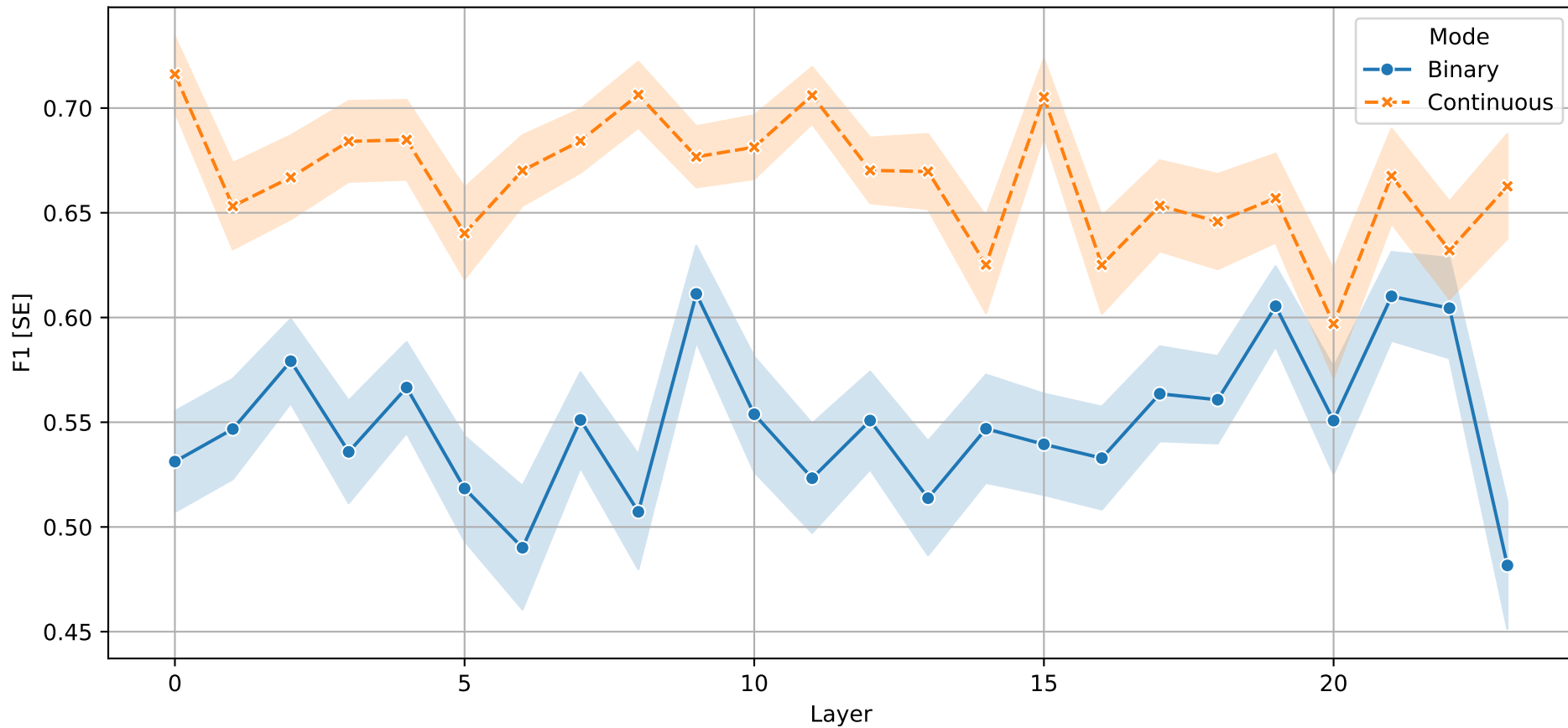
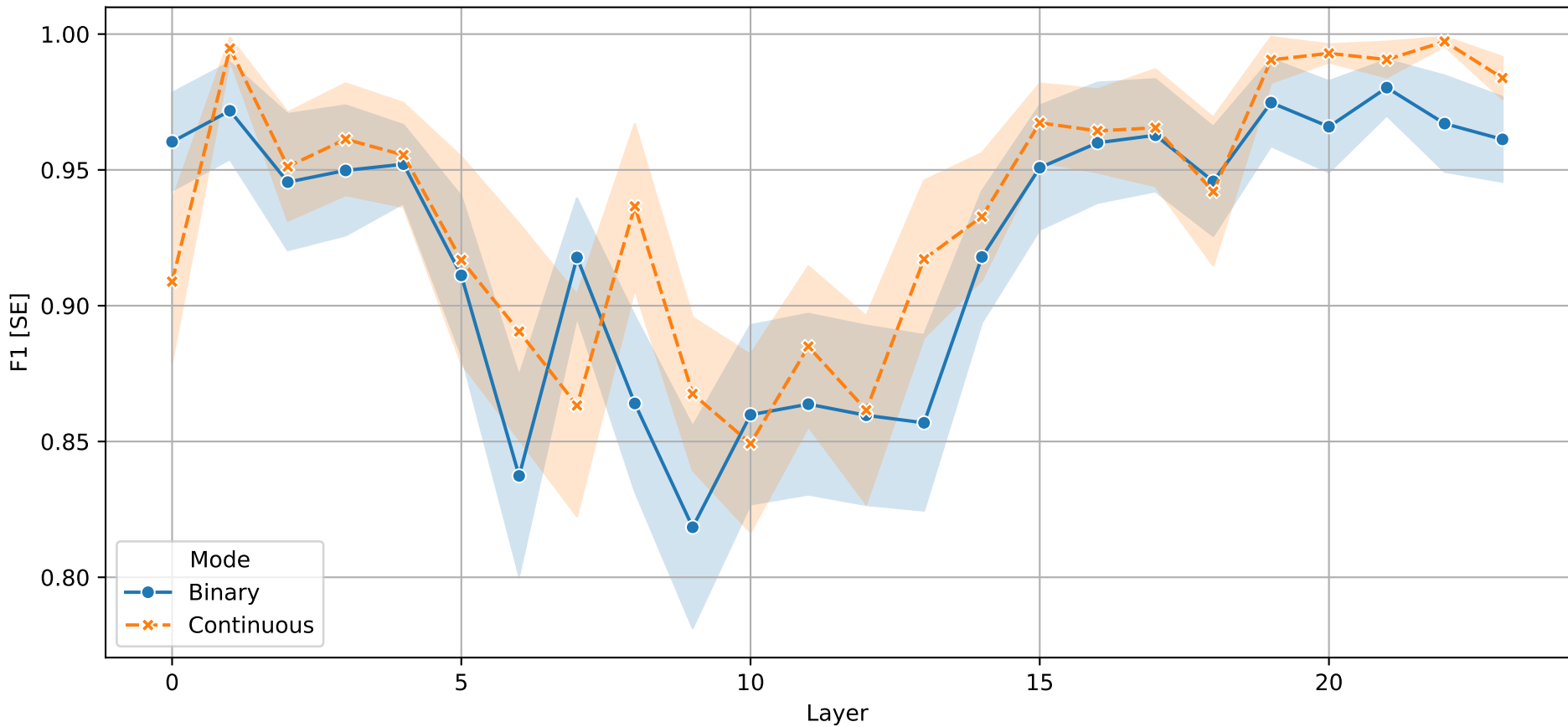


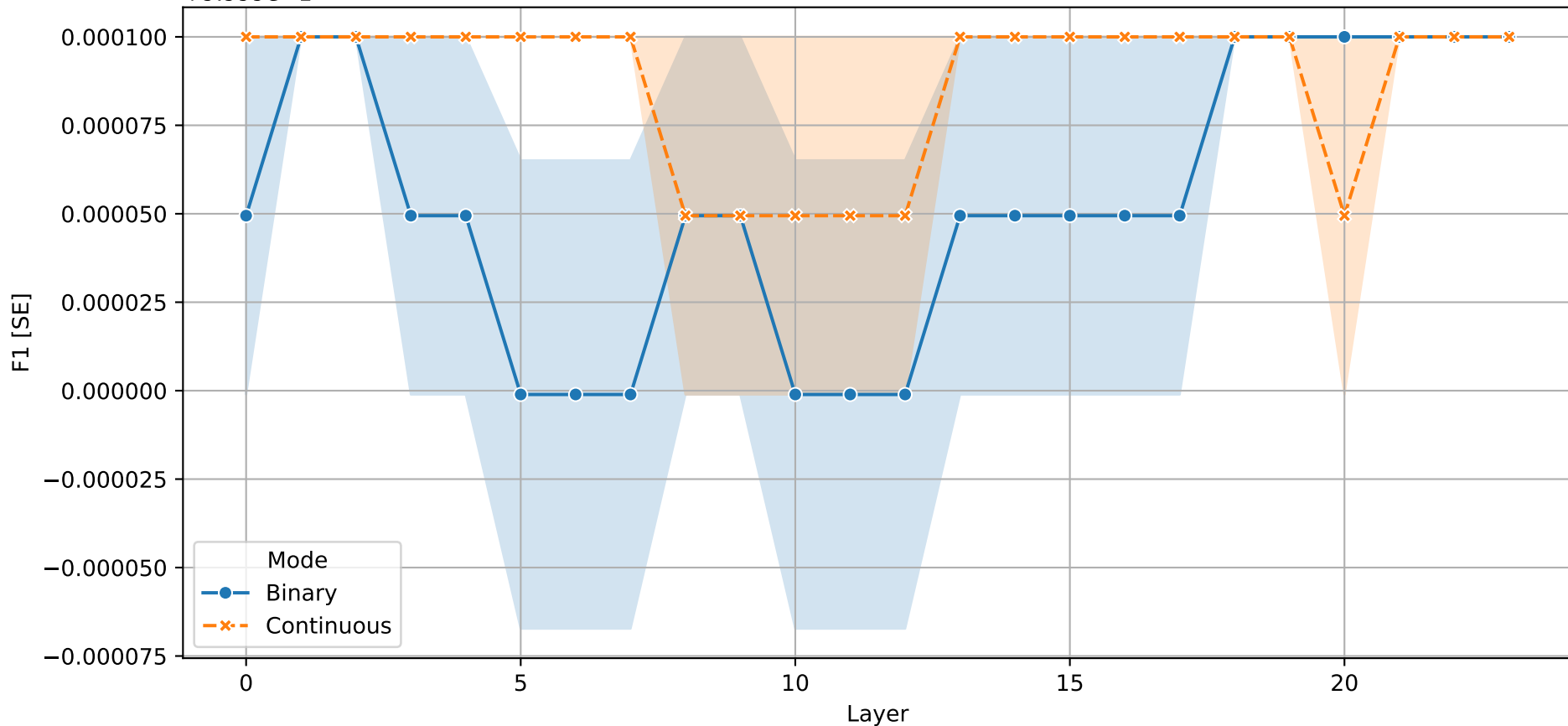
F1 per Layer - Single Neuron Probing



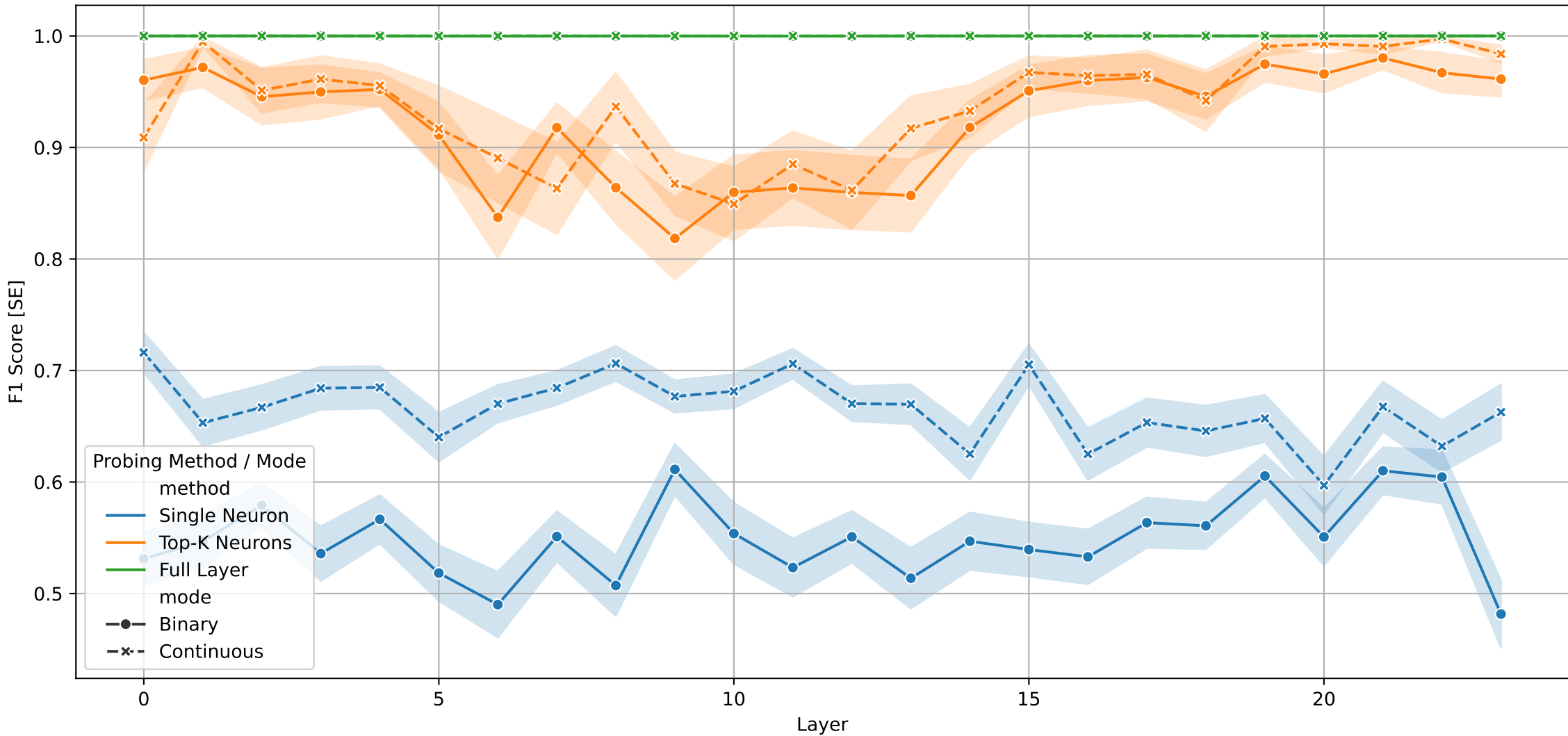
# F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing

 $+9.999e-1$ 

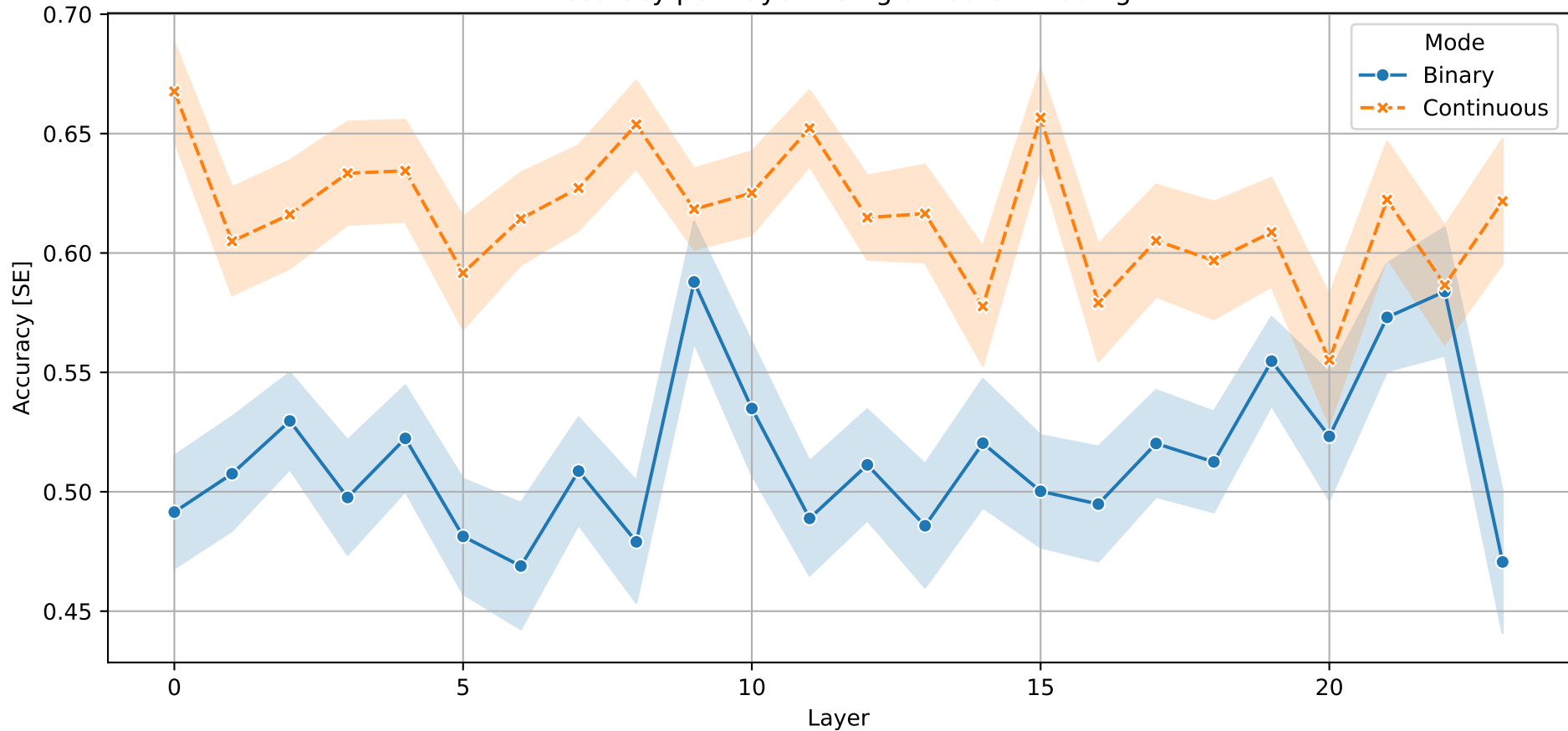
Overall F1 per Layer - All Methods



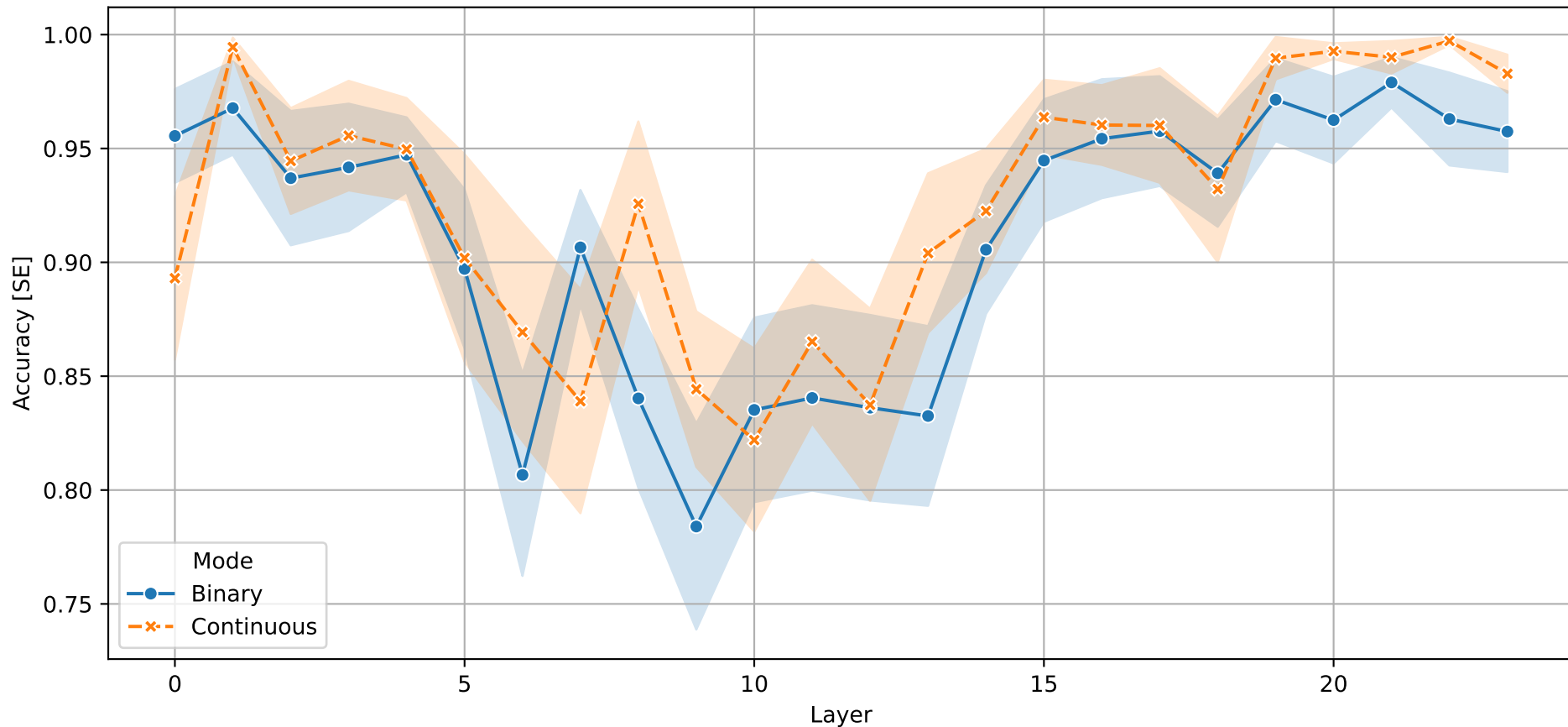
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	1.0	0.0
Full Layer	f1_max	1.0	1.0
Full Layer	f1_mean	1.0	1.0
Full Layer	f1_std	0.0001	0.0001
Single Neuron	f1_best_layer	9.0	0.0
Single Neuron	f1_max	1.0	1.0
Single Neuron	f1_mean	0.549	0.6659
Single Neuron	f1_std	0.2202	0.1789
Top-K Neurons	f1_best_layer	21.0	22.0
Top-K Neurons	f1_max	1.0	1.0
Top-K Neurons	f1_mean	0.9231	0.9369
Top-K Neurons	f1_std	0.083	0.0807

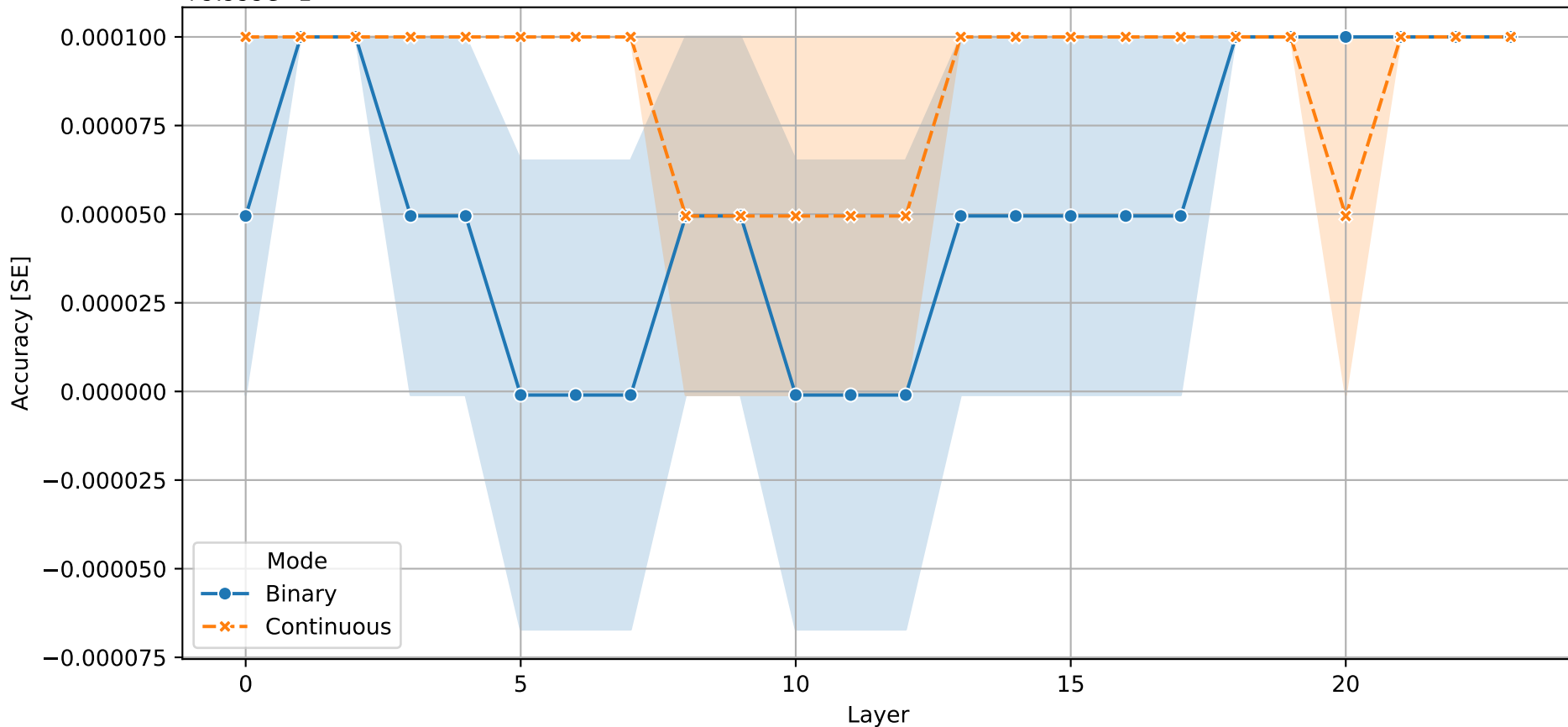
Accuracy per Layer - Single Neuron Probing



Accuracy per Layer - Top-K Neurons Probing

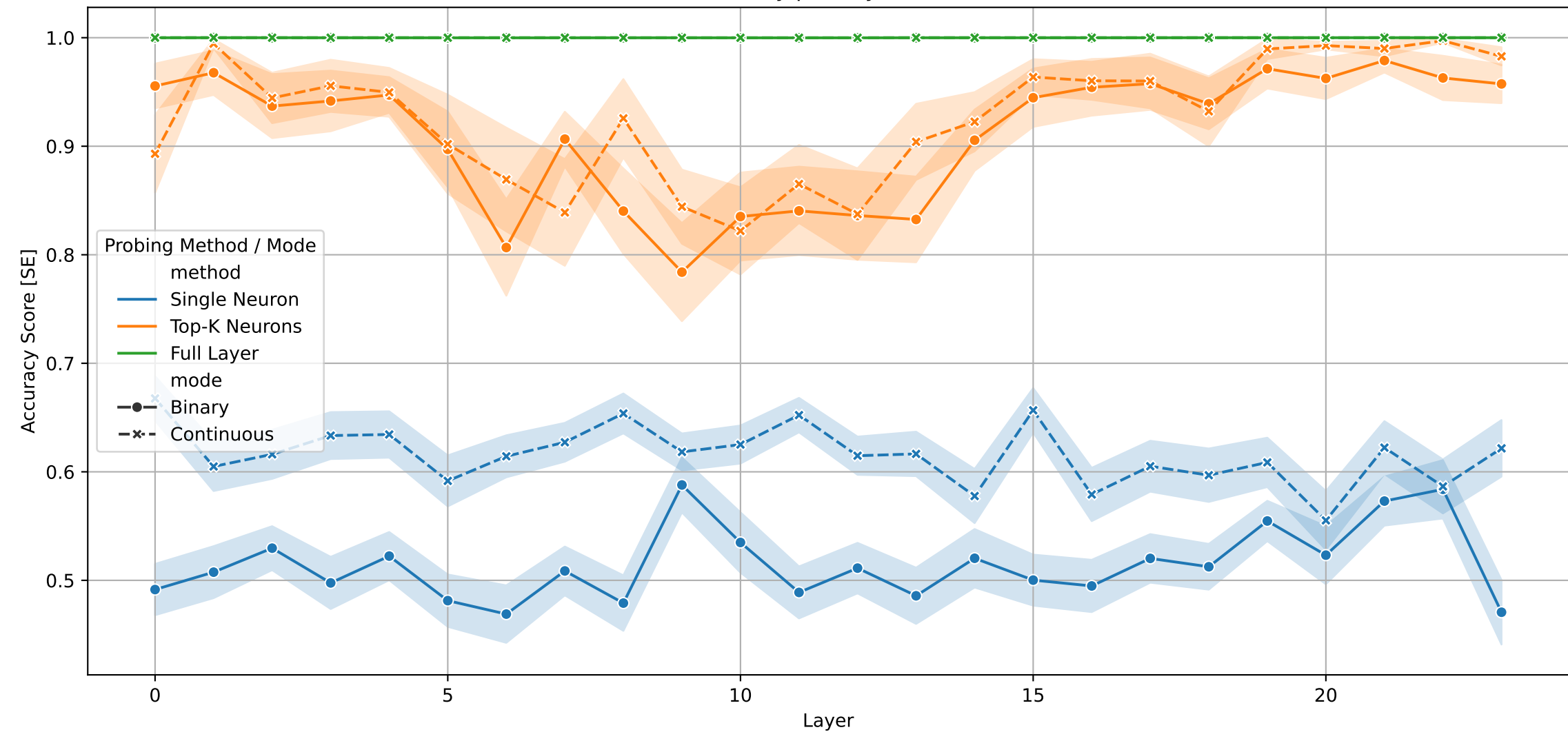


Accuracy per Layer - Full Layer Probing

 $+9.999e-1$ 



Overall Accuracy per Layer - All Methods



Accuracy Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	accuracy_best_layer	1.0	0.0
Full Layer	accuracy_max	1.0	1.0
Full Layer	accuracy_mean	1.0	1.0
Full Layer	accuracy_std	0.0001	0.0001
Single Neuron	accuracy_best_layer	9.0	0.0
Single Neuron	accuracy_max	1.0	1.0
Single Neuron	accuracy_mean	0.5145	0.6159
Single Neuron	accuracy_std	0.2187	0.1959
Top-K Neurons	accuracy_best_layer	21.0	22.0
Top-K Neurons	accuracy_max	1.0	1.0
Top-K Neurons	accuracy_mean	0.9109	0.9266
Top-K Neurons	accuracy_std	0.0993	0.0964