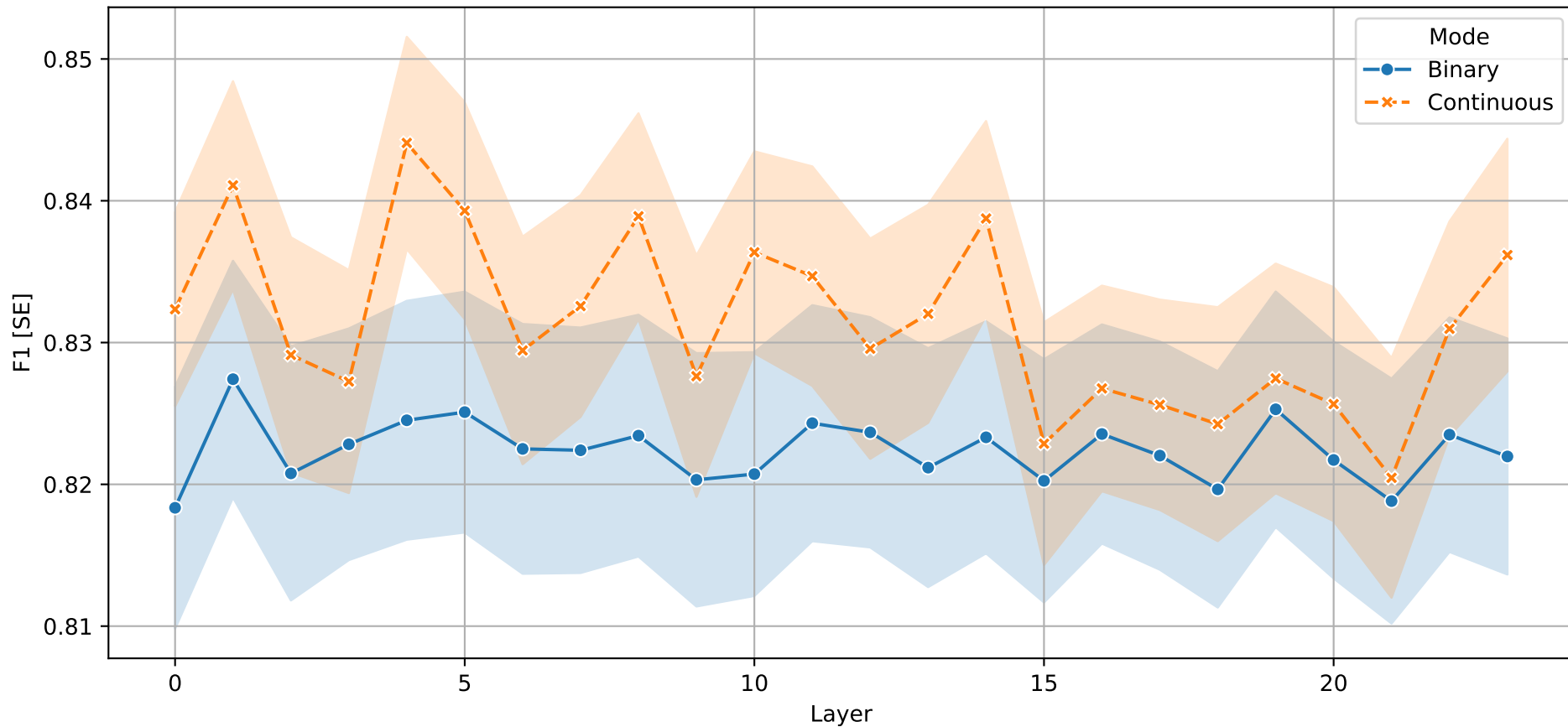
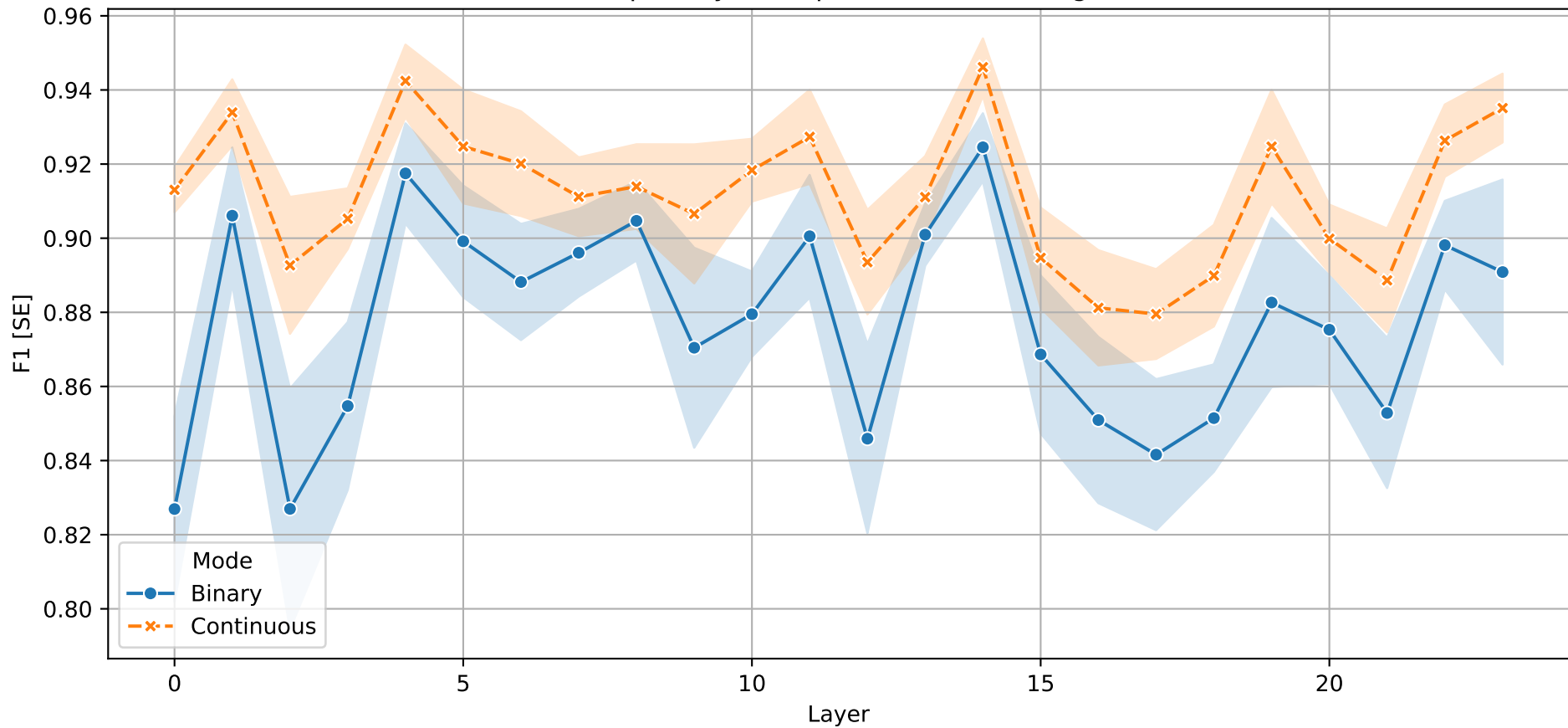


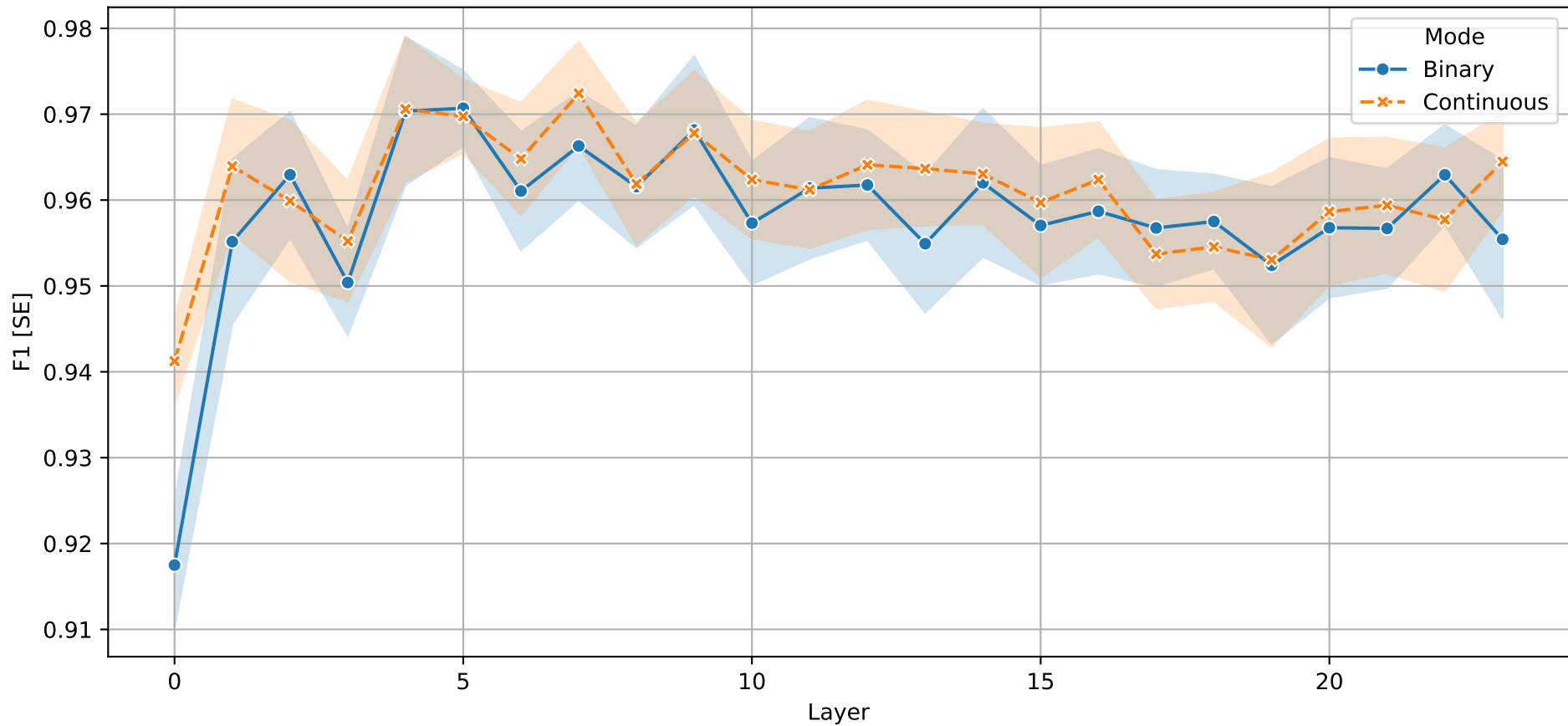
# F1 per Layer - Single Neuron Probing



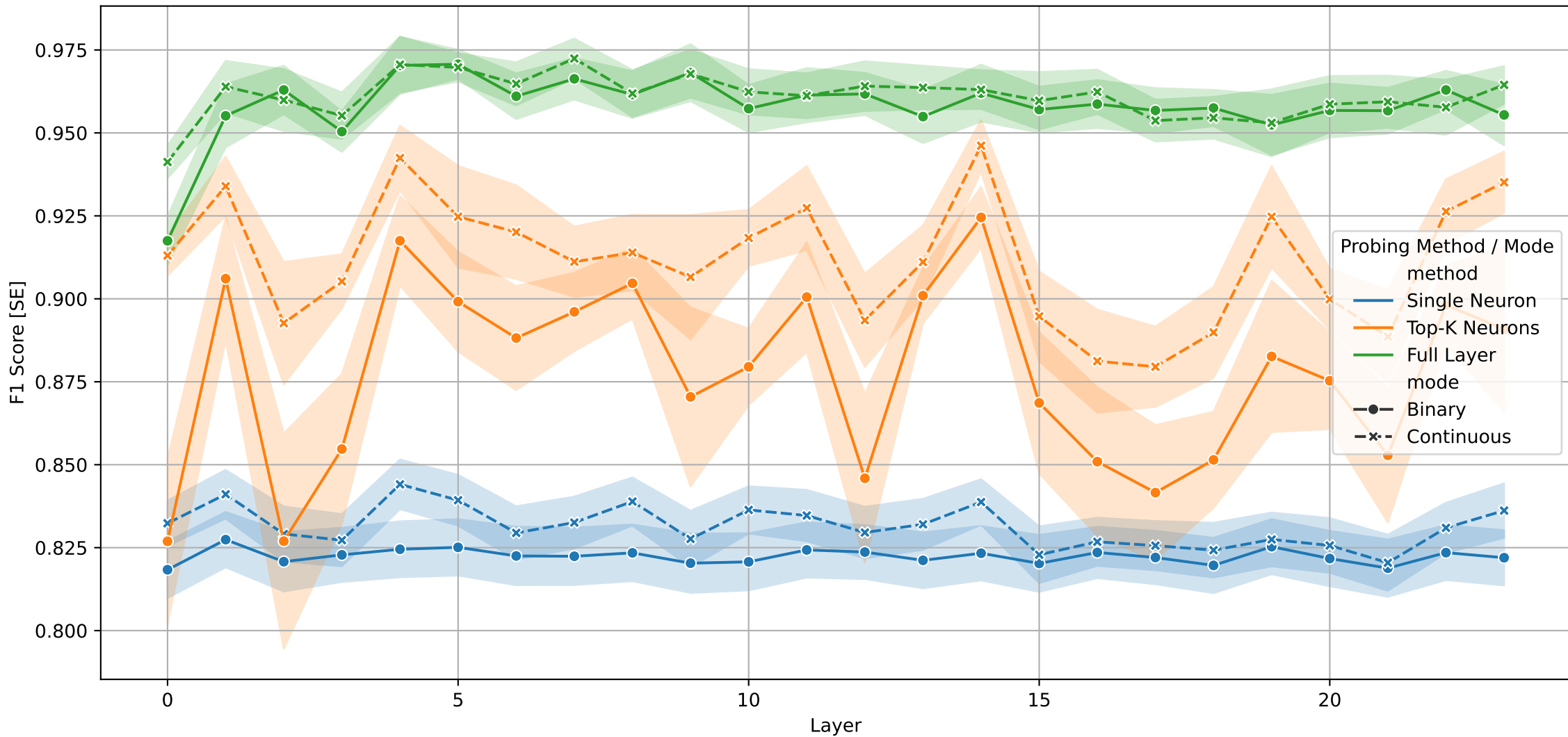
# F1 per Layer - Top-K Neurons Probing



# F1 per Layer - Full Layer Probing



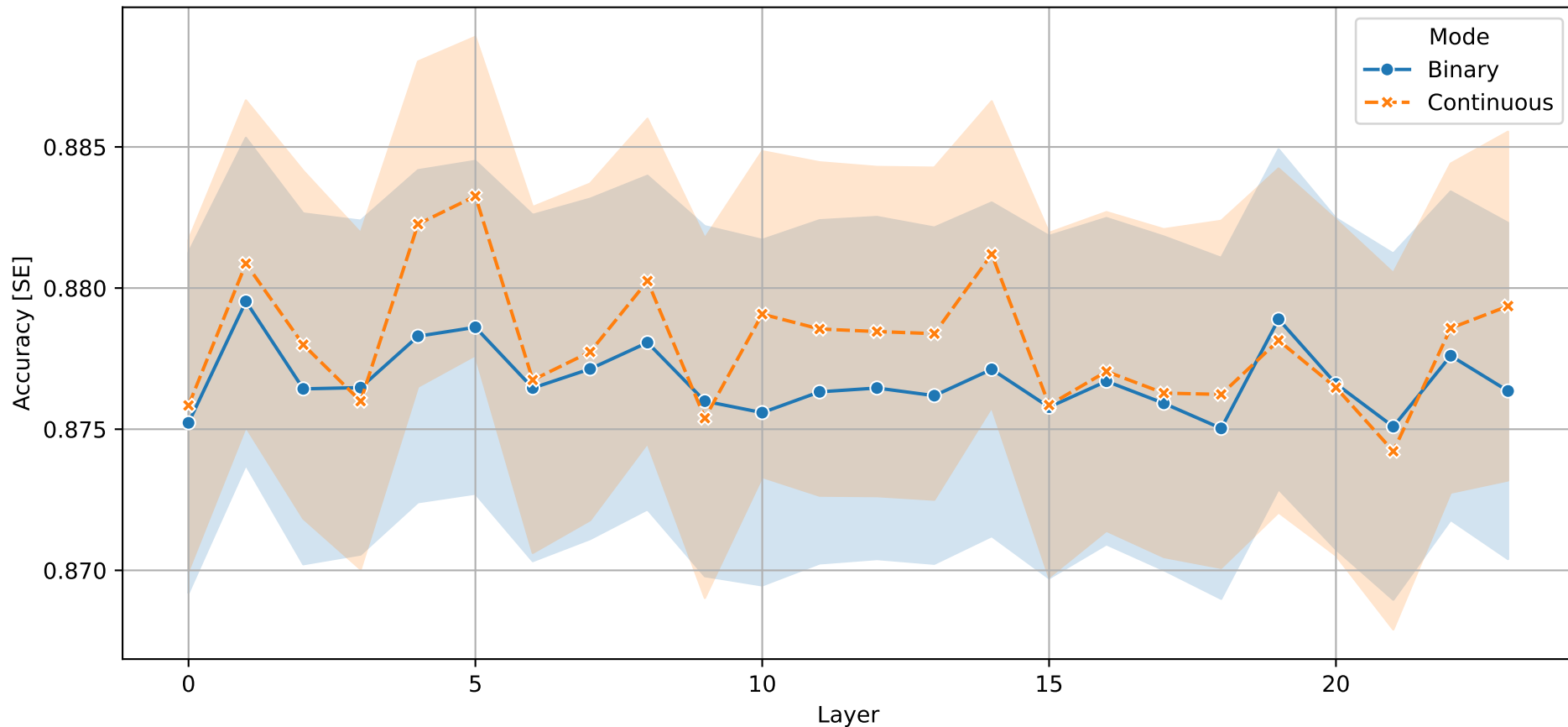
Overall F1 per Layer - All Methods



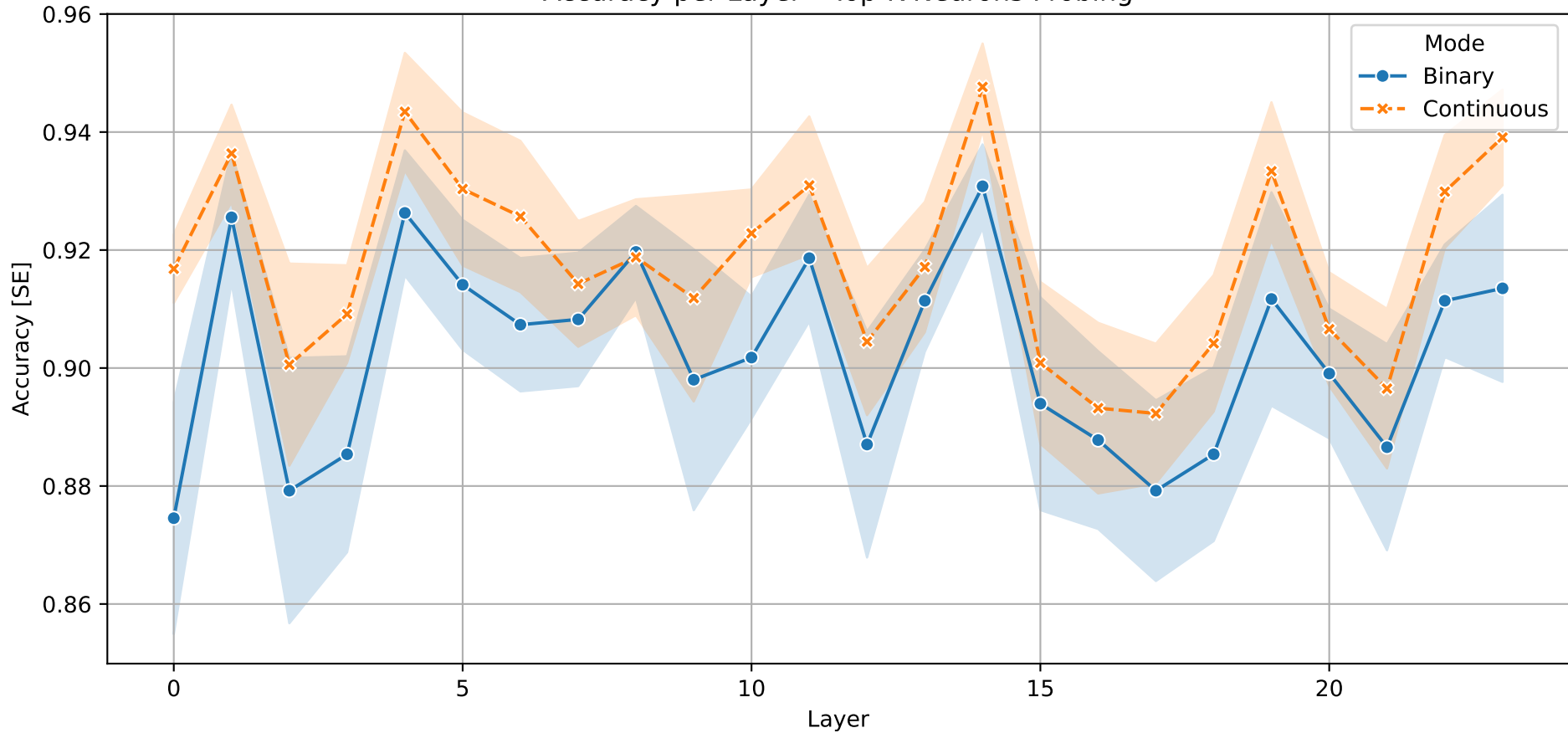
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	5.0	7.0
Full Layer	f1_max	0.9988	0.9988
Full Layer	f1_mean	0.9582	0.9611
Full Layer	f1_std	0.022	0.0202
Single Neuron	f1_best_layer	1.0	4.0
Single Neuron	f1_max	0.9976	0.9988
Single Neuron	f1_mean	0.8224	0.8314
Single Neuron	f1_std	0.0749	0.0694
Top-K Neurons	f1_best_layer	14.0	14.0
Top-K Neurons	f1_max	0.9976	0.9988
Top-K Neurons	f1_mean	0.8773	0.9117
Top-K Neurons	f1_std	0.0579	0.0379

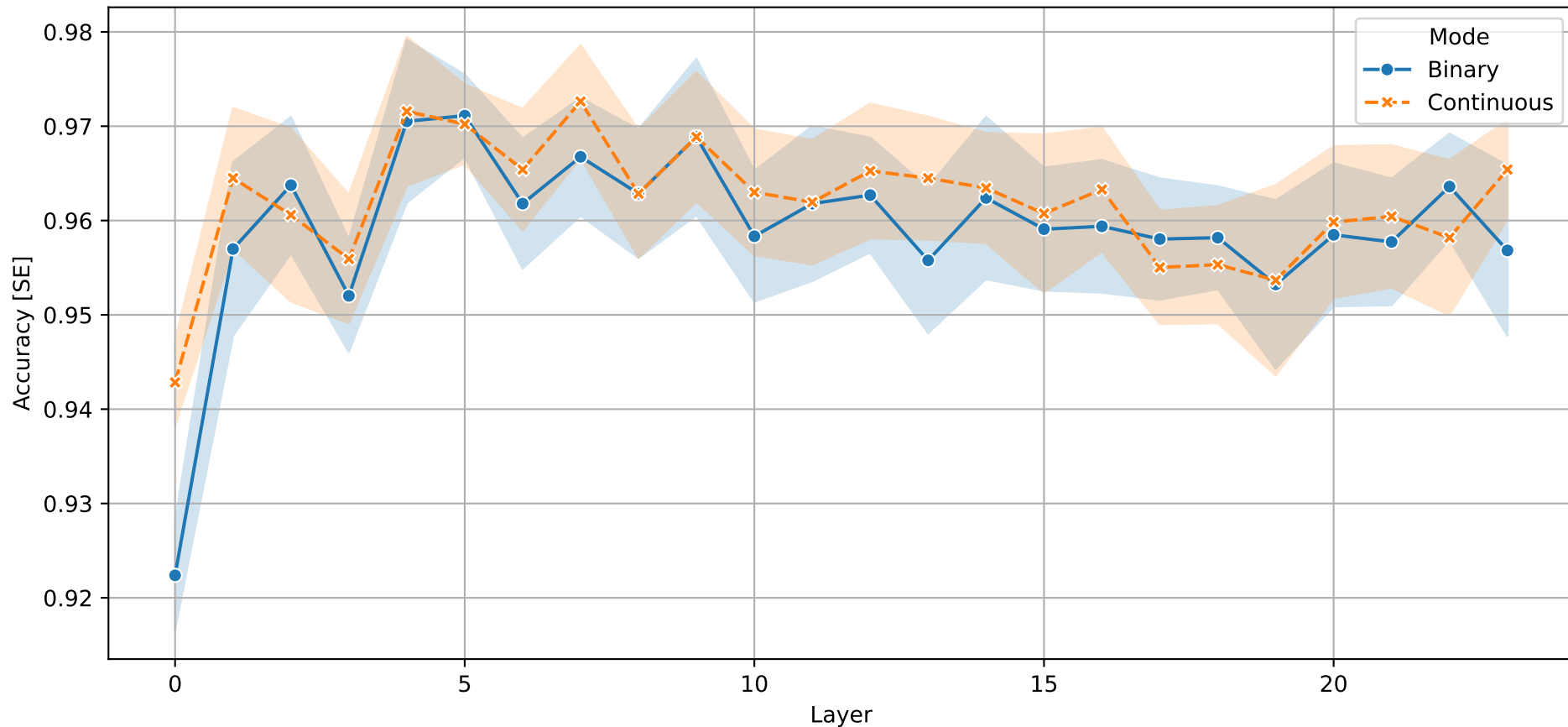
Accuracy per Layer - Single Neuron Probing



Accuracy per Layer - Top-K Neurons Probing

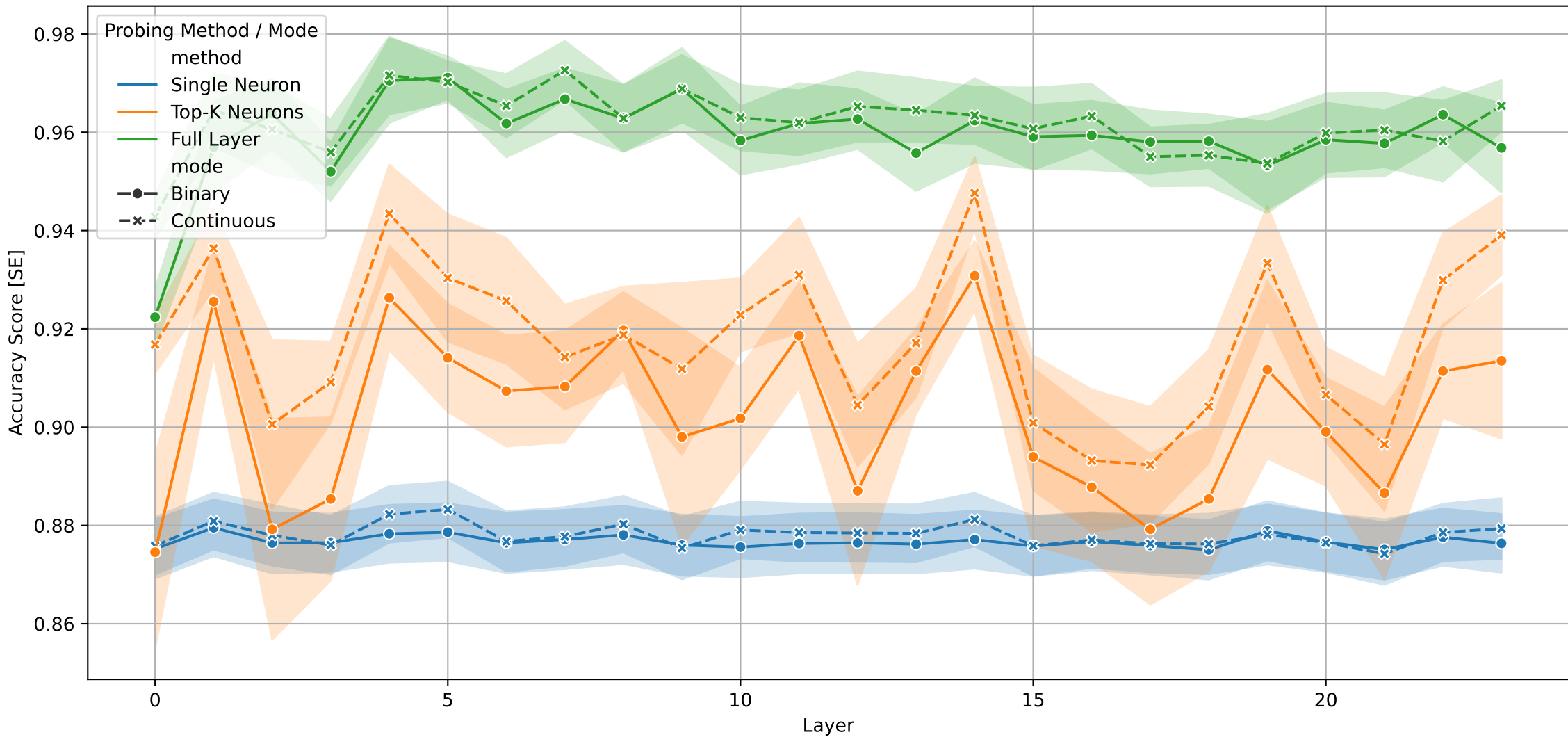


Accuracy per Layer - Full Layer Probing





Overall Accuracy per Layer - All Methods



Accuracy Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	accuracy_best_layer	5.0	7.0
Full Layer	accuracy_max	0.9988	0.9988
Full Layer	accuracy_mean	0.9593	0.9619
Full Layer	accuracy_std	0.0211	0.0196
Single Neuron	accuracy_best_layer	1.0	5.0
Single Neuron	accuracy_max	0.9976	0.9988
Single Neuron	accuracy_mean	0.8767	0.8781
Single Neuron	accuracy_std	0.0532	0.0528
Top-K Neurons	accuracy_best_layer	14.0	14.0
Top-K Neurons	accuracy_max	0.9976	0.9988
Top-K Neurons	accuracy_mean	0.9024	0.9178
Top-K Neurons	accuracy_std	0.0421	0.0344