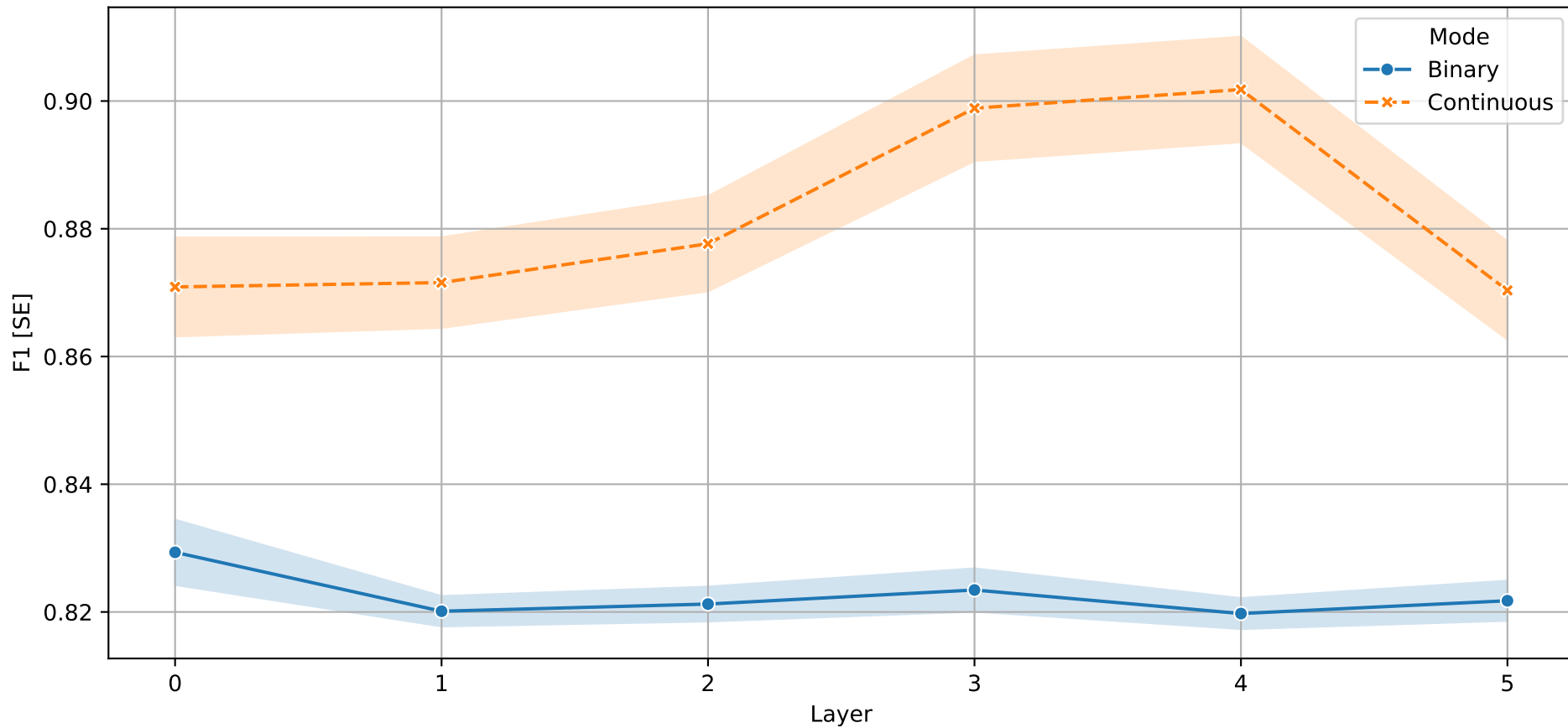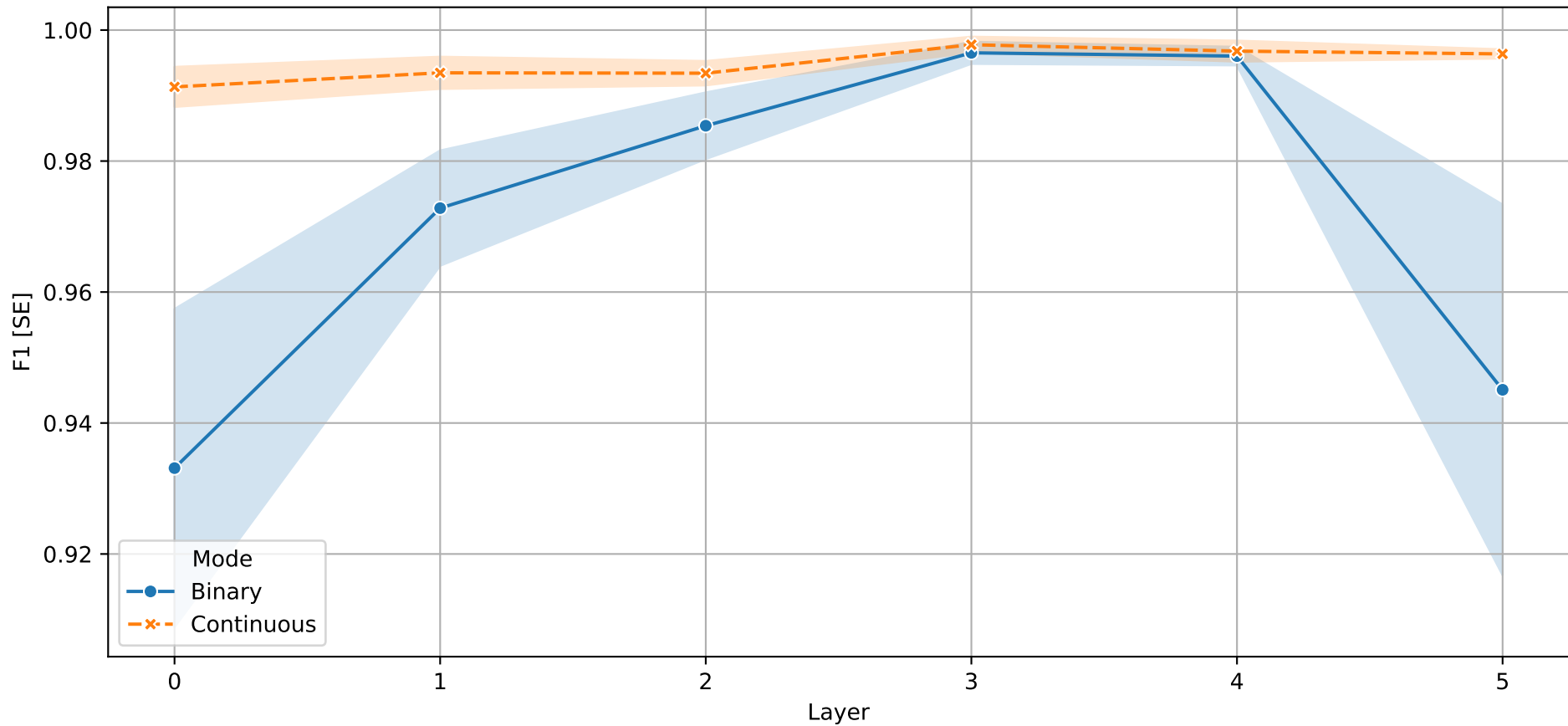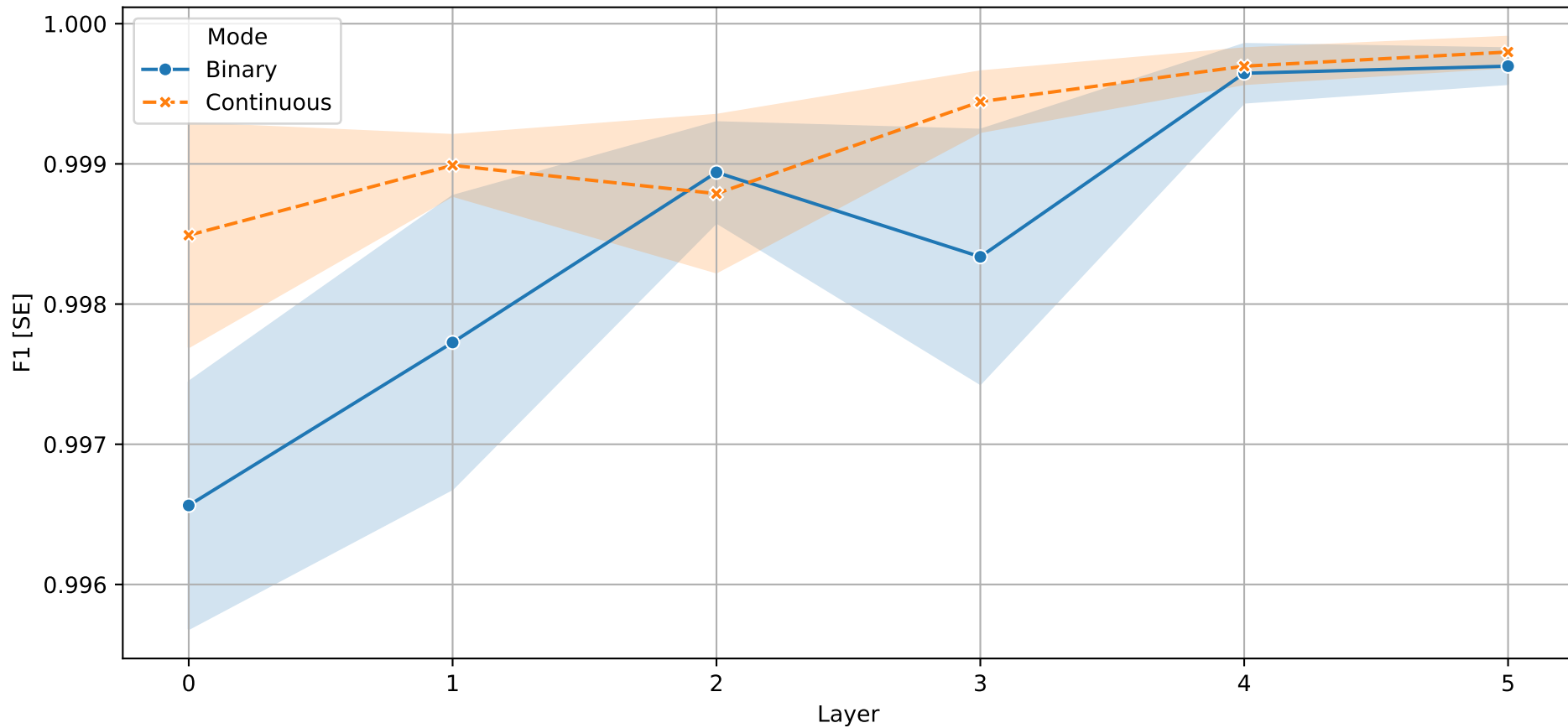F1 per Layer – Single Neuron Probing
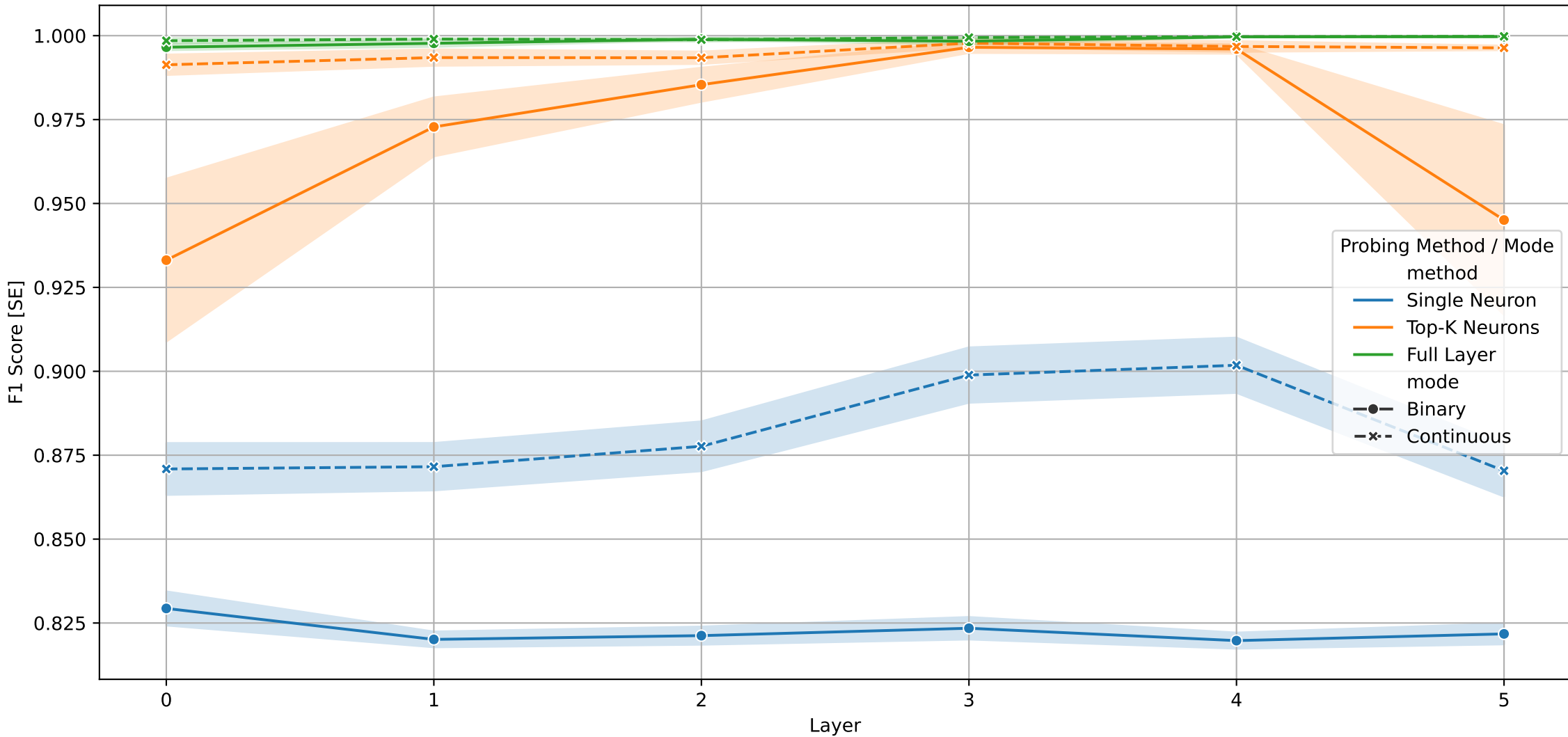
F1 per Layer – Top-K Neurons Probing

F1 per Layer – Full Layer Probing

Overall F1 per Layer – All Methods

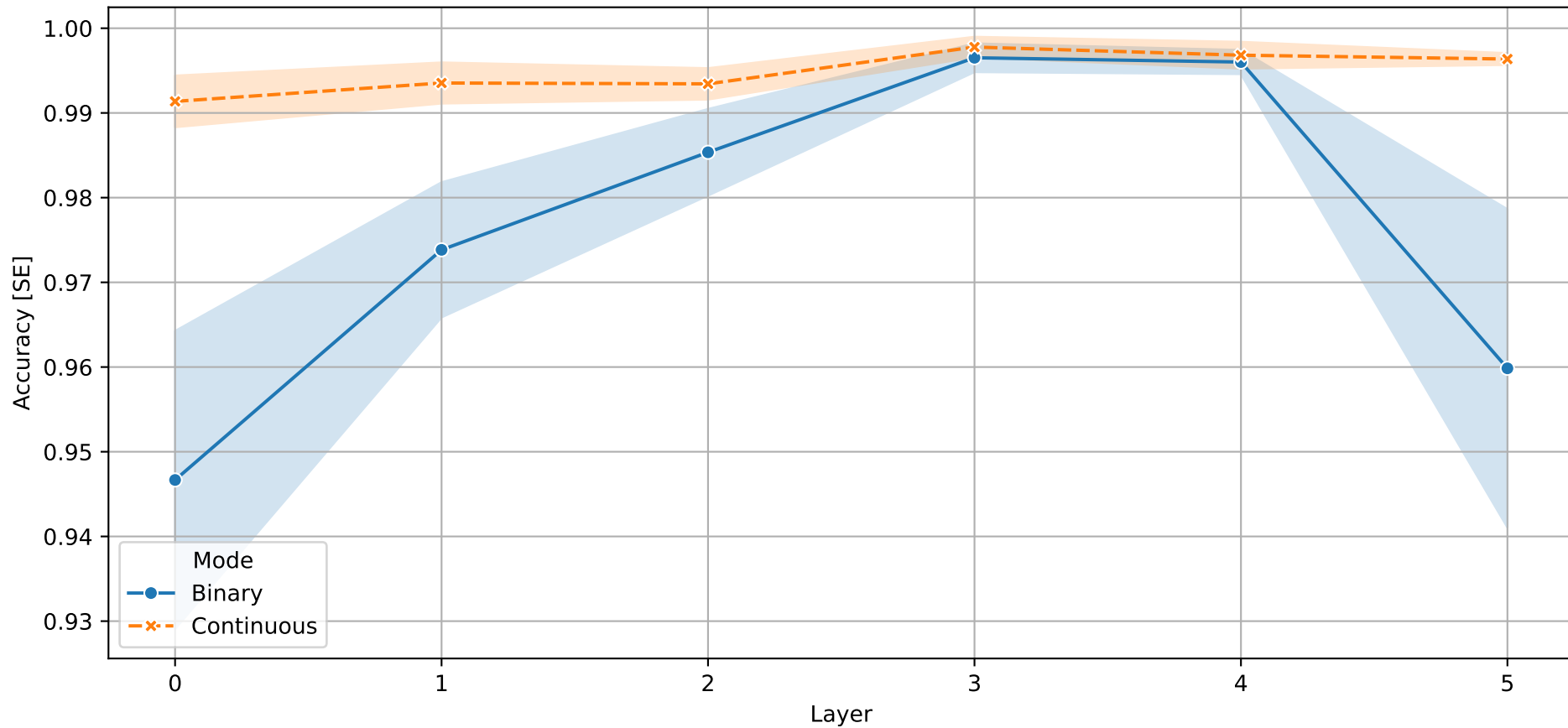## F1 Score Summary by Probing Method

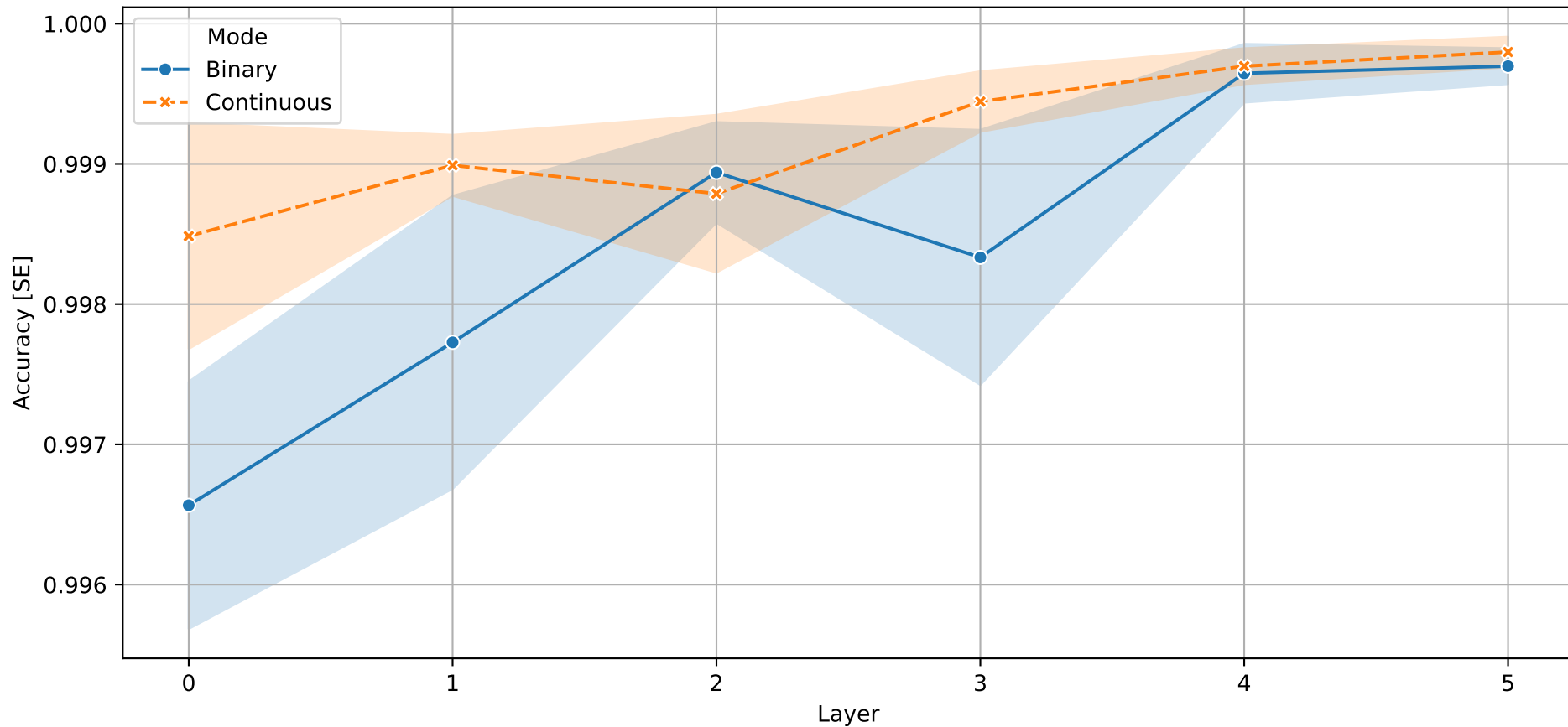| Method | Metric | Binary | Continuous |
|---|---|---|---|
| Full Layer | f1_best_layer | 5.0 | 5.0 |
| Full Layer | f1_max | 1.0 | 1.0 |
| Full Layer | f1_mean | 0.9985 | 0.9992 |
| Full Layer | f1_std | 0.0022 | 0.0012 |
| Single Neuron | f1_best_layer | 0.0 | 4.0 |
| Single Neuron | f1_max | 1.0 | 1.0 |
| Single Neuron | f1_mean | 0.8226 | 0.8819 |
| Single Neuron | f1_std | 0.0296 | 0.0702 |
| Top-K Neurons | f1_best_layer | 3.0 | 3.0 |
| Top-K Neurons | f1_max | 1.0 | 1.0 |
| Top-K Neurons | f1_mean | 0.9715 | 0.9949 |
| Top-K Neurons | f1_std | 0.0489 | 0.0057 |

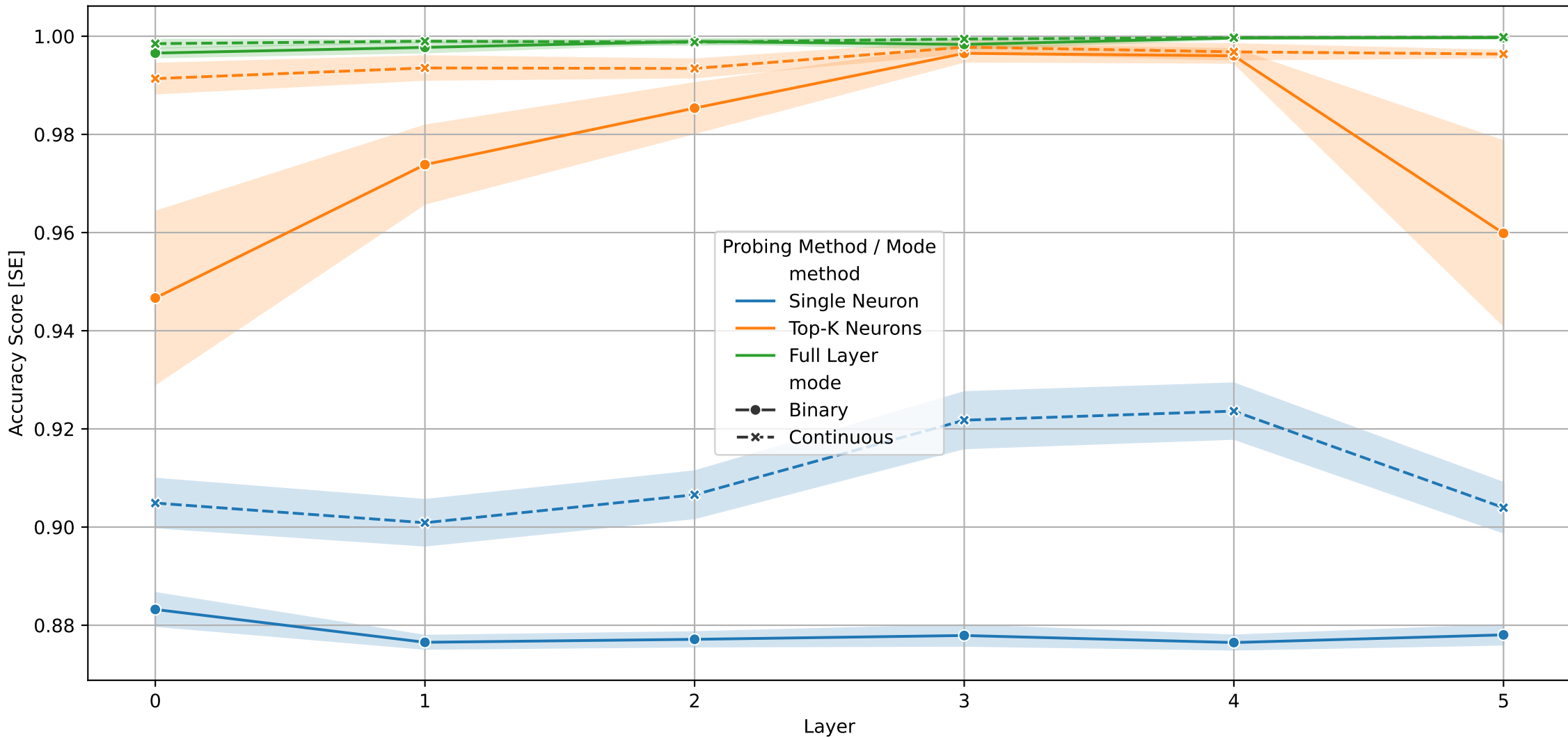Accuracy per Layer – Single Neuron Probing

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 | 5.0 | 5.0 |
| Full Layer | accuracy_max | 1.0 | 1.0 |
| Full Layer | accuracy_mean | 0.9985 | 0.9992 |
| Full Layer | accuracy_std | 0.0022 | 0.0012 |
| Single Neuron | accuracy_best_layer | 0.0 | 4.0 |
| Single Neuron | accuracy_max | 1.0 | 1.0 |
| Single Neuron | accuracy_mean | 0.8782 | 0.9103 |
| Single Neuron | accuracy_std | 0.0186 | 0.0468 |
| Top-K Neurons | accuracy_best_layer | 3.0 | 3.0 |
| Top-K Neurons | accuracy_max | 1.0 | 1.0 |
| Top-K Neurons | accuracy_mean | 0.9764 | 0.9949 |
| Top-K Neurons | accuracy_std | 0.0353 | 0.0057 |