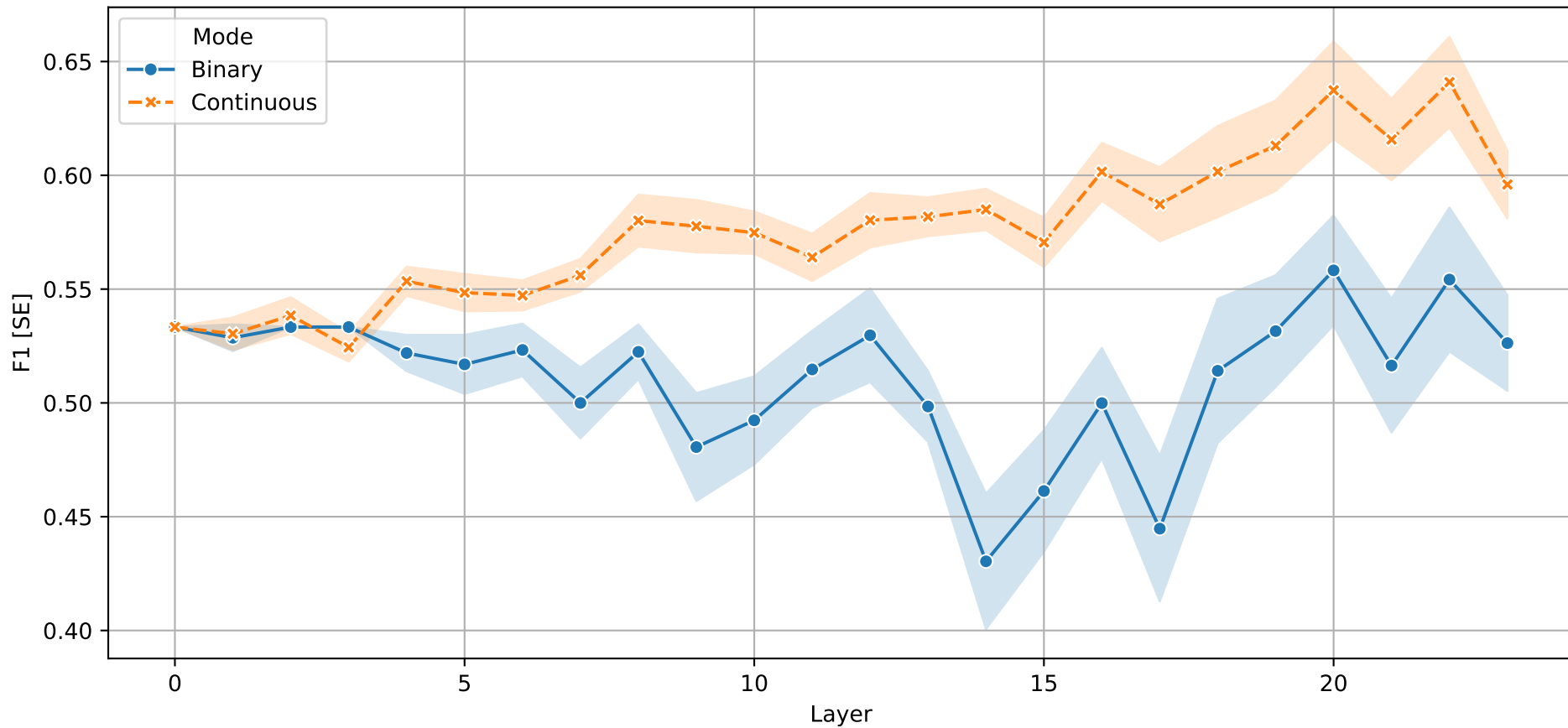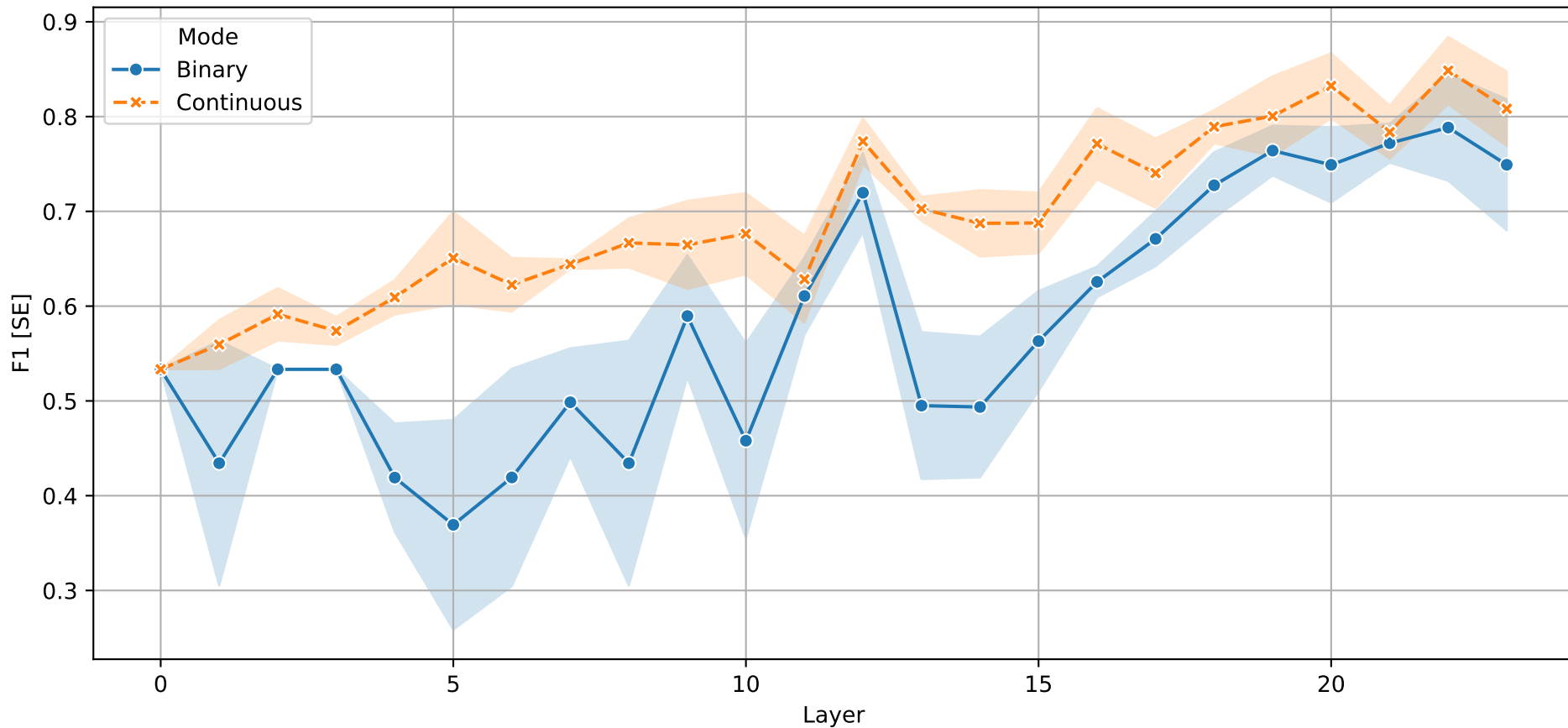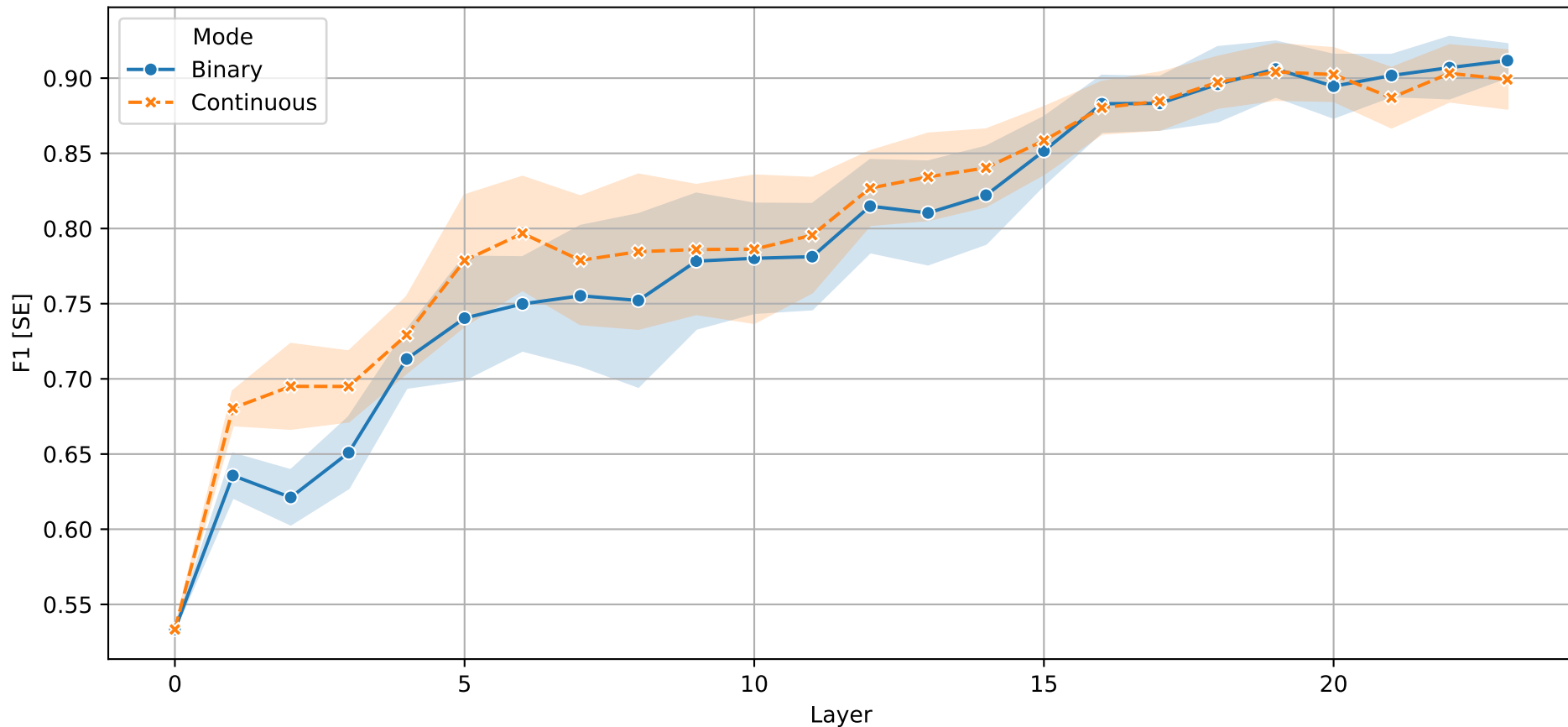F1 per Layer – Single Neuron Probing
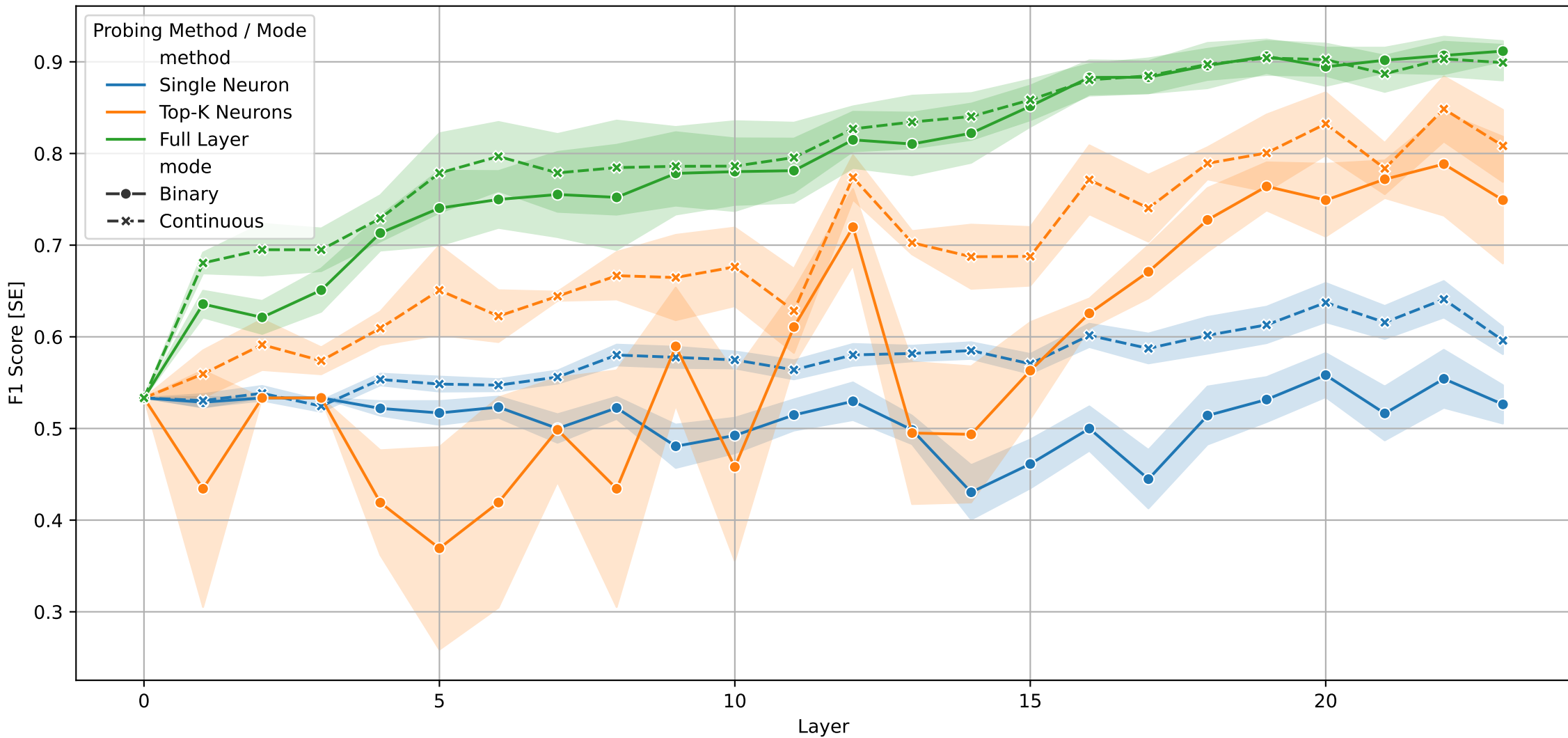
F1 per Layer – Top-K Neurons Probing

F1 per Layer – Full Layer Probing

Overall F1 per Layer – All Methods

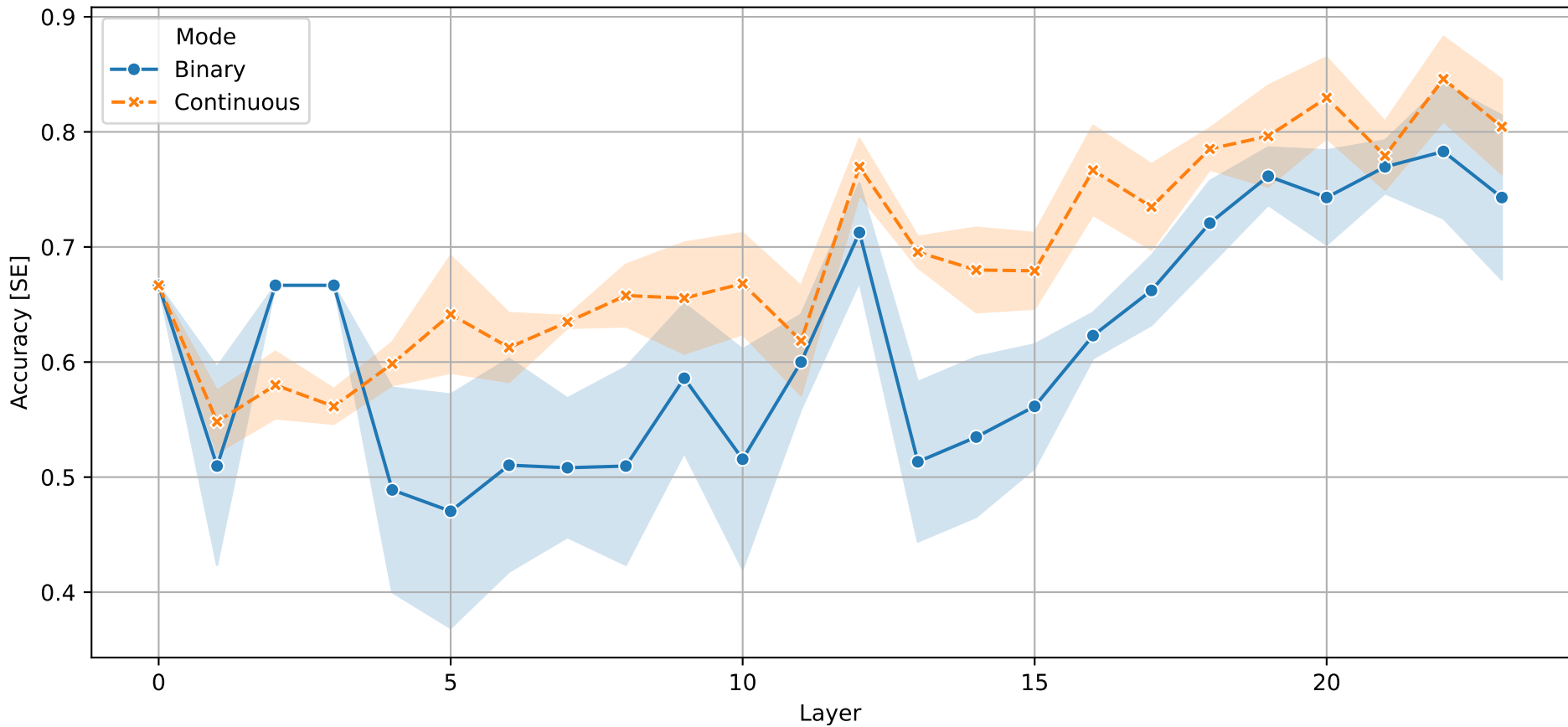## F1 Score Summary by Probing Method

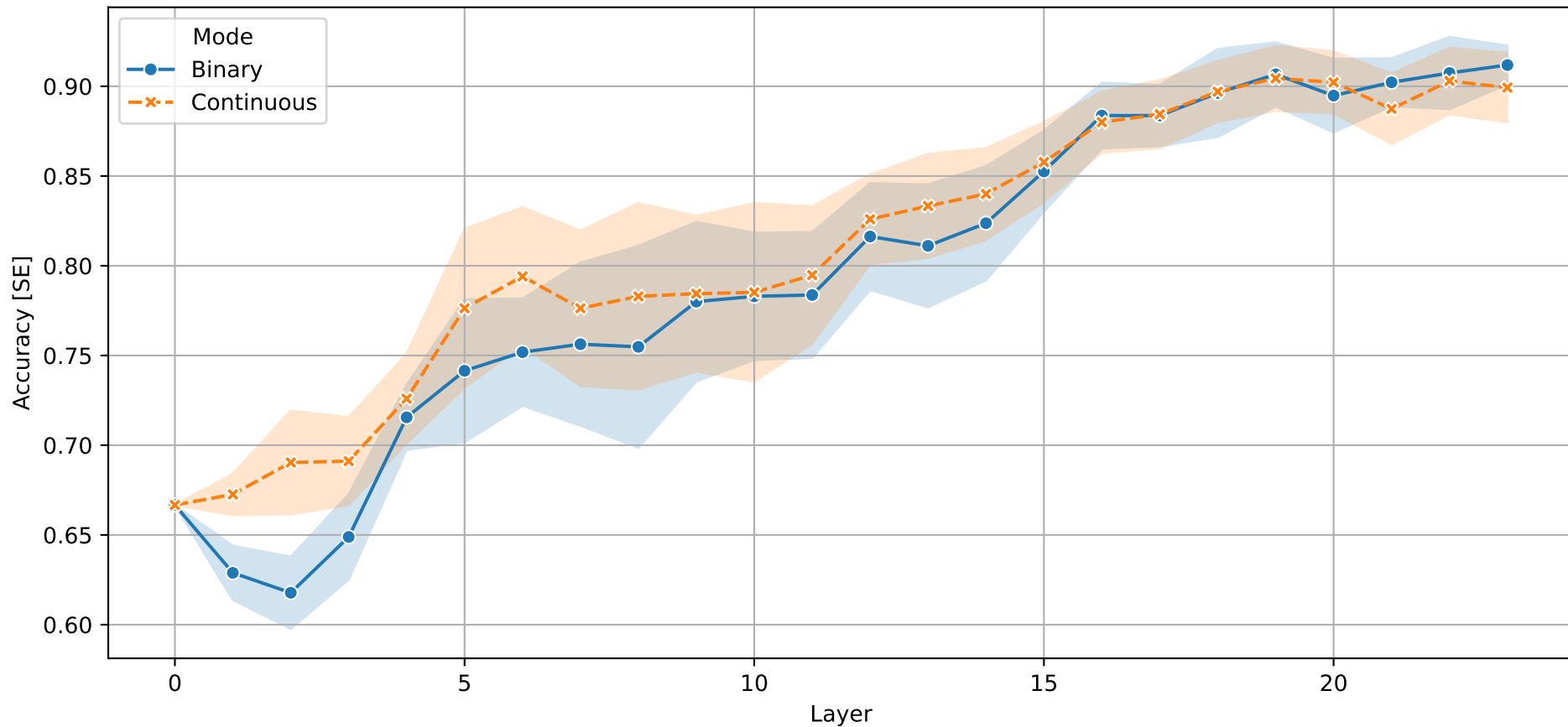| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | f1_best_layer | 23.0 | 19.0 |
| Full Layer | f1_max | 0.9468 | 0.9405 |
| Full Layer | f1_mean | 0.7906 | 0.8066 |
| Full Layer | f1_std | 0.1109 | 0.0994 |
| Single Neuron | f1_best_layer | 20.0 | 22.0 |
| Single Neuron | f1_max | 0.8844 | 0.9151 |
| Single Neuron | f1_mean | 0.5111 | 0.5766 |
| Single Neuron | f1_std | 0.1159 | 0.075 |
| Top-K Neurons | f1_best_layer | 22.0 | 22.0 |
| Top-K Neurons | f1_max | 0.8839 | 0.916 |
| Top-K Neurons | f1_mean | 0.5813 | 0.6936 |
| Top-K Neurons | f1_std | 0.1627 | 0.1007 |

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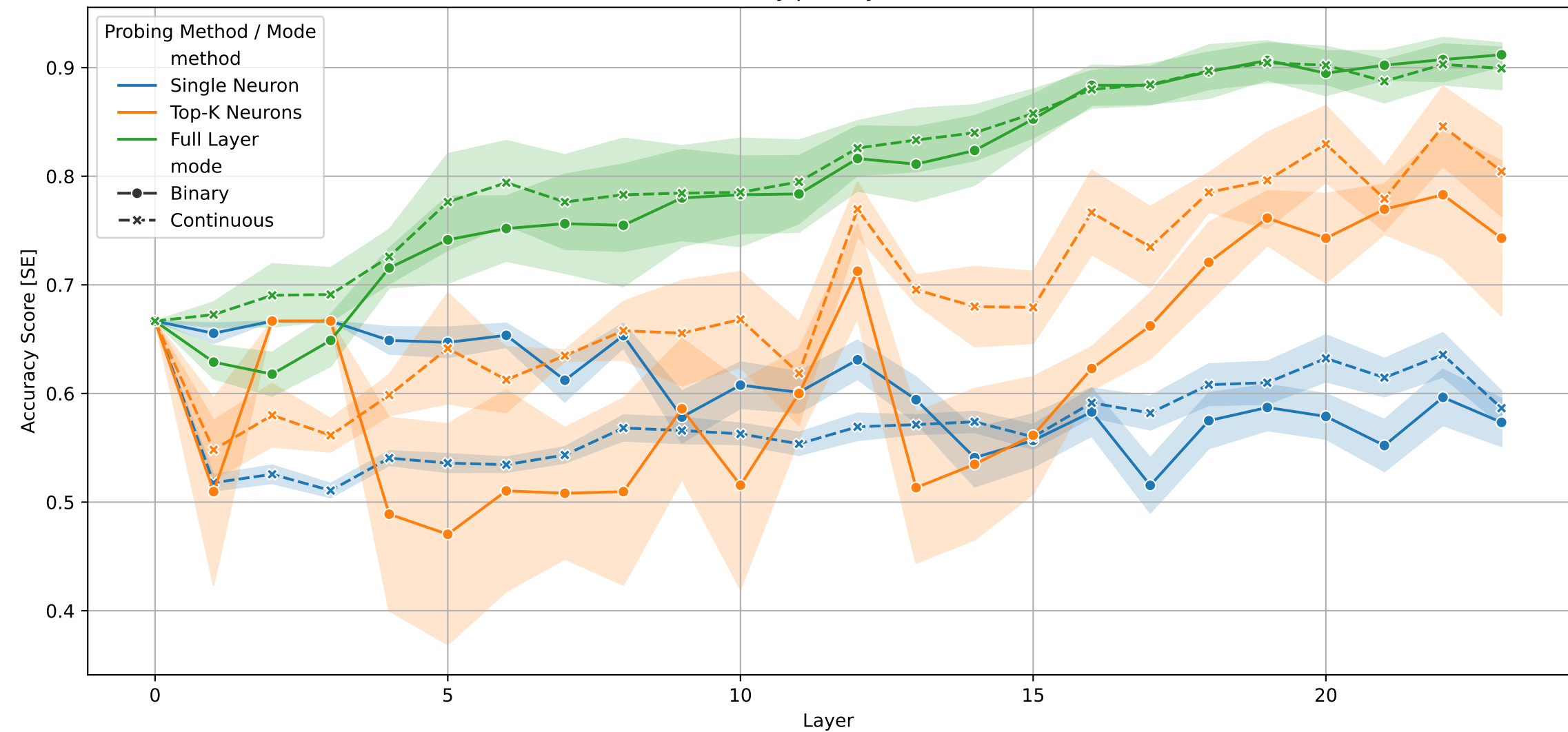
Accuracy per Layer – Top-K Neurons Probing

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 | 23.0 | 19.0 |
| Full Layer | accuracy_max | 0.9467 | 0.94 |
| Full Layer | accuracy_mean | 0.7966 | 0.8106 |
| Full Layer | accuracy_std | 0.1012 | 0.0882 |
| Single Neuron | accuracy_best_layer | 0.0 | 0.0 |
| Single Neuron | accuracy_max | 0.8822 | 0.9156 |
| Single Neuron | accuracy_mean | 0.6059 | 0.5734 |
| Single Neuron | accuracy_std | 0.112 | 0.0785 |
| Top-K Neurons | accuracy_best_layer | 22.0 | 22.0 |
| Top-K Neurons | accuracy_max | 0.8822 | 0.9156 |
| Top-K Neurons | accuracy_mean | 0.6178 | 0.6921 |
| Top-K Neurons | accuracy_std | 0.1342 | 0.0984 |