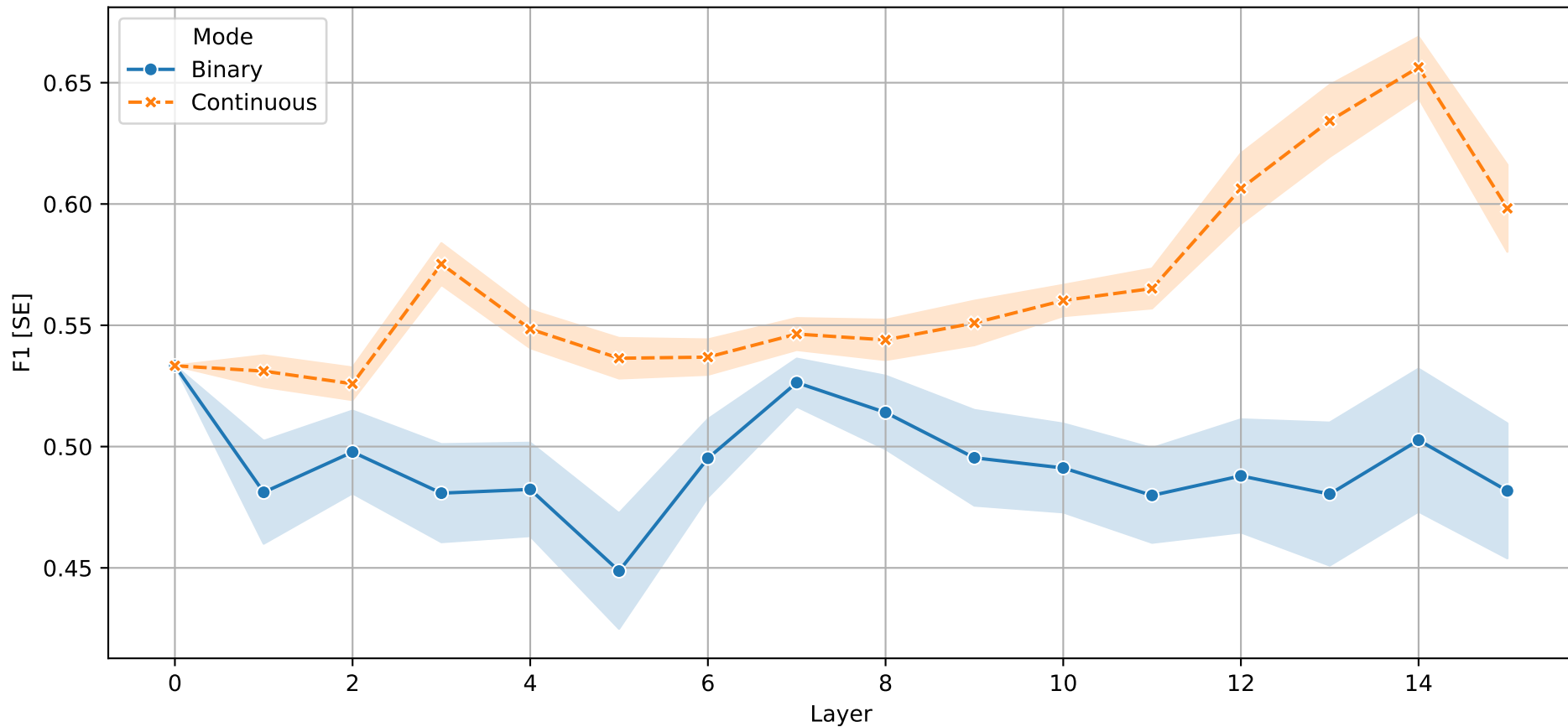
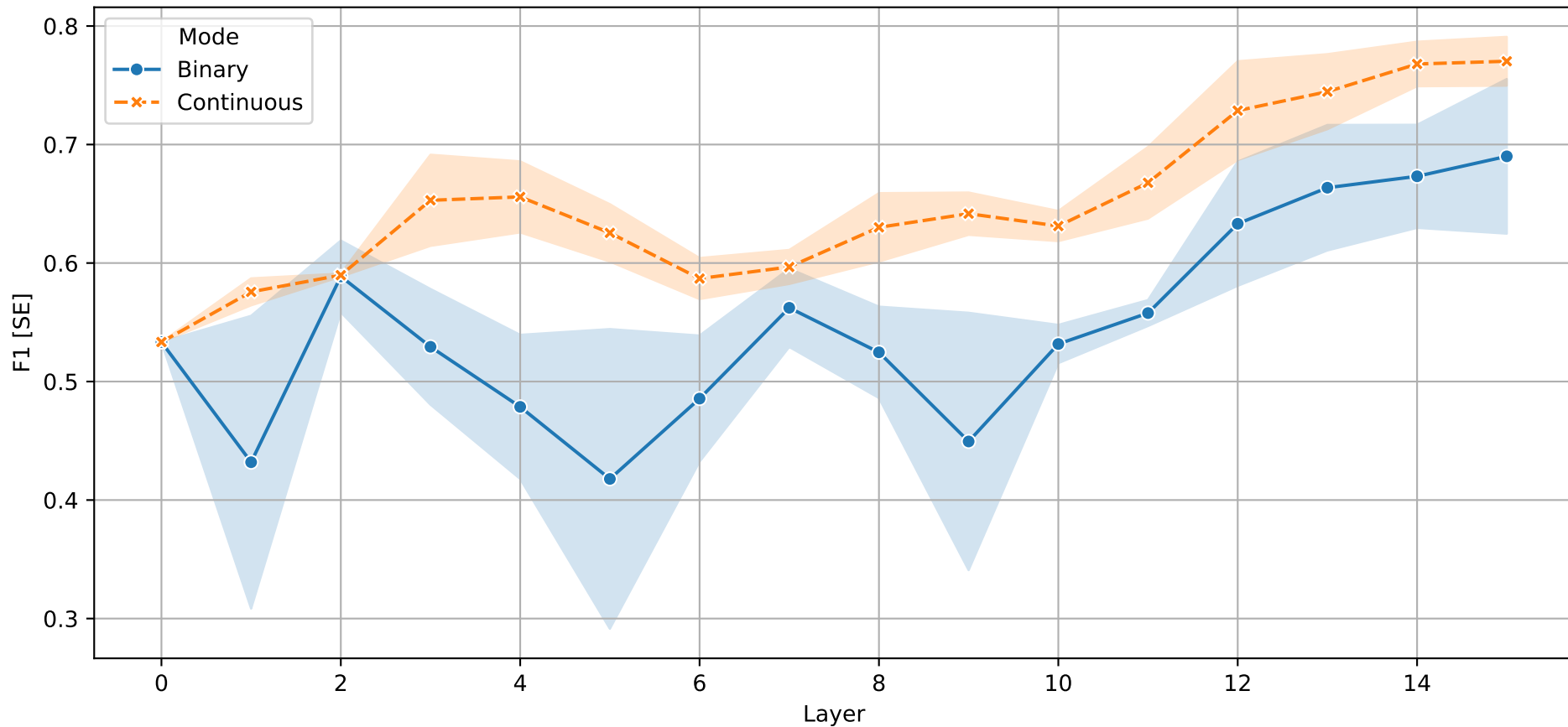


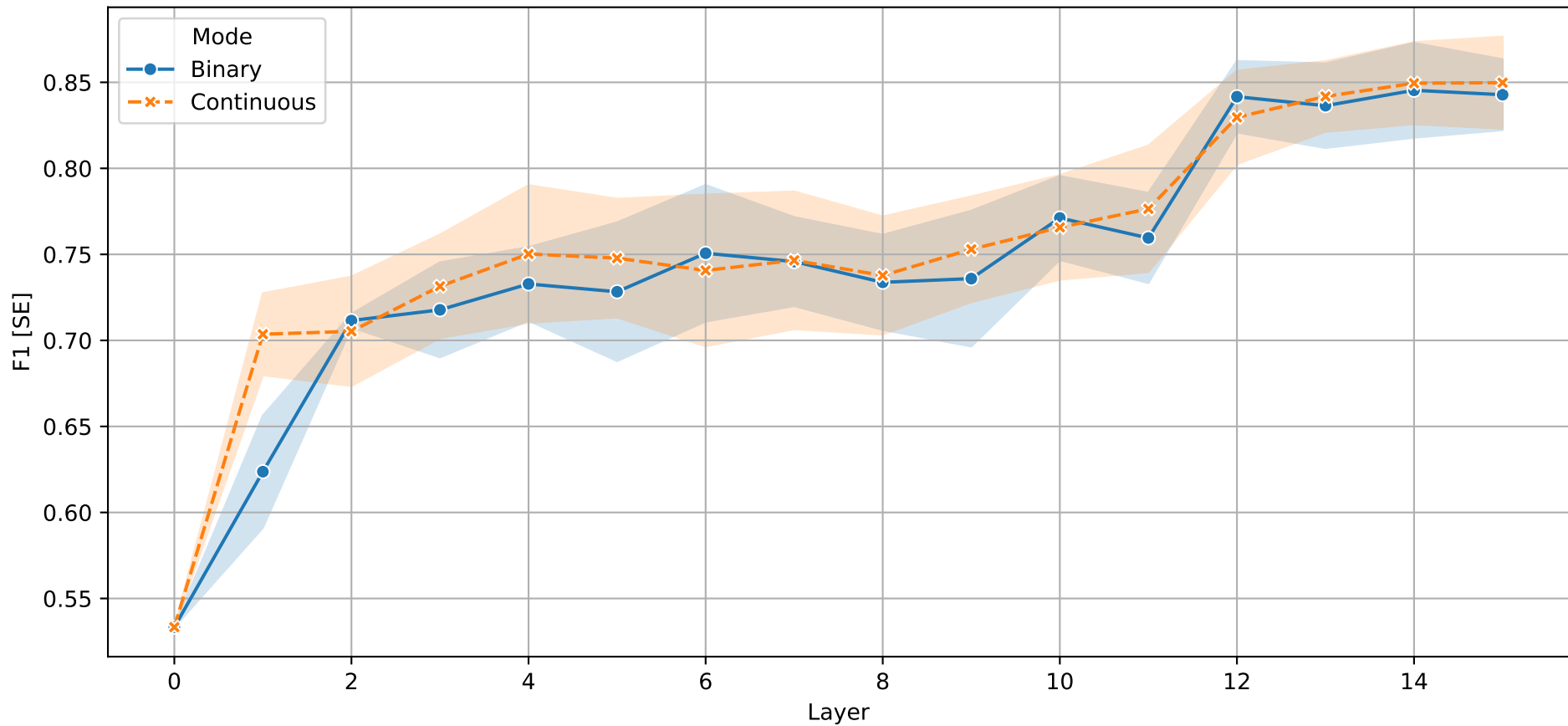
F1 per Layer - Single Neuron Probing



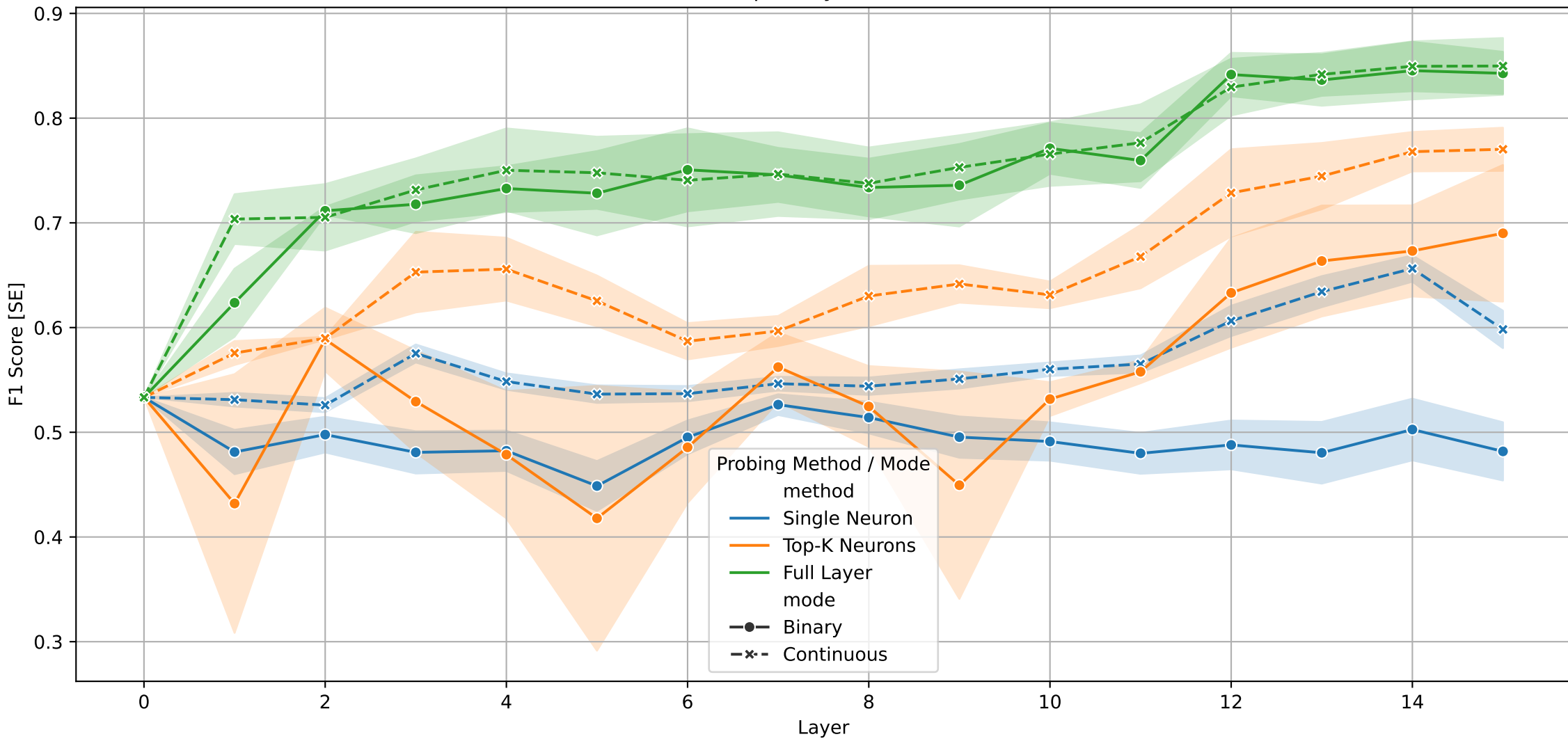
F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



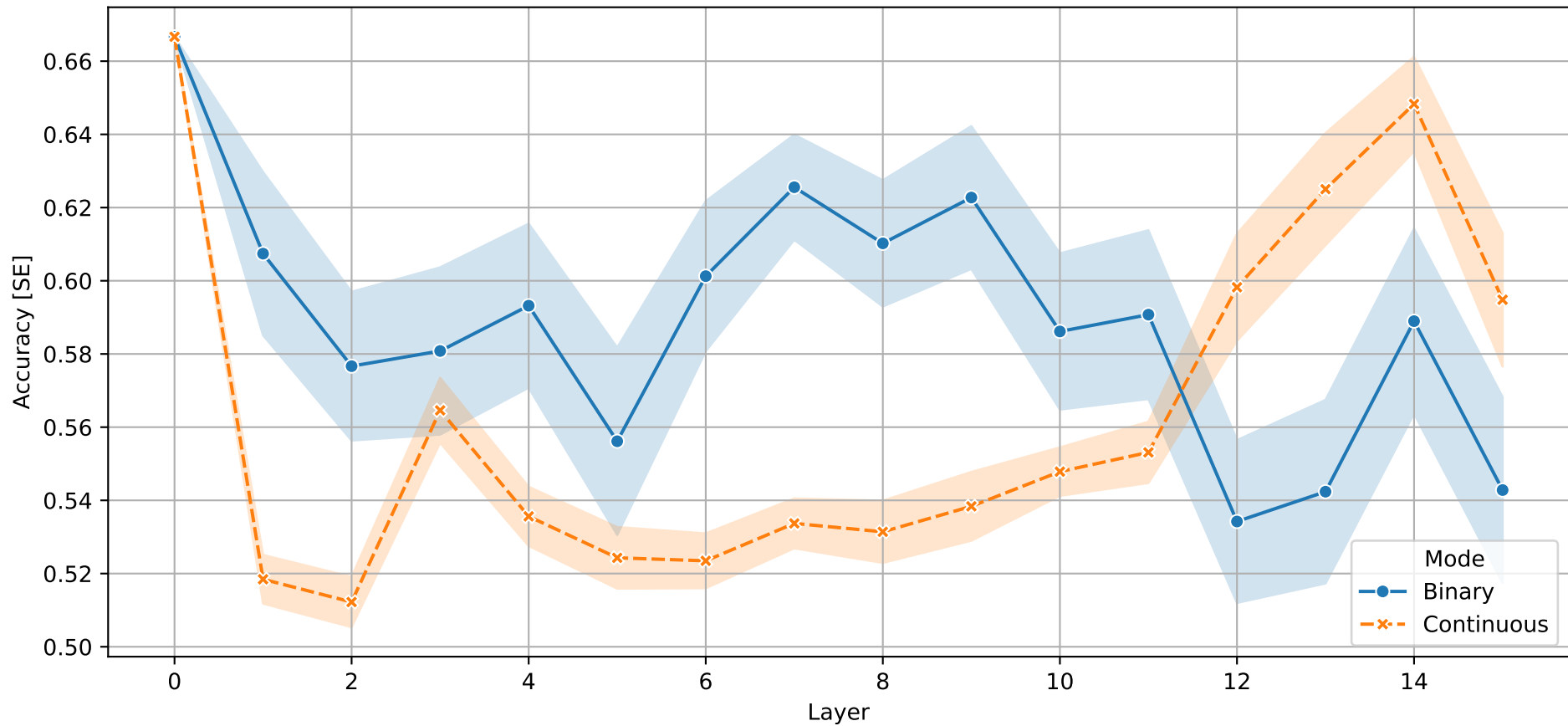
Overall F1 per Layer - All Methods



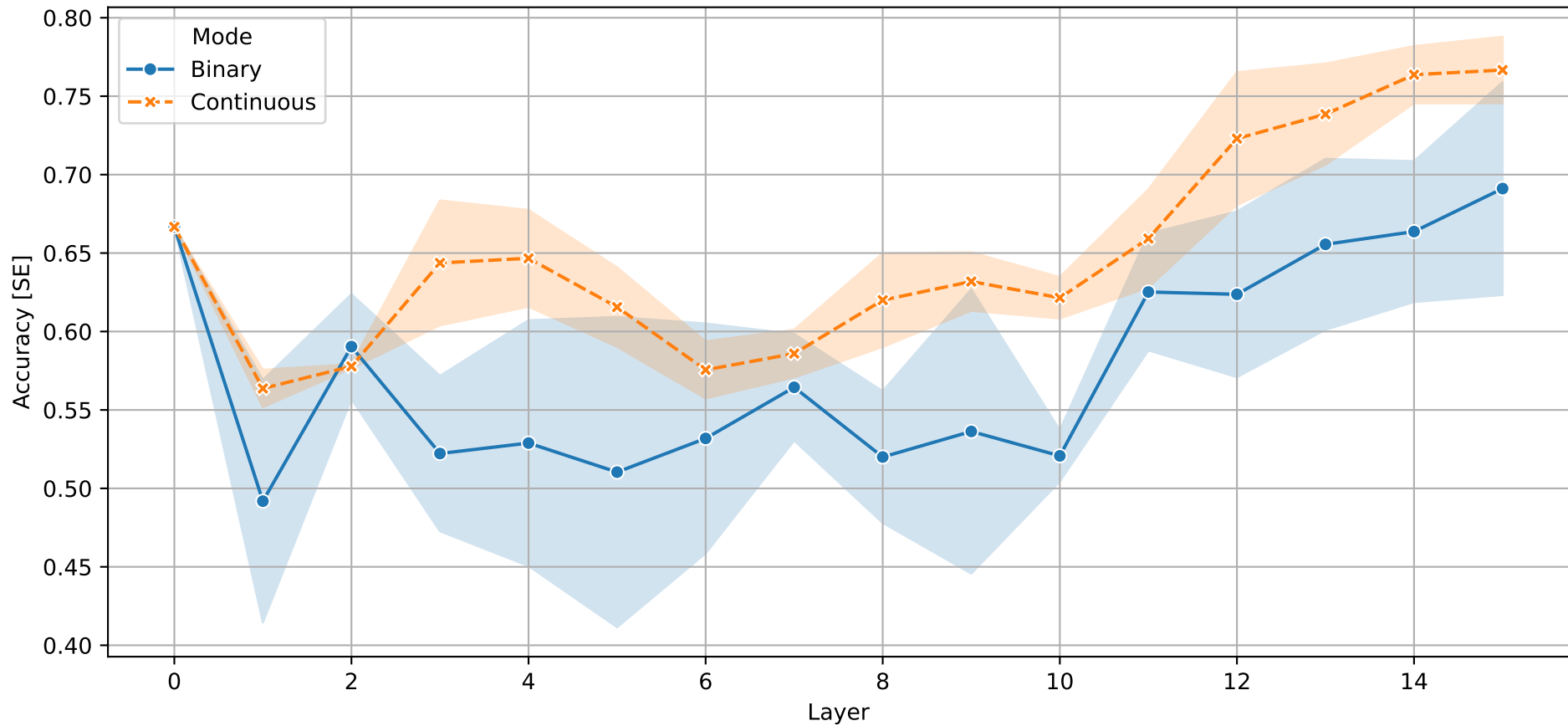
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	14.0	15.0
Full Layer	f1_max	0.8956	0.894
Full Layer	f1_mean	0.7444	0.7539
Full Layer	f1_std	0.089	0.0869
Single Neuron	f1_best_layer	0.0	14.0
Single Neuron	f1_max	0.7776	0.8091
Single Neuron	f1_mean	0.4924	0.5656
Single Neuron	f1_std	0.1129	0.0651
Top-K Neurons	f1_best_layer	15.0	15.0
Top-K Neurons	f1_max	0.8085	0.8111
Top-K Neurons	f1_mean	0.5469	0.6499
Top-K Neurons	f1_std	0.1247	0.0773

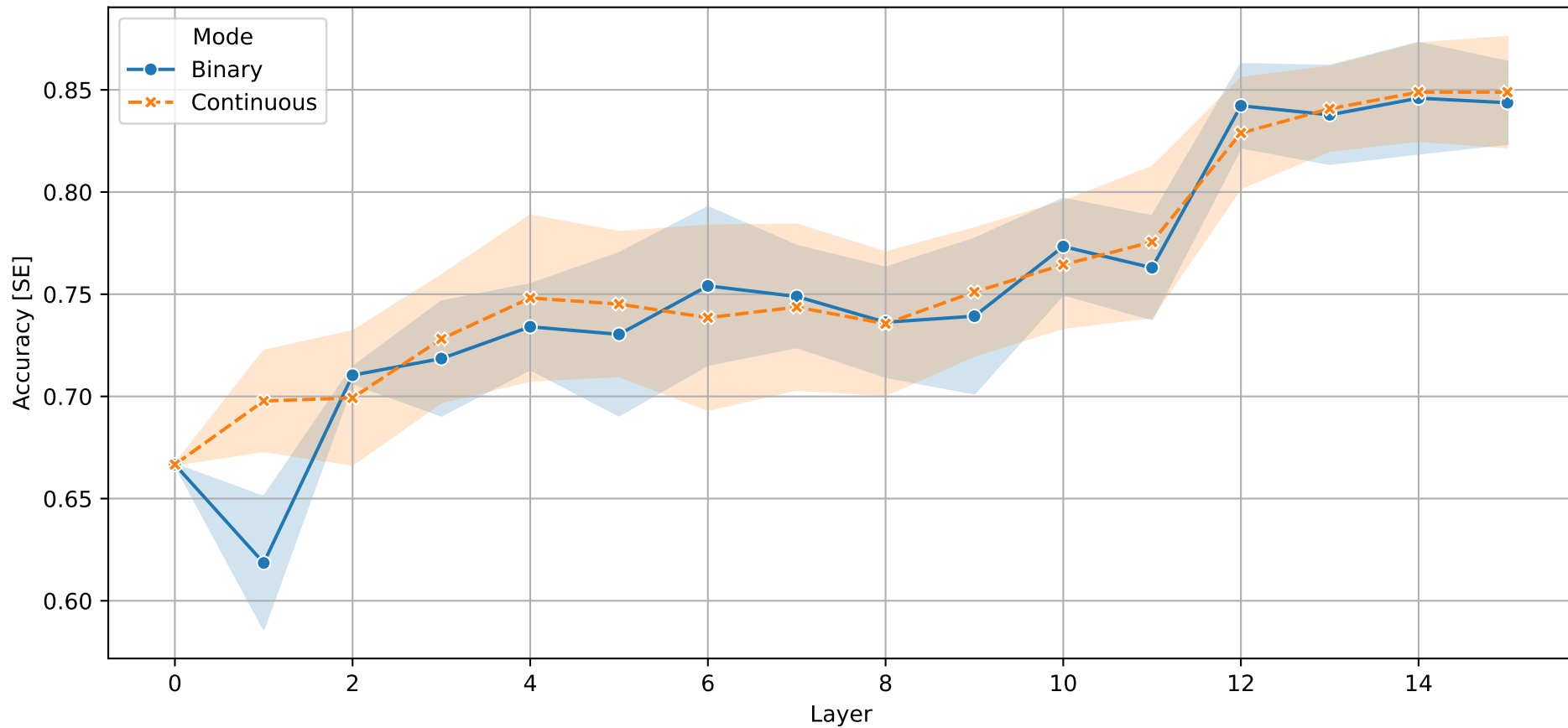
Accuracy per Layer - Single Neuron Probing



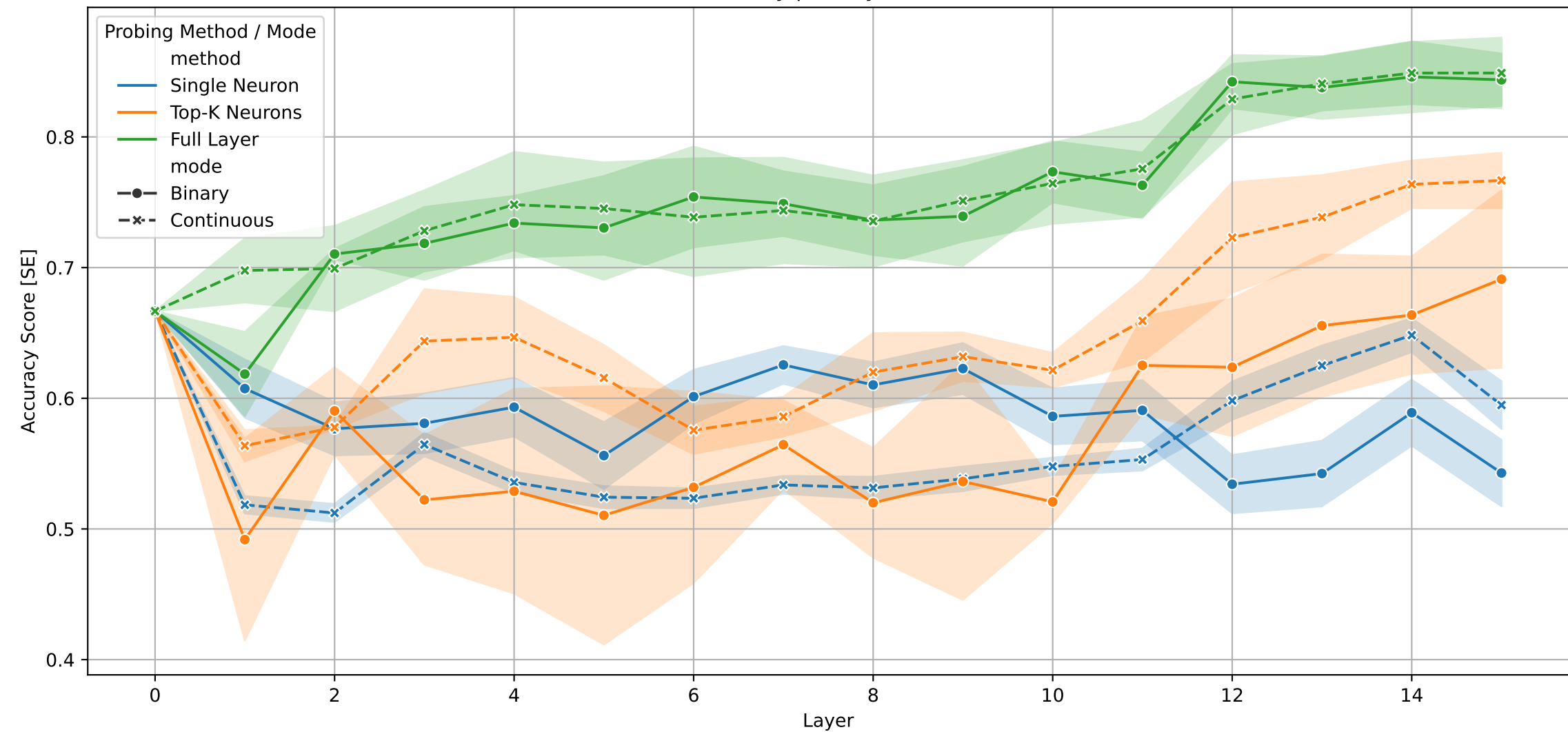
Accuracy per Layer - Top-K Neurons Probing



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	14.0	14.0
Full Layer	accuracy_max	0.8956	0.8933
Full Layer	accuracy_mean	0.7539	0.7601
Full Layer	accuracy_std	0.0737	0.0708
Single Neuron	accuracy_best_layer	0.0	0.0
Single Neuron	accuracy_max	0.7911	0.8044
Single Neuron	accuracy_mean	0.5891	0.5635
Single Neuron	accuracy_std	0.1199	0.0719
Top-K Neurons	accuracy_best_layer	15.0	15.0
Top-K Neurons	accuracy_max	0.8156	0.8067
Top-K Neurons	accuracy_mean	0.5777	0.65
Top-K Neurons	accuracy_std	0.1067	0.0738