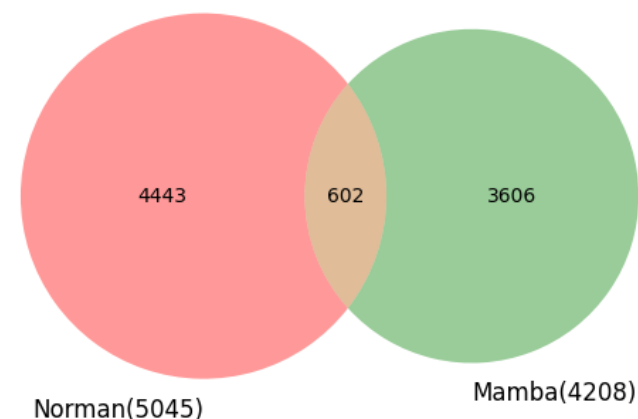# Perturbation prediction tasks

Dataset: Perturb-seq dataset (Norman et al.) containing 131 two-gene perturbations.

```
Local copy of split is detected. Loading...
Simulation split test composition:
combo_seen0:9
combo_seen1:43
combo_seen2:19
unseen_single:36
Done!
Creating dataloaders....
Done!
```
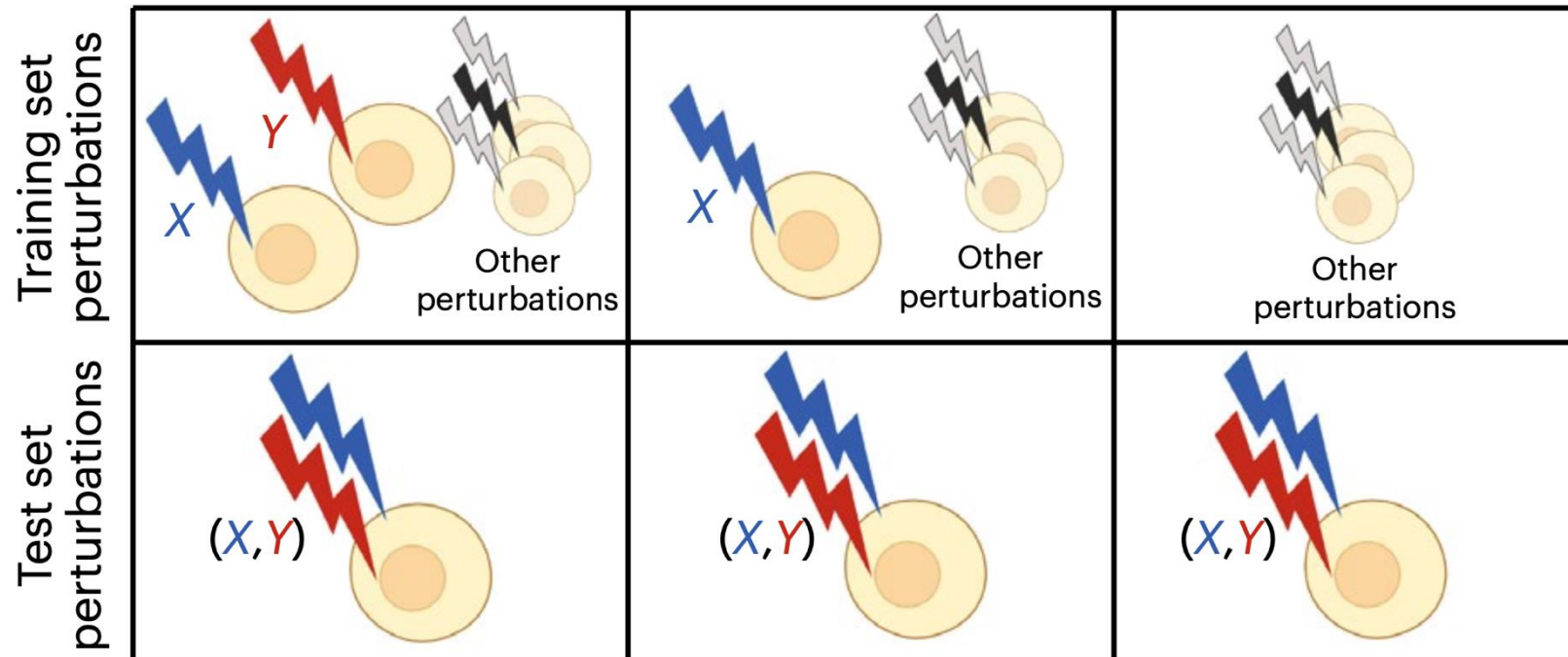
Model: GeneMamba + GEARS



Venn Diagram(number of genes): Norman vs Mamba

## One-gene perturbation

Training set perturbations

Other perturbations

Test set perturbations

$Y$

(1 unseen of 1)

```
'DUSP9+ctrl',
'BCORL1+ctrl',
'MEIS1+ctrl',
'CBL+ctrl',
'SLC4A1+ctrl',
'COL2A1+ctrl',
'S1PR2+ctrl',
```

## Two-gene perturbation

**Increasing difficulty of generalization**

Training set perturbations

$X$ $Y$ Other perturbations

$X$ Other perturbations

Other perturbations

Test set perturbations

$(X,Y)$

$(X,Y)$

$(X,Y)$

(0 unseen of 2)

(1 unseen of 2)
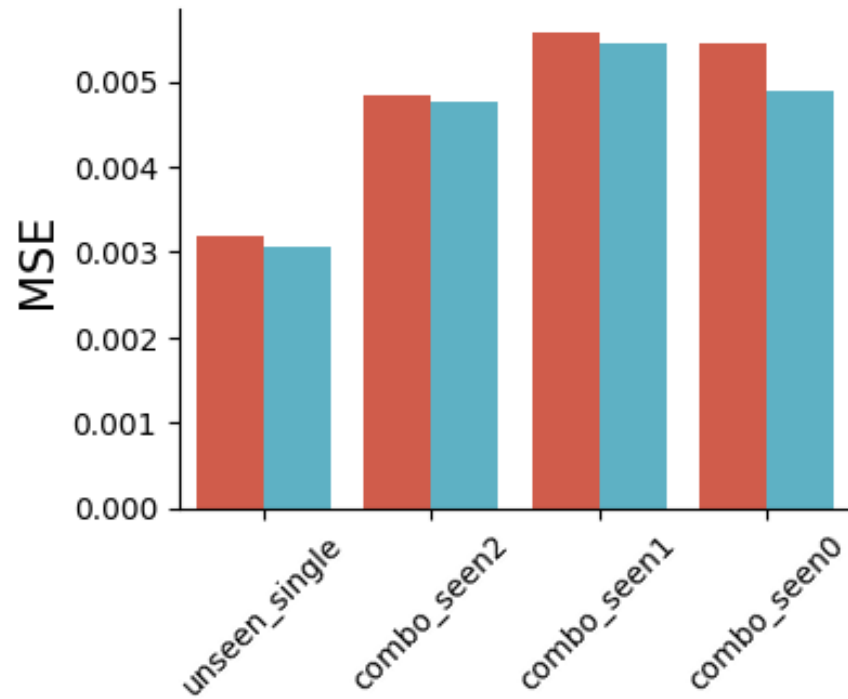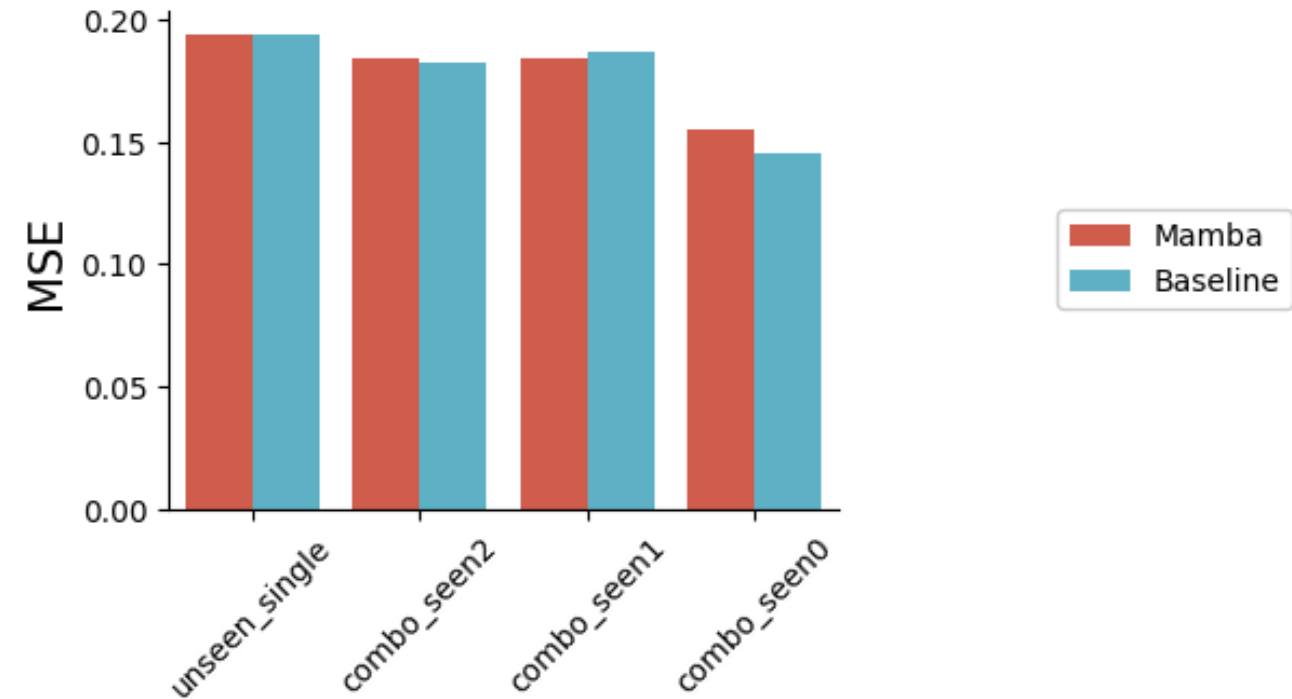
(2 unseen of 2)

```
'AHR+KLF1',
'CEBPE+CNN1',
'CEBPE+KLF1',
'CNN1+MAPK1',
'ETS2+CEBPE',
'ETS2+CNN1',
'ETS2+MAPK1',
'FEV+ISL2',
'FOSB+CEBPE',
```

Two-gene perturbation dataset: Norman et al.[9]

Since there was no single-cell-level ground truth in the perturbed data, we computed the averaged mean square error (MSE) of the top 20 differentially expressed (DE) genes <mark>between pre- and post-gene expression profiles</mark> for evaluation
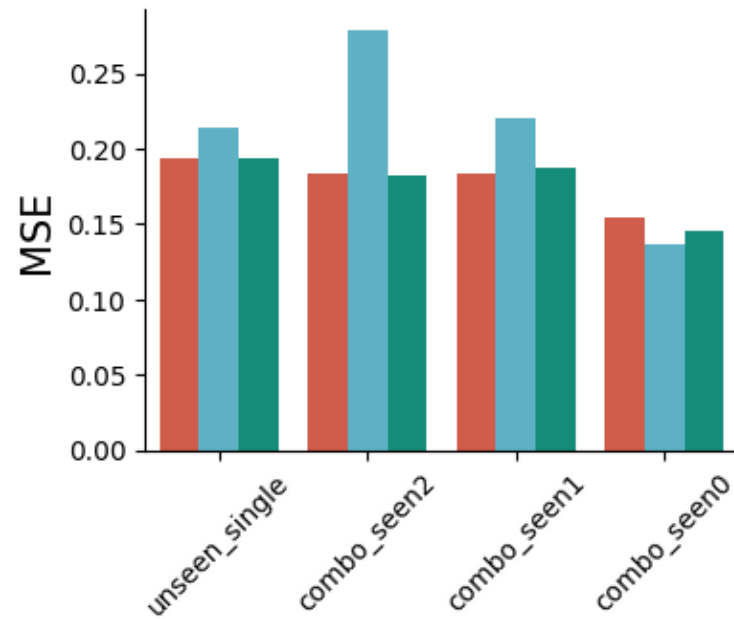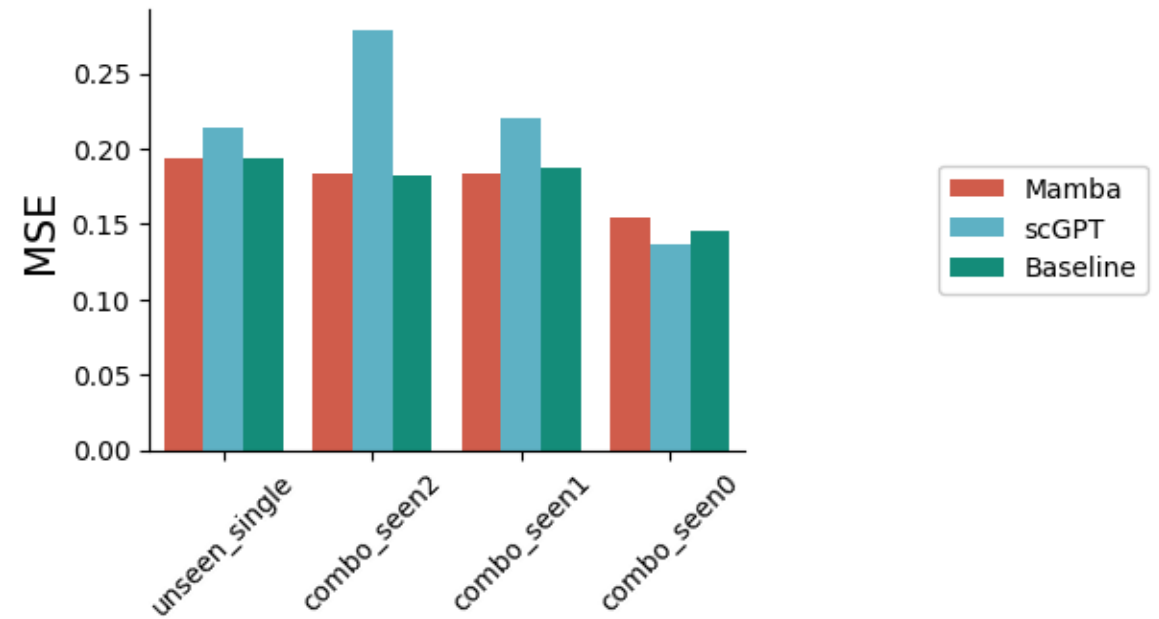


All Genes

Top 20 DE genes

Since there was no single-cell-level ground truth in the perturbed data, we computed the averaged mean square error (MSE) of the top 20 differentially expressed (DE) genes ==between pre- and post-gene expression profiles== for evaluation
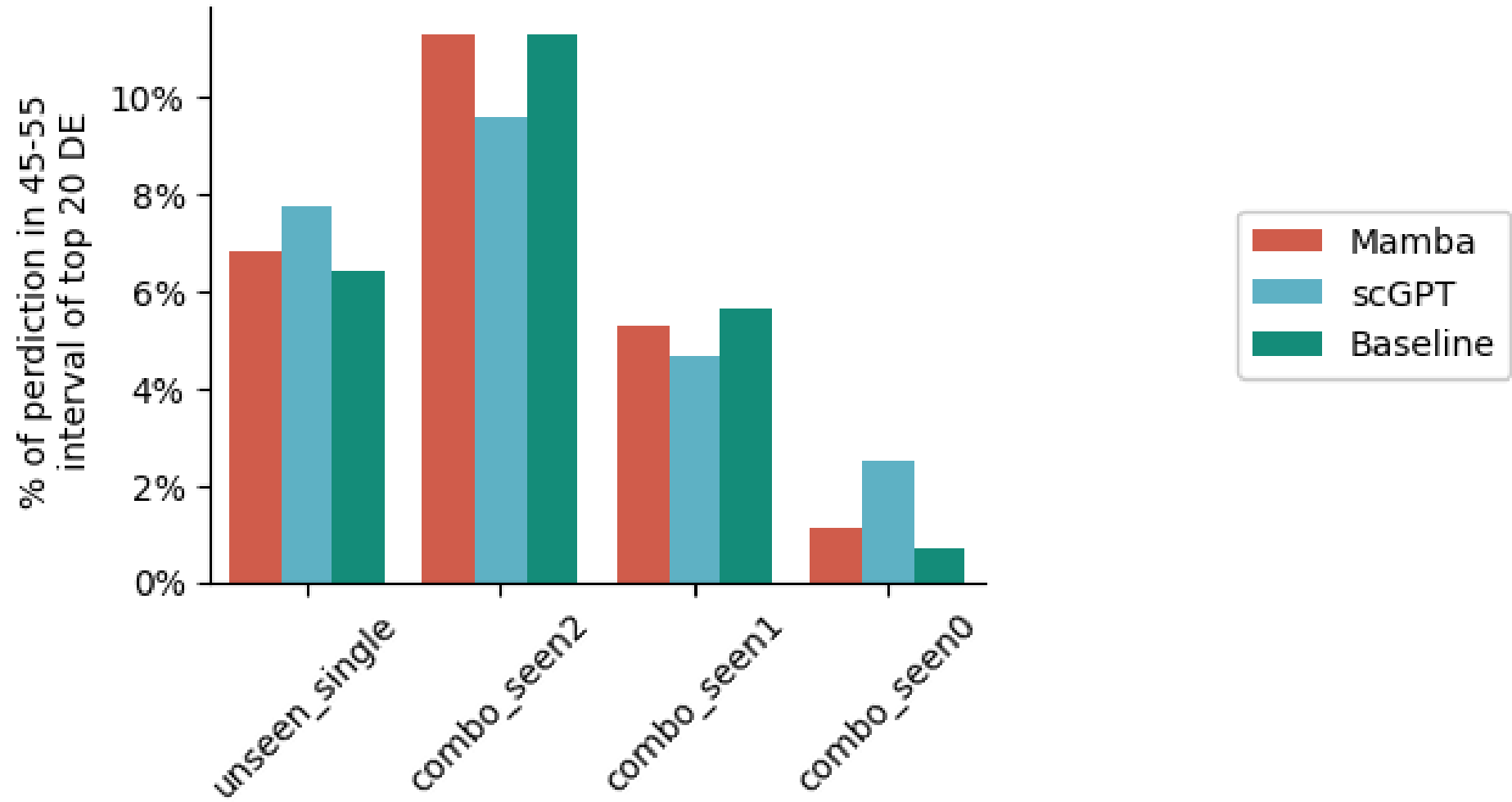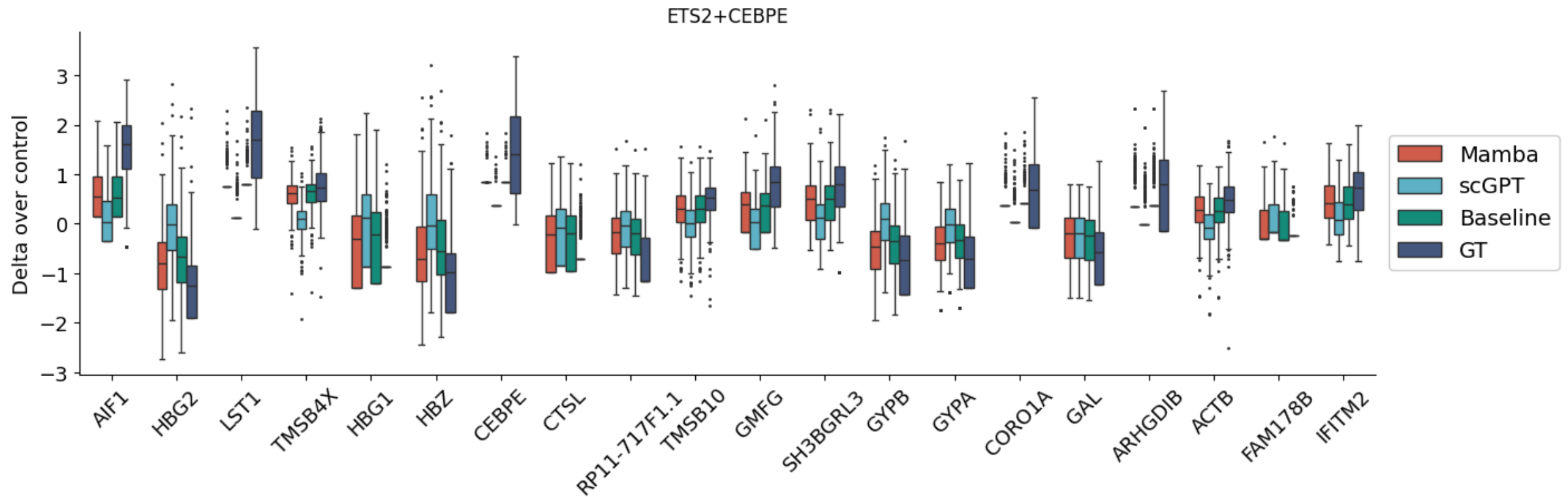


All Genes

Top 20 DE genes

For each two-gene perturbation in the test set, we further examined the proportion of the top 20 DE genes with mean predicted values falling in the 45–55% quantile of the true expression distribution interval.
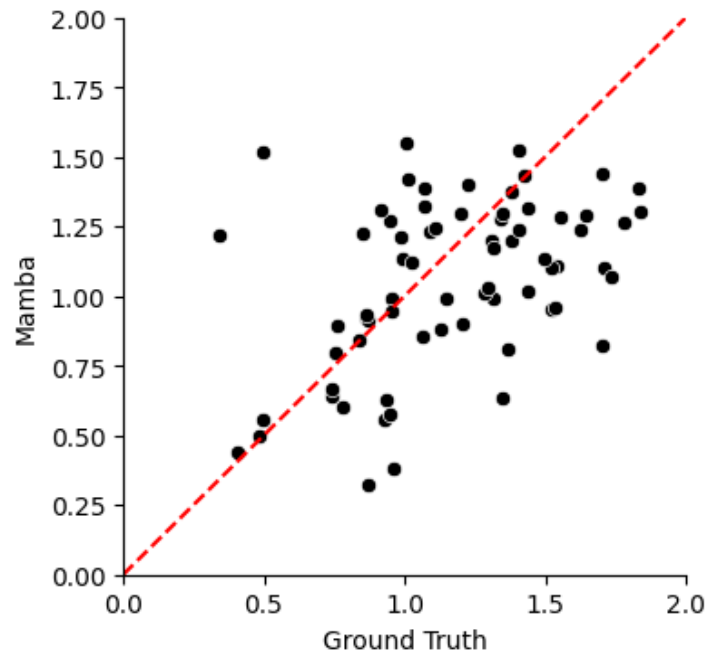
The predicted gene expression over control for the top 20 most DE genes after a combinatorial perturbation (ETS2 + CEBPE). The red and blue boxes indicate gene prediction results by the Mamba-based GEARS model and the baseline GEARS model, respectively. The purple boxes represent the ground truth post-gene distribution.
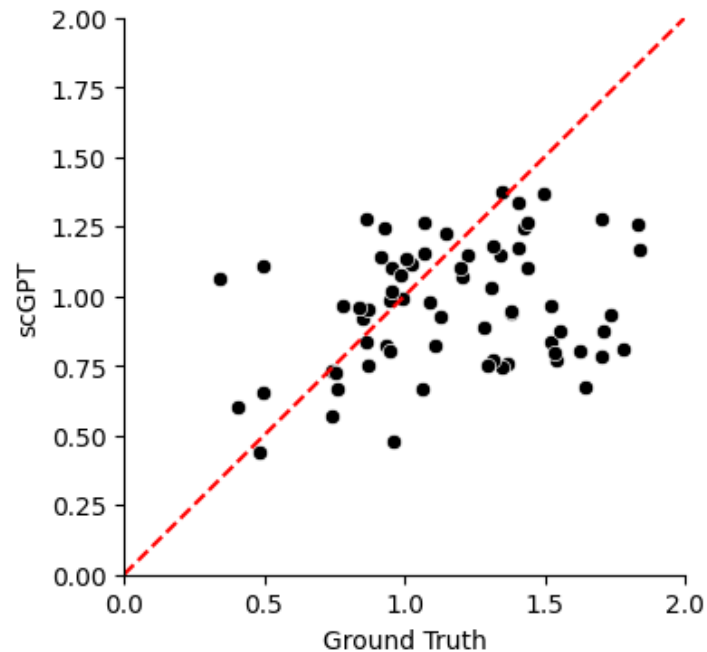
Magnitude scores computed for all test perturbing combinations on the Norman dataset. Each dot represents a specific perturbing combination. The y axis shows the magnitude score computed from the prediction results, while the x axis represents the ground truth magnitude score computed using real post-gene expression.
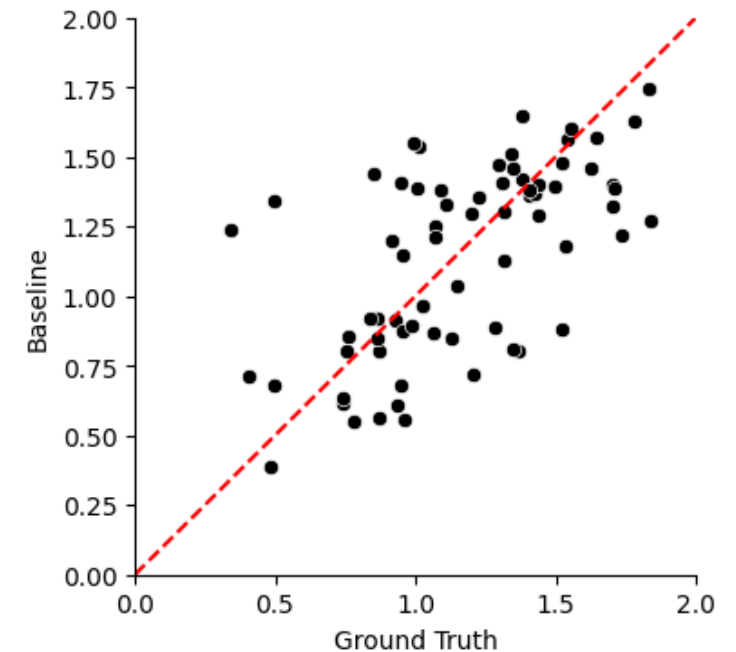
**Magnitude scores**



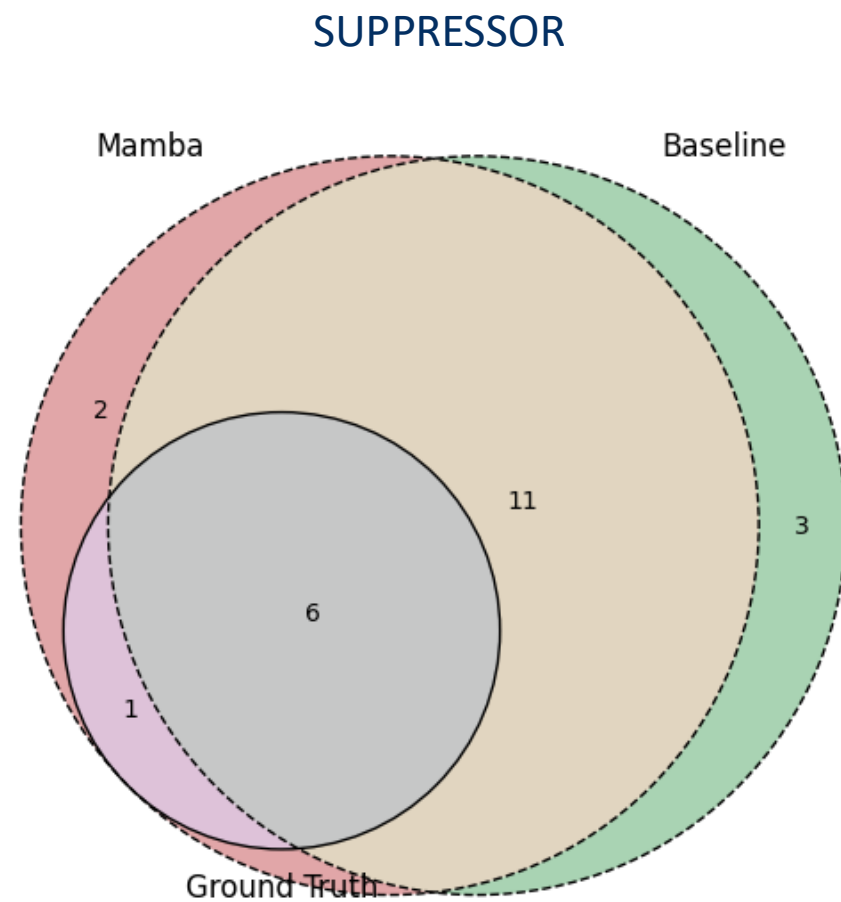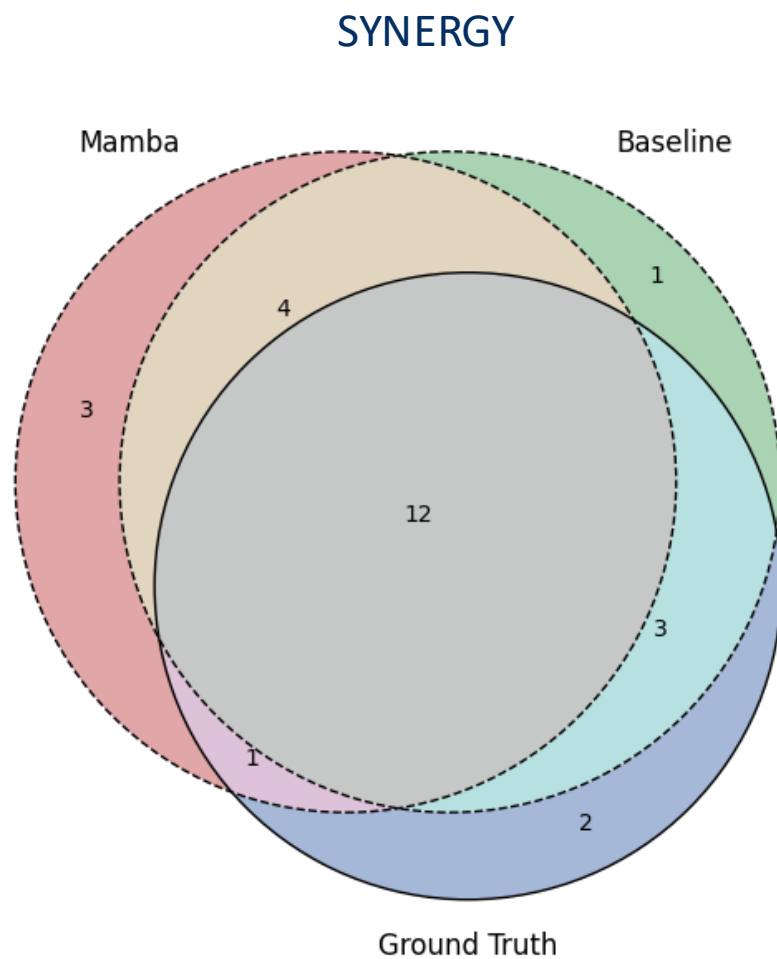PearsonRResult(statistic=0.411691900021 76964, pvalue=0.000361129511 6439501)

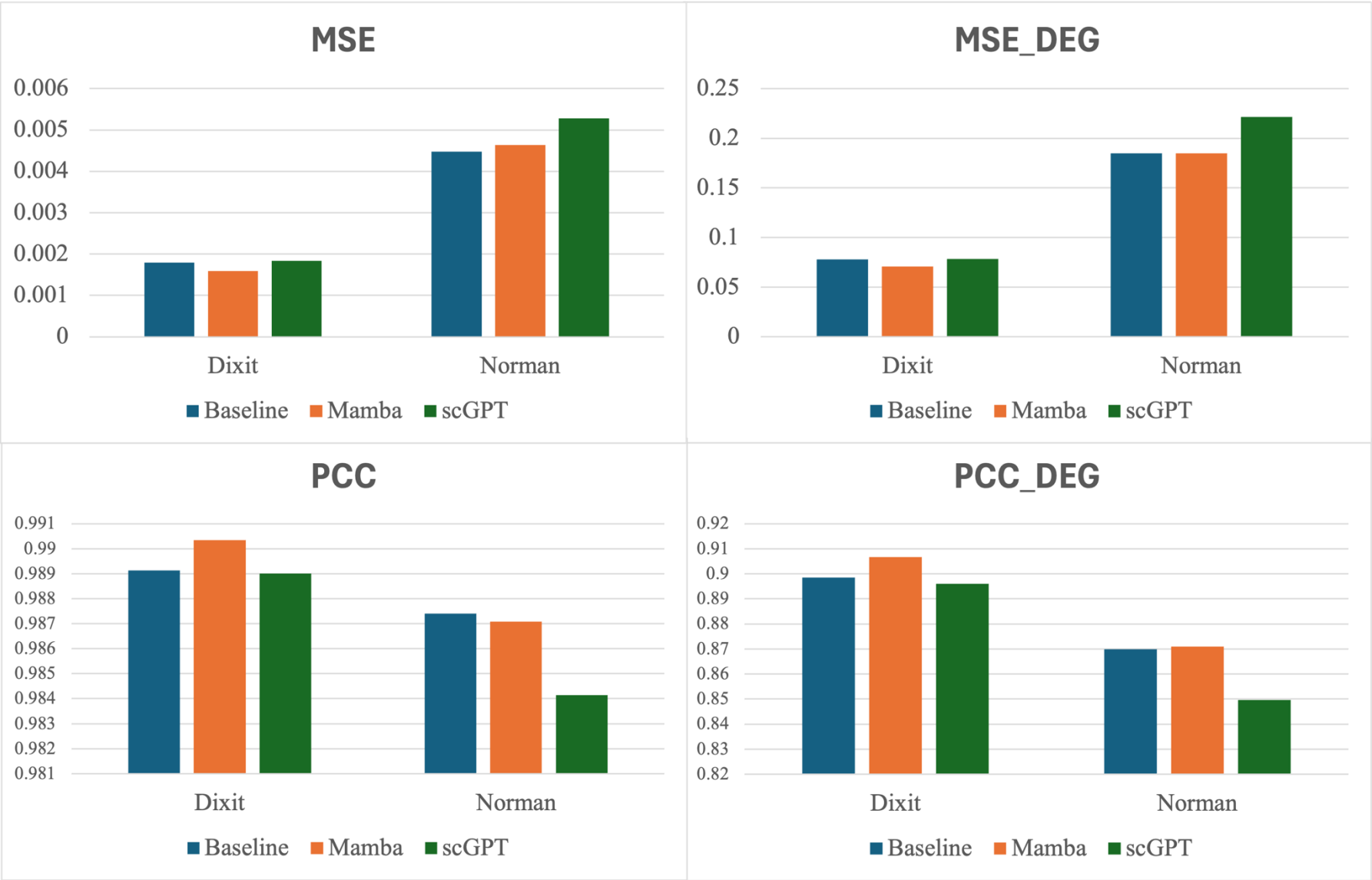PearsonRResult(statistic=0.221630443 1984232, pvalue=0.0632376566721118)

PearsonRResult(statistic=0.56855930073 20305, pvalue=2.3106232323576483e-07)

Top 20 perturbations with synergistic and suppressor gene interaction types identified using Mamba and baseline methods. The Venn plot illustrates the relationship between the identified perturbation set and the verified perturbation set.

Scores of gene expression prediction using perturb-seq datasets based on GEARS

# scELMo: Embeddings from Language Models are Good Learners for Single-cell Data Analysis