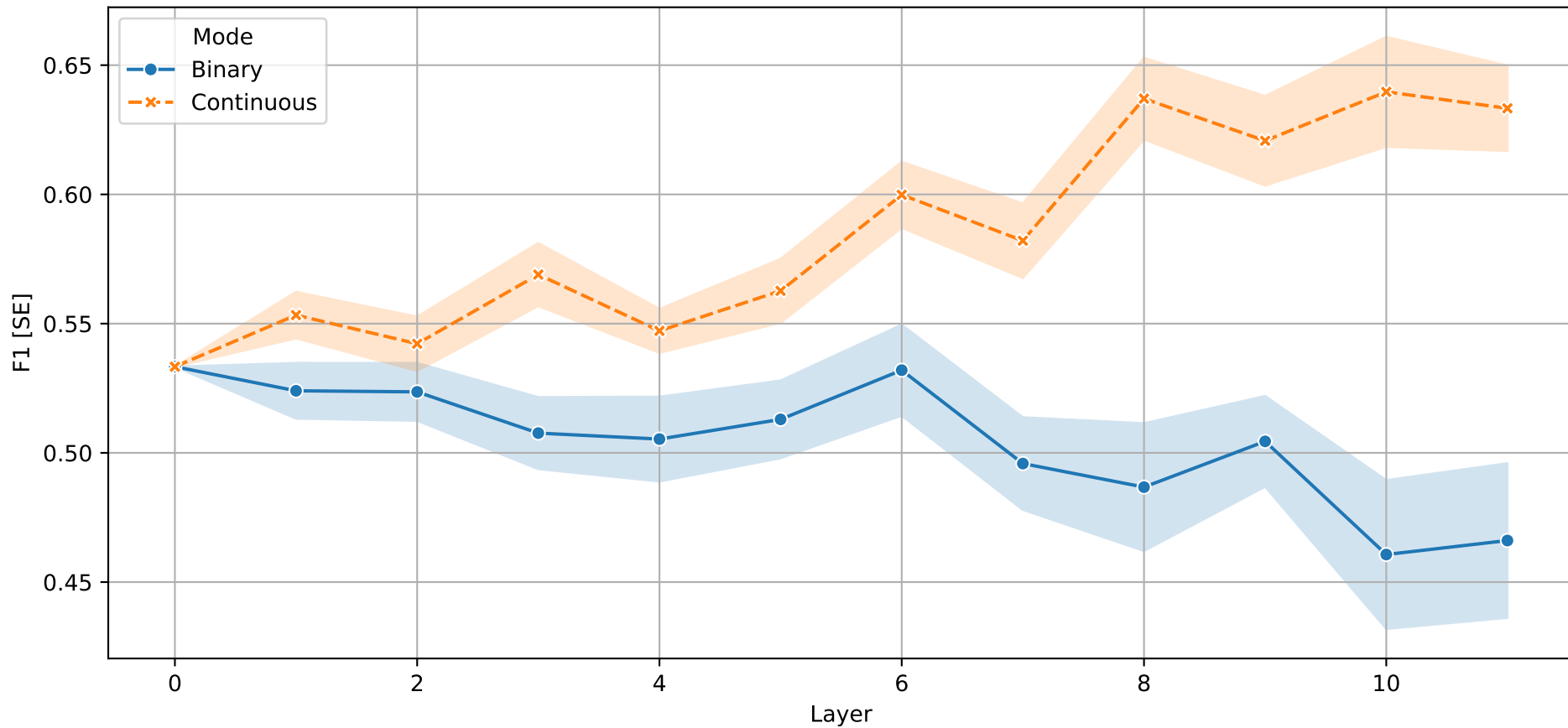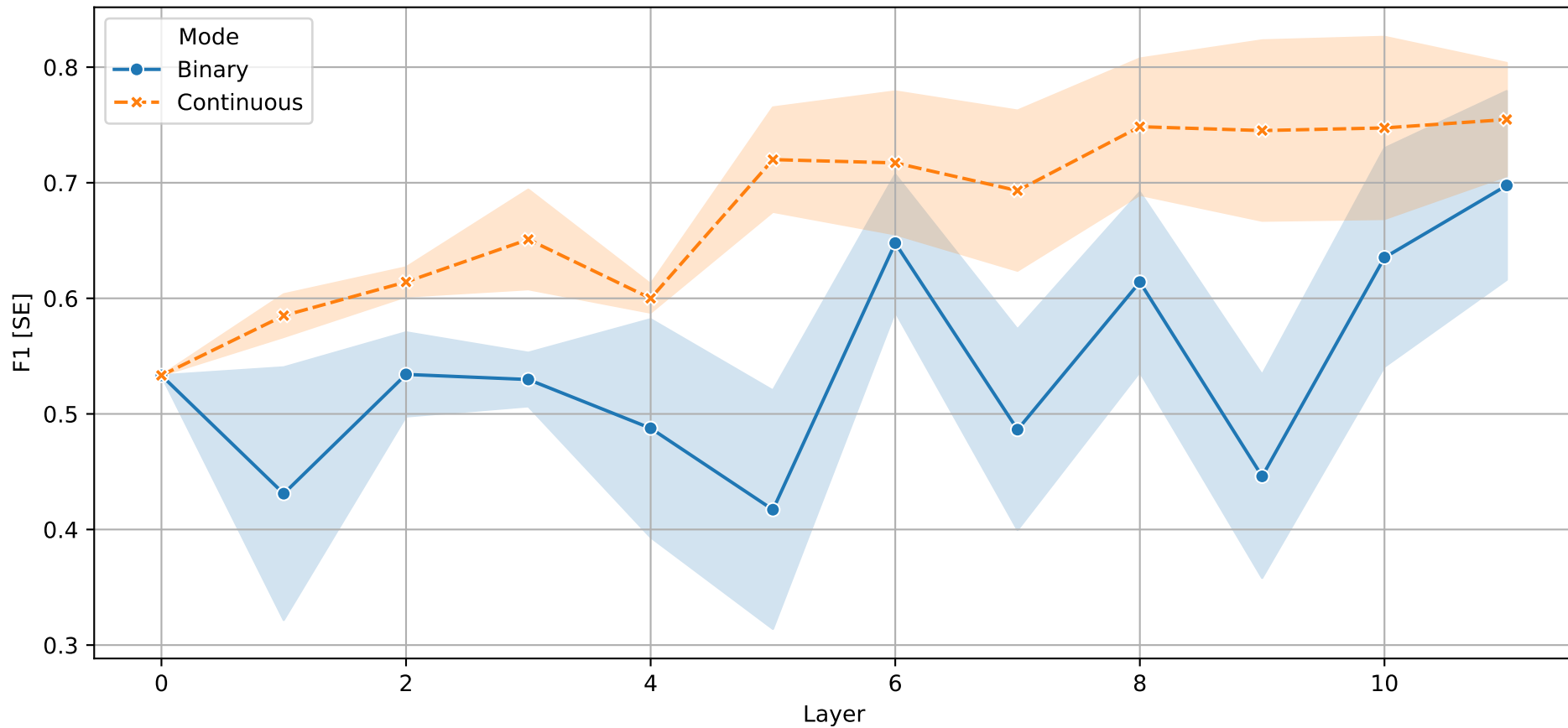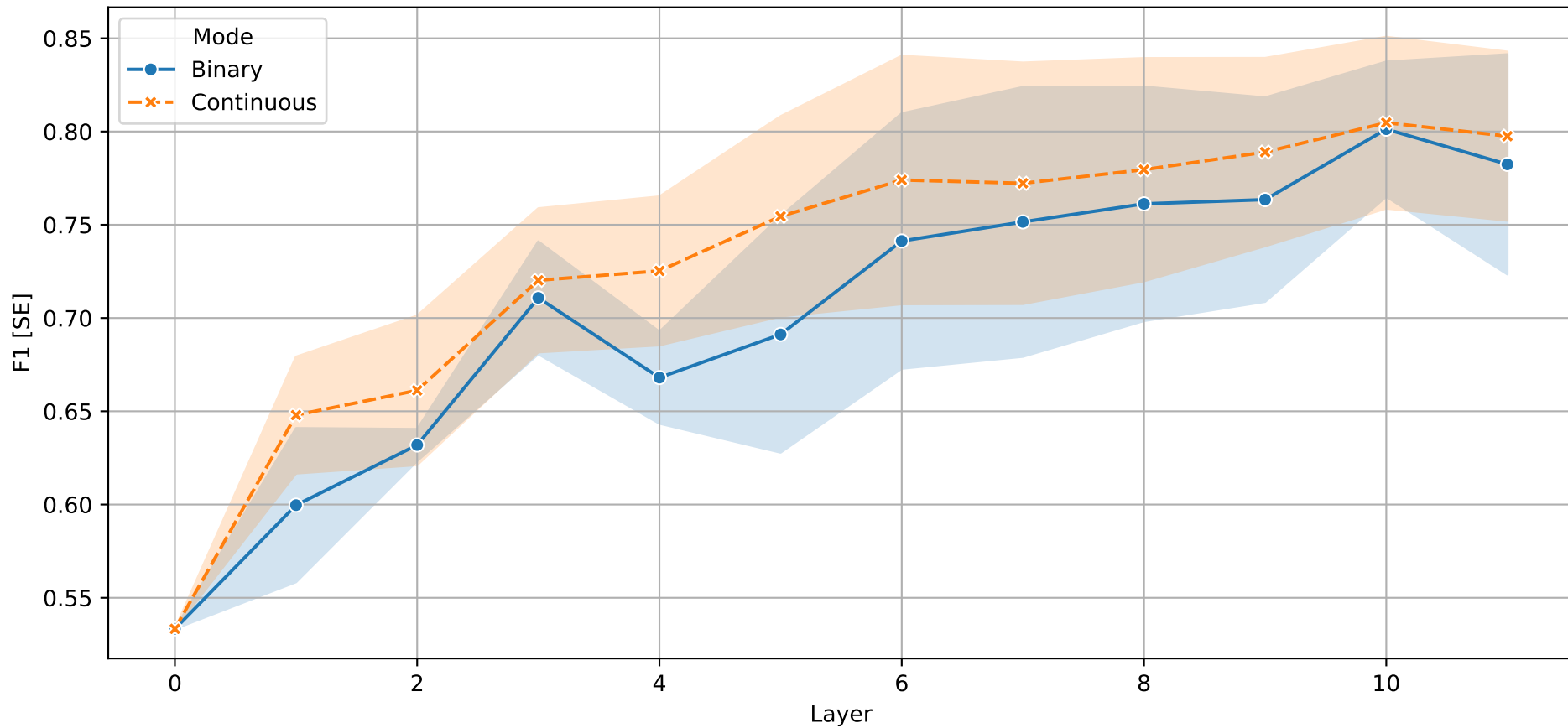F1 per Layer – Single Neuron Probing
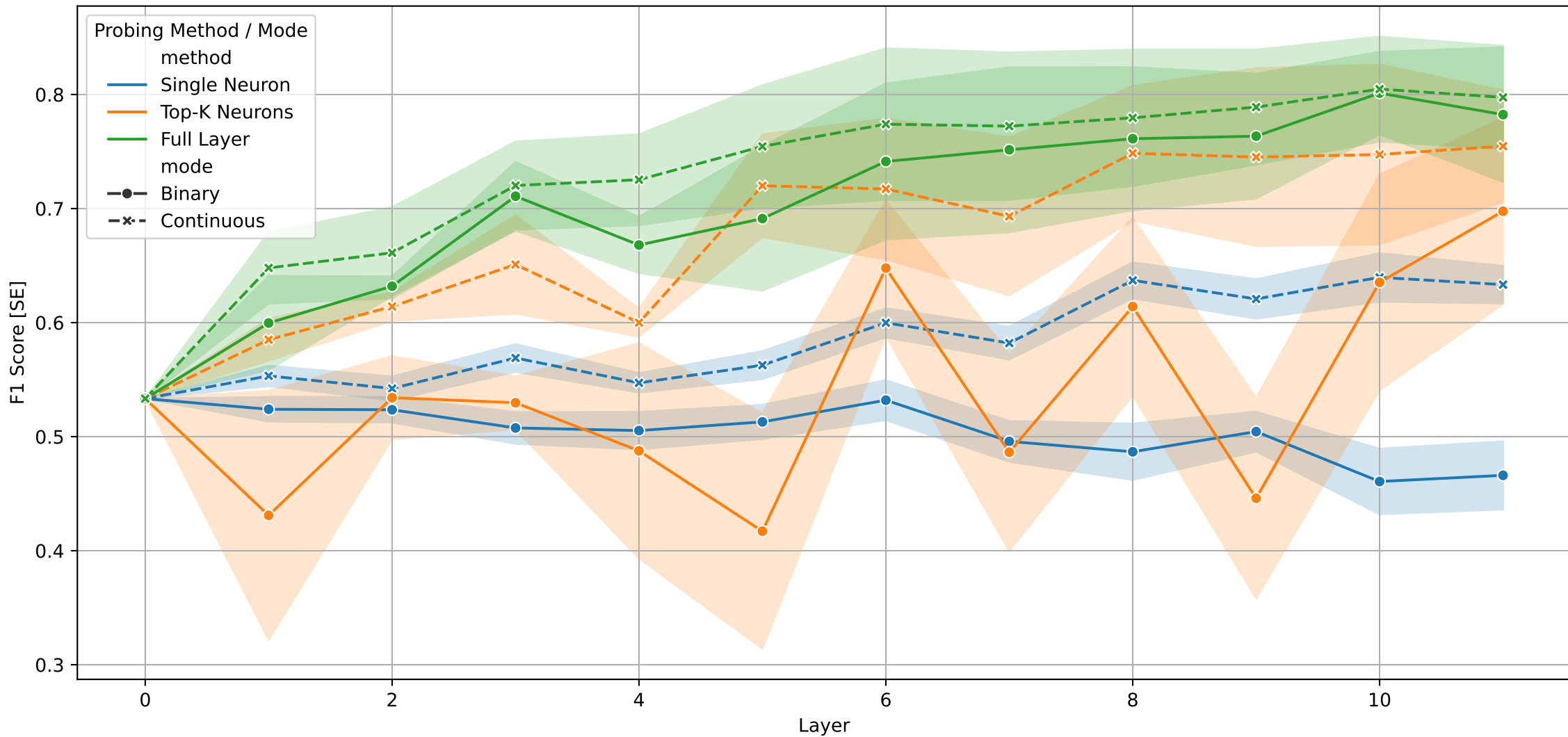
F1 per Layer – Top-K Neurons Probing

F1 per Layer – Full Layer Probing

Overall F1 per Layer – All Methods

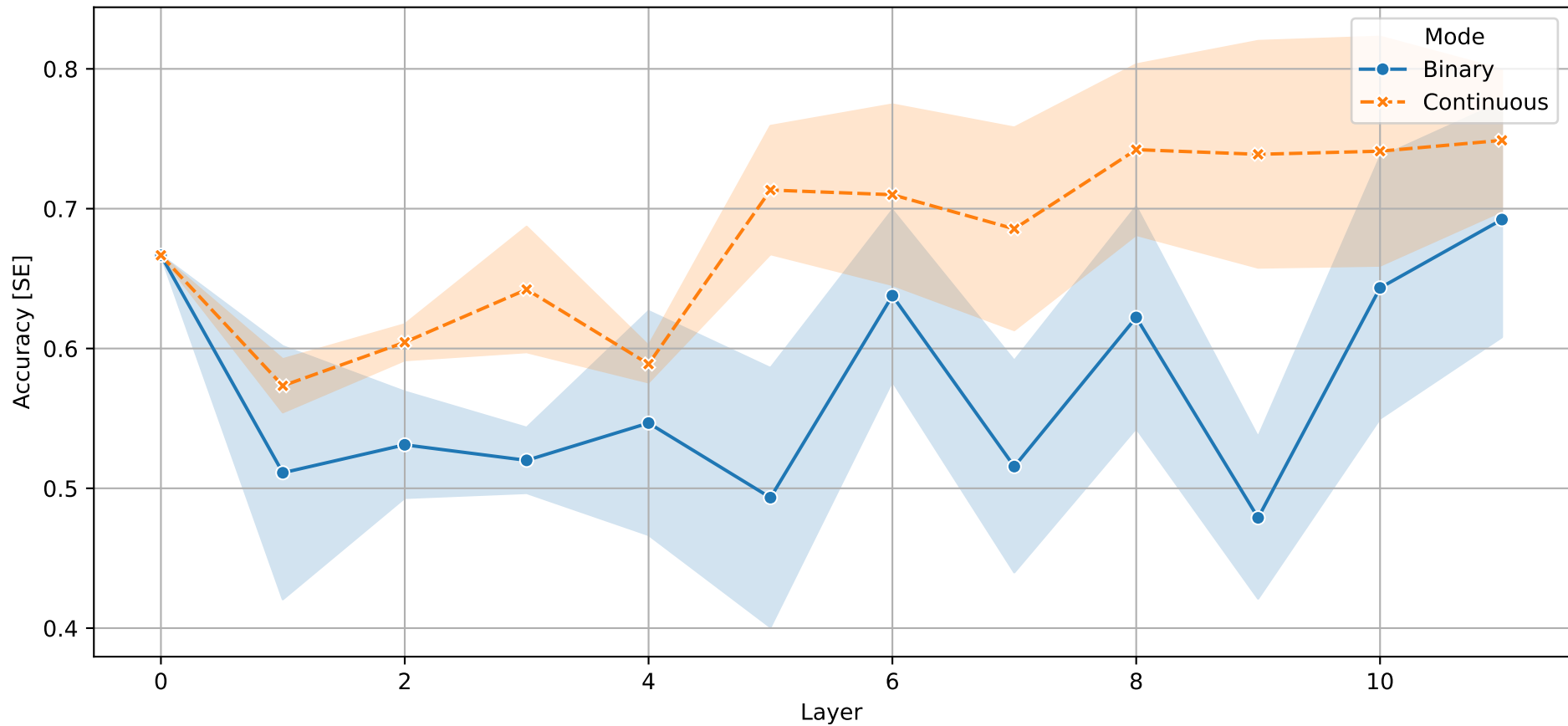## F1 Score Summary by Probing Method

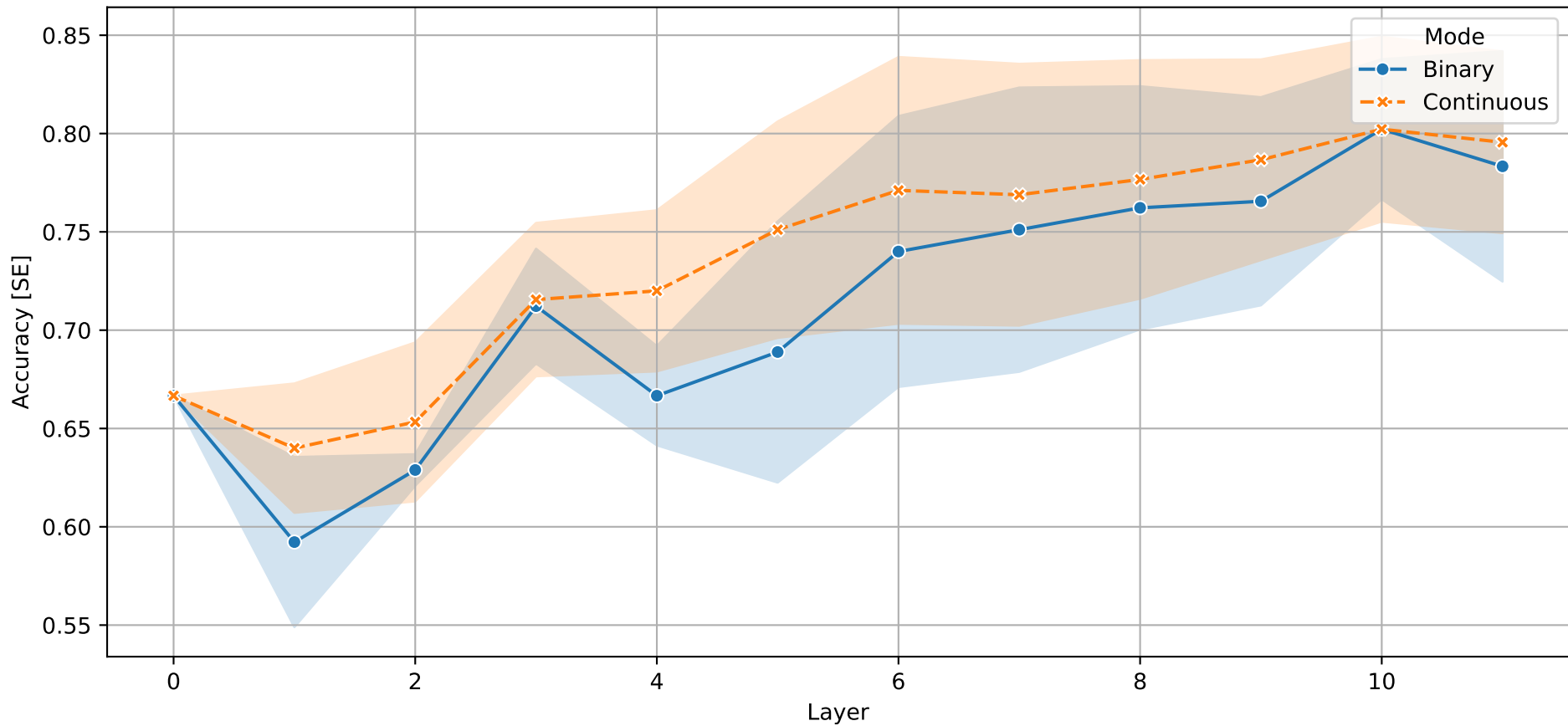| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | f1_best_layer | 10.0 | 10.0 |
| Full Layer | f1_max | 0.8931 | 0.8938 |
| Full Layer | f1_mean | 0.703 | 0.73 |
| Full Layer | f1_std | 0.1058 | 0.1037 |
| Single Neuron | f1_best_layer | 0.0 | 10.0 |
| Single Neuron | f1_max | 0.8549 | 0.862 |
| Single Neuron | f1_mean | 0.5044 | 0.585 |
| Single Neuron | f1_std | 0.1034 | 0.0824 |
| Top-K Neurons | f1_best_layer | 11.0 | 11.0 |
| Top-K Neurons | f1_max | 0.8549 | 0.8808 |
| Top-K Neurons | f1_mean | 0.5384 | 0.6758 |
| Top-K Neurons | f1_std | 0.1435 | 0.1042 |

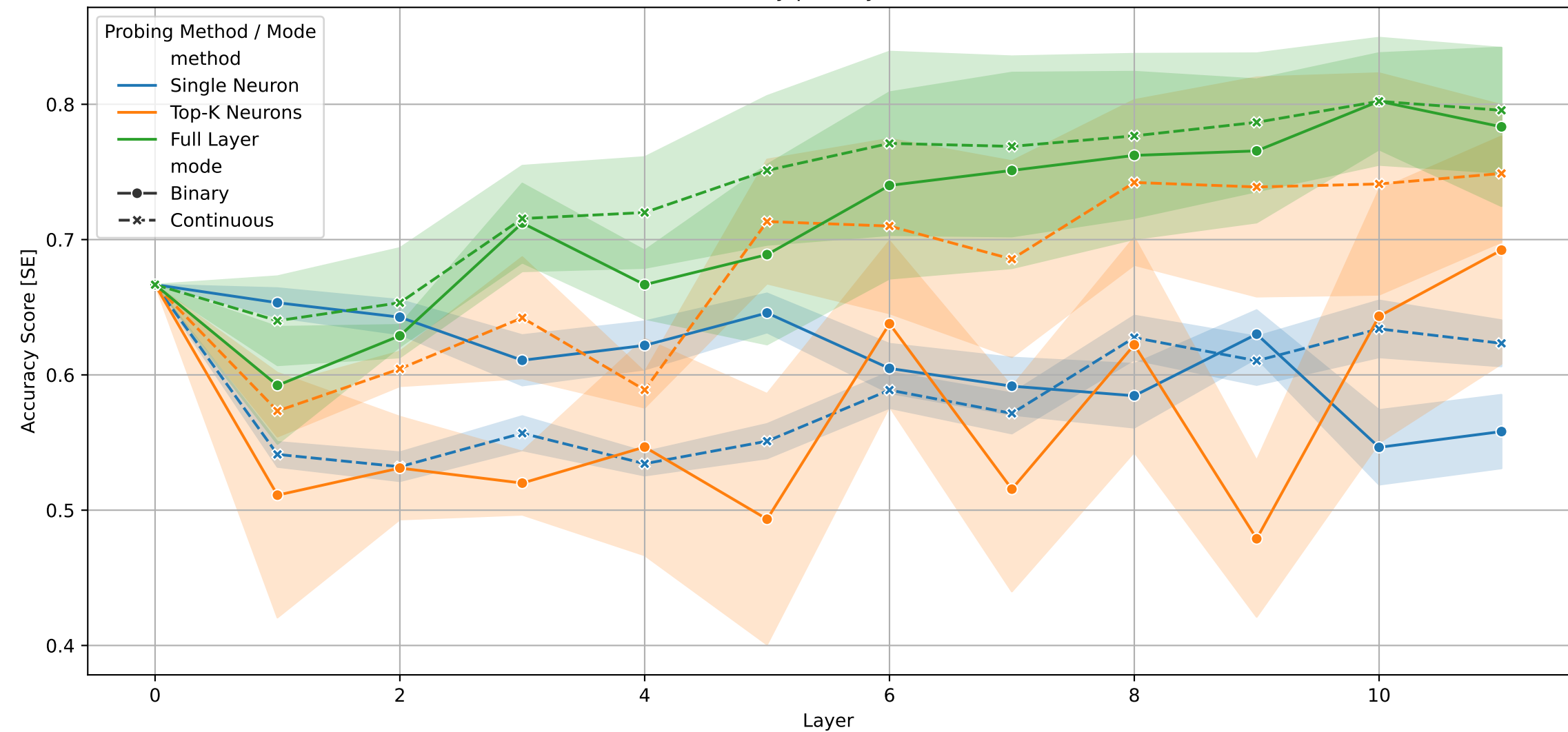Accuracy per Layer – Single Neuron Probing

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 | 0.8933 | 0.8933 |
| Full Layer | accuracy_mean | 0.7133 | 0.7373 |
| Full Layer | accuracy_std | 0.0949 | 0.0895 |
| Single Neuron | accuracy_best_layer | 0.0 | 0.0 |
| Single Neuron | accuracy_max | 0.8533 | 0.86 |
| Single Neuron | accuracy_mean | 0.6131 | 0.5865 |
| Single Neuron | accuracy_std | 0.1082 | 0.0864 |
| Top-K Neurons | accuracy_best_layer | 11.0 | 11.0 |
| Top-K Neurons | accuracy_max | 0.8533 | 0.88 |
| Top-K Neurons | accuracy_mean | 0.5716 | 0.6796 |
| Top-K Neurons | accuracy_std | 0.125 | 0.0981 |