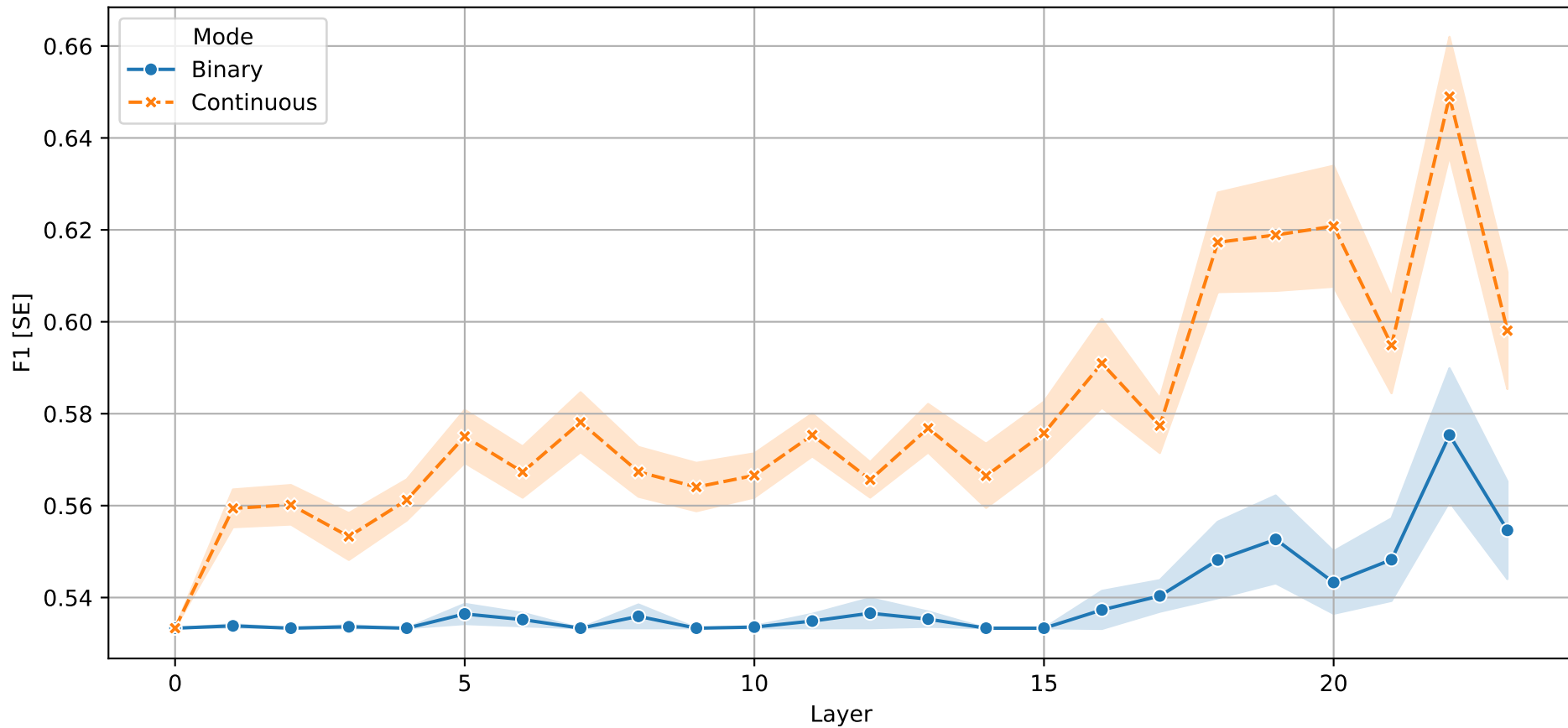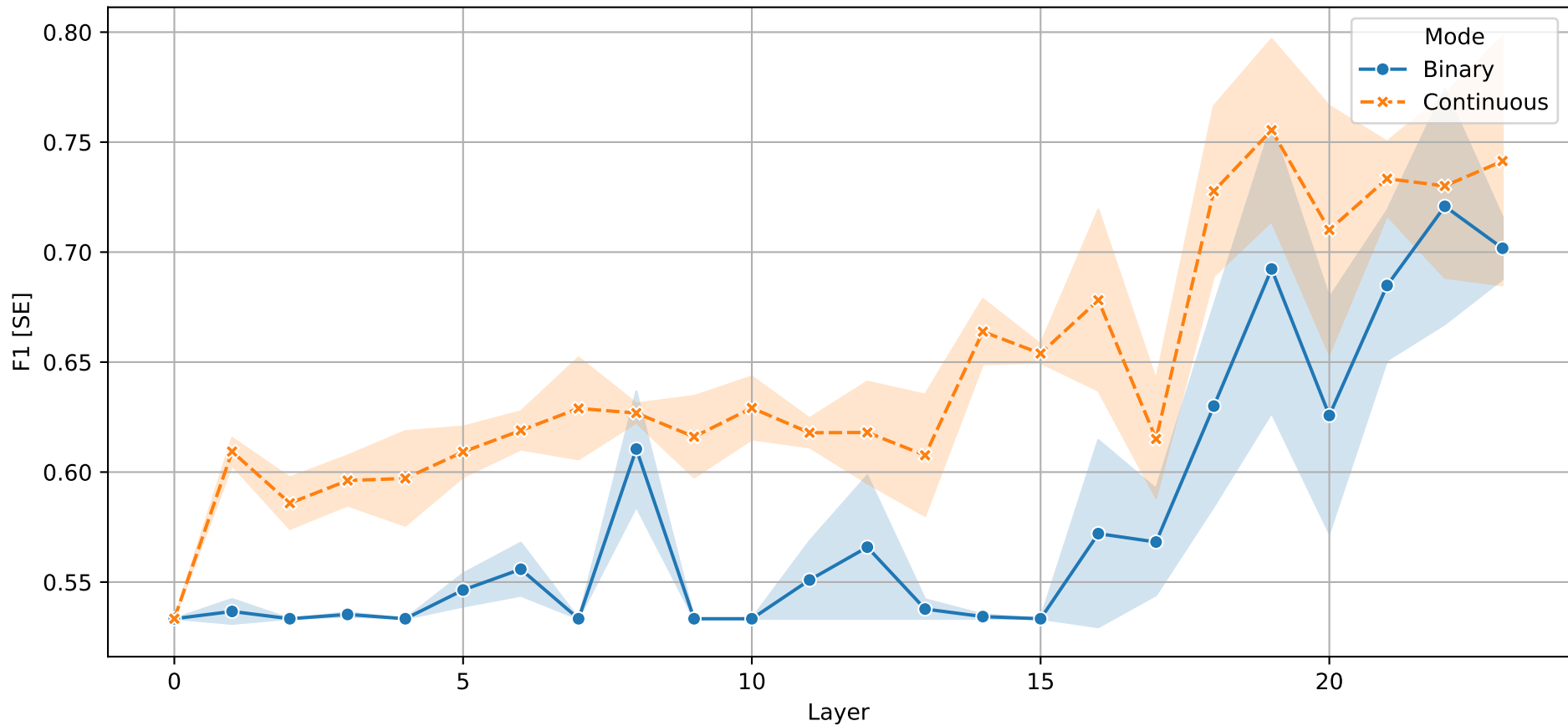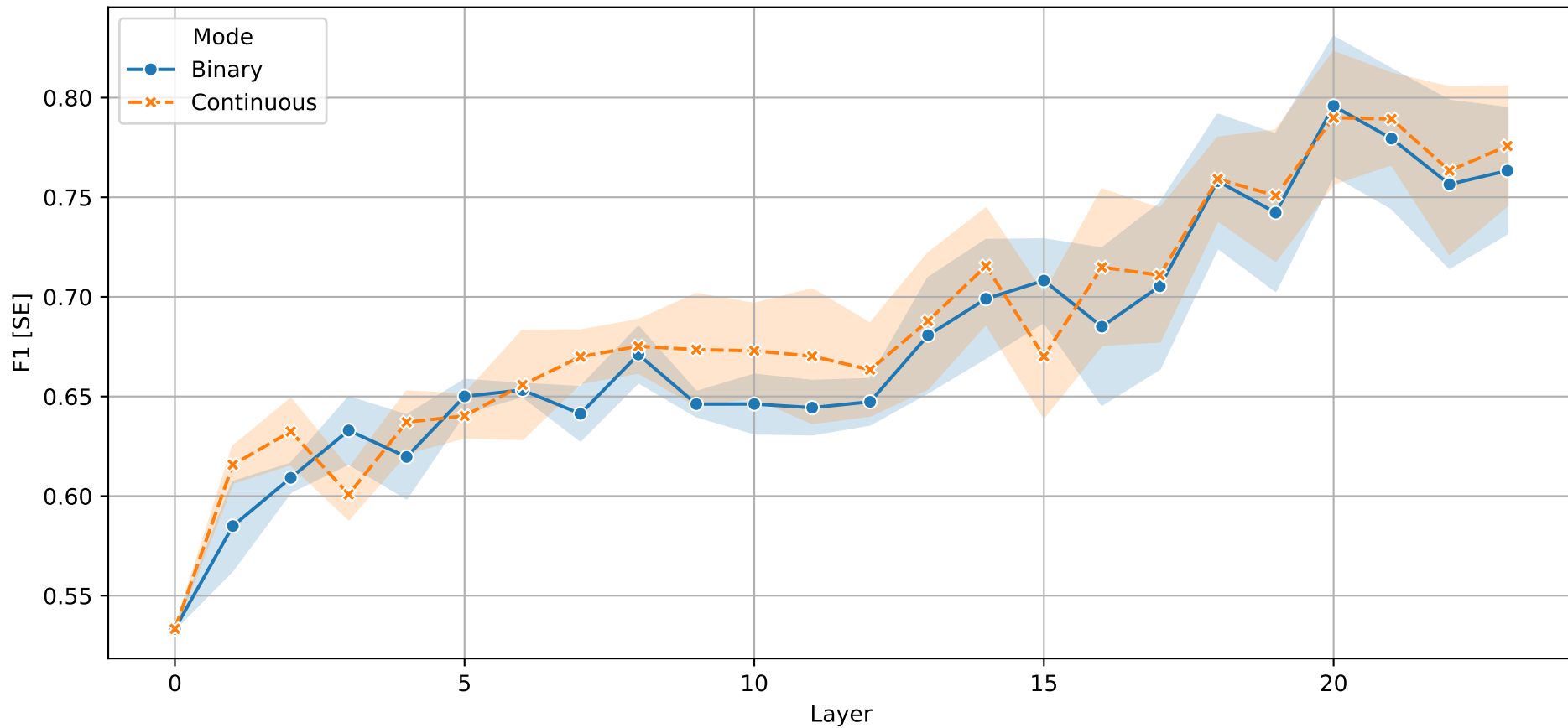F1 per Layer – Single Neuron Probing
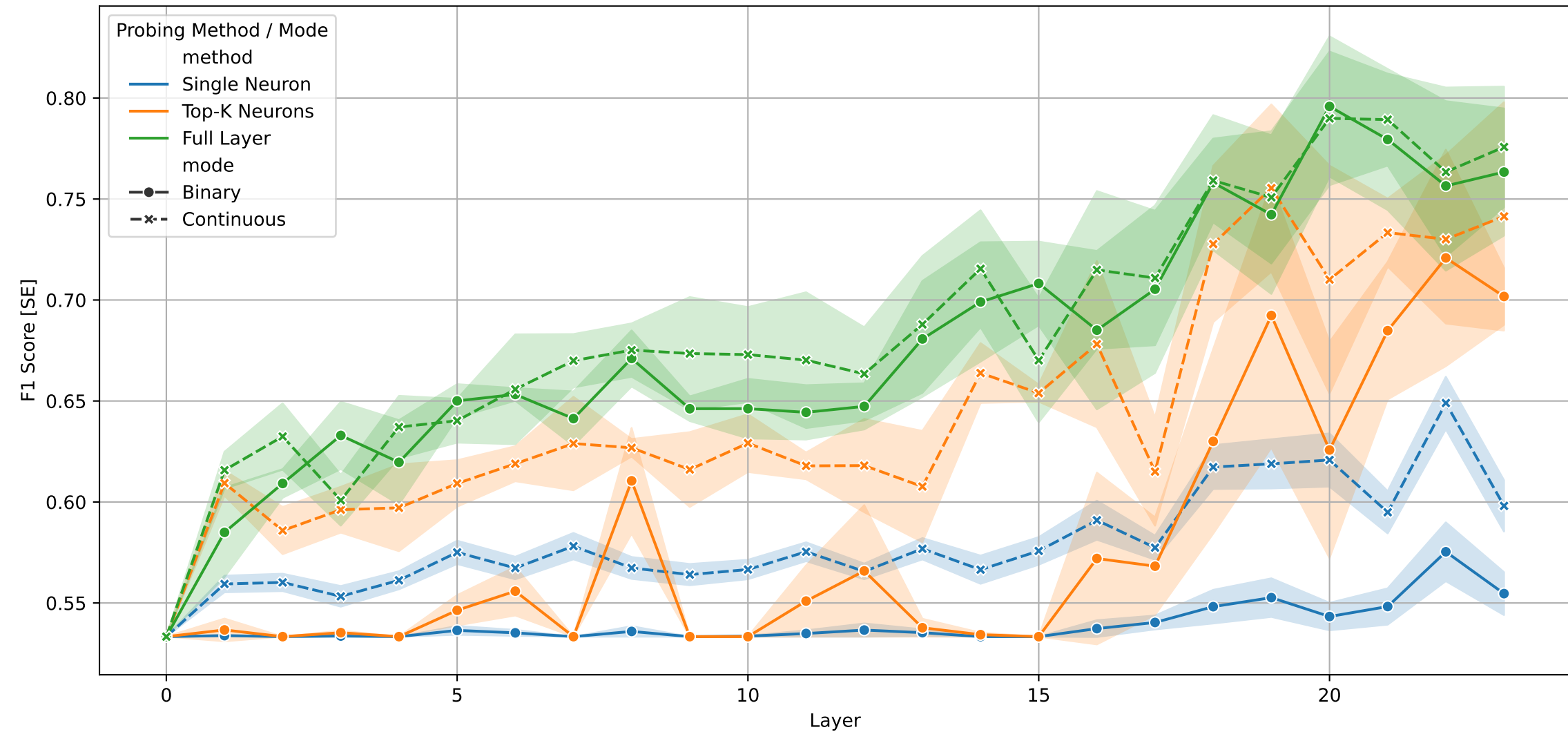
F1 per Layer – Top-K Neurons Probing

F1 per Layer – Full Layer Probing

Overall F1 per Layer – All Methods

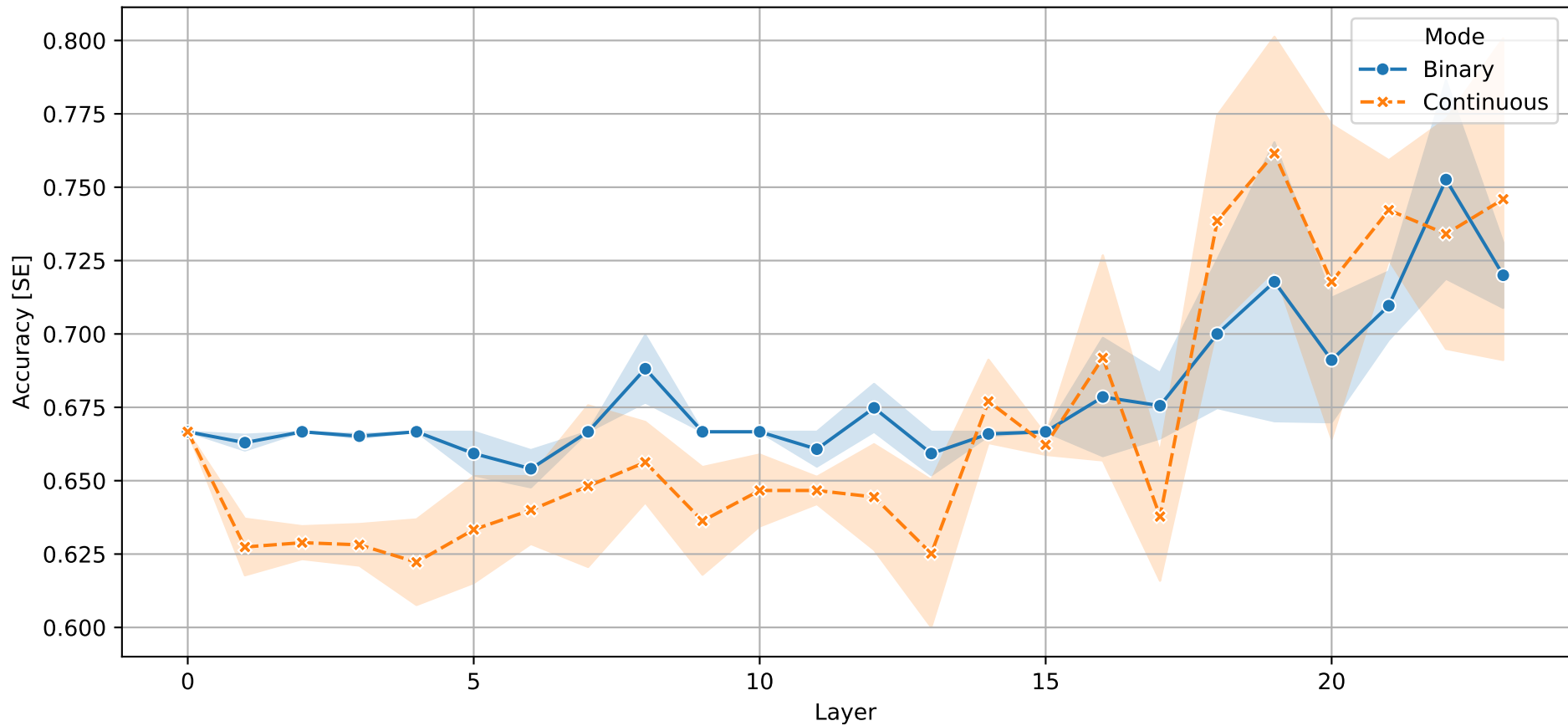## F1 Score Summary by Probing Method

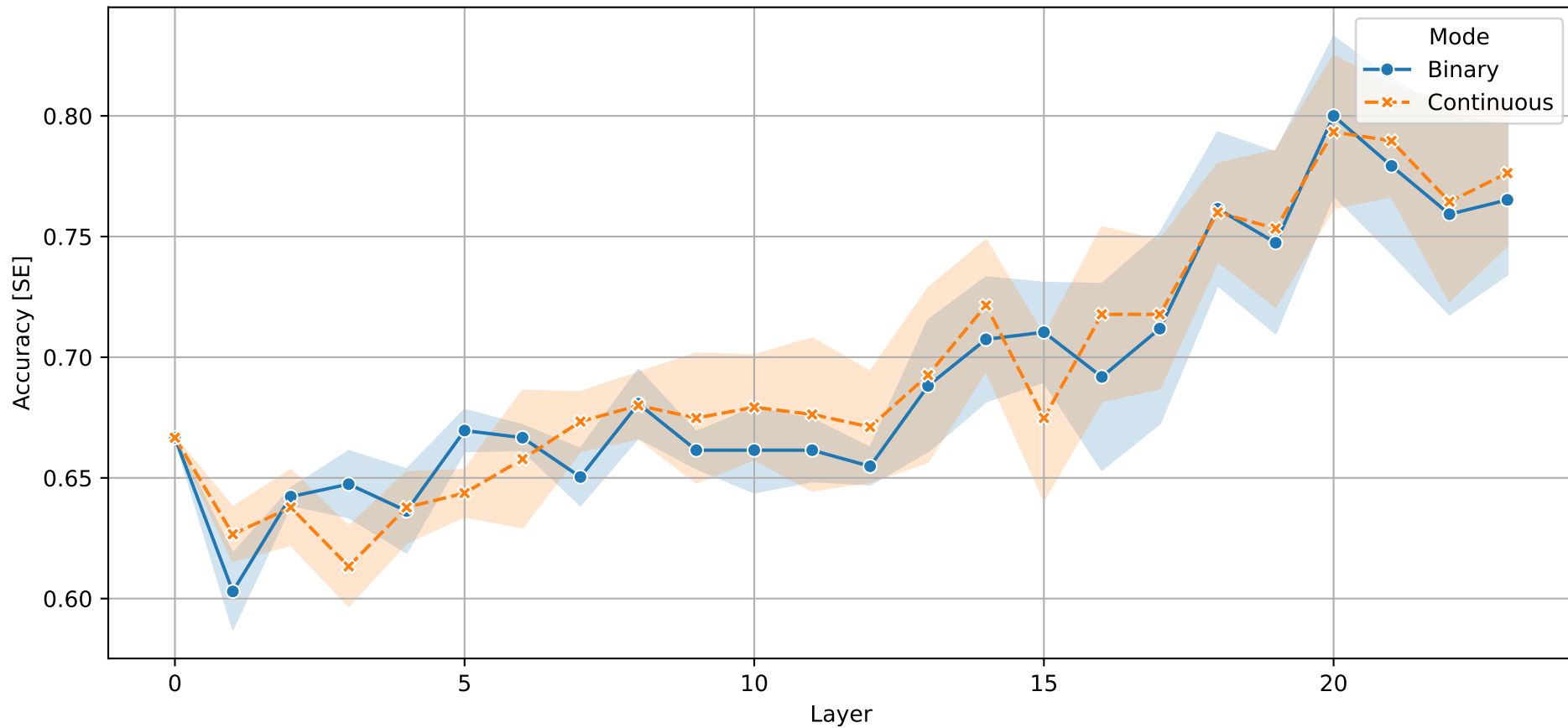| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | f1_best_layer | 20.0 | 20.0 |
| Full Layer | f1_max | 0.8641 | 0.8517 |
| Full Layer | f1_mean | 0.6772 | 0.6862 |
| Full Layer | f1_std | 0.0735 | 0.0728 |
| Single Neuron | f1_best_layer | 22.0 | 22.0 |
| Single Neuron | f1_max | 0.7915 | 0.8365 |
| Single Neuron | f1_mean | 0.5395 | 0.5797 |
| Single Neuron | f1_std | 0.0301 | 0.0487 |
| Top-K Neurons | f1_best_layer | 22.0 | 19.0 |
| Top-K Neurons | f1_max | 0.8117 | 0.8528 |
| Top-K Neurons | f1_mean | 0.5793 | 0.646 |
| Top-K Neurons | f1_std | 0.073 | 0.0693 |

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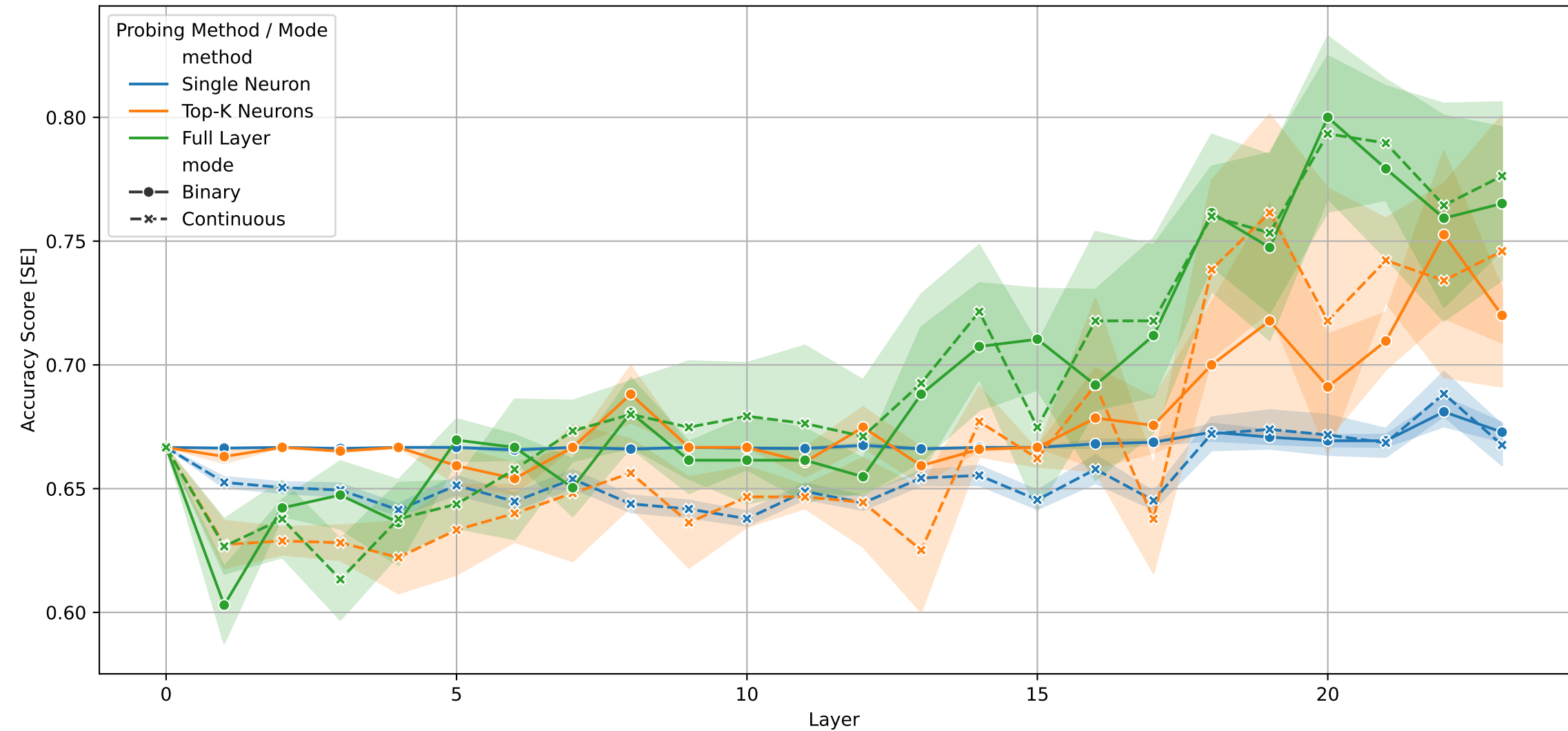
Accuracy per Layer – Top-K Neurons Probing

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 | 20.0 | 20.0 |
| Full Layer | accuracy_max | 0.8644 | 0.8533 |
| Full Layer | accuracy_mean | 0.6927 | 0.6958 |
| Full Layer | accuracy_std | 0.0615 | 0.0638 |
| Single Neuron | accuracy_best_layer | 22.0 | 22.0 |
| Single Neuron | accuracy_max | 0.7889 | 0.8378 |
| Single Neuron | accuracy_mean | 0.6682 | 0.6553 |
| Single Neuron | accuracy_std | 0.0109 | 0.0288 |
| Top-K Neurons | accuracy_best_layer | 22.0 | 19.0 |
| Top-K Neurons | accuracy_max | 0.8156 | 0.8533 |
| Top-K Neurons | accuracy_mean | 0.6793 | 0.6691 |
| Top-K Neurons | accuracy_std | 0.0325 | 0.0577 |