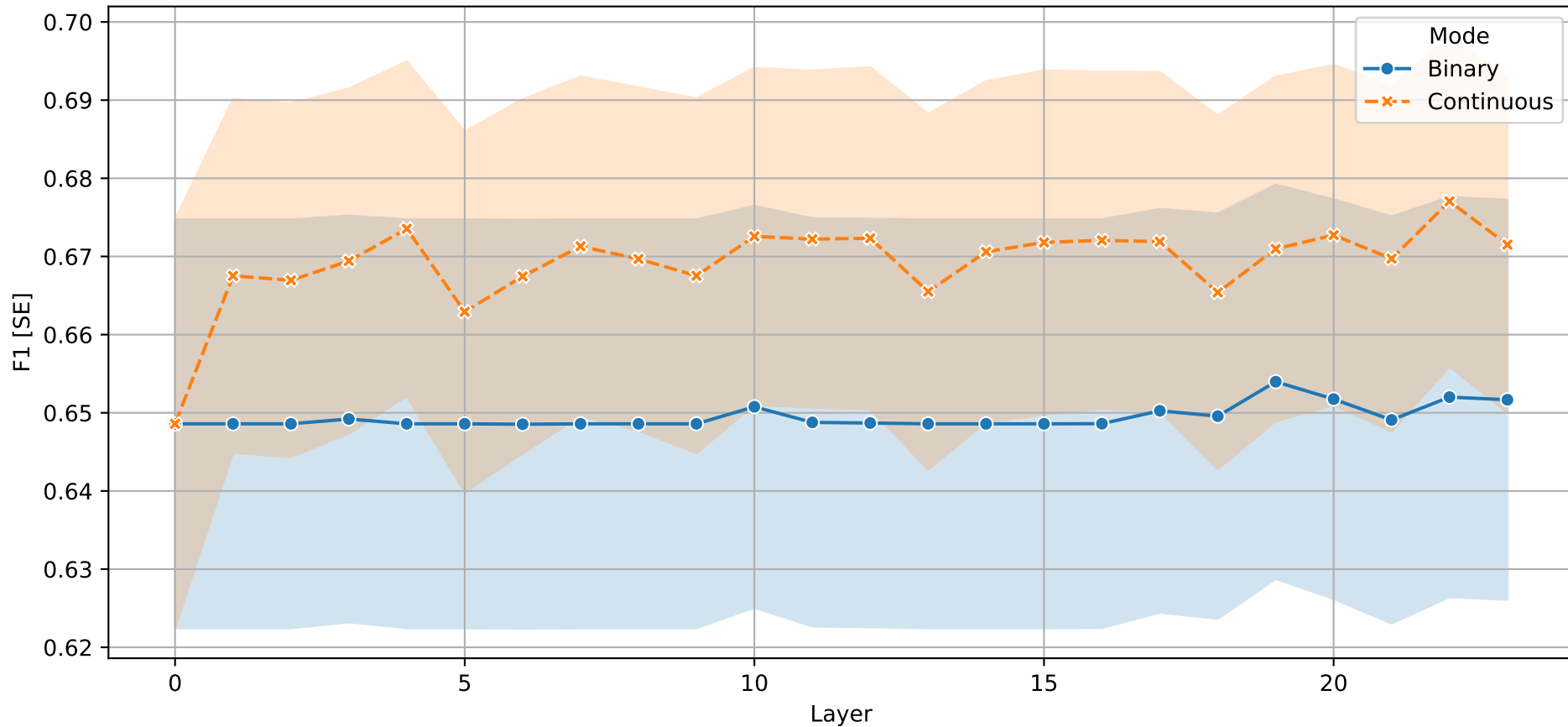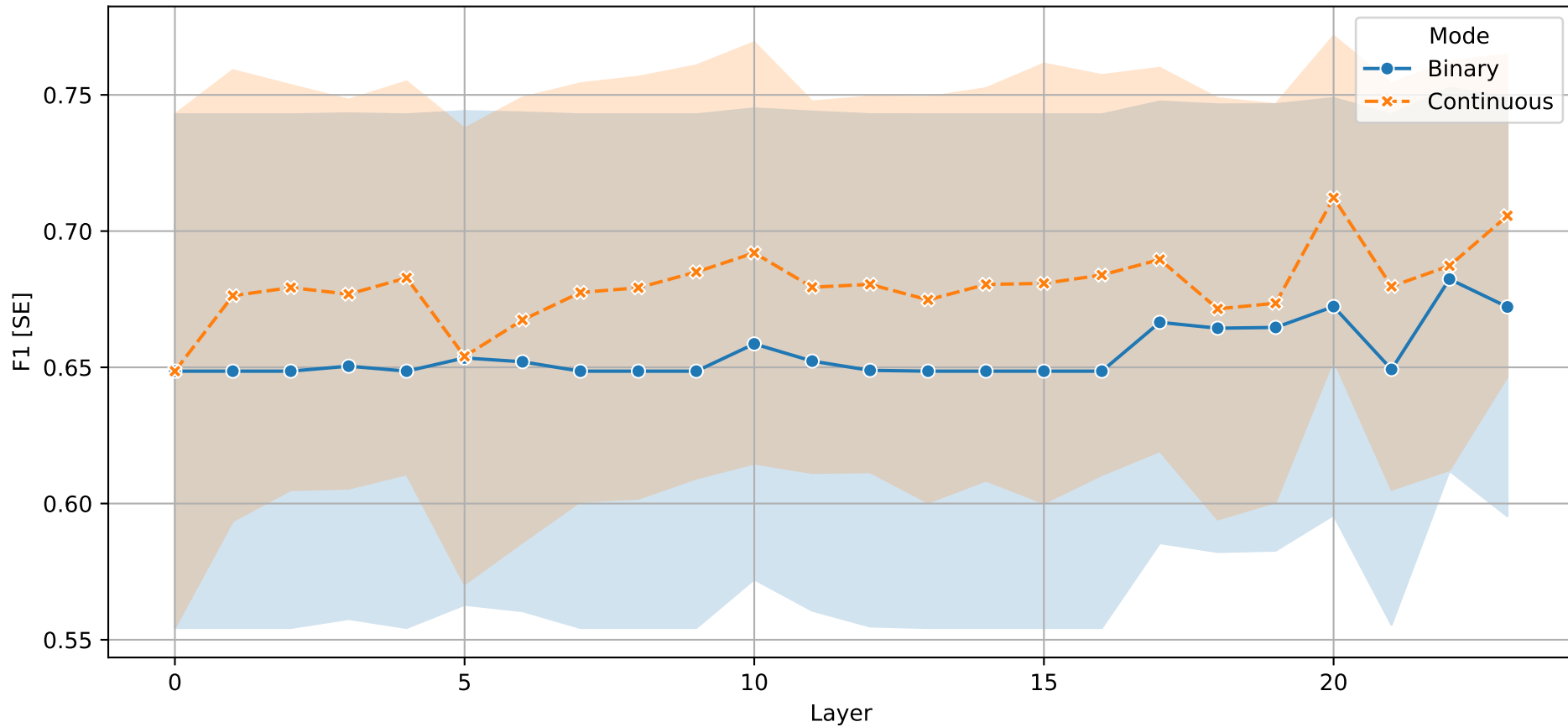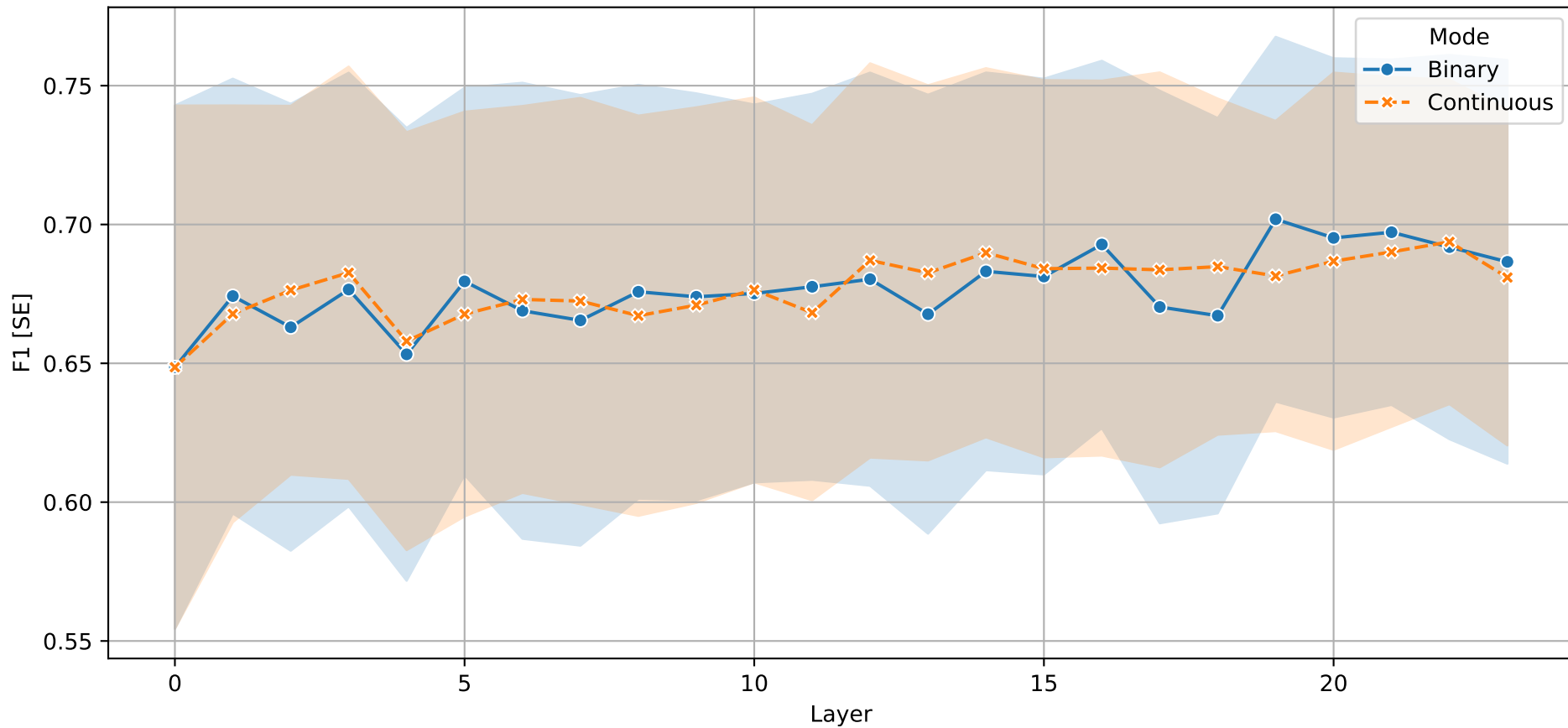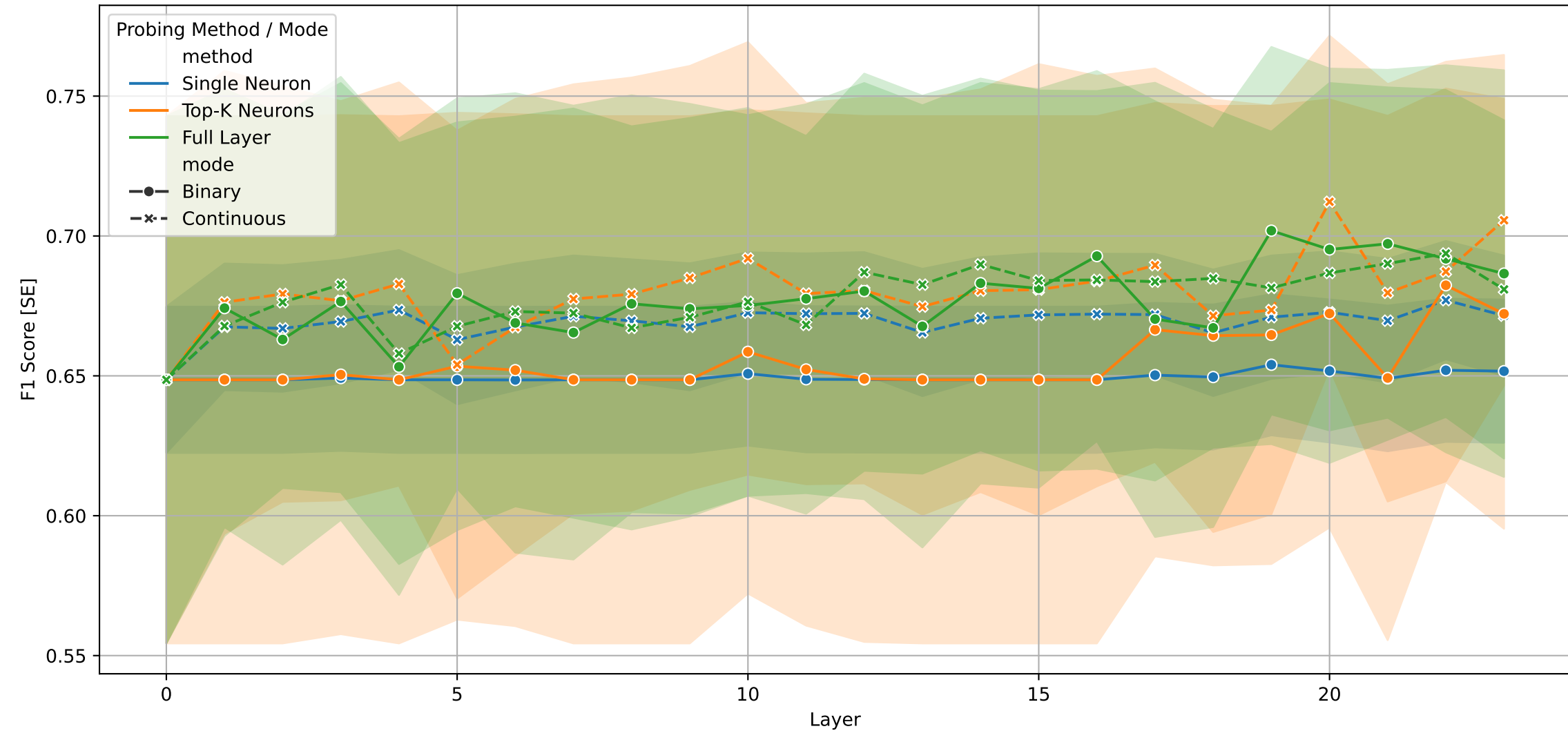F1 per Layer – Single Neuron Probing

F1 per Layer – Top-K Neurons Probing

F1 per Layer – Full Layer Probing

Overall F1 per Layer – All Methods

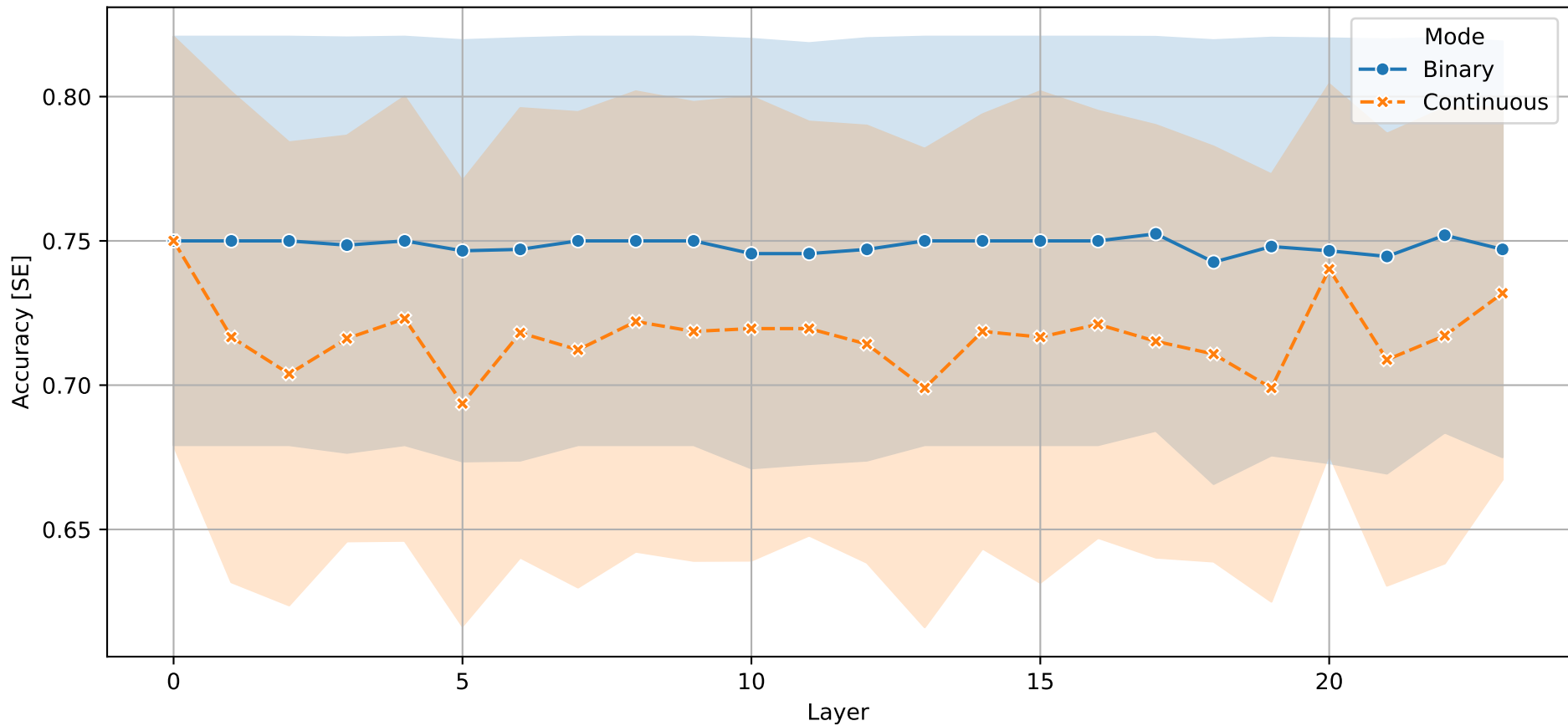## F1 Score Summary by Probing Method

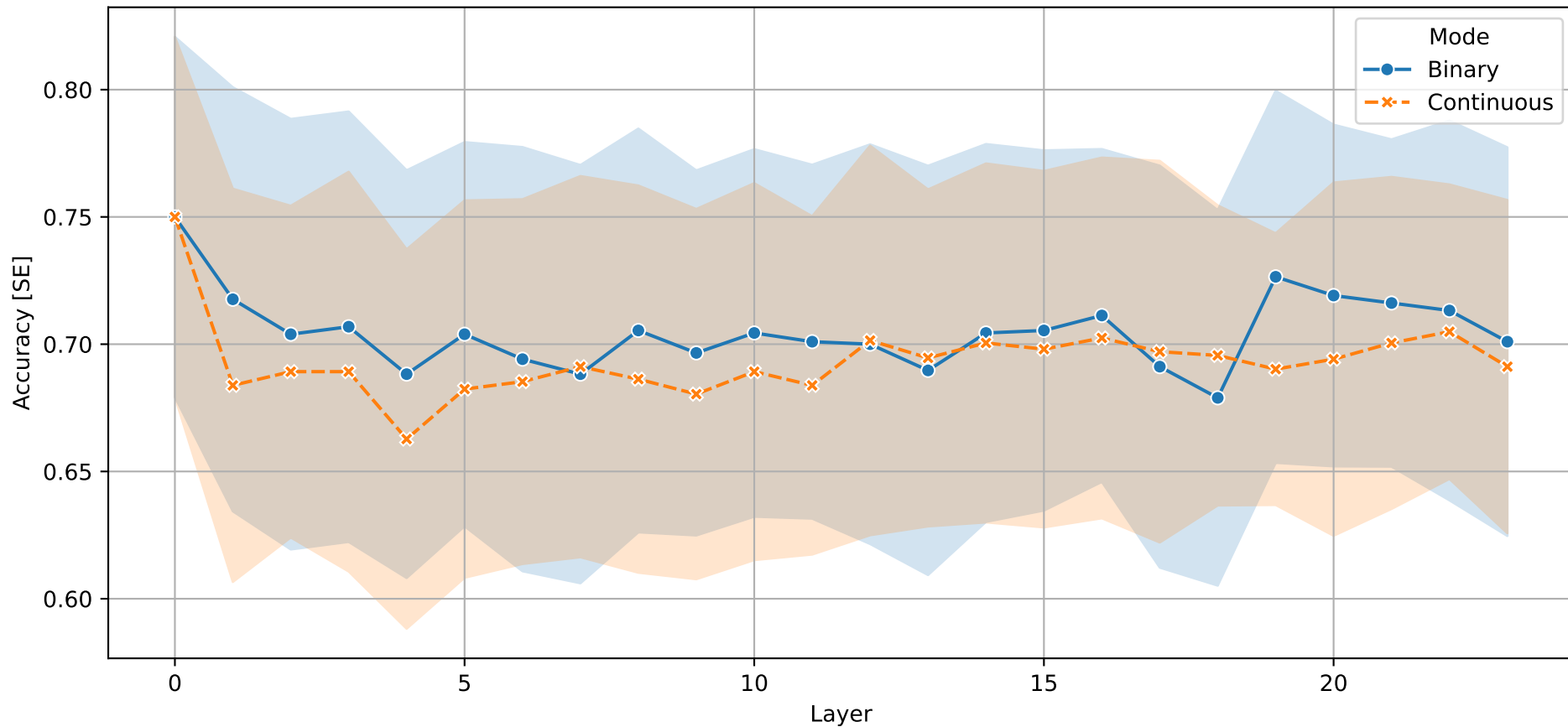| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | f1_best_layer | 19.0 | 22.0 |
| Full Layer | f1_max | 0.8559 | 0.8651 |
| Full Layer | f1_mean | 0.677 | 0.6774 |
| Full Layer | f1_std | 0.1301 | 0.1215 |
| Single Neuron | f1_best_layer | 19.0 | 22.0 |
| Single Neuron | f1_max | 0.8526 | 0.8526 |
| Single Neuron | f1_mean | 0.6495 | 0.6692 |
| Single Neuron | f1_std | 0.1623 | 0.139 |
| Top-K Neurons | f1_best_layer | 22.0 | 20.0 |
| Top-K Neurons | f1_max | 0.8526 | 0.8784 |
| Top-K Neurons | f1_mean | 0.6551 | 0.6799 |
| Top-K Neurons | f1_std | 0.1565 | 0.1313 |

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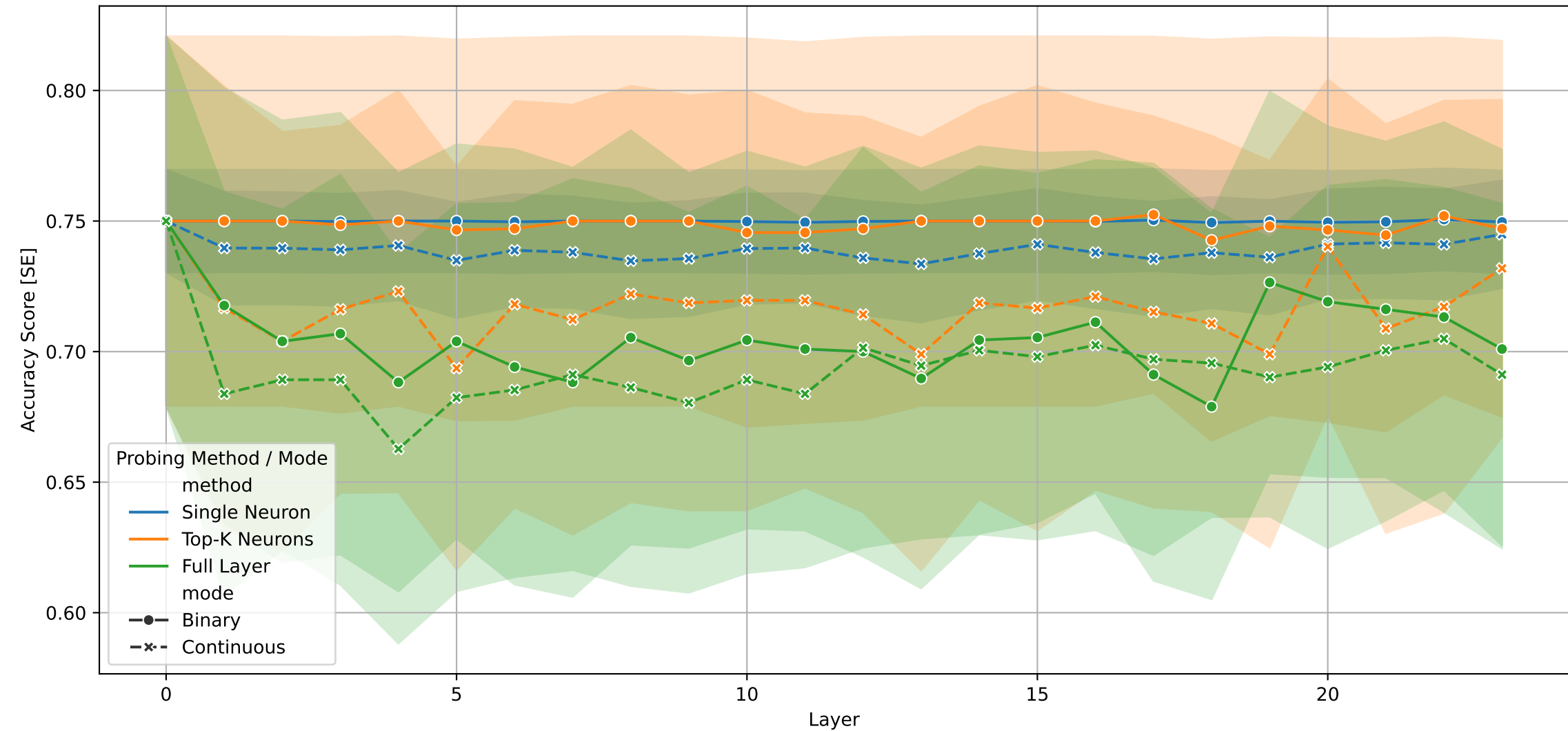
Accuracy per Layer – Top-K Neurons Probing

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 | 0.0 | 0.0 |
| Full Layer | accuracy_max | 0.9 | 0.9 |
| Full Layer | accuracy_mean | 0.7049 | 0.6935 |
| Full Layer | accuracy_std | 0.1328 | 0.1229 |
| Single Neuron | accuracy_best_layer | 22.0 | 0.0 |
| Single Neuron | accuracy_max | 0.9 | 0.9 |
| Single Neuron | accuracy_mean | 0.7499 | 0.7389 |
| Single Neuron | accuracy_std | 0.1227 | 0.1337 |
| Top-K Neurons | accuracy_best_layer | 17.0 | 0.0 |
| Top-K Neurons | accuracy_max | 0.9 | 0.9 |
| Top-K Neurons | accuracy_mean | 0.7485 | 0.7169 |
| Top-K Neurons | accuracy_std | 0.1251 | 0.1337 |