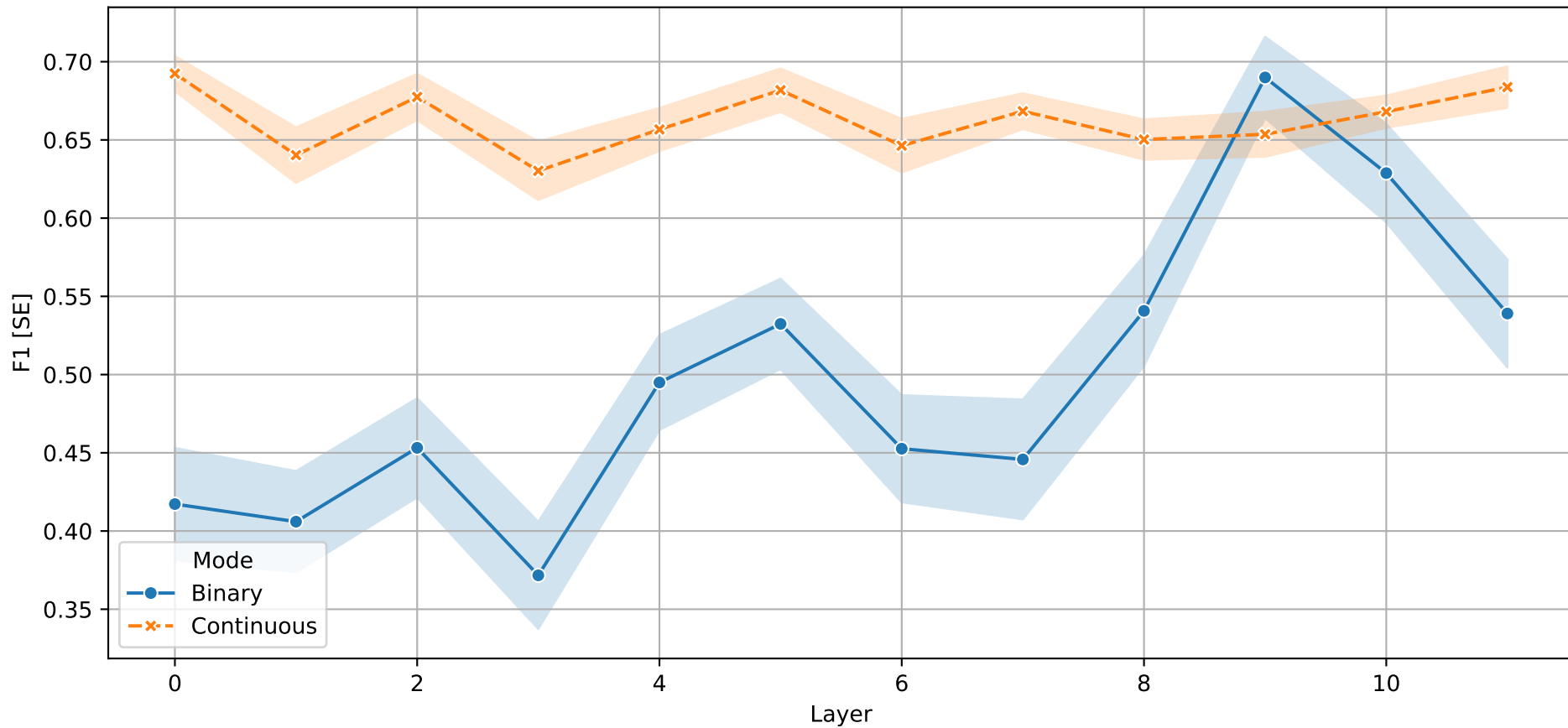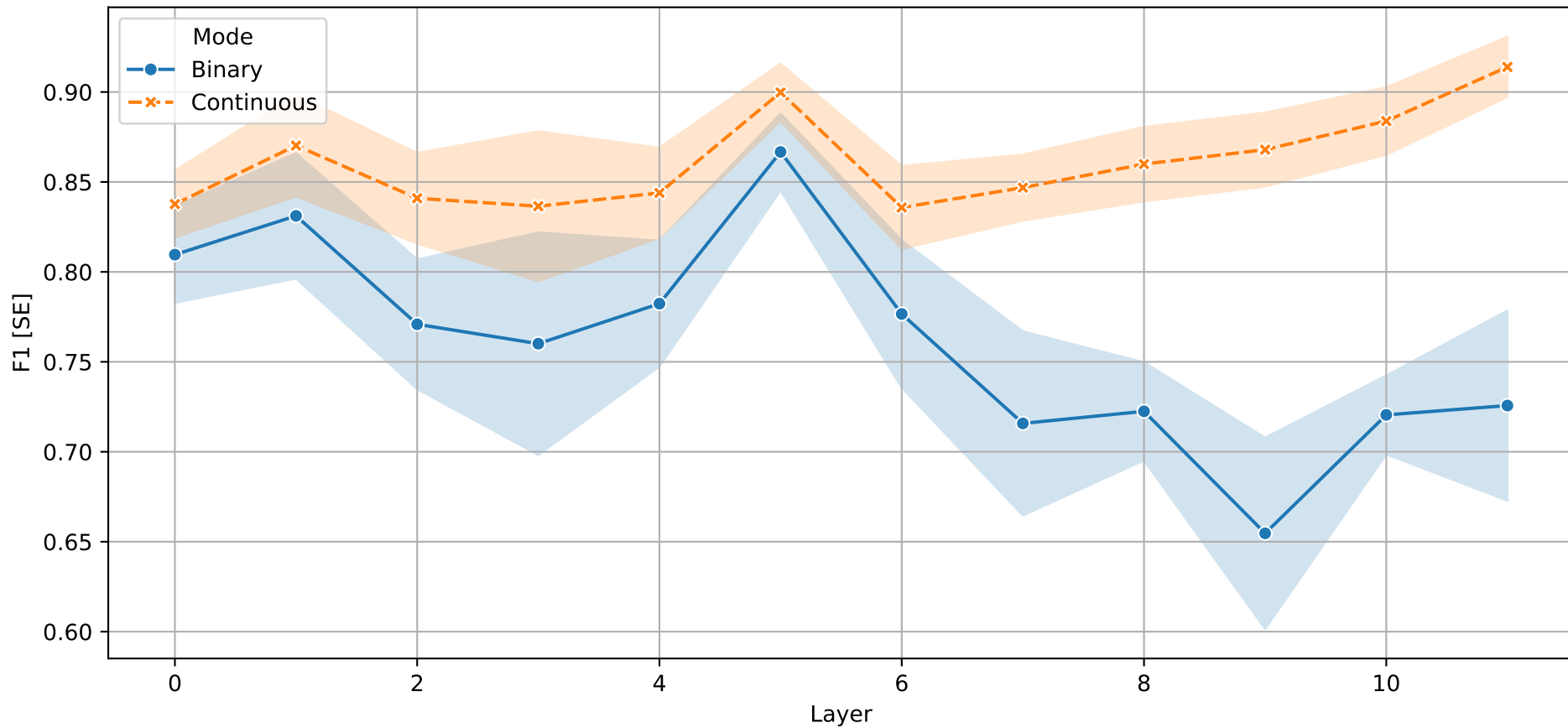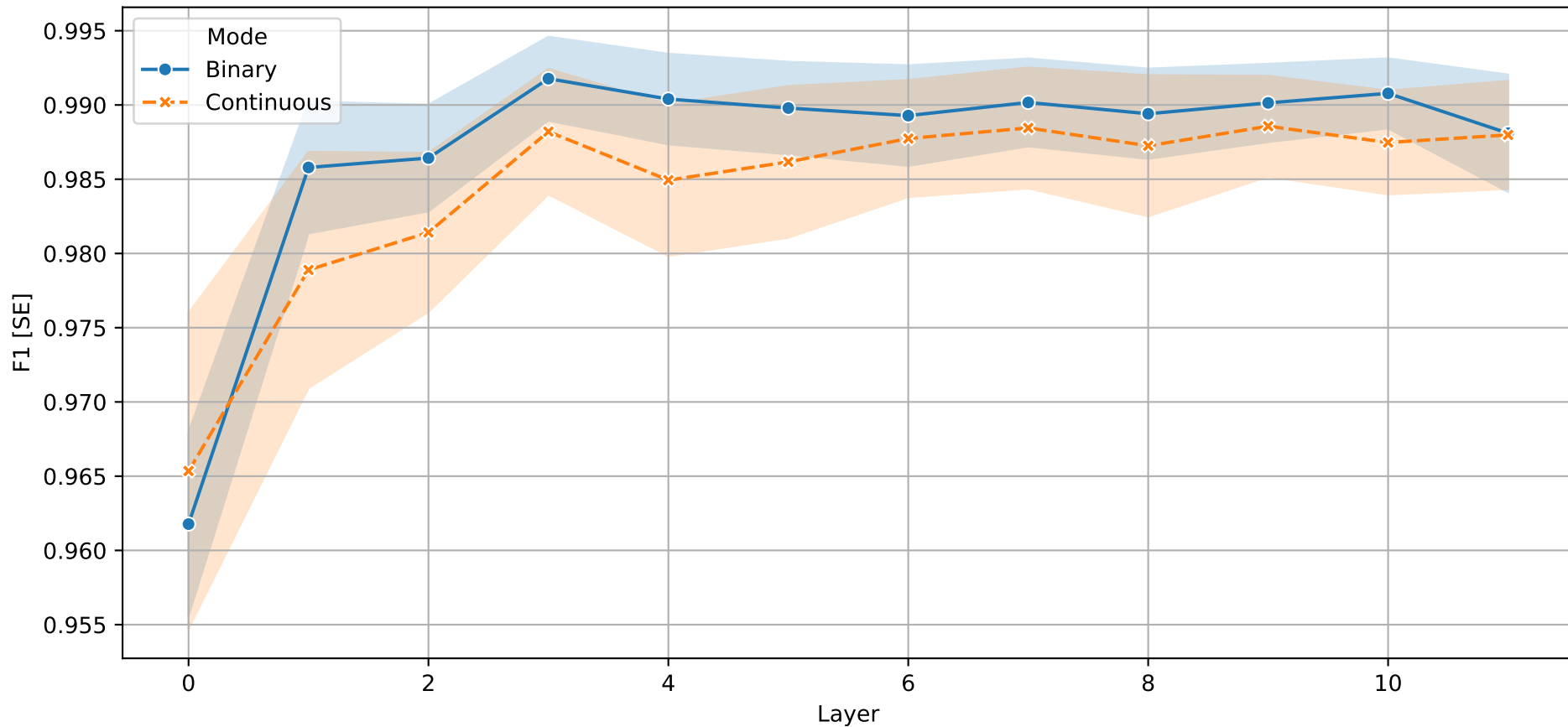F1 per Layer – Single Neuron Probing
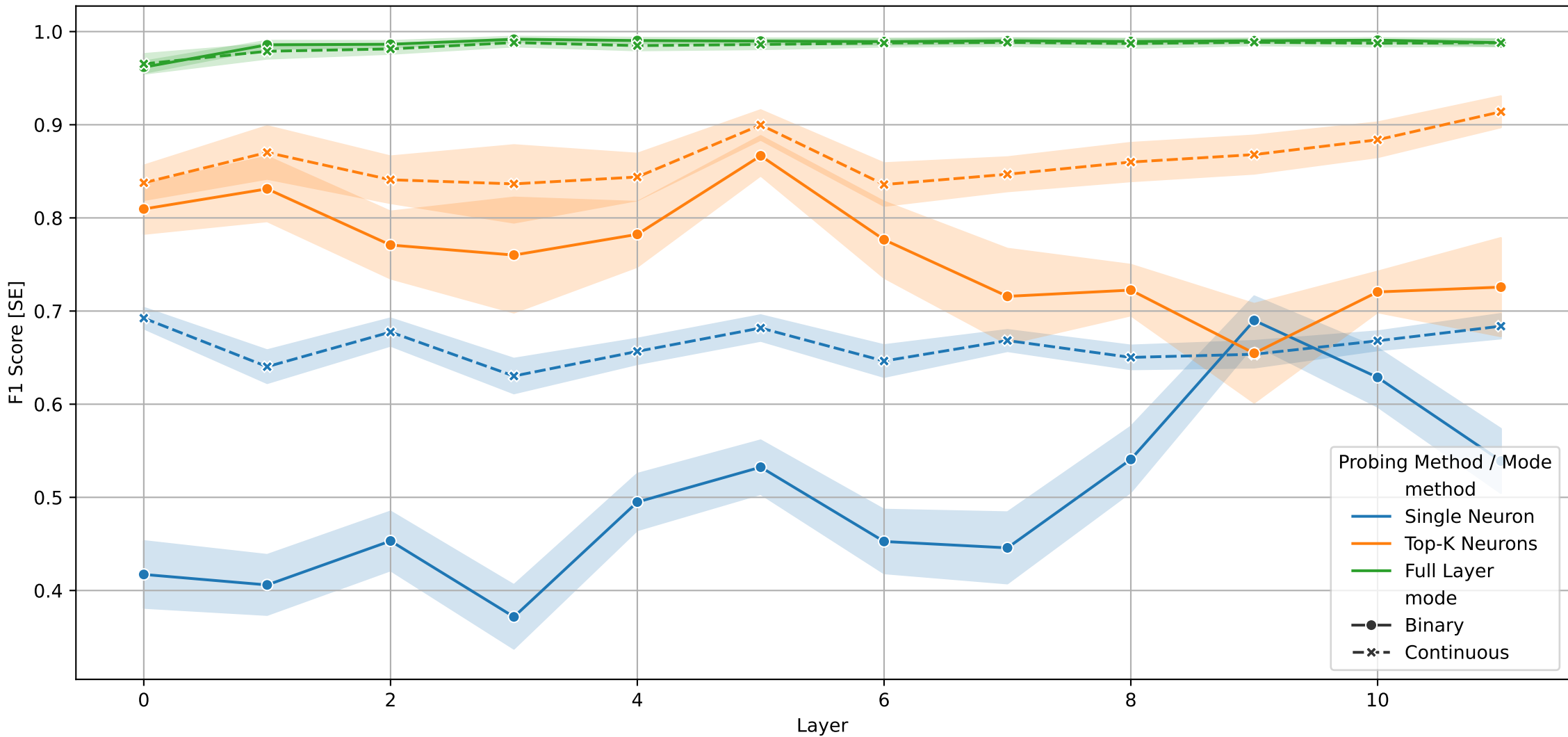
F1 per Layer – Top-K Neurons Probing

F1 per Layer – Full Layer Probing

Overall F1 per Layer – All Methods

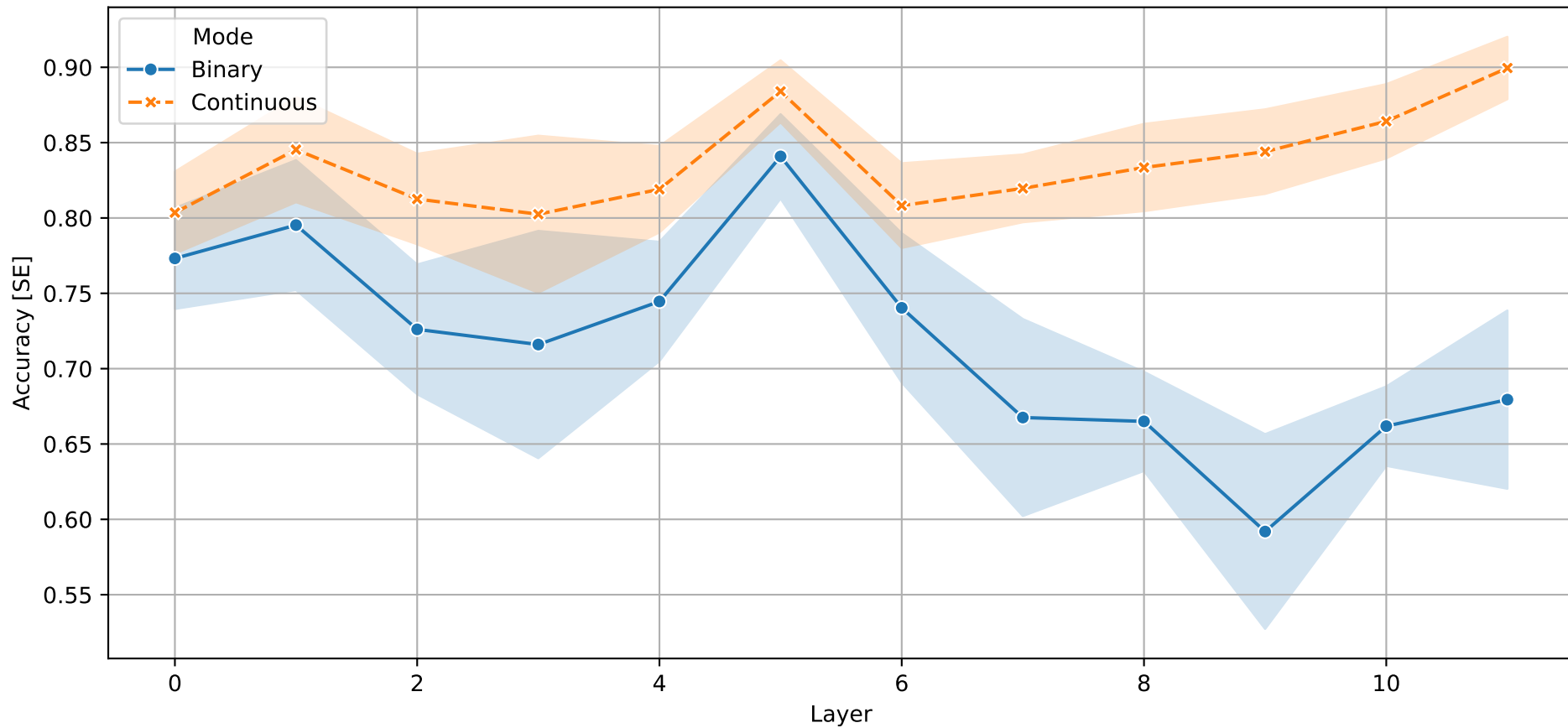## F1 Score Summary by Probing Method

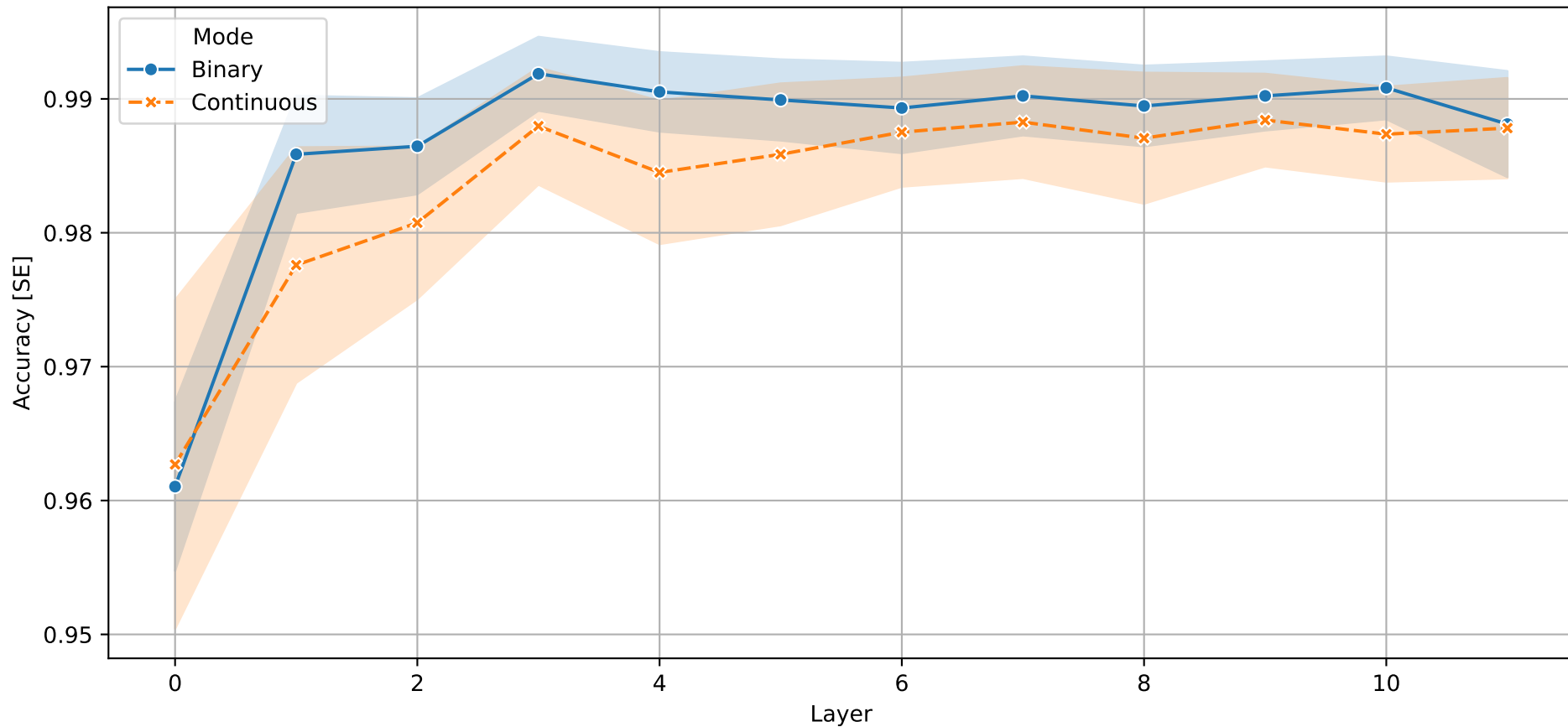| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | f1_best_layer | 3.0 | 9.0 |
| Full Layer | f1_max | 0.9976 | 0.9976 |
| Full Layer | f1_mean | 0.987 | 0.9844 |
| Full Layer | f1_std | 0.0123 | 0.016 |
| Single Neuron | f1_best_layer | 9.0 | 0.0 |
| Single Neuron | f1_max | 0.9761 | 0.9748 |
| Single Neuron | f1_mean | 0.4977 | 0.6624 |
| Single Neuron | f1_std | 0.3058 | 0.1284 |
| Top-K Neurons | f1_best_layer | 5.0 | 11.0 |
| Top-K Neurons | f1_max | 0.9787 | 0.9916 |
| Top-K Neurons | f1_mean | 0.7614 | 0.8615 |
| Top-K Neurons | f1_std | 0.1217 | 0.0679 |

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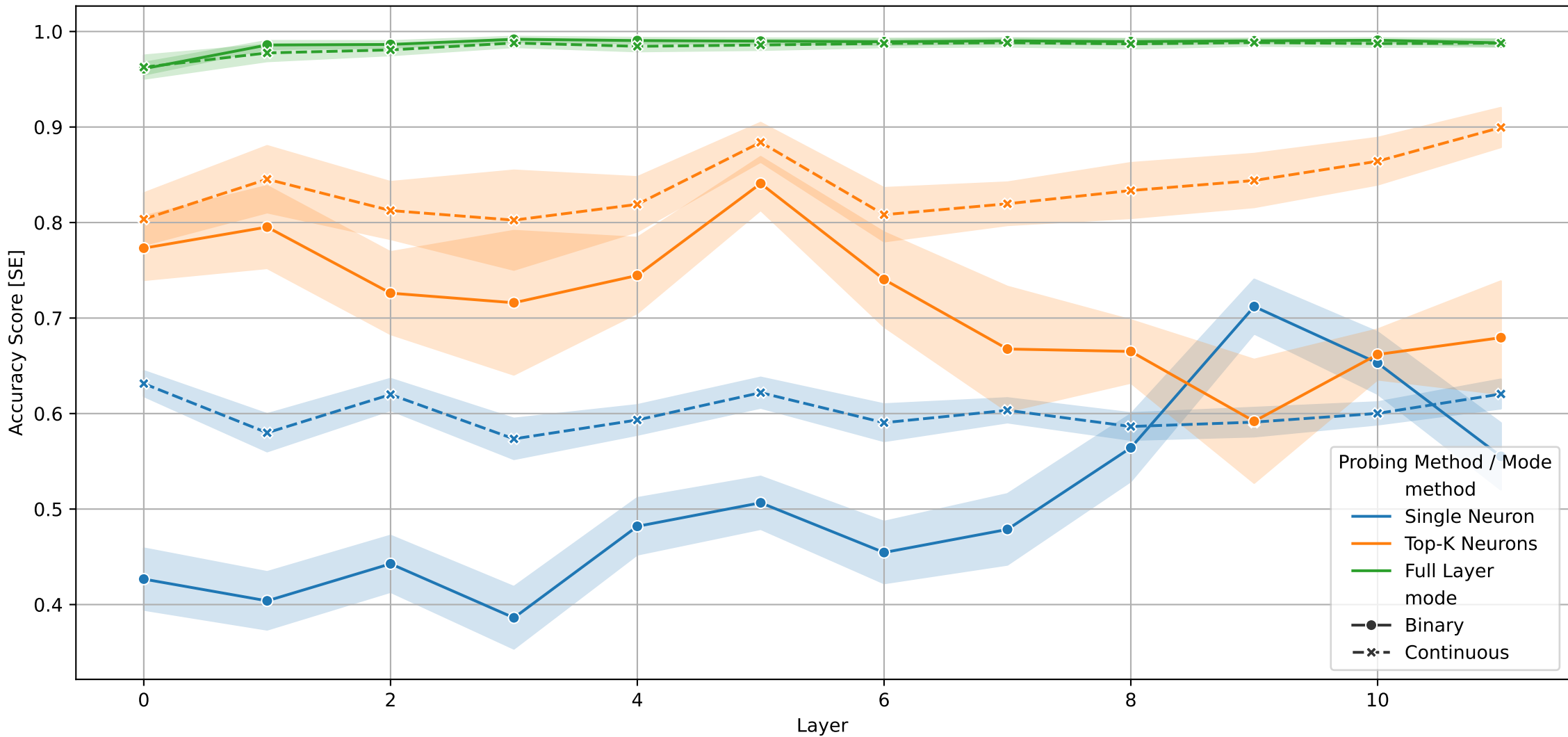
Accuracy per Layer – Top-K Neurons Probing

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 | 3.0 | 9.0 |
| Full Layer | accuracy_max | 0.9976 | 0.9976 |
| Full Layer | accuracy_mean | 0.987 | 0.9838 |
| Full Layer | accuracy_std | 0.0125 | 0.0175 |
| Single Neuron | accuracy_best_layer | 9.0 | 0.0 |
| Single Neuron | accuracy_max | 0.9771 | 0.9759 |
| Single Neuron | accuracy_mean | 0.5055 | 0.601 |
| Single Neuron | accuracy_std | 0.2995 | 0.1445 |
| Top-K Neurons | accuracy_best_layer | 5.0 | 11.0 |
| Top-K Neurons | accuracy_max | 0.9783 | 0.9916 |
| Top-K Neurons | accuracy_mean | 0.7168 | 0.8363 |
| Top-K Neurons | accuracy_std | 0.1464 | 0.0854 |