**  **Cell mixture**  **Cell expression**  **Cell manifold**



What are the genes differentiating cell types?

# Caveats of conventional scRNA-seq differentiation analysis workflow



**Caveats:**
- ❏ The **uncertainty induced by the nonlinear embedding** is ignored
- ❏ The **stochasticity in cluster assignments** is ignored
- ❏ The **dependency among genes** is ignored
- ❏ **Only genes enriched in single clusters are highlighted** as the signature

# Adversarial Clustering Explanation (ACE) overcomes limitations of existing methods
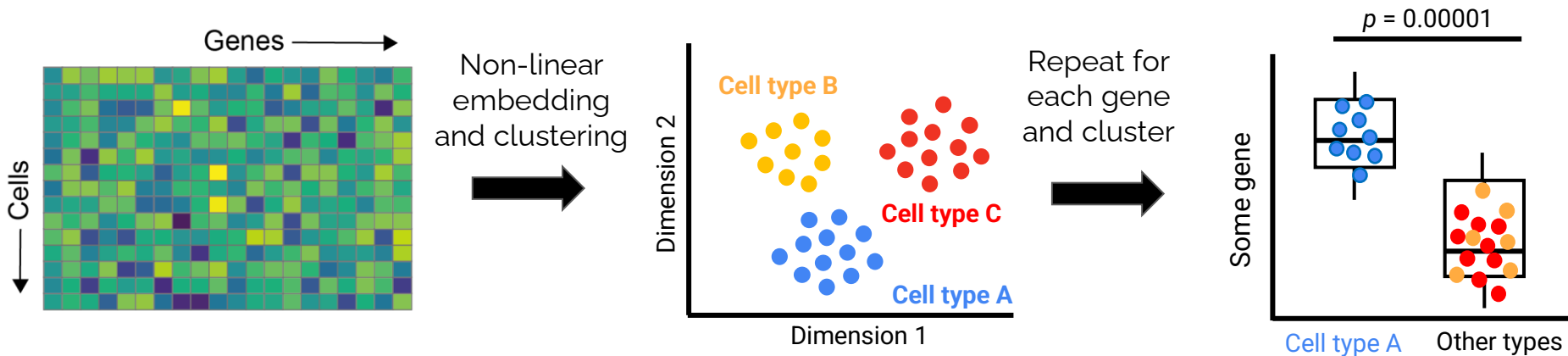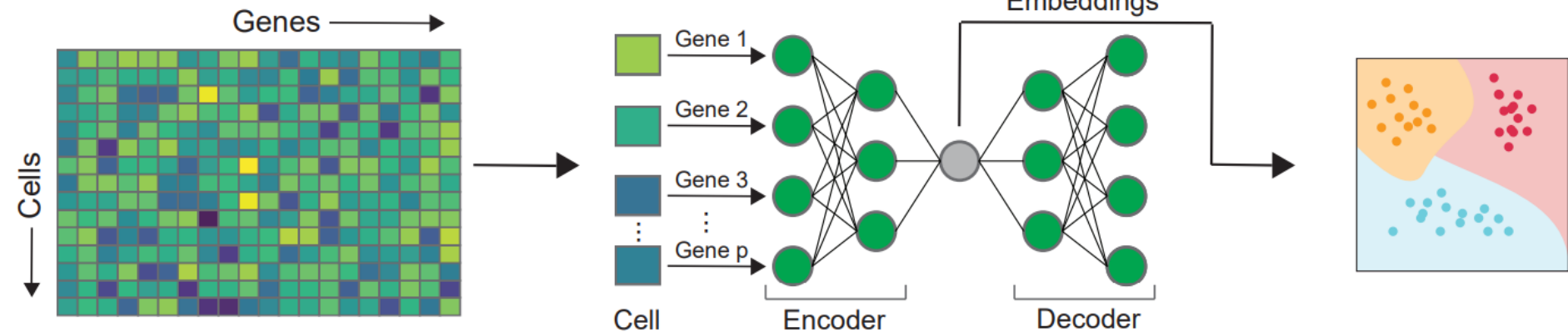
| | ACE | DESeq2 | Jensen-Shannon Distance (Monocle 3) | Global Counterfactual Explanation | Gene Relevance Score |
|---|---|---|---|---|---|
| Dependency among genes | ✓ | ✗ | ✗ | ✓ | ✓ |
| Uncertainty by the nonlinear embedding | ✓ | ✗ | ✗ | ✗ | ✓ |
| Stochasticity in cluster assignments | ✓ | ✗ | ✗ | ✗ | ✗ |
| Does not limit to enriched genes | ✓ | ✗ | ✗ | ✓ | ✓ |
| Additional limitations | NA | NA | NA | Limited to linear transformation | Limited to diffusion maps |

*Love et al, Genome Biology (2014)*
*Cao et al, Nature (2019)*
*Plumb et al, ICML (2020)*
*Angerer et al, Bioinformatics (2020)*

# ACE aims to jointly explain the embedding and clustering



1. Input: gene expression matrix
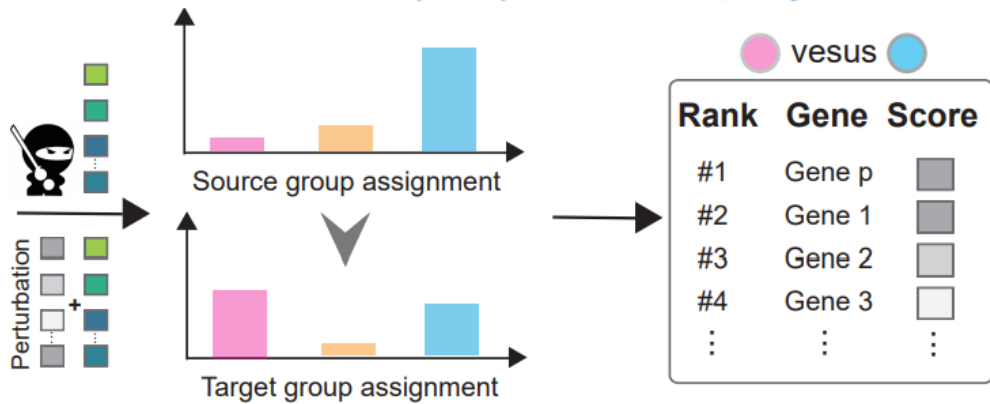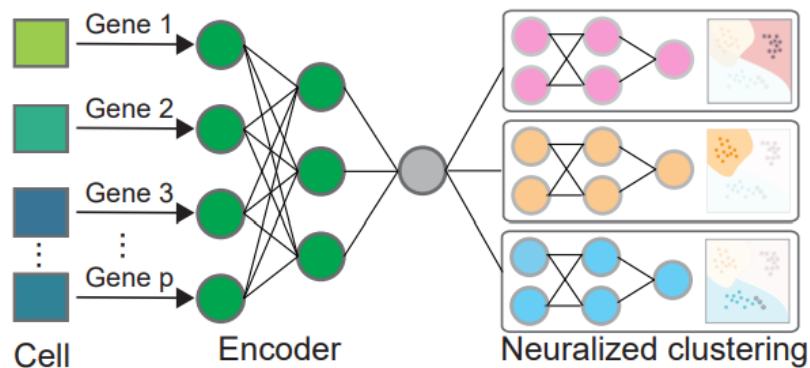
2. Deep autoencoder learns low-dimensional representation

3. Embedding clustering

Embeddings

Cell | Encoder | Decoder

4. Clustering is neuralized and concatenated with the encoder

5. Differentiation analysis by ACE

6. Output: gene relevance

Cell | Encoder | Neuralized clustering

Perturbation

Source group assignment

Target group assignment

vesus

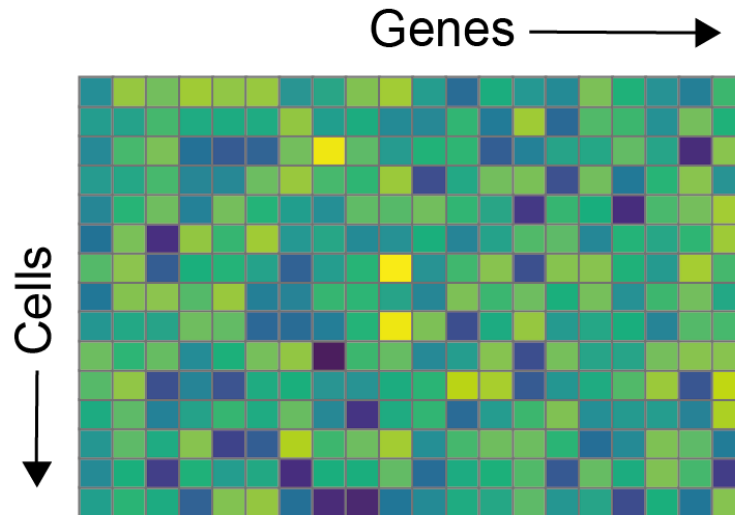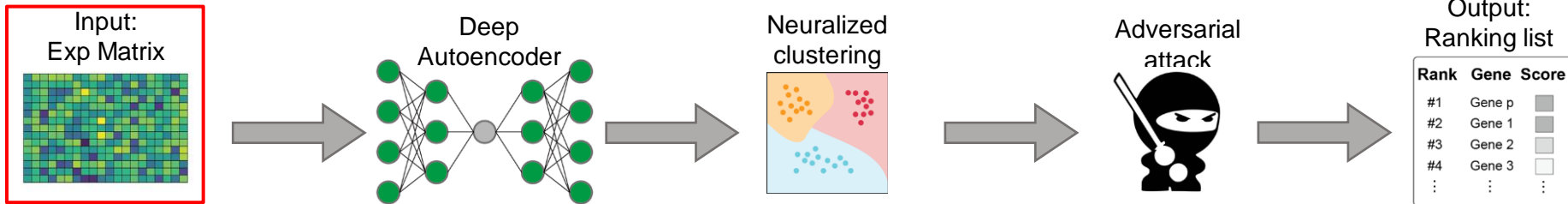| Rank | Gene | Score |
|------|--------|-------|
| #1 | Gene p | |
| #2 | Gene 1 | |
| #3 | Gene 2 | |
| #4 | Gene 3 | |
| ⋮ | ⋮ | ⋮ |

# ACE takes as input the expression matrix and a pre-specified number of clusters



**Input:**

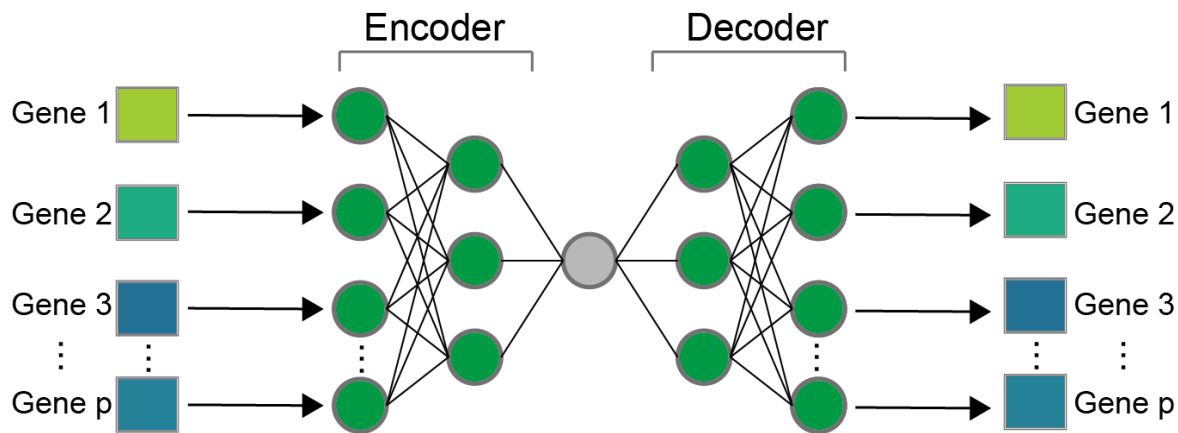❏ The scRNA-seq expression matrix
❏ The specified cluster number **k**

Genes →

Cells →

**At a high level:**

Input: Exp Matrix

Deep Autoencoder

Neuralized clustering

Adversarial attack

Output: Ranking list

| Rank | Gene | Score |
|------|--------|-------|
| #1 | Gene p | |
| #2 | Gene 1 | |
| #3 | Gene 2 | |
| #4 | Gene 3 | |
| ⋮ | ⋮ | ⋮ |

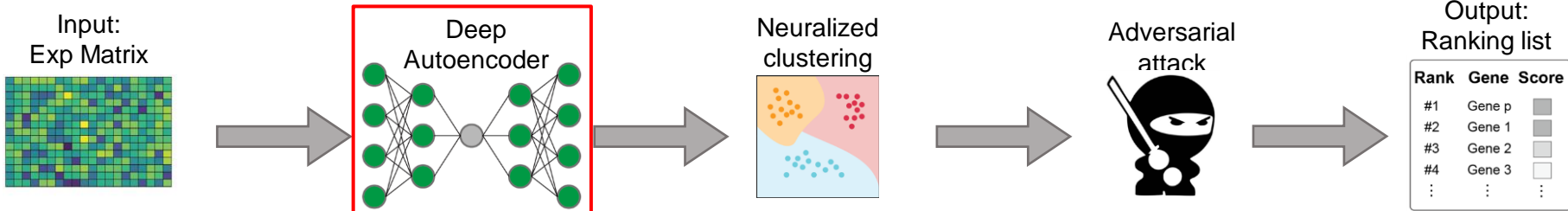# ACE projects the expression data into a low-dimensional embedding using a deep autoencoder

**Embedding by Autoencoder:**

❏ Dimension reduction similar to UMAP, t-SNE, or PCA

❏ Batch correction

❏ Applicable to any scRNA-seq embedding method
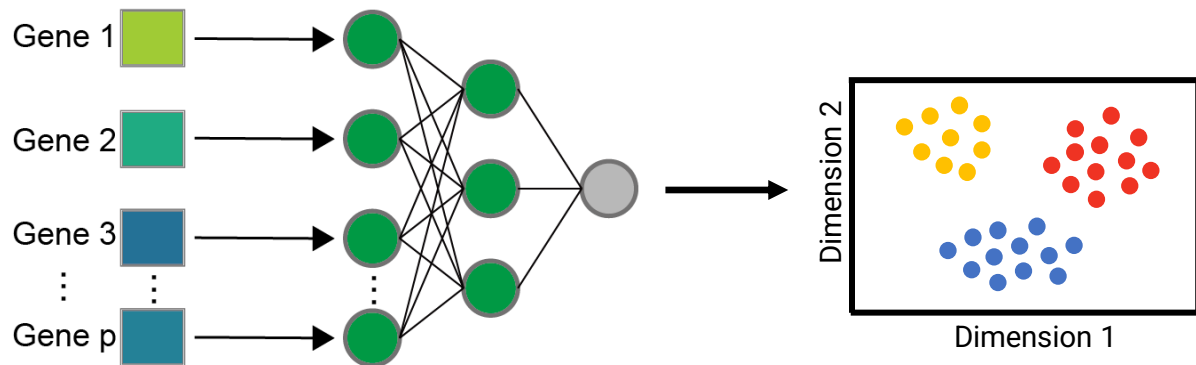


*Amodio et al, Nature Methods (2019)*
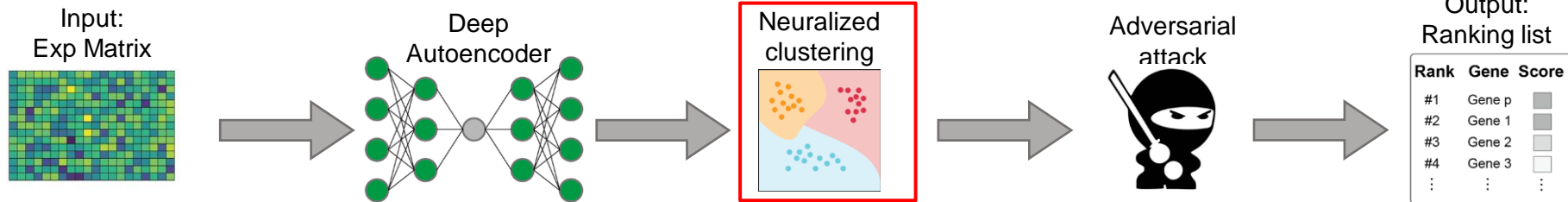
**At a high level:**

# ACE performs k-means clustering in the learned embedding space

**Clustering:**

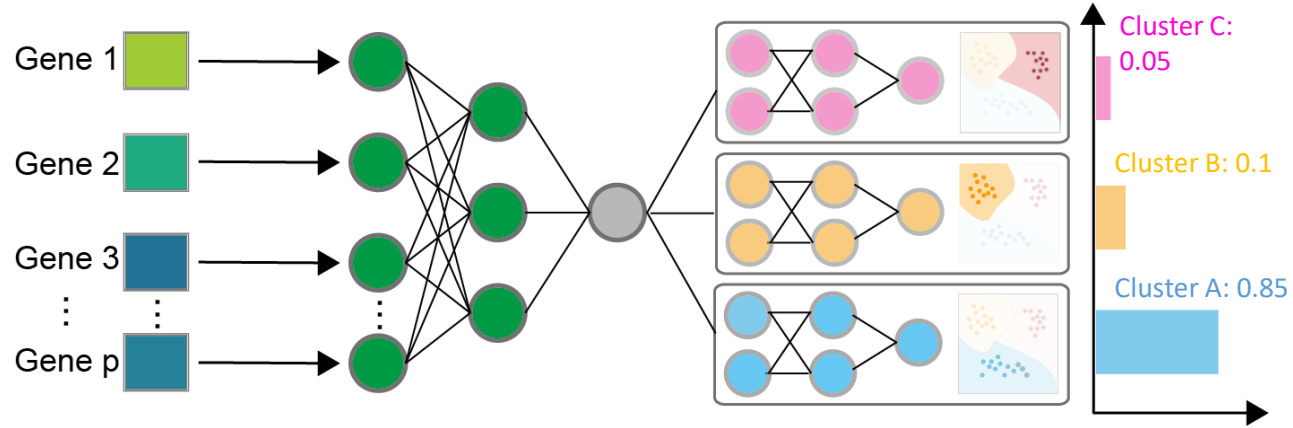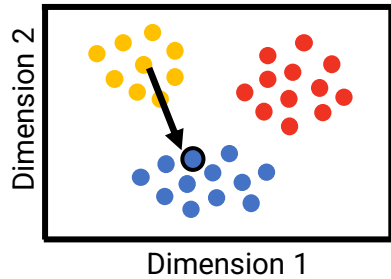❑ K-means clustering in the embedding space



**At a high level:**

# ACE reformulates the k-means clustering as a functionally equivalent multi-layer neural network
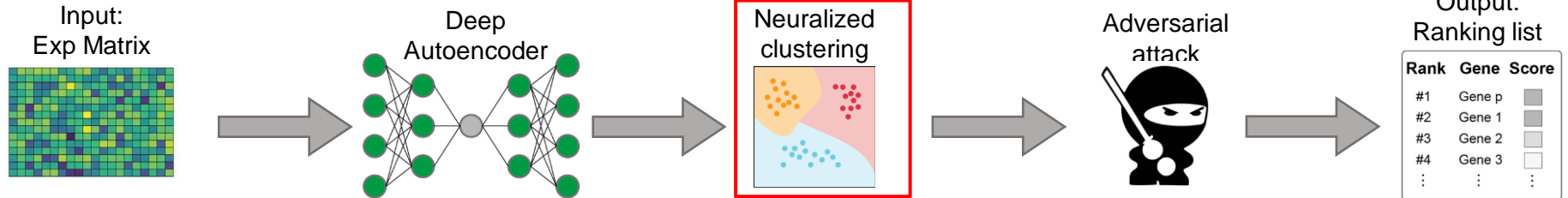
**Cluster neuralization:**
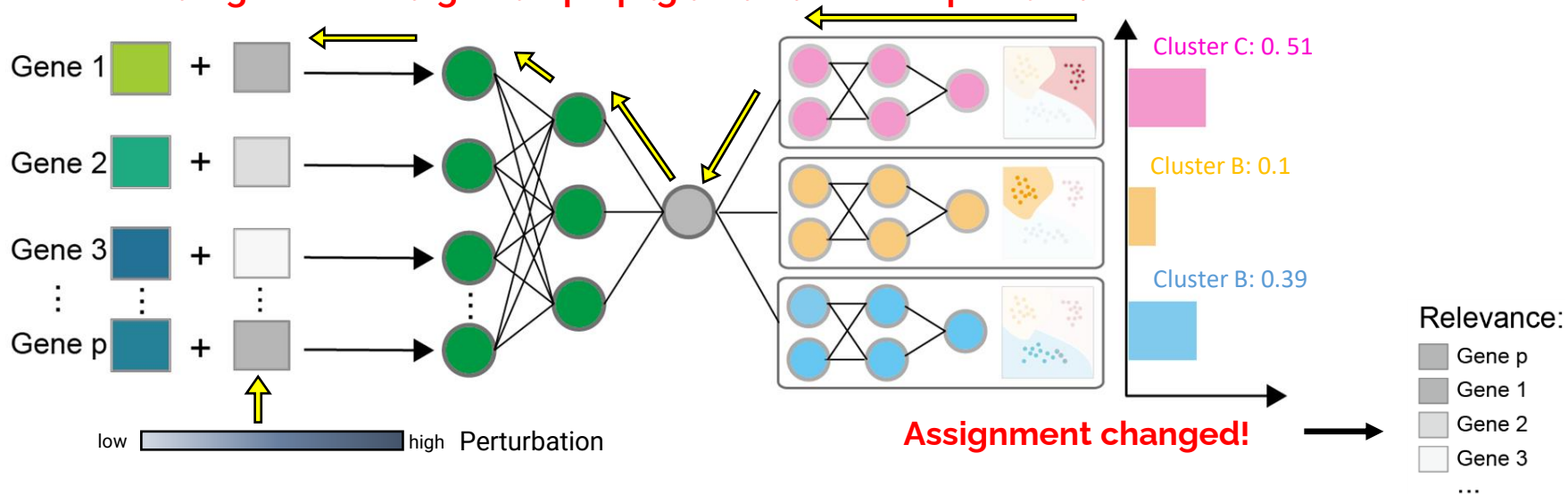
❏ The cluster assignment is preserved for each cell



*Kauffmann et al, arXiv:1906.07633 (2019)*
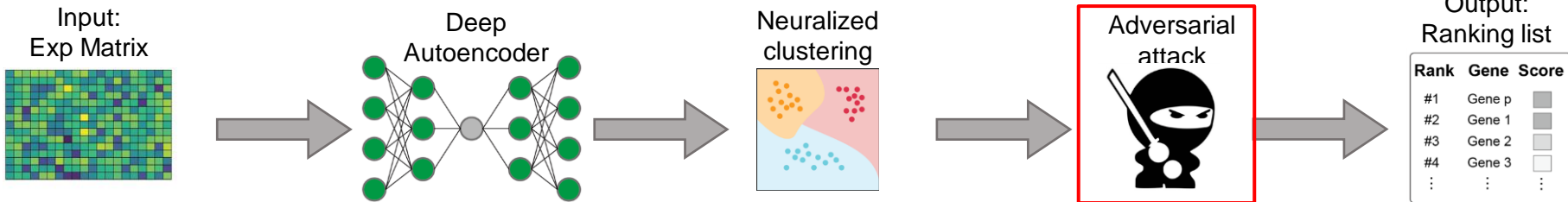
**At a high level:**

# ACE finds the **minimal perturbation** to a cell that causes the clustering assignment to change



**The assignment changes are propagated back to the perturbation**

Gene 1 +
Gene 2 +
Gene 3 +
⋮
Gene p +

low ▭ high  Perturbation

Cluster C: 0. 51

Cluster B: 0.1

Cluster B: 0.39

**Assignment changed!**

Relevance:
▭ Gene p
▭ Gene 1
▭ Gene 2
▭ Gene 3
…

**At a high level:**

Input:
Exp Matrix

Deep
Autoencoder

Neuralized
clustering

Adversarial
attack

Output:
Ranking list

| Rank | Gene | Score |
|------|--------|-------|
| #1 | Gene p | ▭ |
| #2 | Gene 1 | ▭ |
| #3 | Gene 2 | ▭ |
| #4 | Gene 3 | ▭ |
| ⋮ | ⋮ | ⋮ |

# Sanity check: we applied ACE to identify digit transitions in a pixel-wise manner



Perturbation range: -1 [color bar] +1    Pixels in initial digit: ▪

# ACE is applied to a simulated dataset with many **redundant** genes



20 causal genes

100 dependent genes

100 noise genes

combined dataset

Both causal and noise genes are simulated for 500 cells by using **SymSim** toolkit.

**Dependent genes are weighted sums of random causal genes with noises:**
$\text{Dep}_1 = w_1\text{causal}_1 + w_2\text{causal}_3$
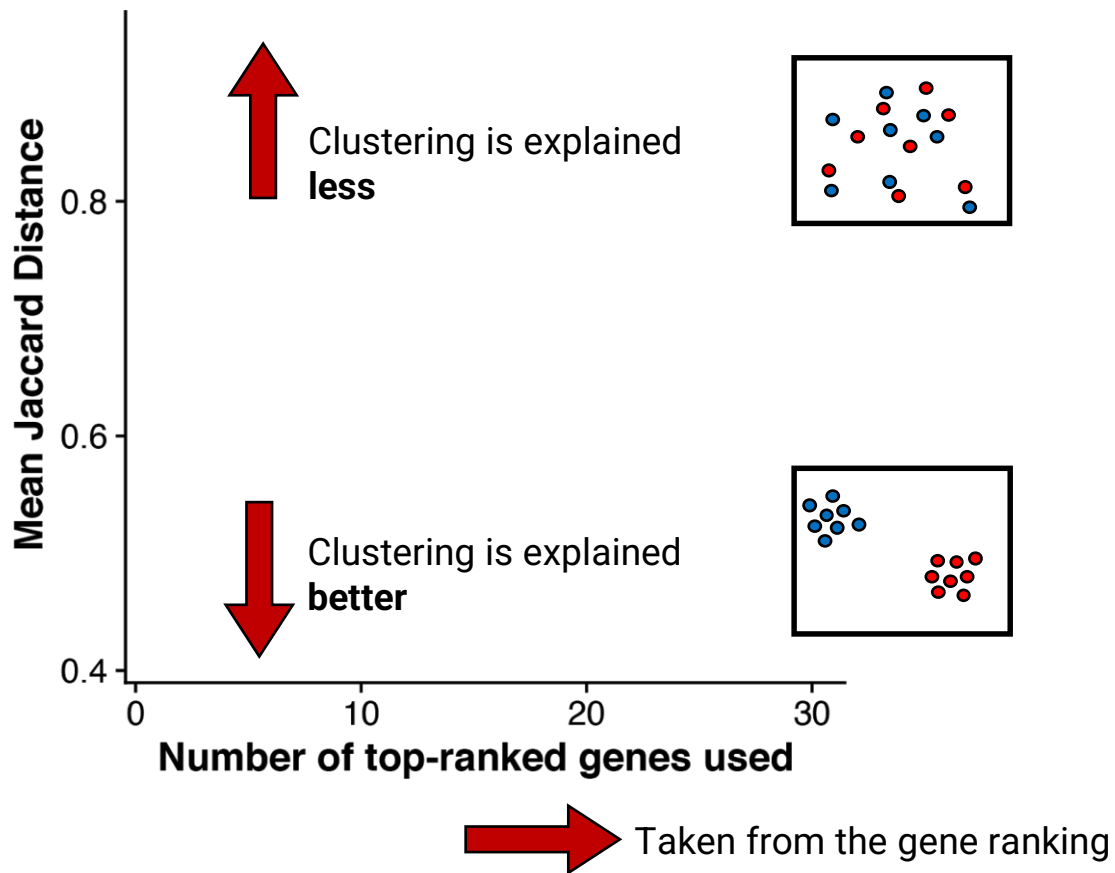$\text{Dep}_2 = w_1\text{causal}_2 + w_2\text{causal}_4 + w_3\text{causal}_6$
$\text{Dep}_3 = w_1\text{causal}_{10}$

…

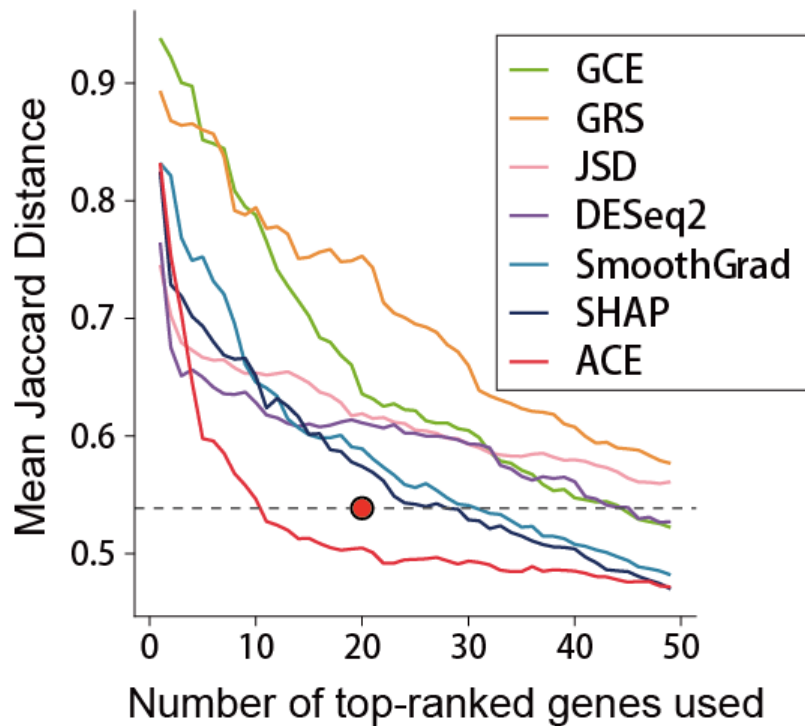**Task:** Rank the genes by relevance

**Metric:** Quantify how well the top $k$ genes in the ranking capture the clustering

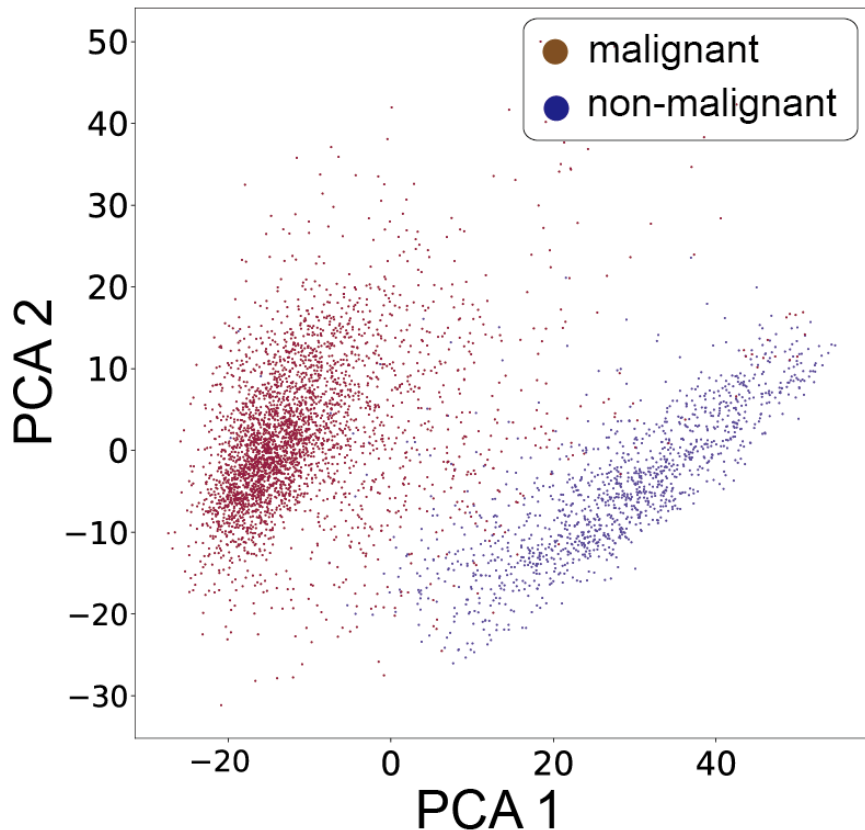*Zhang et al, Nature Communications (2019)*

# The **Mean Jaccard Distance** is a metric for how well a subset of features capture the cluster structure

# ACE is **competitive** against existing methods
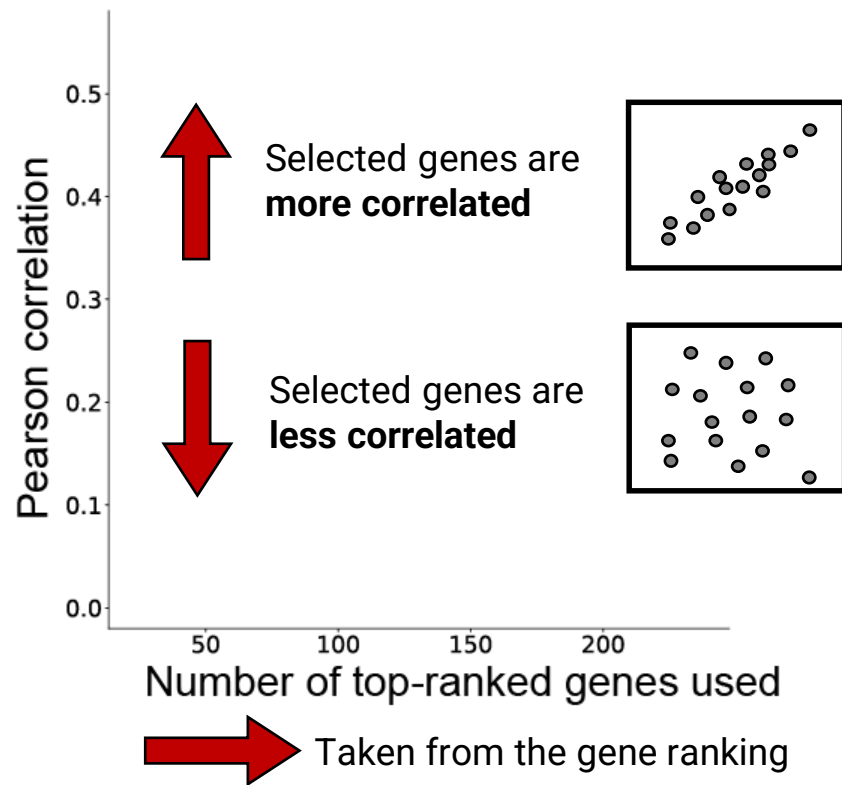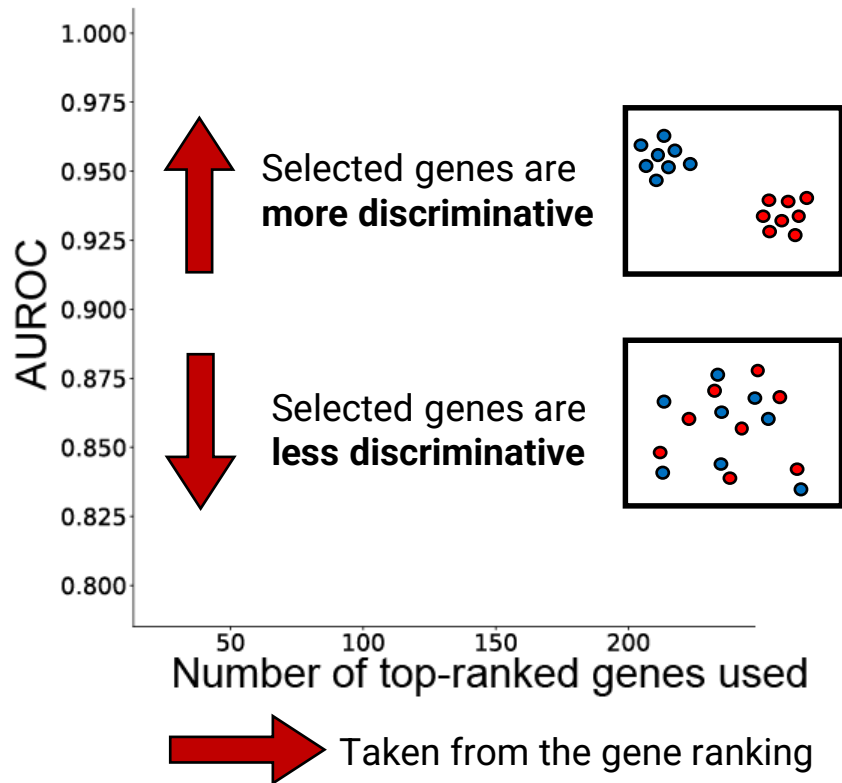
# ACE is applied to a real melanoma dataset



2 cell types (malignant vs. non-malignant)
4513 cells (1257 malignant and 3256 non-malignant)
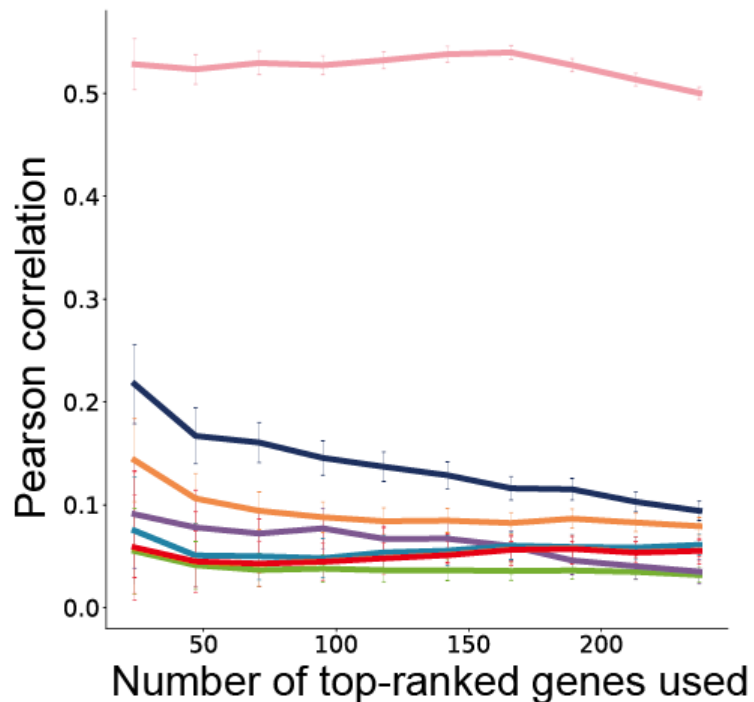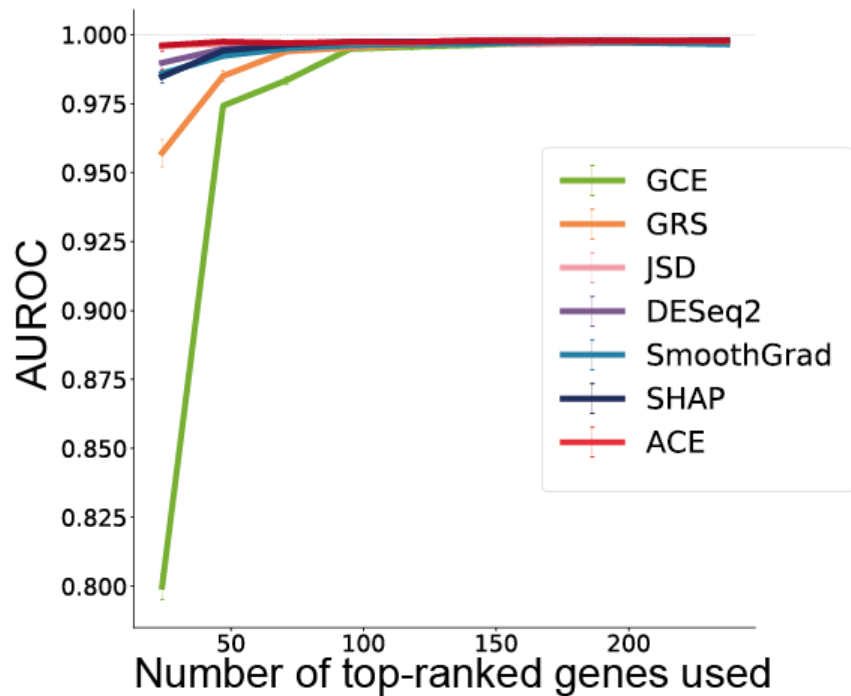23686 genes

**Task:** Rank the genes by relevance

**Metric:**
❏ Quantify how well the top *k* genes in ranking discriminate malignant cells
❏ Quantify how non-redundant/diverse are the top *k* genes in the ranking

*Tiroshi et al, Science (2016)*

# Ideally we want the selected top-ranked genes to be both **highly discriminative** and **non-redundant**
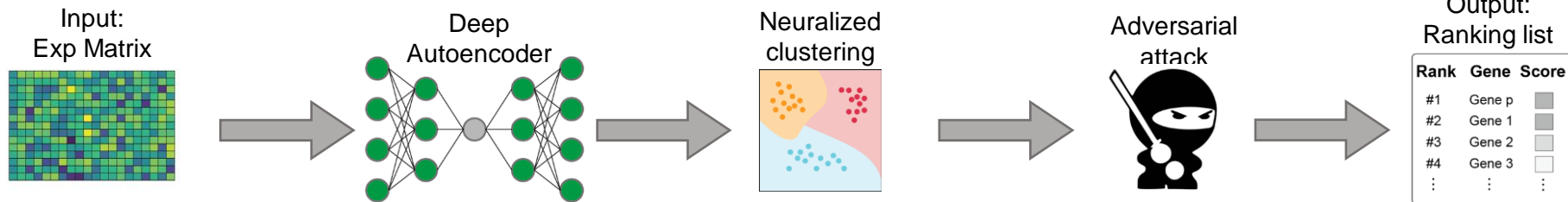
# ACE is **competitive** against existing methods in both discriminative power and minimum redundancy

# Conclusions

❑ ACE finds the minimal set of genes that best explain clustering and is competitive against existing methods.

❑ The selected highly-discriminative genes can be both enriched and depleted.

❑ ACE is potentially useful in domains beyond biology.

❑ Open-source code availability: https://bitbucket.org/noblelab/ace

**At a high level:**



Input: Exp Matrix → Deep Autoencoder → Neuralized clustering → Adversarial attack → Output: Ranking list

| Rank | Gene | Score |
|------|------|-------|
| #1 | Gene p | |
| #2 | Gene 1 | |
| #3 | Gene 2 | |
| #4 | Gene 3 | |
| ⋮ | ⋮ | ⋮ |

# Acknowledgements

- Noble lab members: