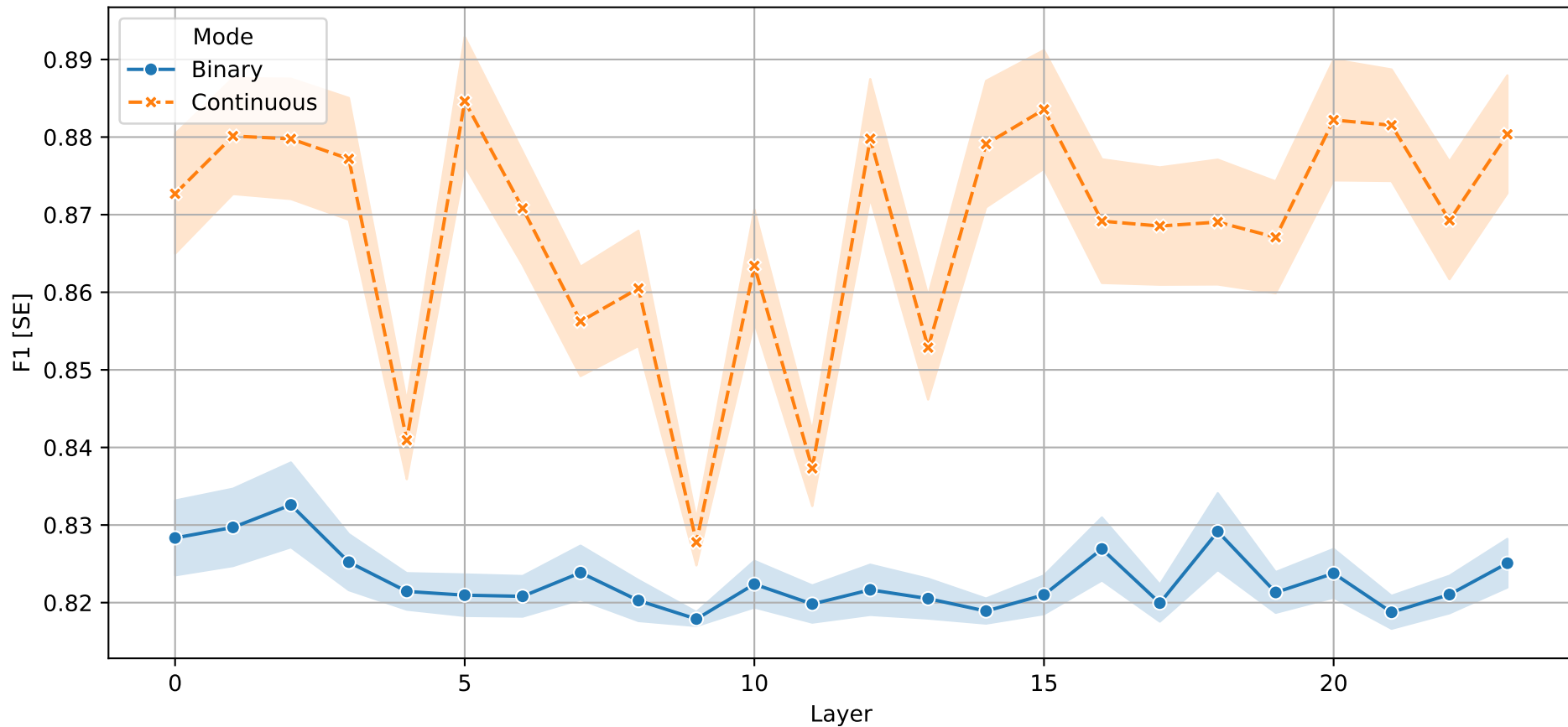
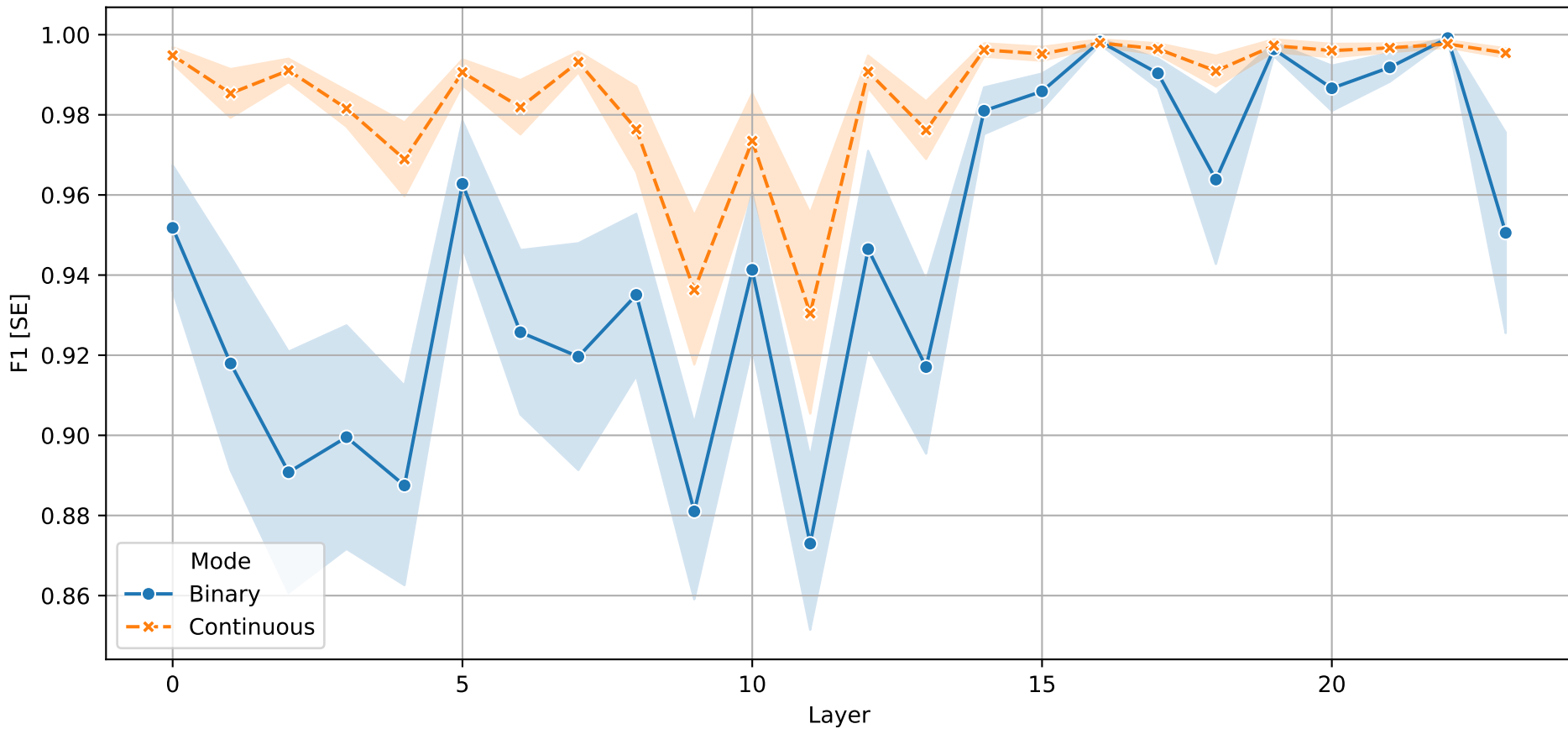


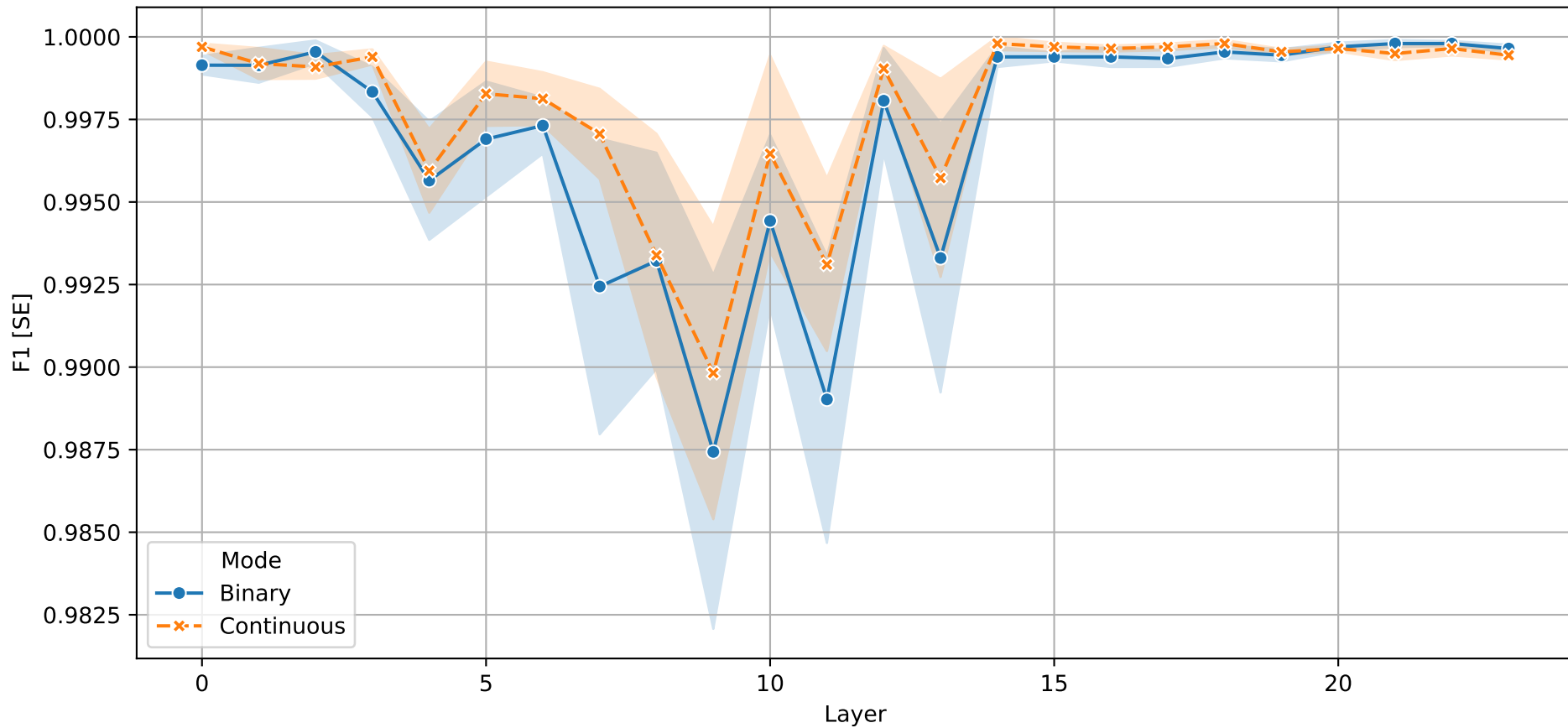
# F1 per Layer - Single Neuron Probing



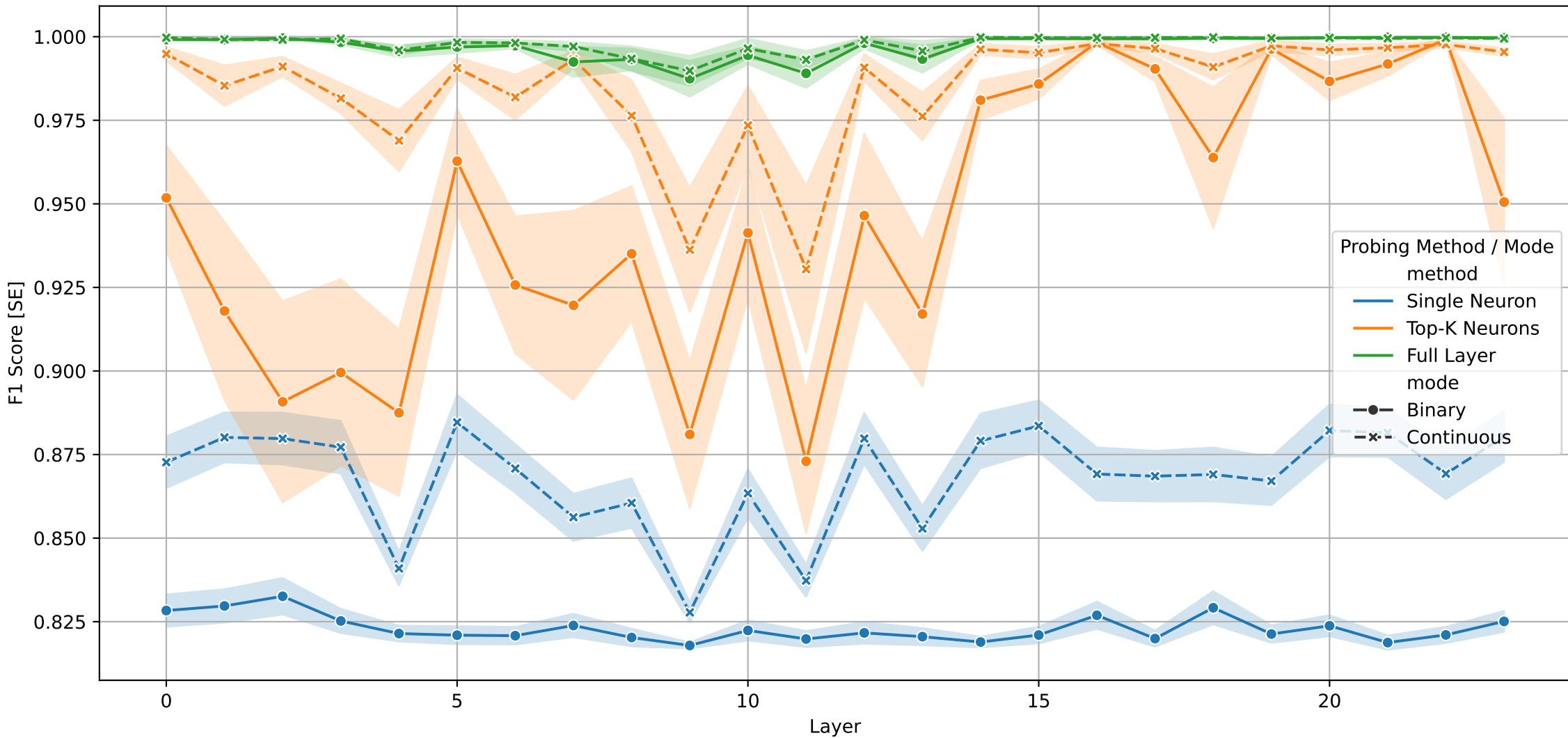
# F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



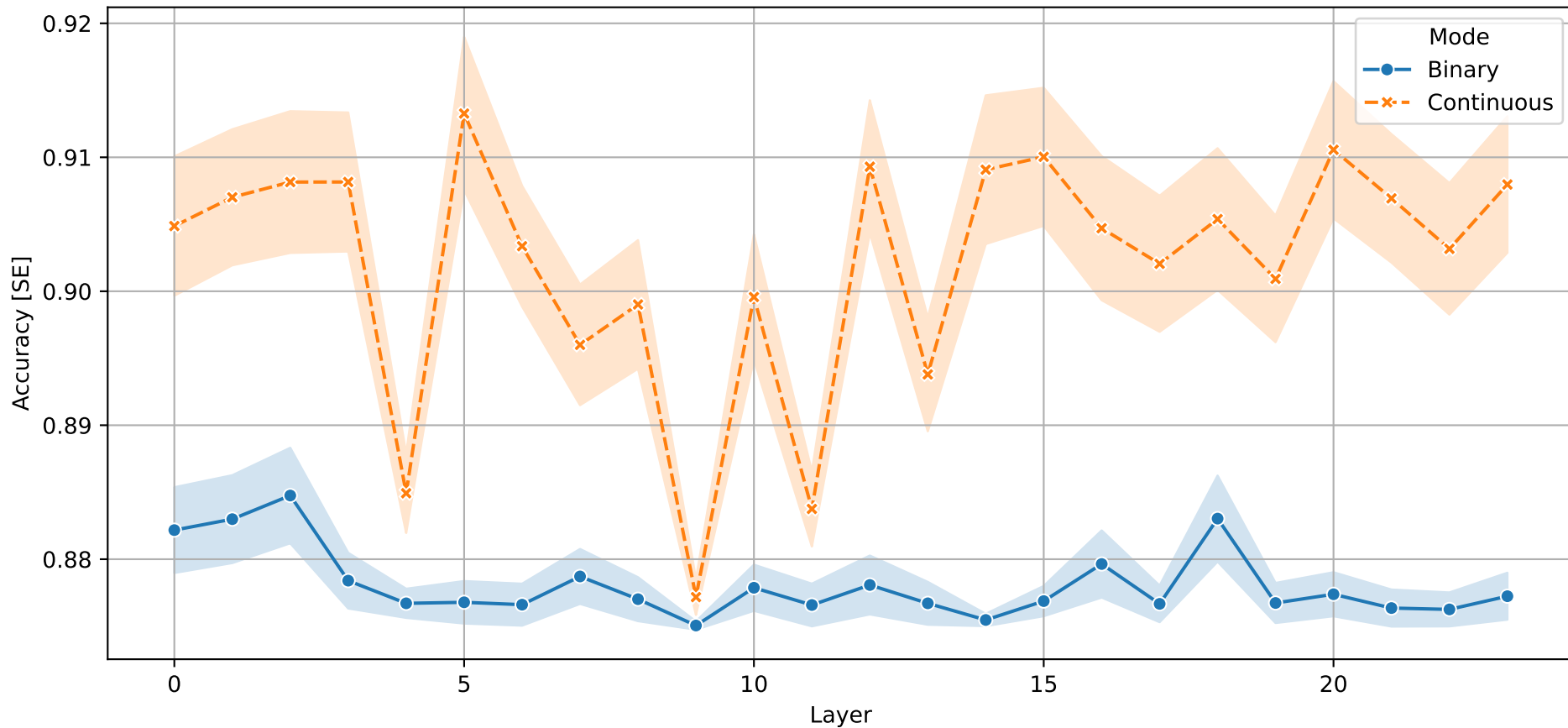
Overall F1 per Layer - All Methods



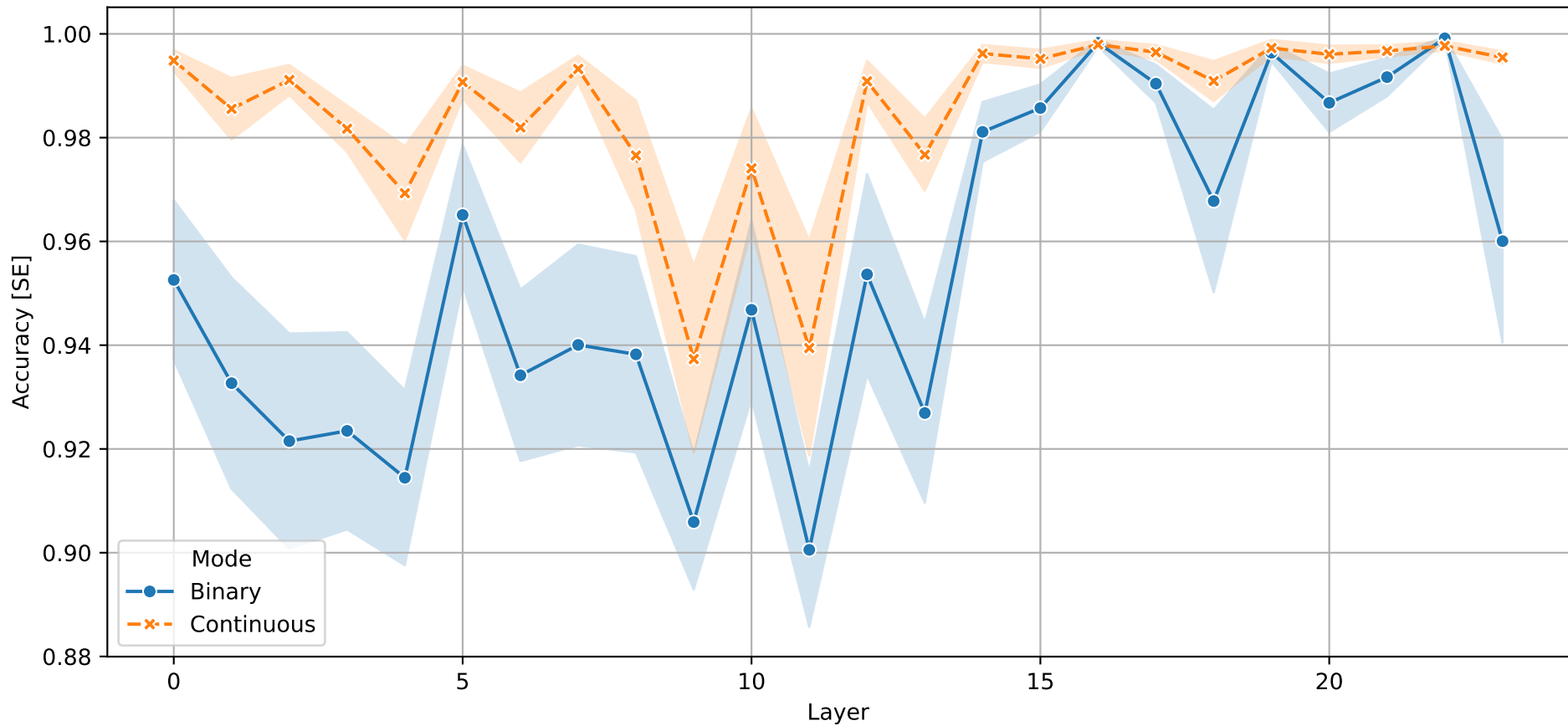
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	21.0	18.0
Full Layer	f1_max	1.0	1.0
Full Layer	f1_mean	0.9971	0.9979
Full Layer	f1_std	0.0067	0.005
Single Neuron	f1_best_layer	2.0	5.0
Single Neuron	f1_max	1.0	1.0
Single Neuron	f1_mean	0.823	0.8681
Single Neuron	f1_std	0.0292	0.0663
Top-K Neurons	f1_best_layer	22.0	16.0
Top-K Neurons	f1_max	1.0	1.0
Top-K Neurons	f1_mean	0.9456	0.9846
Top-K Neurons	f1_std	0.0645	0.0276

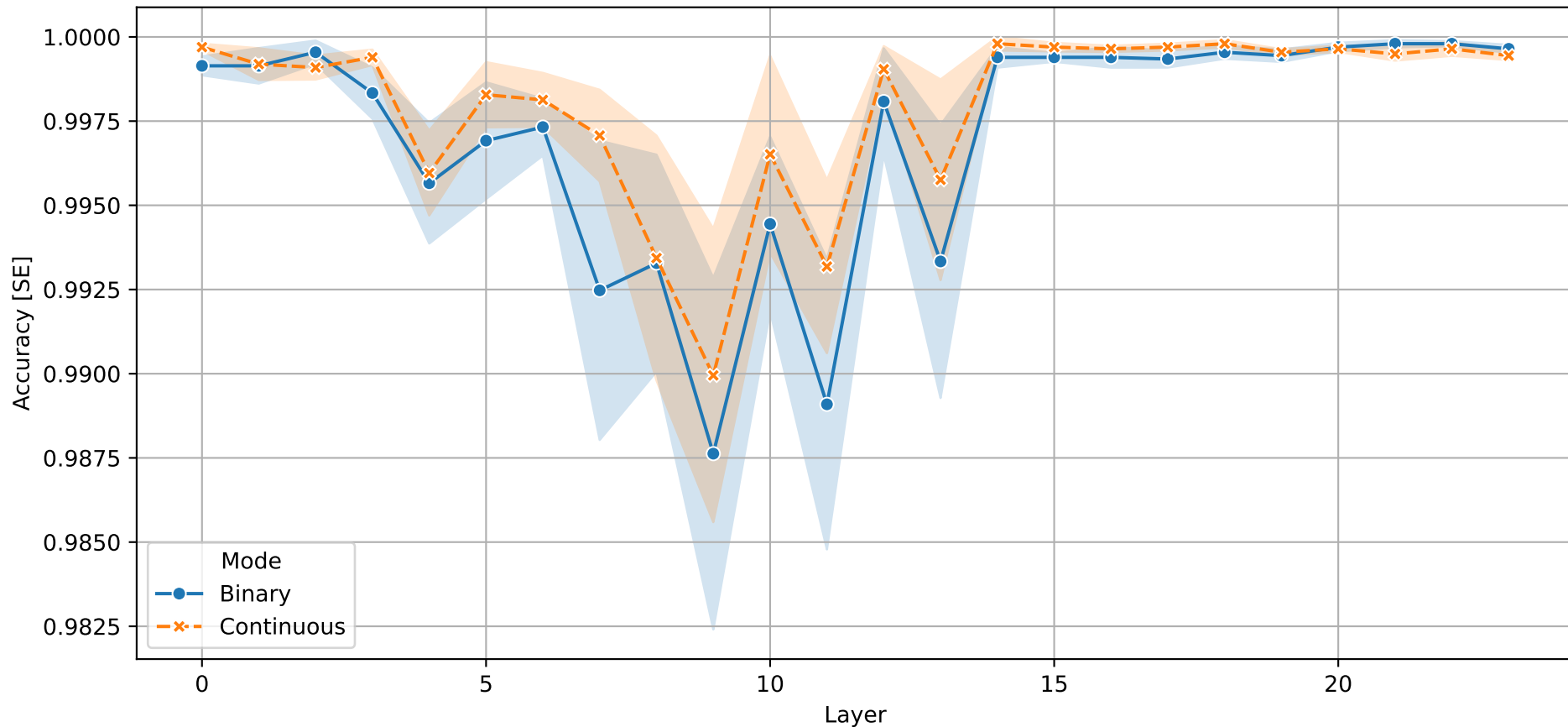
Accuracy per Layer - Single Neuron Probing



Accuracy per Layer - Top-K Neurons Probing

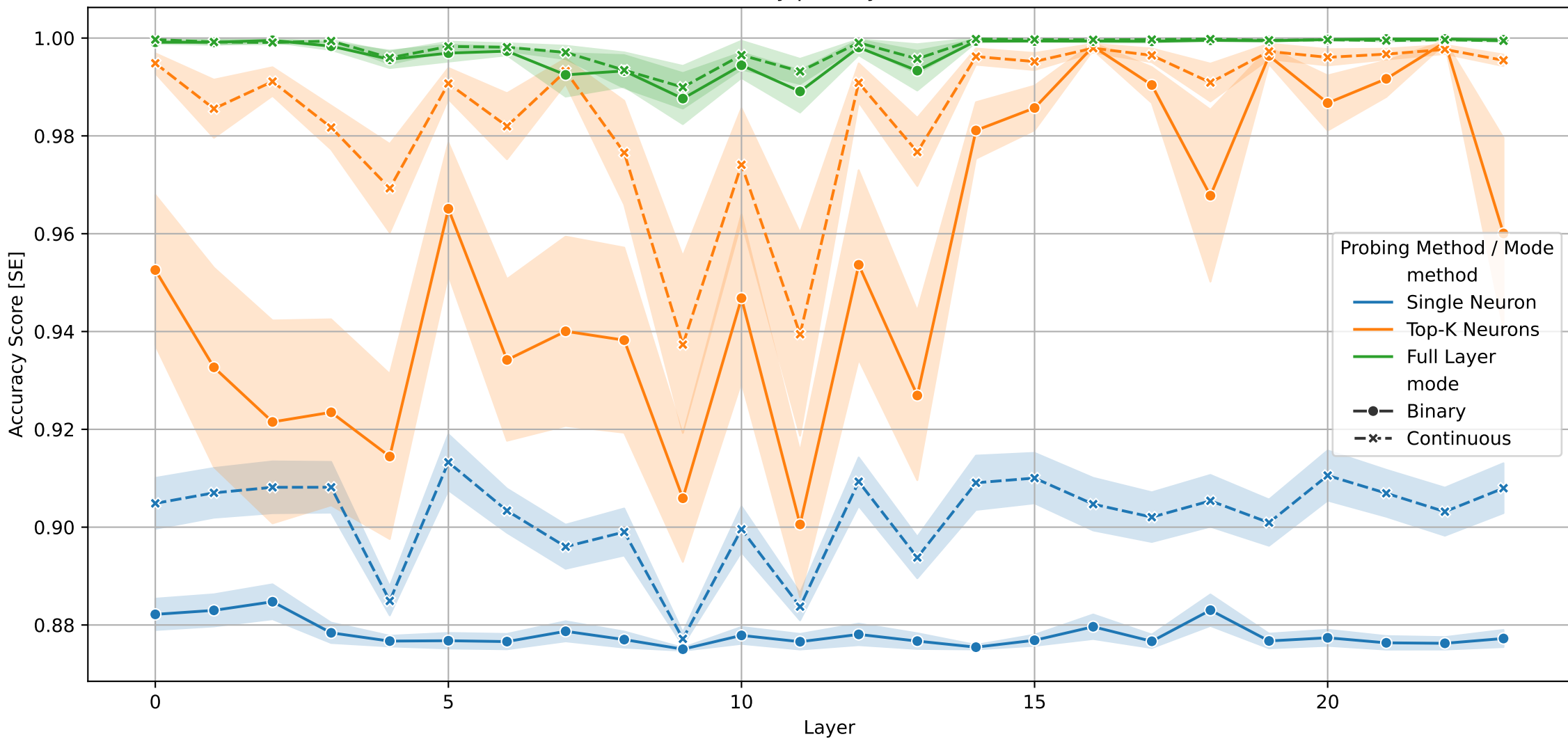


Accuracy per Layer - Full Layer Probing





Overall Accuracy per Layer - All Methods



Accuracy Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	accuracy_best_layer	21.0	14.0
Full Layer	accuracy_max	1.0	1.0
Full Layer	accuracy_mean	0.9971	0.998
Full Layer	accuracy_std	0.0067	0.005
Single Neuron	accuracy_best_layer	2.0	5.0
Single Neuron	accuracy_max	1.0	1.0
Single Neuron	accuracy_mean	0.8781	0.902
Single Neuron	accuracy_std	0.0179	0.0434
Top-K Neurons	accuracy_best_layer	22.0	16.0
Top-K Neurons	accuracy_max	1.0	1.0
Top-K Neurons	accuracy_mean	0.9547	0.9851
Top-K Neurons	accuracy_std	0.0493	0.0255