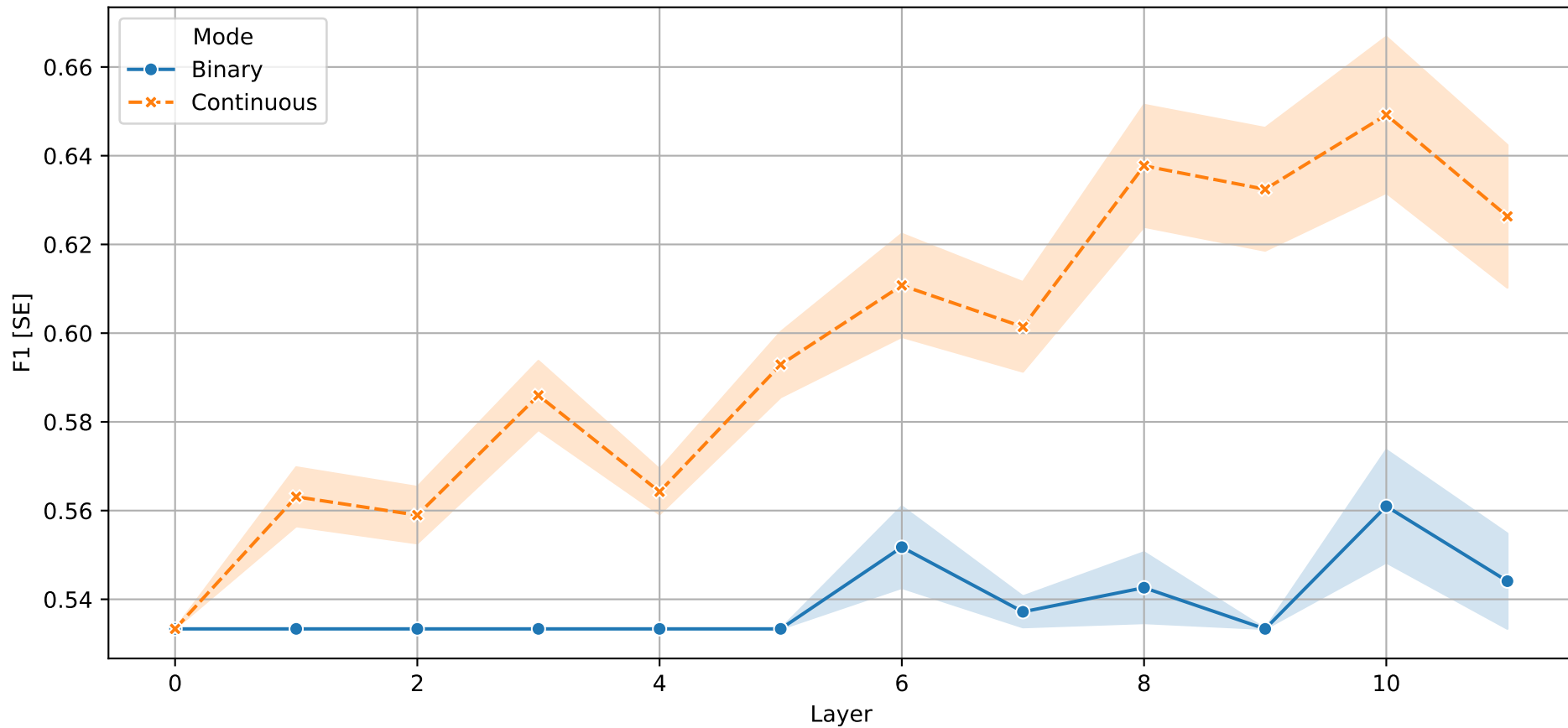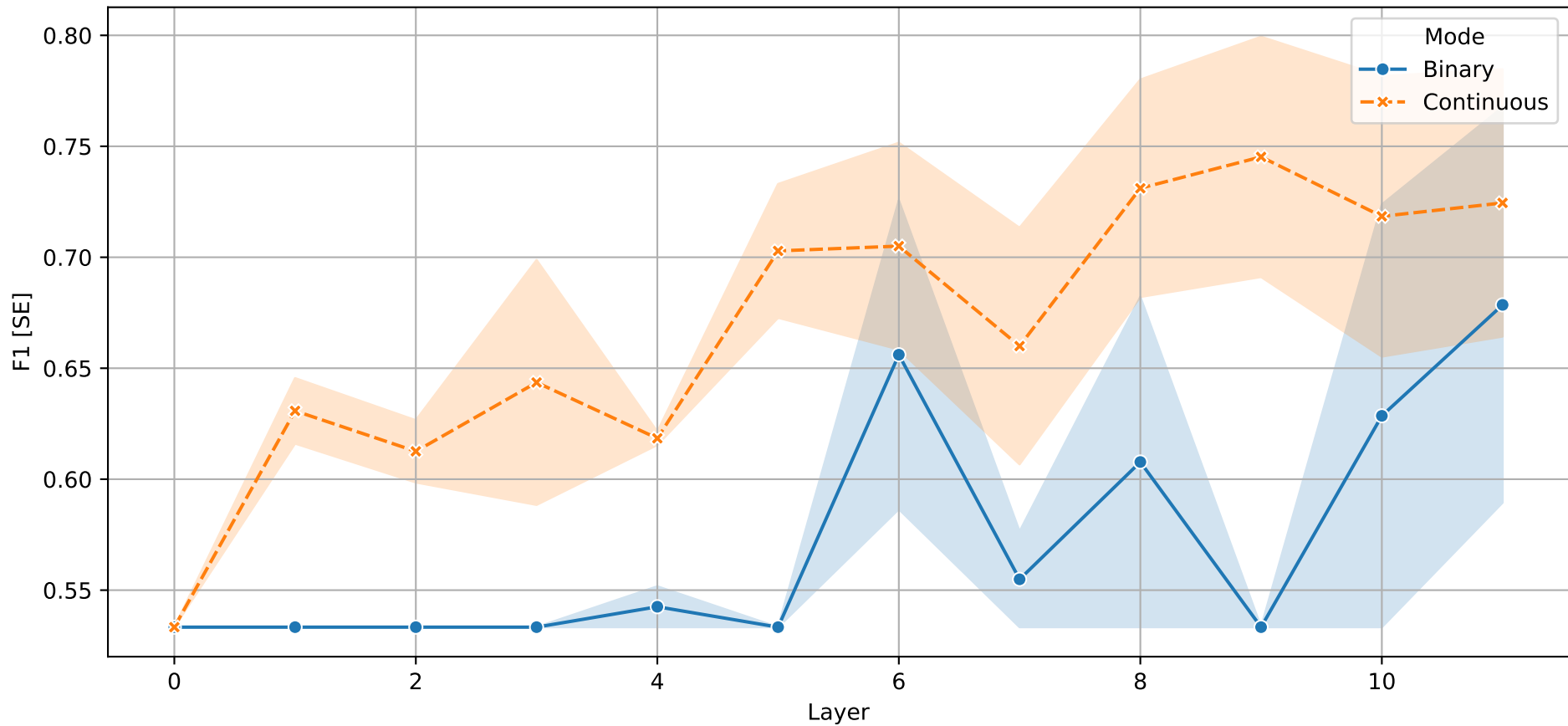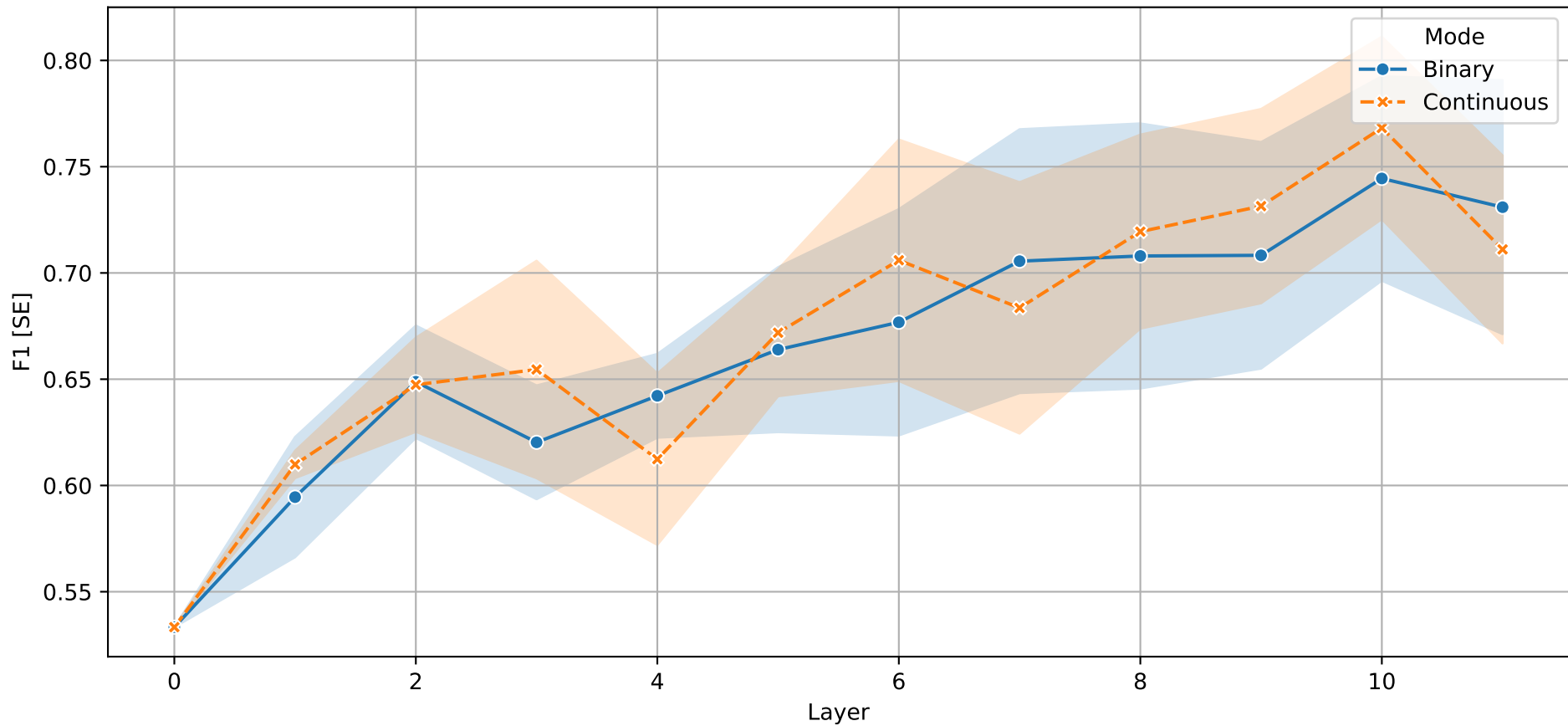F1 per Layer – Single Neuron Probing
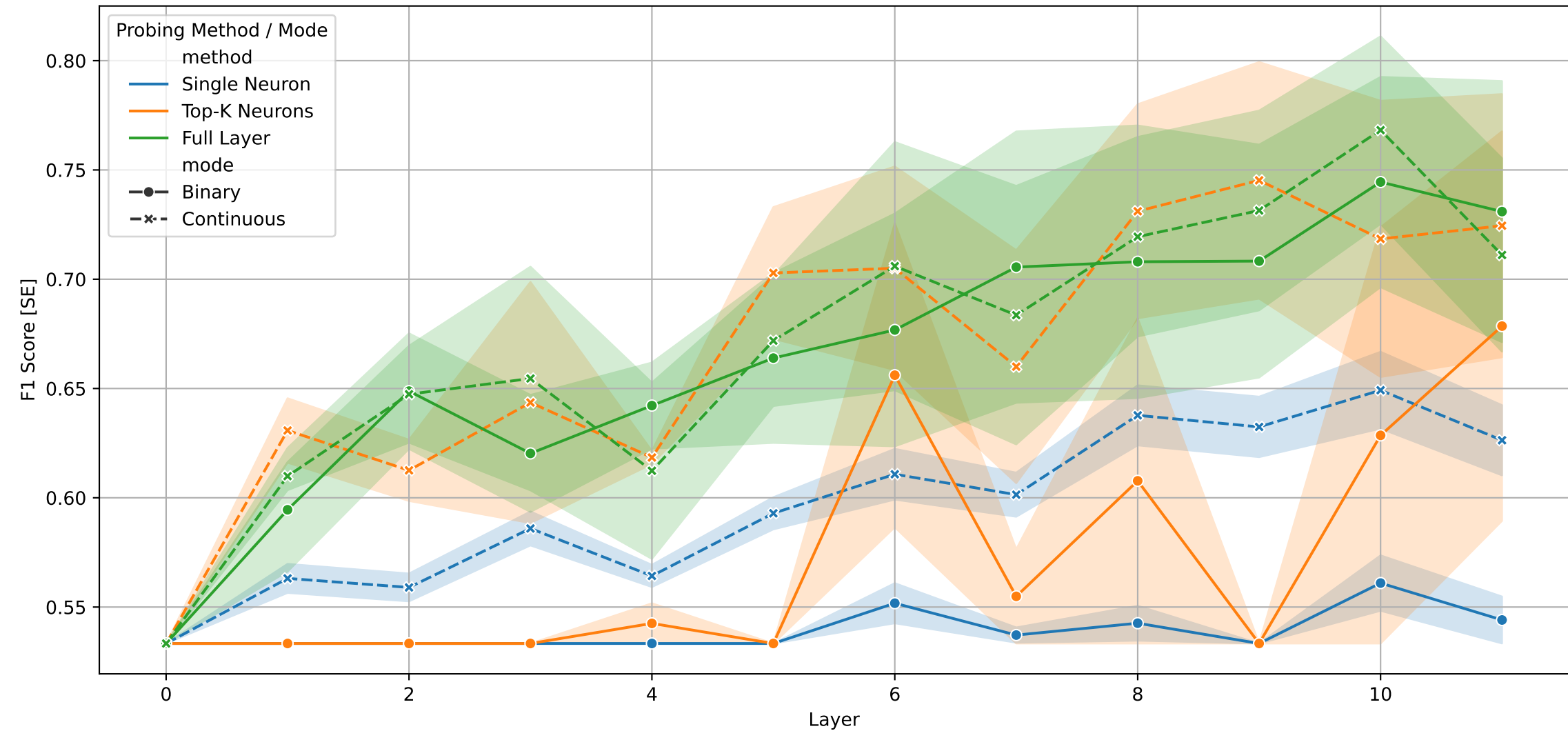
F1 per Layer – Top-K Neurons Probing

F1 per Layer – Full Layer Probing

Overall F1 per Layer – All Methods

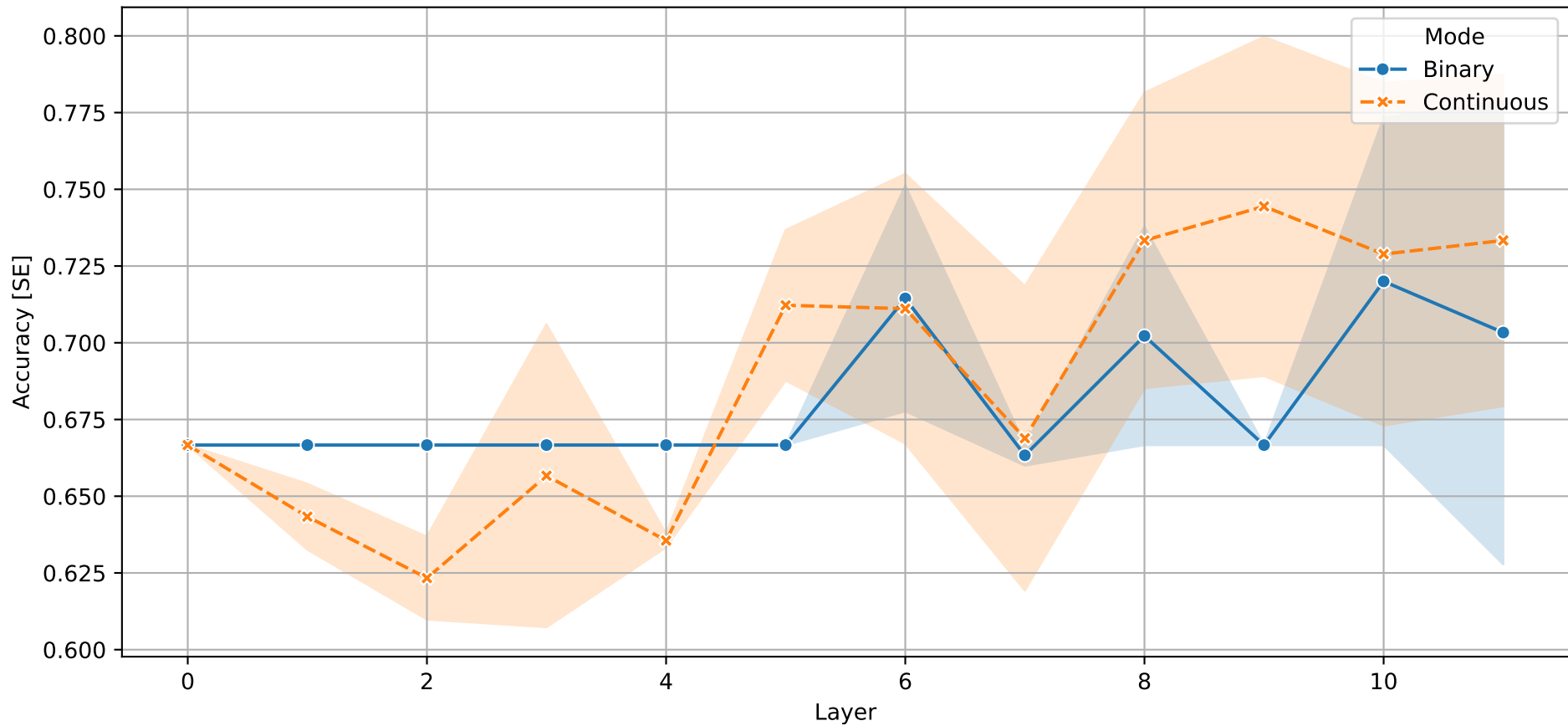## F1 Score Summary by Probing Method

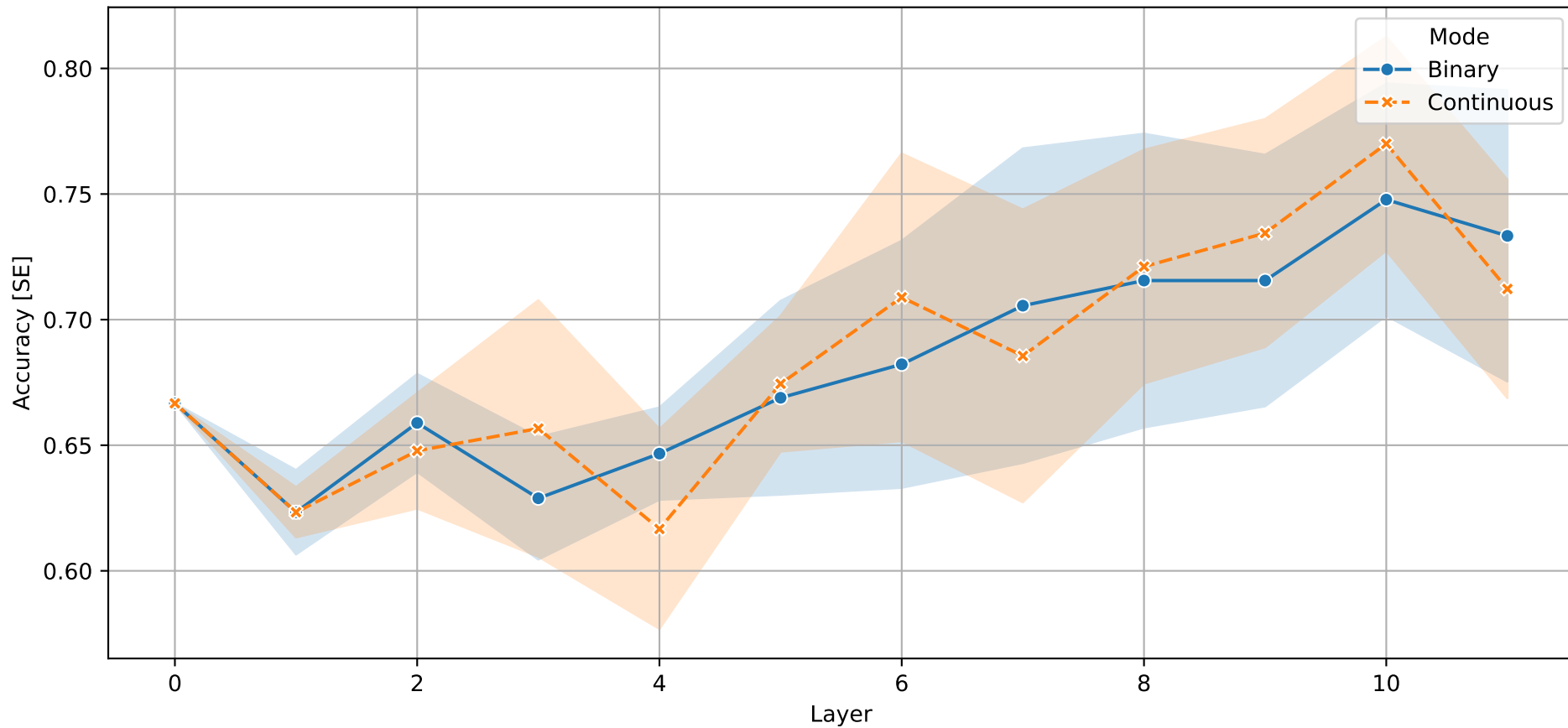| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | f1_best_layer | 10.0 | 10.0 |
| Full Layer | f1_max | 0.8431 | 0.8517 |
| Full Layer | f1_mean | 0.6648 | 0.6708 |
| Full Layer | f1_std | 0.0869 | 0.0859 |
| Single Neuron | f1_best_layer | 10.0 | 10.0 |
| Single Neuron | f1_max | 0.8549 | 0.861 |
| Single Neuron | f1_mean | 0.5392 | 0.5964 |
| Single Neuron | f1_std | 0.0336 | 0.068 |
| Top-K Neurons | f1_best_layer | 11.0 | 9.0 |
| Top-K Neurons | f1_max | 0.8549 | 0.849 |
| Top-K Neurons | f1_mean | 0.5724 | 0.6688 |
| Top-K Neurons | f1_std | 0.0873 | 0.0872 |

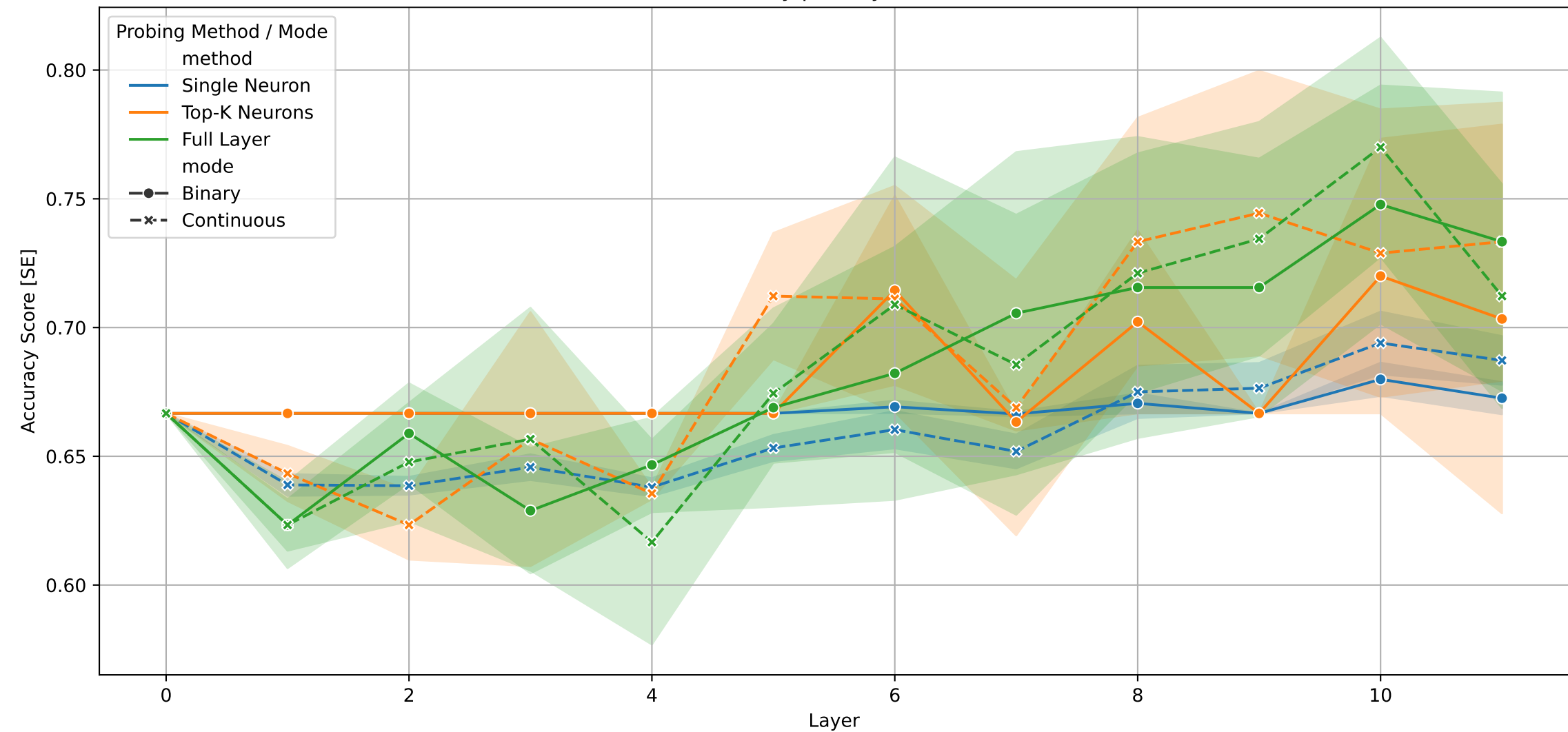Accuracy per Layer – Single Neuron Probing

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 | 10.0 | 10.0 |
| Full Layer | accuracy_max | 0.8433 | 0.8533 |
| Full Layer | accuracy_mean | 0.6828 | 0.6848 |
| Full Layer | accuracy_std | 0.0716 | 0.0739 |
| Single Neuron | accuracy_best_layer | 10.0 | 10.0 |
| Single Neuron | accuracy_max | 0.8533 | 0.86 |
| Single Neuron | accuracy_mean | 0.6688 | 0.6605 |
| Single Neuron | accuracy_std | 0.0162 | 0.0431 |
| Top-K Neurons | accuracy_best_layer | 10.0 | 9.0 |
| Top-K Neurons | accuracy_max | 0.8533 | 0.85 |
| Top-K Neurons | accuracy_mean | 0.6808 | 0.6881 |
| Top-K Neurons | accuracy_std | 0.0487 | 0.0712 |