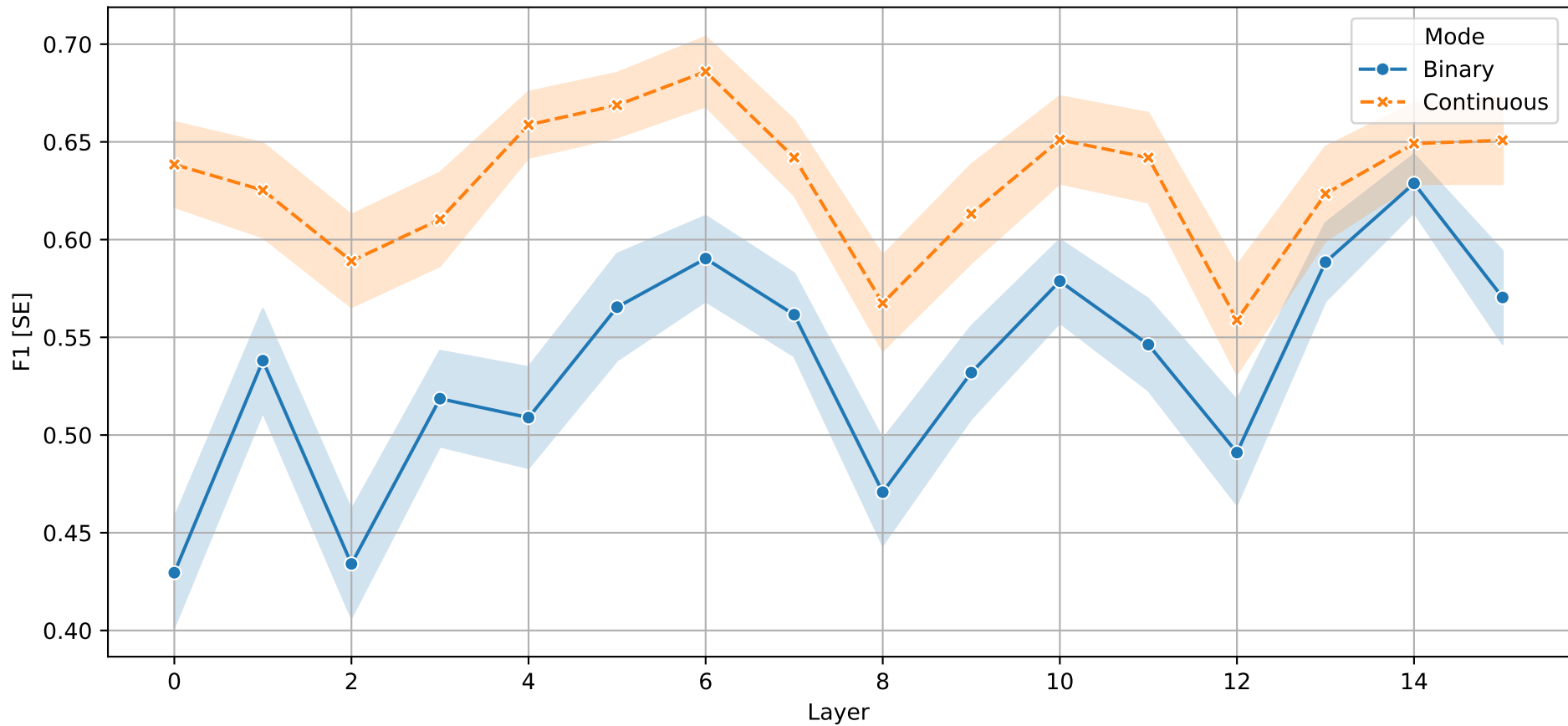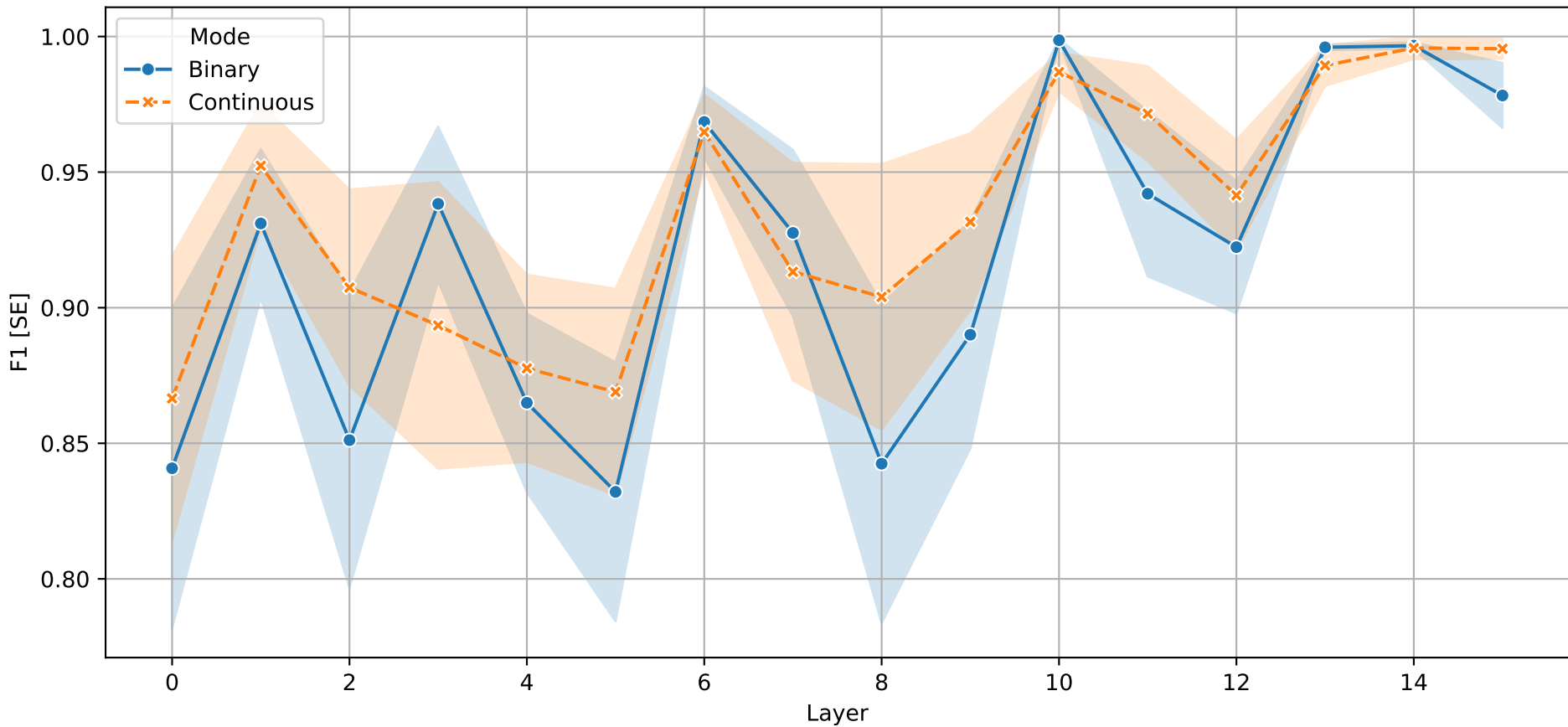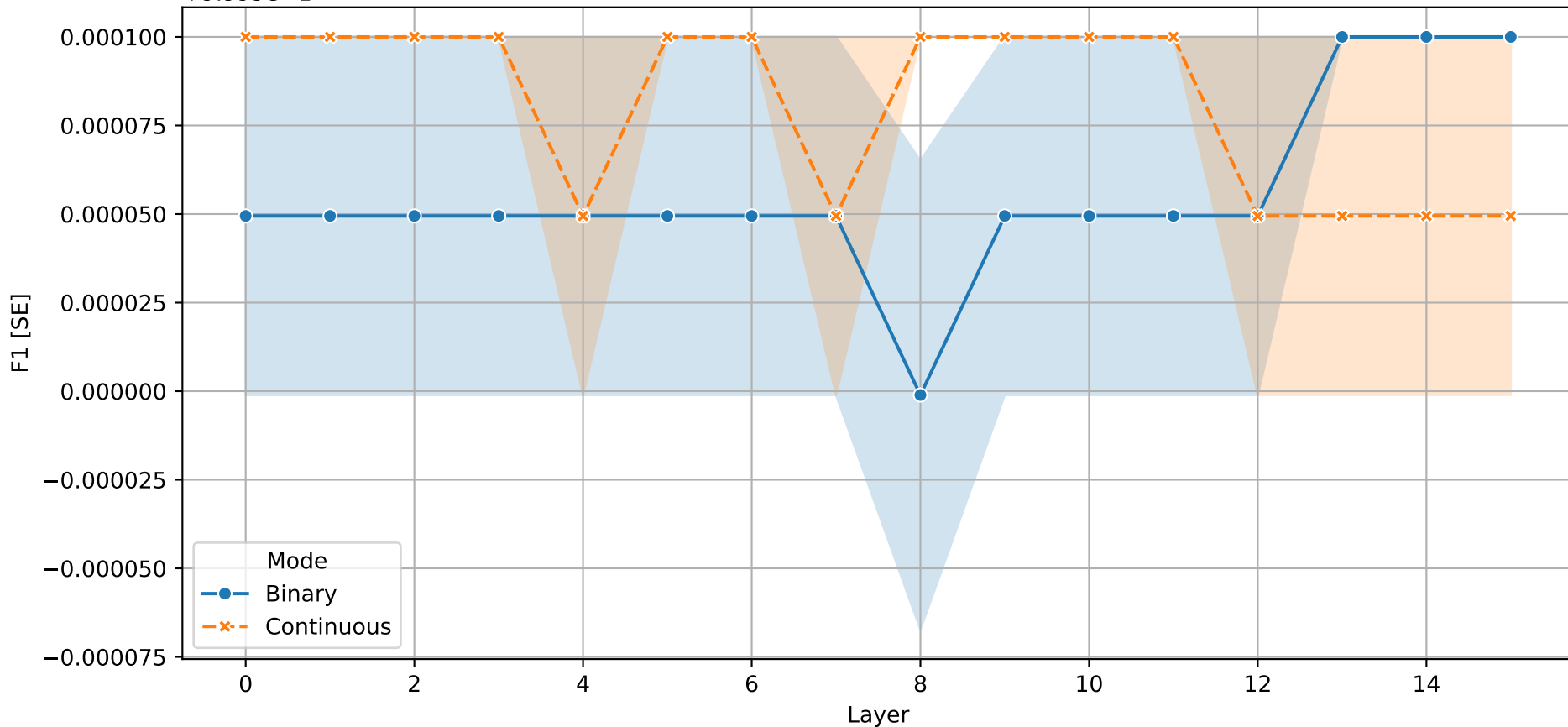F1 per Layer – Single Neuron Probing
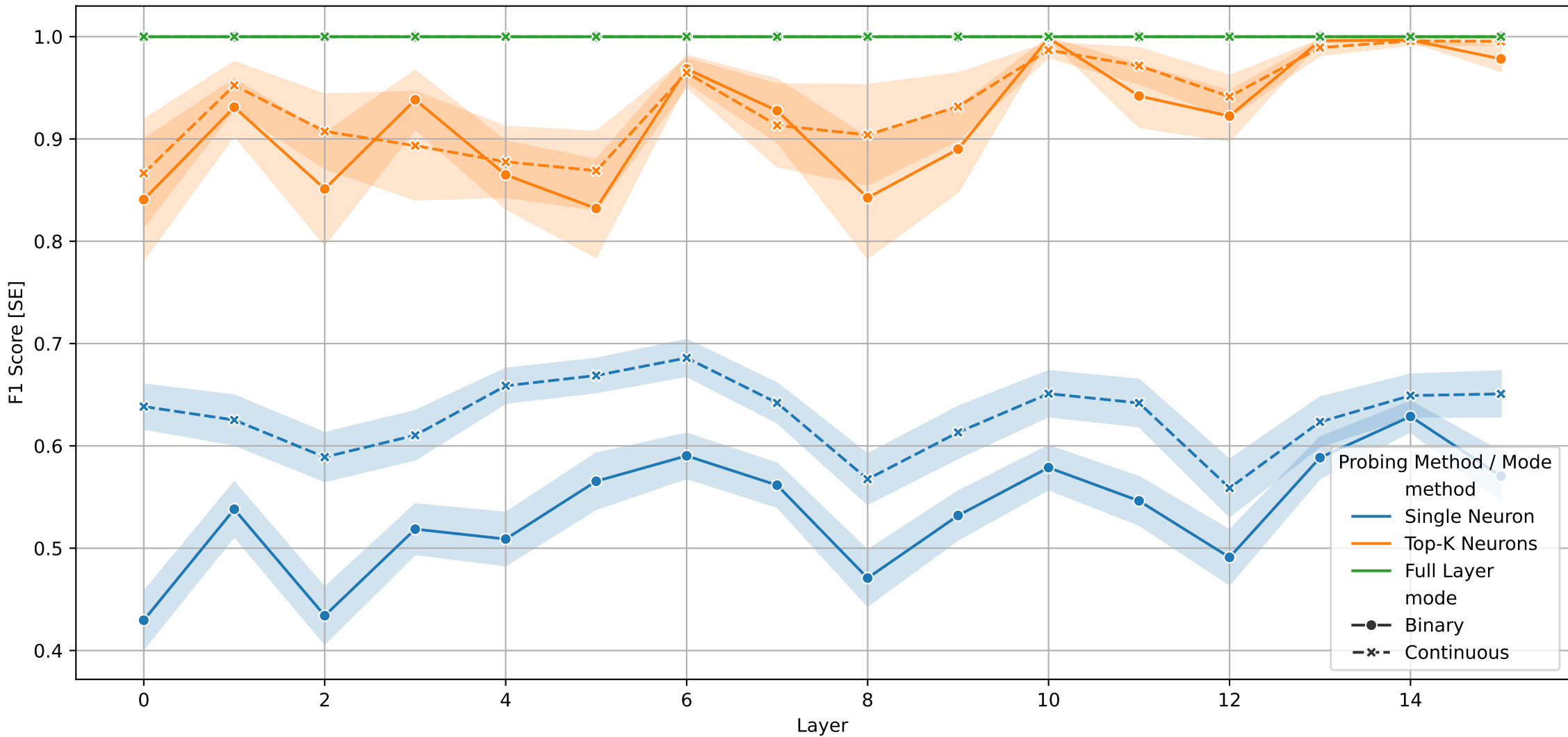
F1 per Layer – Top-K Neurons Probing

F1 per Layer – Full Layer Probing

Overall F1 per Layer – All Methods

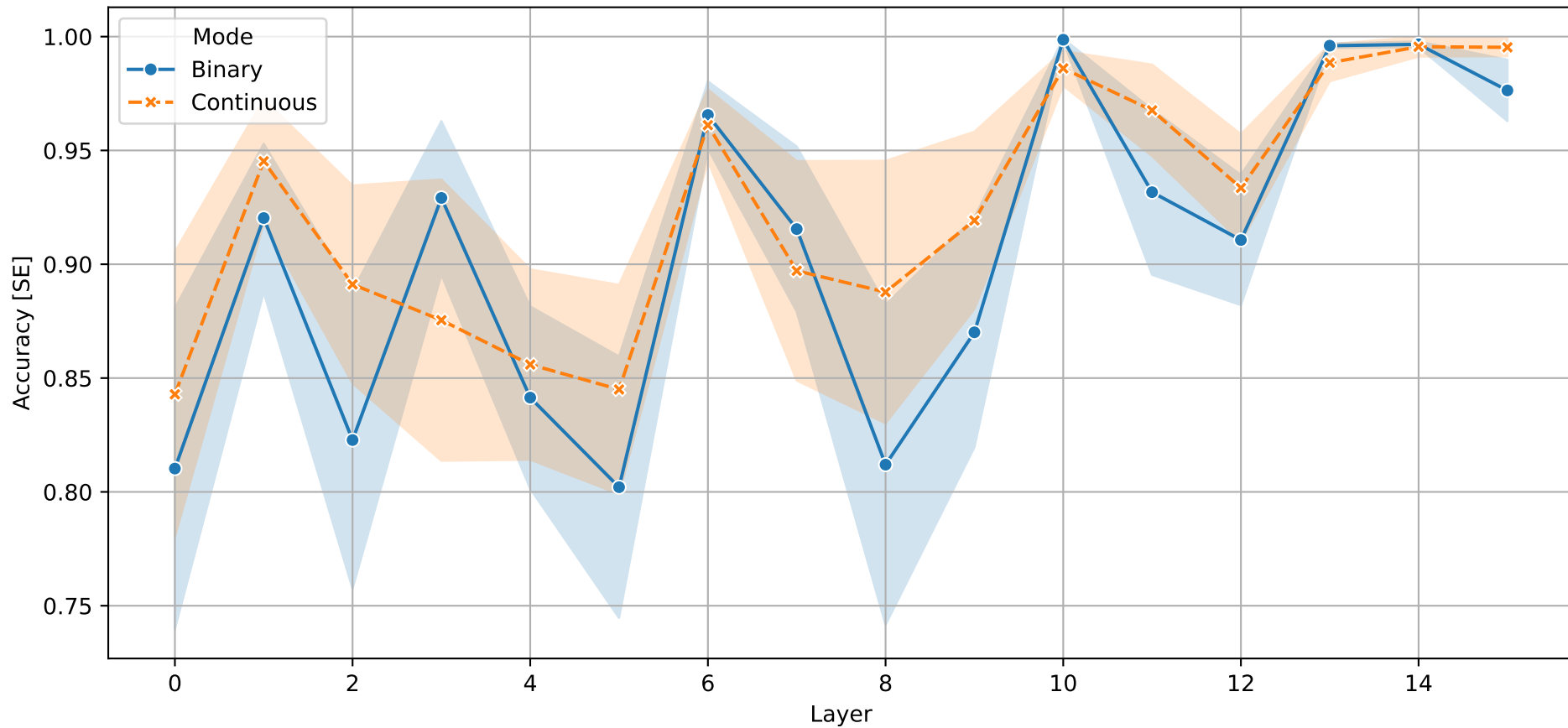## F1 Score Summary by Probing Method

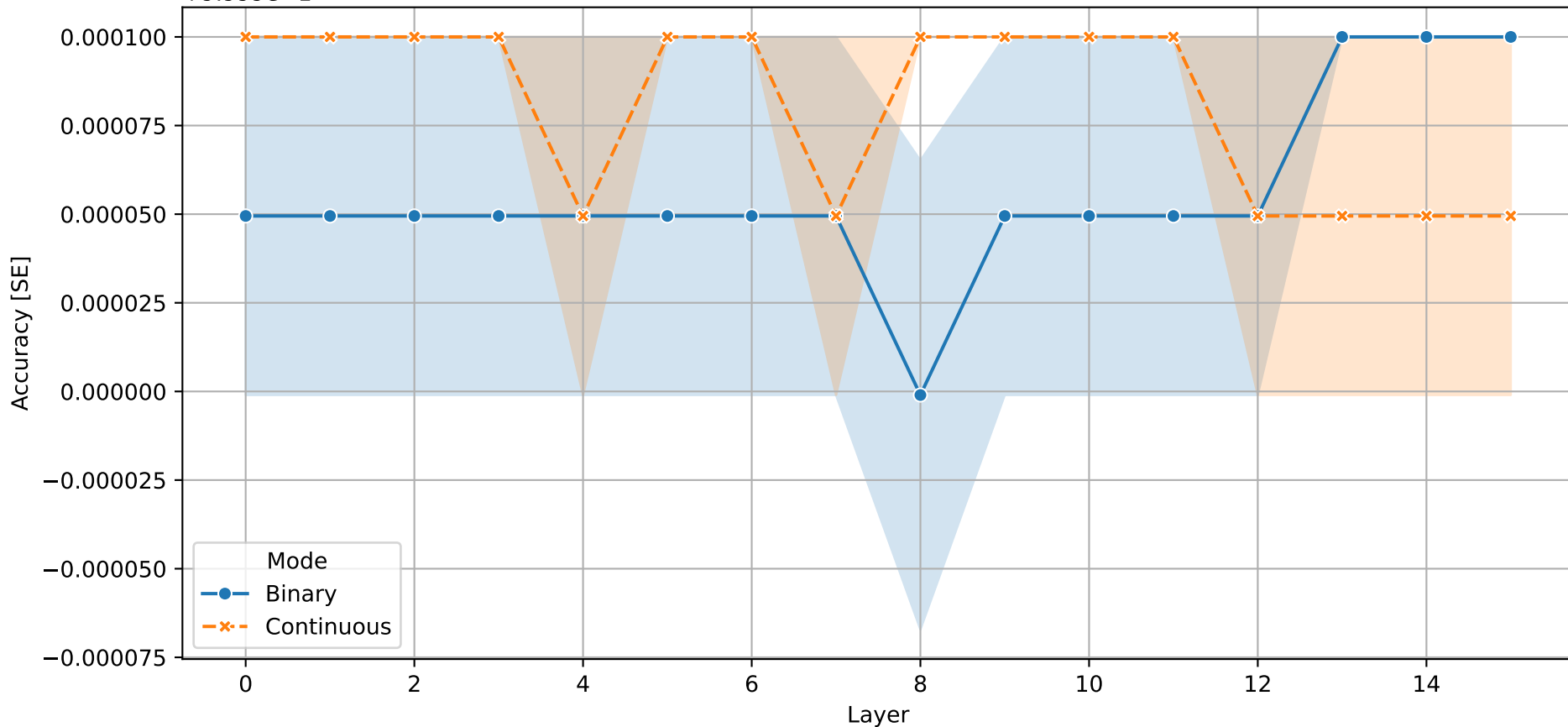| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | f1_best_layer | 13.0 | 0.0 |
| Full Layer | f1_max | 1.0 | 1.0 |
| Full Layer | f1_mean | 1.0 | 1.0 |
| Full Layer | f1_std | 0.0001 | 0.0001 |
| Single Neuron | f1_best_layer | 14.0 | 6.0 |
| Single Neuron | f1_max | 1.0 | 1.0 |
| Single Neuron | f1_mean | 0.5346 | 0.6297 |
| Single Neuron | f1_std | 0.2218 | 0.201 |
| Top-K Neurons | f1_best_layer | 10.0 | 14.0 |
| Top-K Neurons | f1_max | 1.0 | 1.0 |
| Top-K Neurons | f1_mean | 0.9201 | 0.935 |
| Top-K Neurons | f1_std | 0.109 | 0.0955 |

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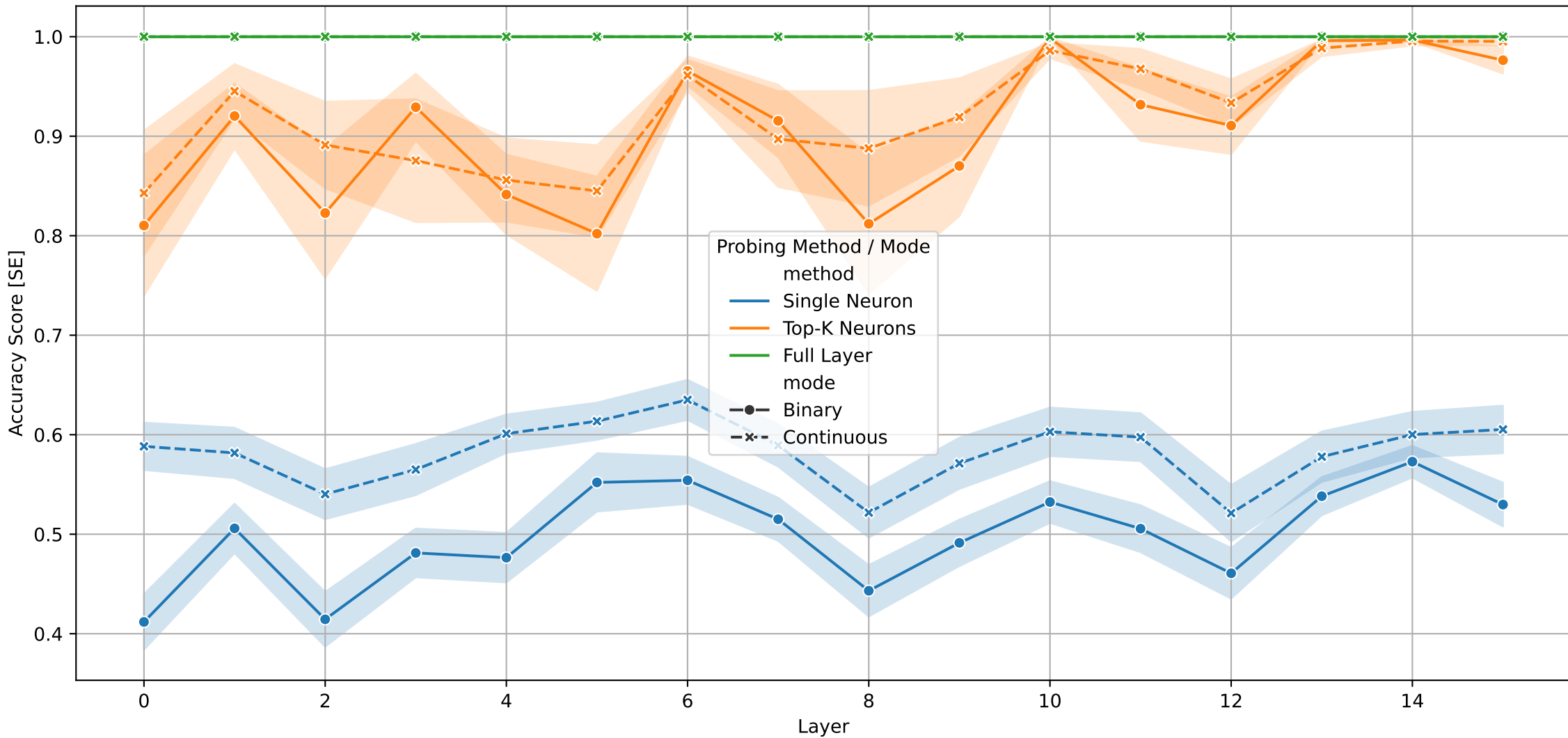
Accuracy per Layer – Top-K Neurons Probing

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 | 13.0 | 0.0 |
| Full Layer | accuracy_max | 1.0 | 1.0 |
| Full Layer | accuracy_mean | 1.0 | 1.0 |
| Full Layer | accuracy_std | 0.0001 | 0.0001 |
| Single Neuron | accuracy_best_layer | 14.0 | 6.0 |
| Single Neuron | accuracy_max | 1.0 | 1.0 |
| Single Neuron | accuracy_mean | 0.4991 | 0.582 |
| Single Neuron | accuracy_std | 0.2184 | 0.214 |
| Top-K Neurons | accuracy_best_layer | 10.0 | 14.0 |
| Top-K Neurons | accuracy_max | 1.0 | 1.0 |
| Top-K Neurons | accuracy_mean | 0.9062 | 0.9242 |
| Top-K Neurons | accuracy_std | 0.1303 | 0.1134 |