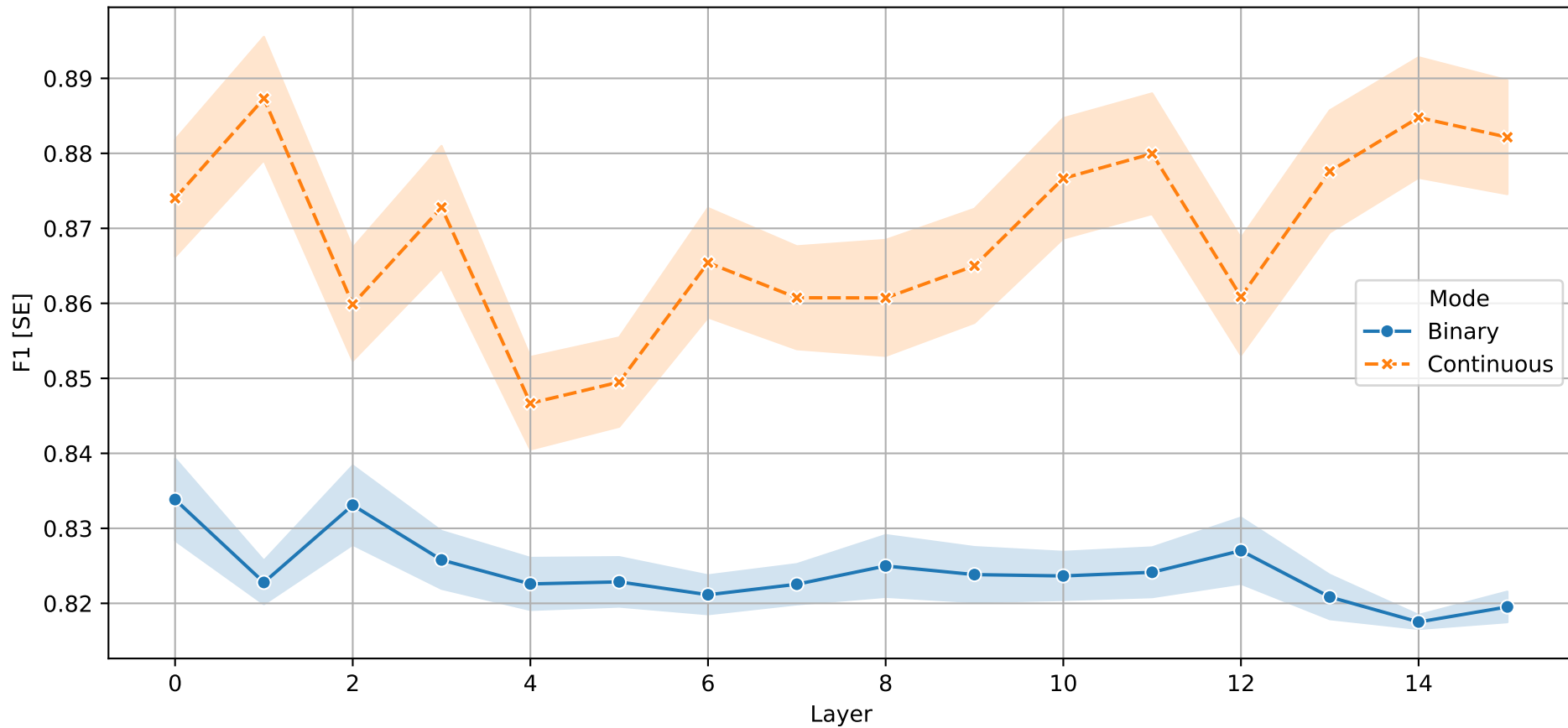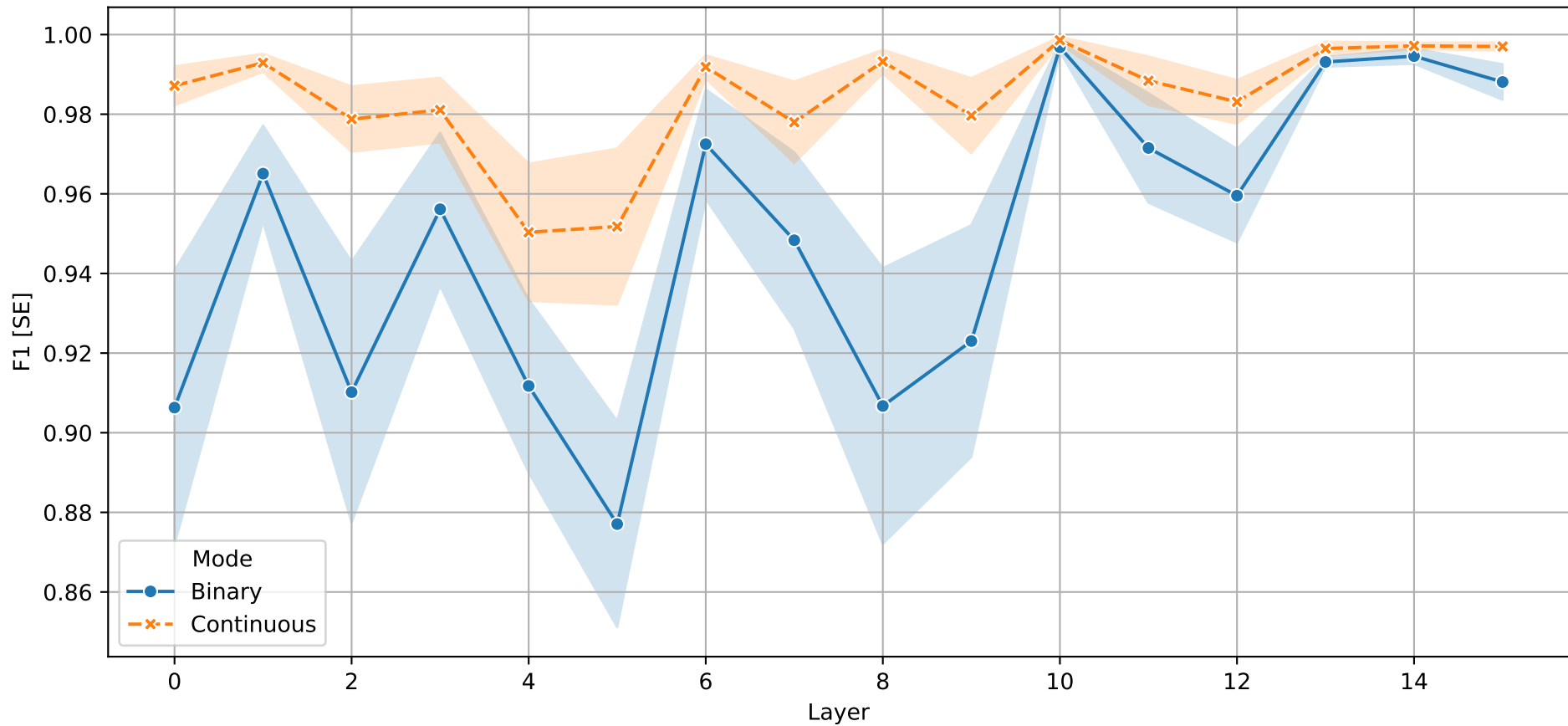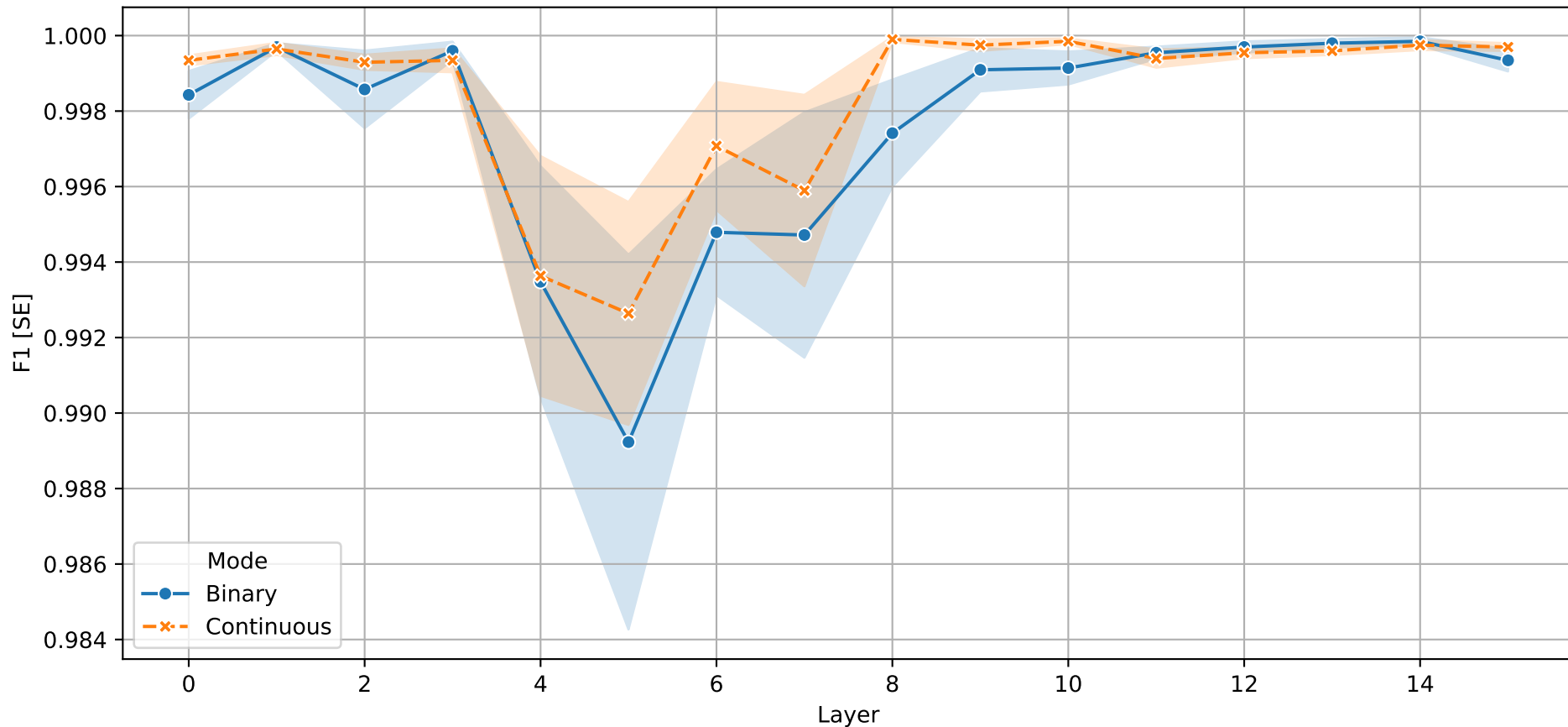F1 per Layer – Single Neuron Probing
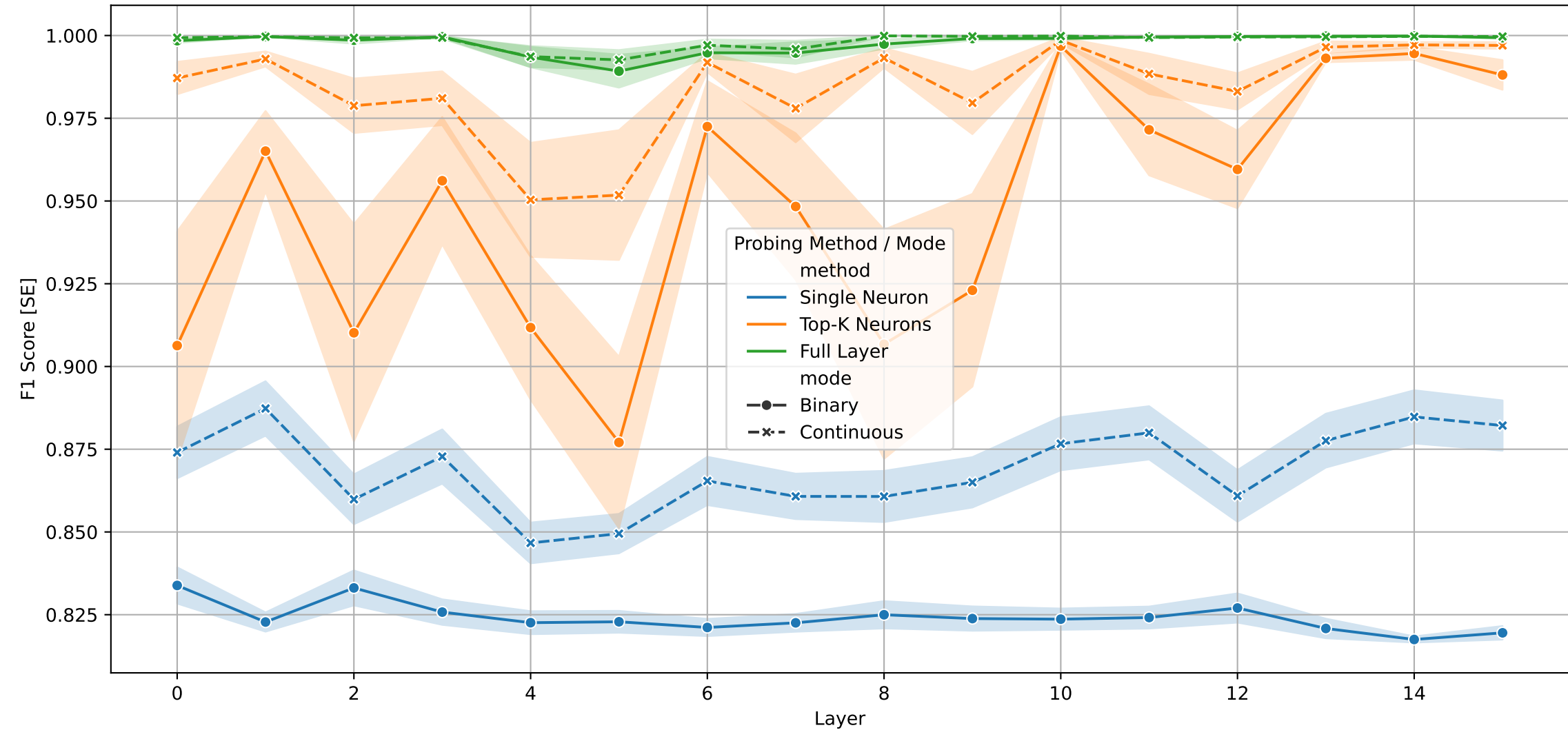
F1 per Layer – Top-K Neurons Probing

F1 per Layer – Full Layer Probing

Overall F1 per Layer – All Methods

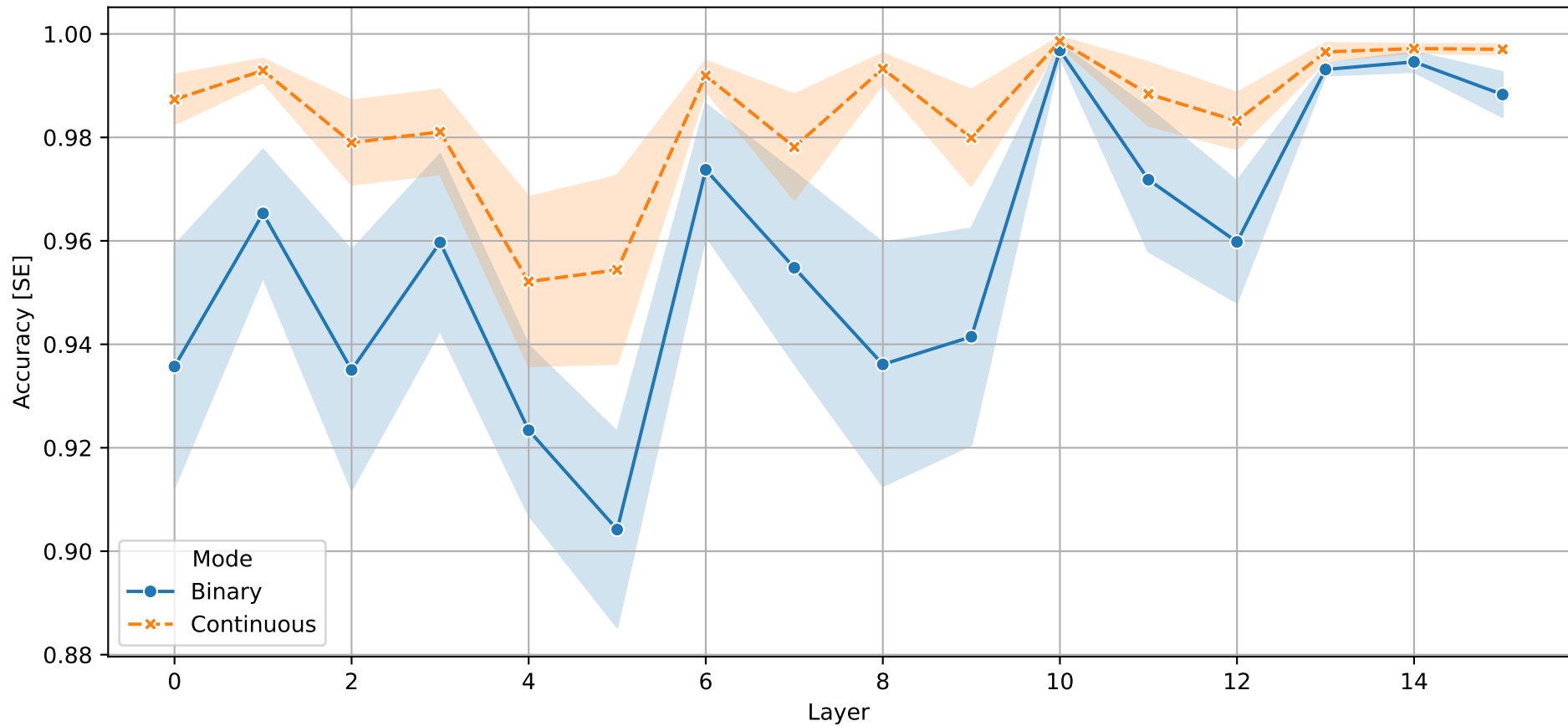## F1 Score Summary by Probing Method

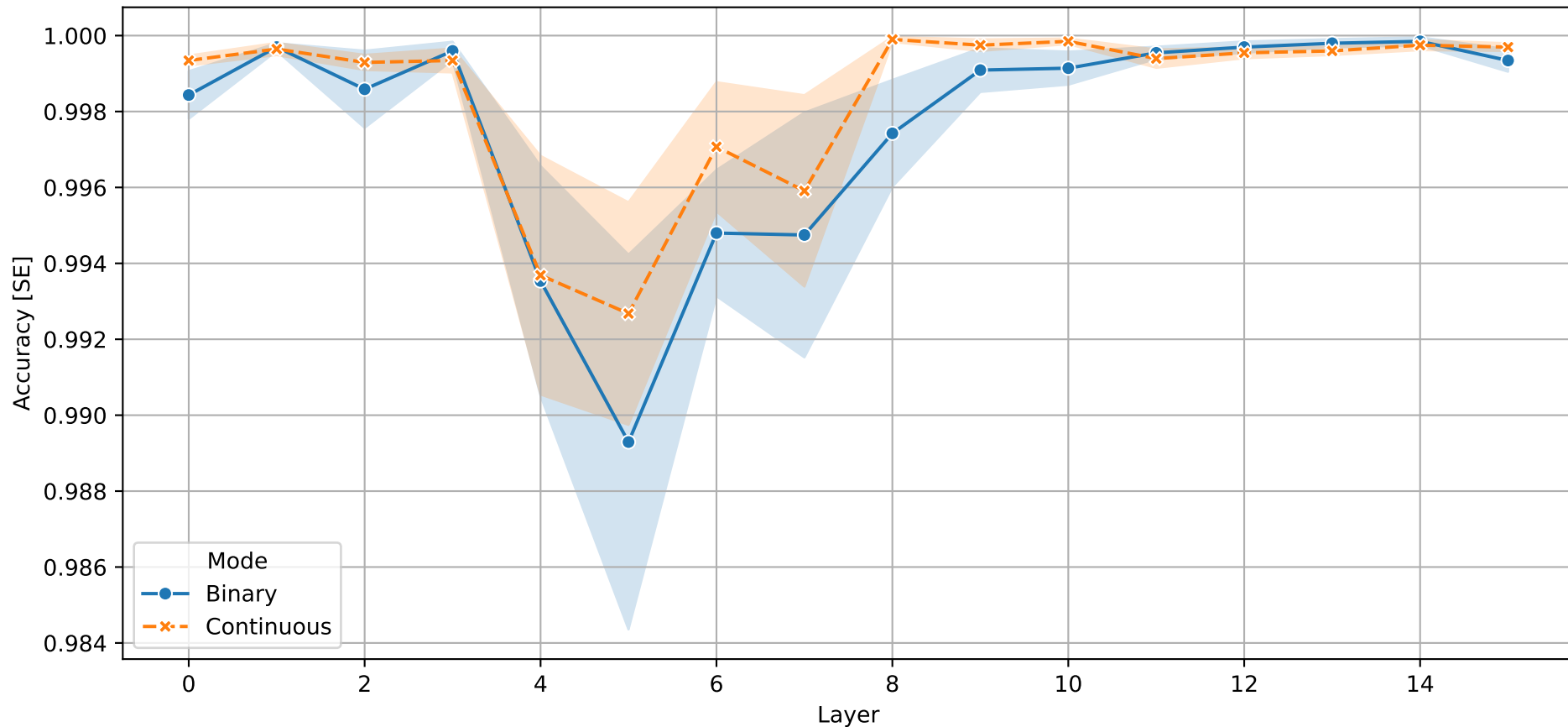| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | f1_best_layer | 14.0 | 8.0 |
| Full Layer | f1_max | 1.0 | 1.0 |
| Full Layer | f1_mean | 0.9976 | 0.9984 |
| Full Layer | f1_std | 0.0056 | 0.0042 |
| Single Neuron | f1_best_layer | 0.0 | 1.0 |
| Single Neuron | f1_max | 1.0 | 1.0 |
| Single Neuron | f1_mean | 0.8241 | 0.869 |
| Single Neuron | f1_std | 0.0319 | 0.0684 |
| Top-K Neurons | f1_best_layer | 10.0 | 10.0 |
| Top-K Neurons | f1_max | 1.0 | 1.0 |
| Top-K Neurons | f1_mean | 0.9488 | 0.9841 |
| Top-K Neurons | f1_std | 0.0667 | 0.0265 |

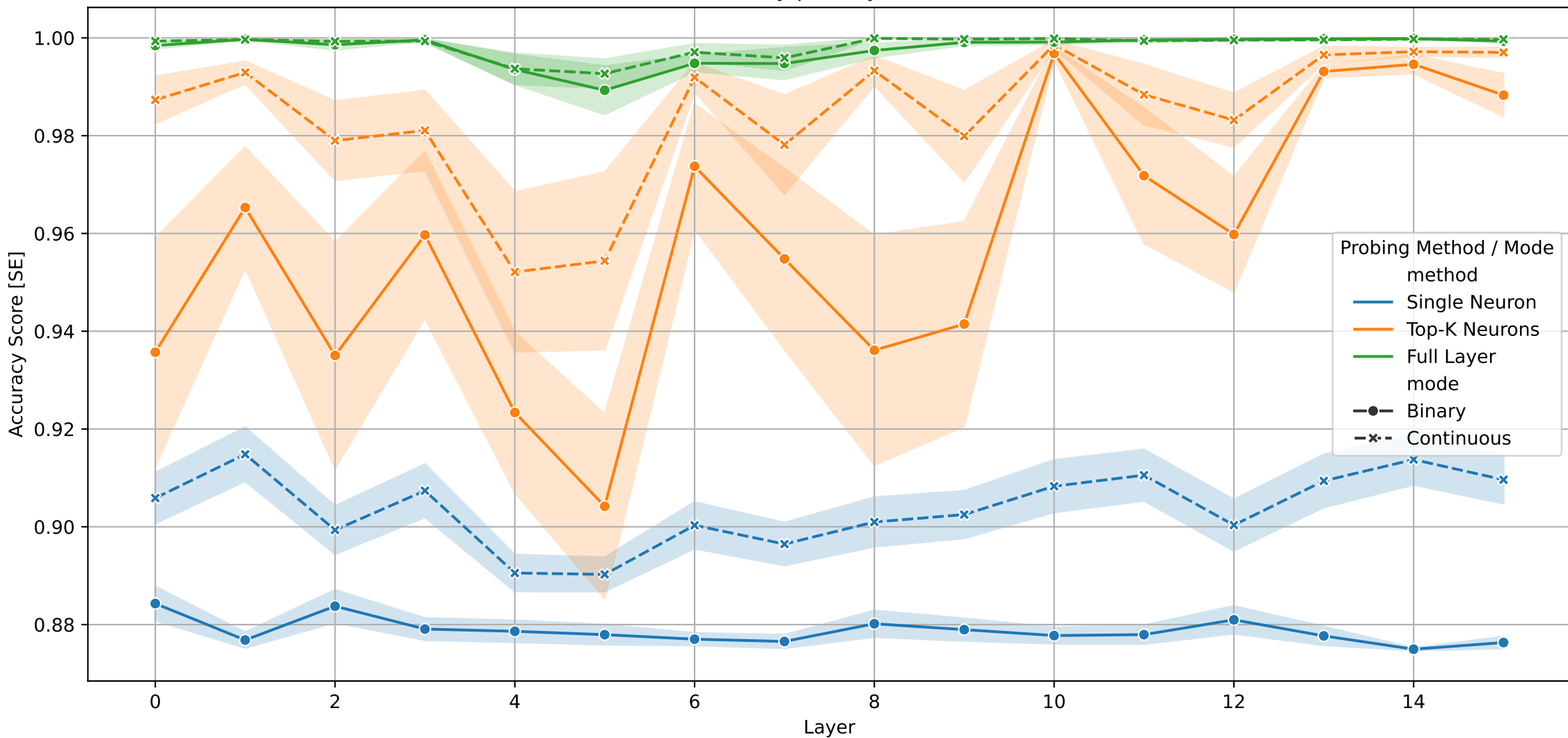Accuracy per Layer – Single Neuron Probing

Accuracy per Layer – Top-K Neurons Probing

Accuracy per Layer – Full Layer Probing

Overall Accuracy per Layer – All Methods

## Accuracy Score Summary by Probing Method

| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | accuracy_best_layer | 14.0 | 8.0 |
| Full Layer | accuracy_max | 1.0 | 1.0 |
| Full Layer | accuracy_mean | 0.9977 | 0.9984 |
| Full Layer | accuracy_std | 0.0056 | 0.0042 |
| Single Neuron | accuracy_best_layer | 0.0 | 1.0 |
| Single Neuron | accuracy_max | 1.0 | 1.0 |
| Single Neuron | accuracy_mean | 0.8787 | 0.9038 |
| Single Neuron | accuracy_std | 0.0195 | 0.0446 |
| Top-K Neurons | accuracy_best_layer | 10.0 | 10.0 |
| Top-K Neurons | accuracy_max | 1.0 | 1.0 |
| Top-K Neurons | accuracy_mean | 0.9584 | 0.9844 |
| Top-K Neurons | accuracy_std | 0.0498 | 0.0253 |