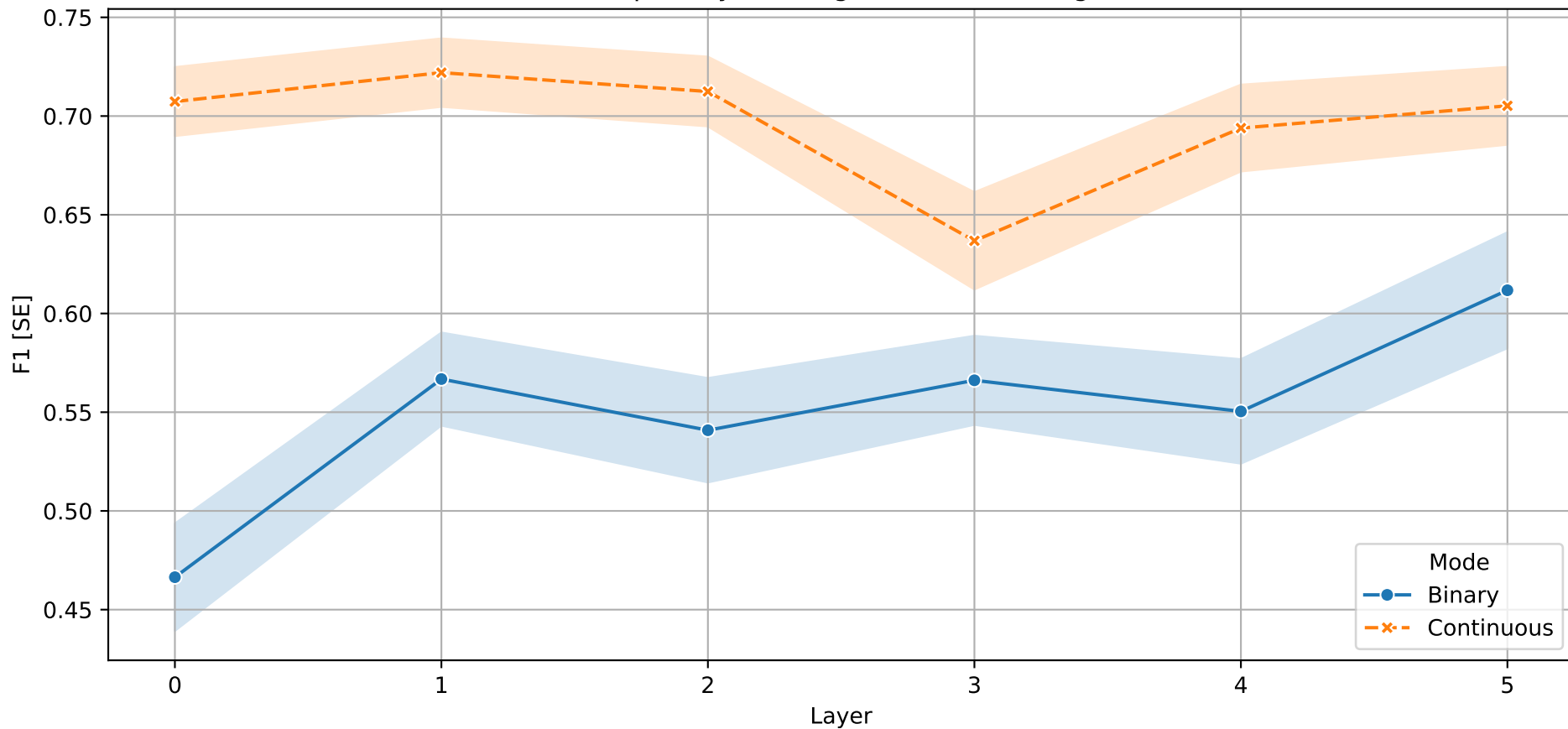
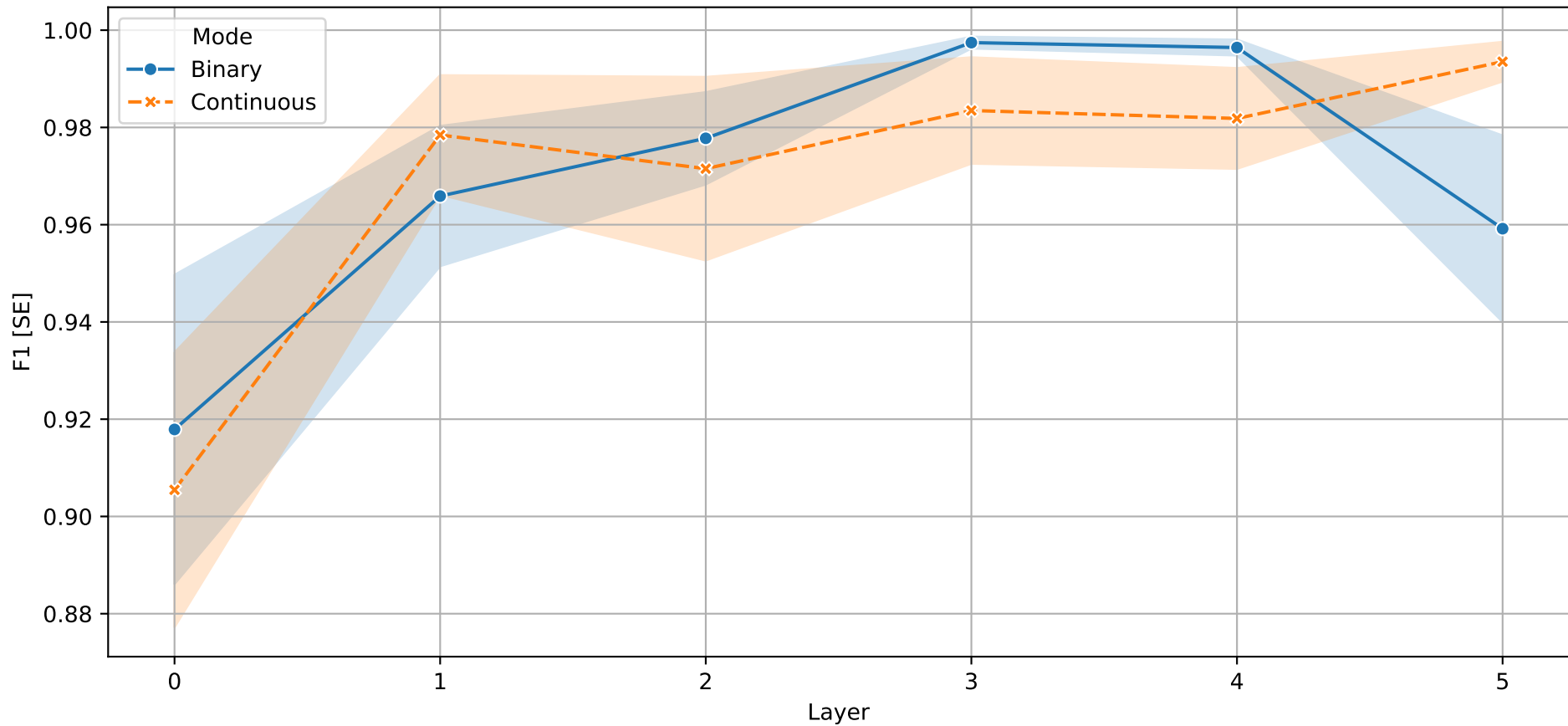


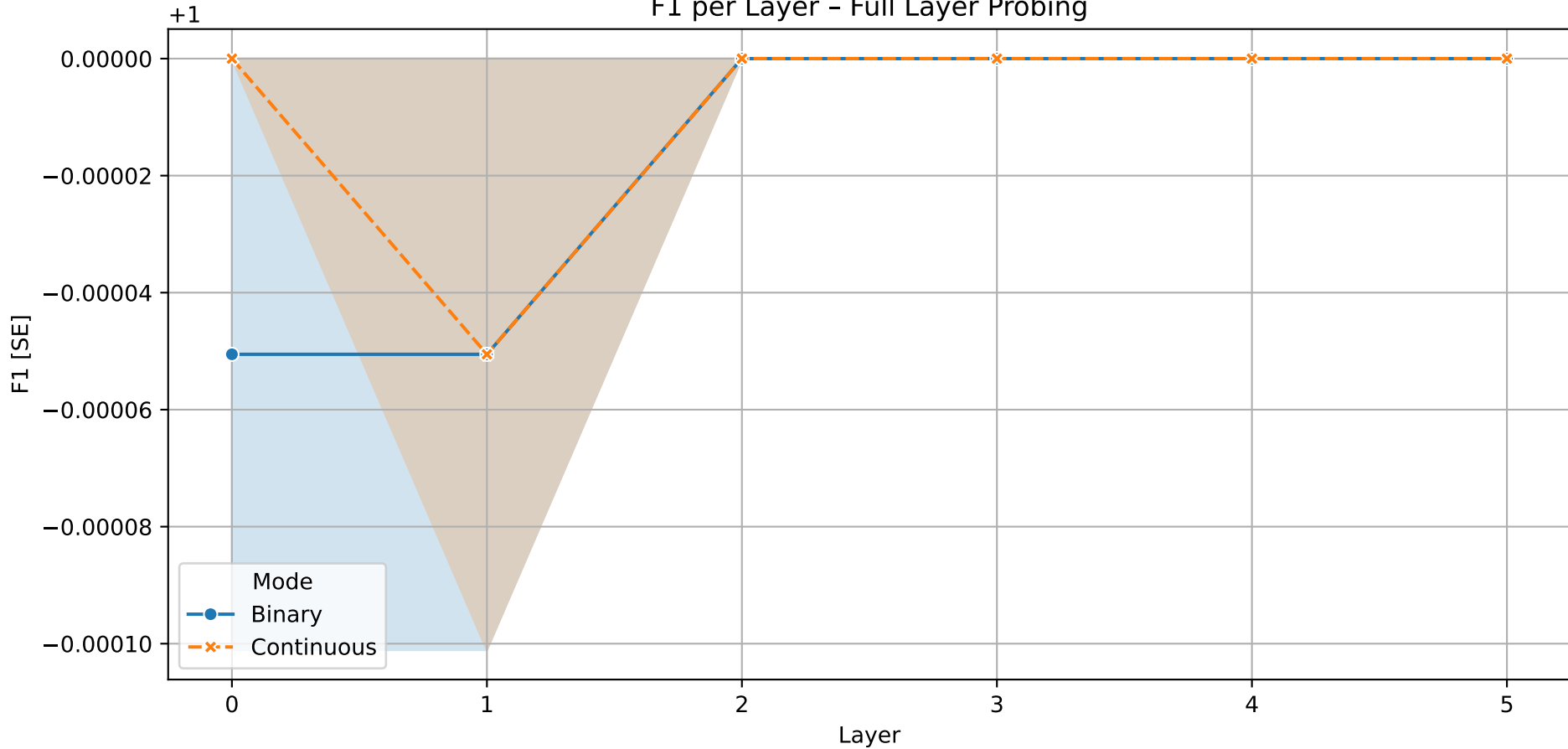
F1 per Layer - Single Neuron Probing



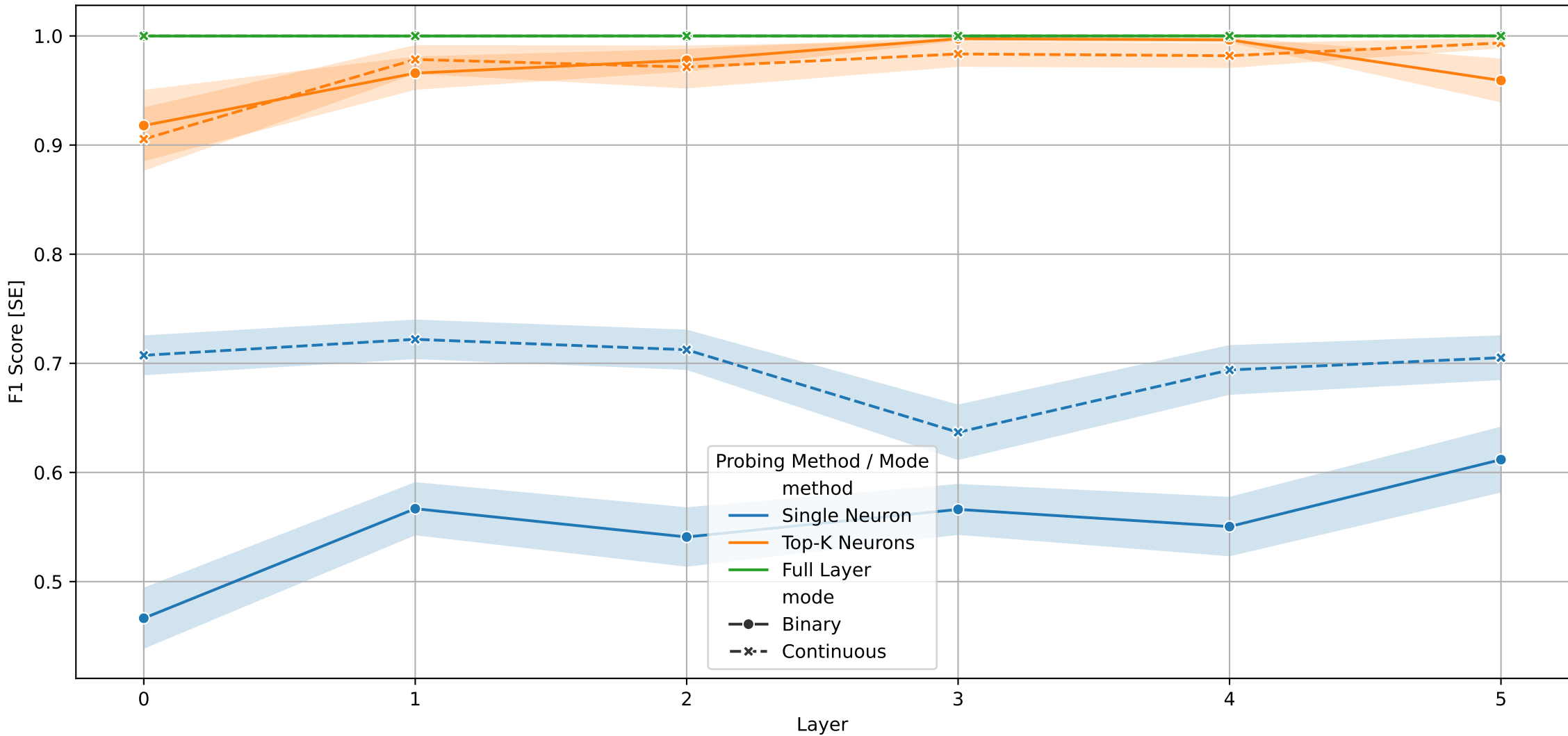
F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



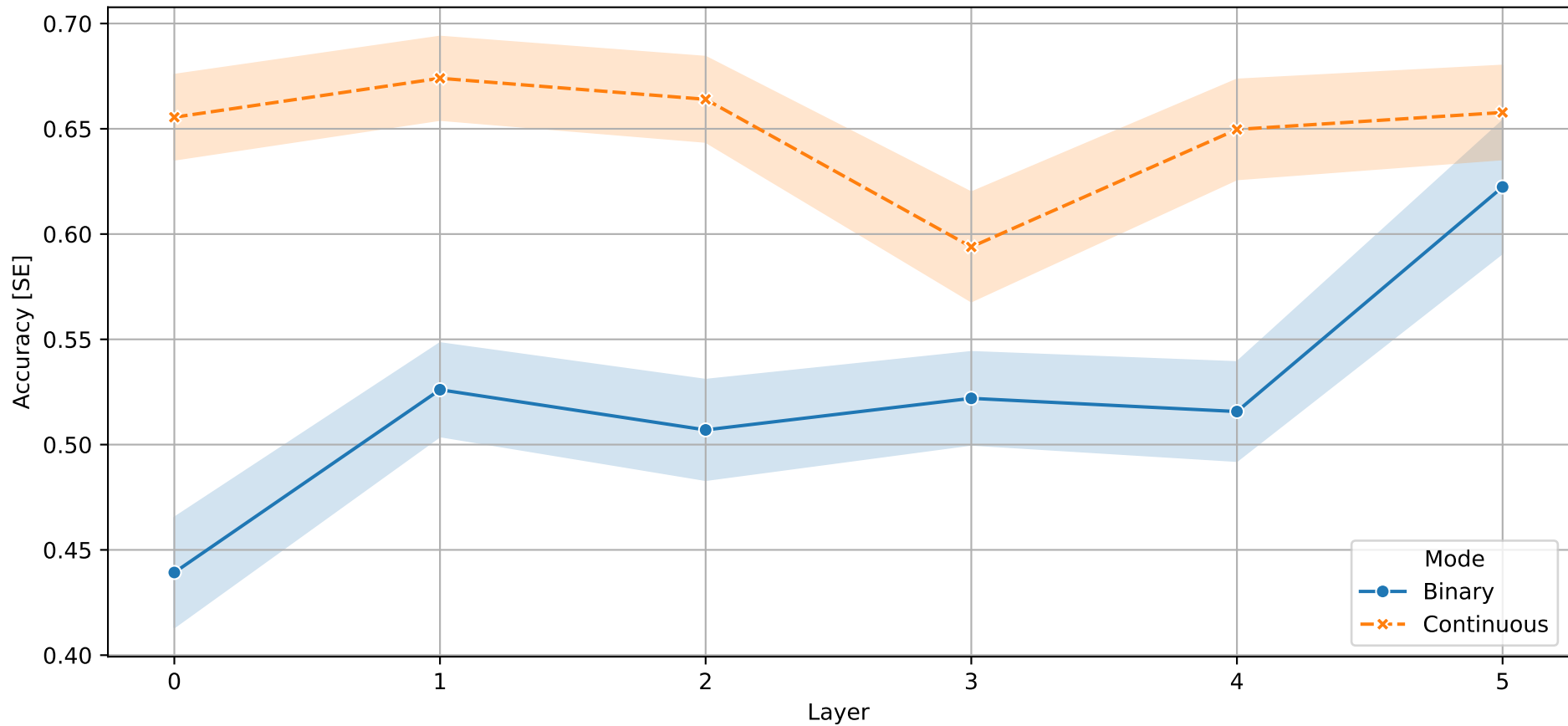
Overall F1 per Layer - All Methods



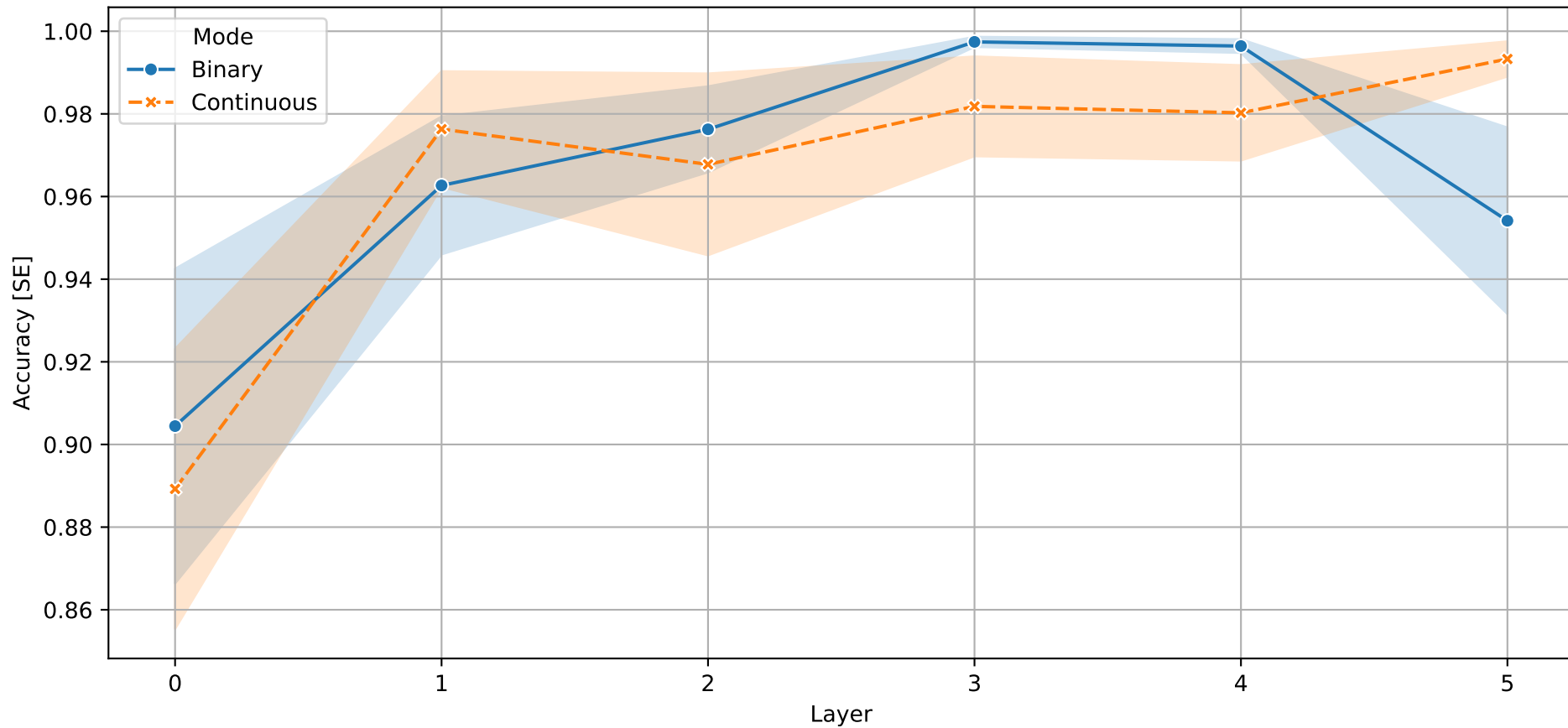
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	2.0	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	5.0	1.0
Single Neuron	f1_max	1.0	1.0
Single Neuron	f1_mean	0.5504	0.6963
Single Neuron	f1_std	0.235	0.1794
Top-K Neurons	f1_best_layer	3.0	5.0
Top-K Neurons	f1_max	1.0	1.0
Top-K Neurons	f1_mean	0.9691	0.969
Top-K Neurons	f1_std	0.0524	0.0521

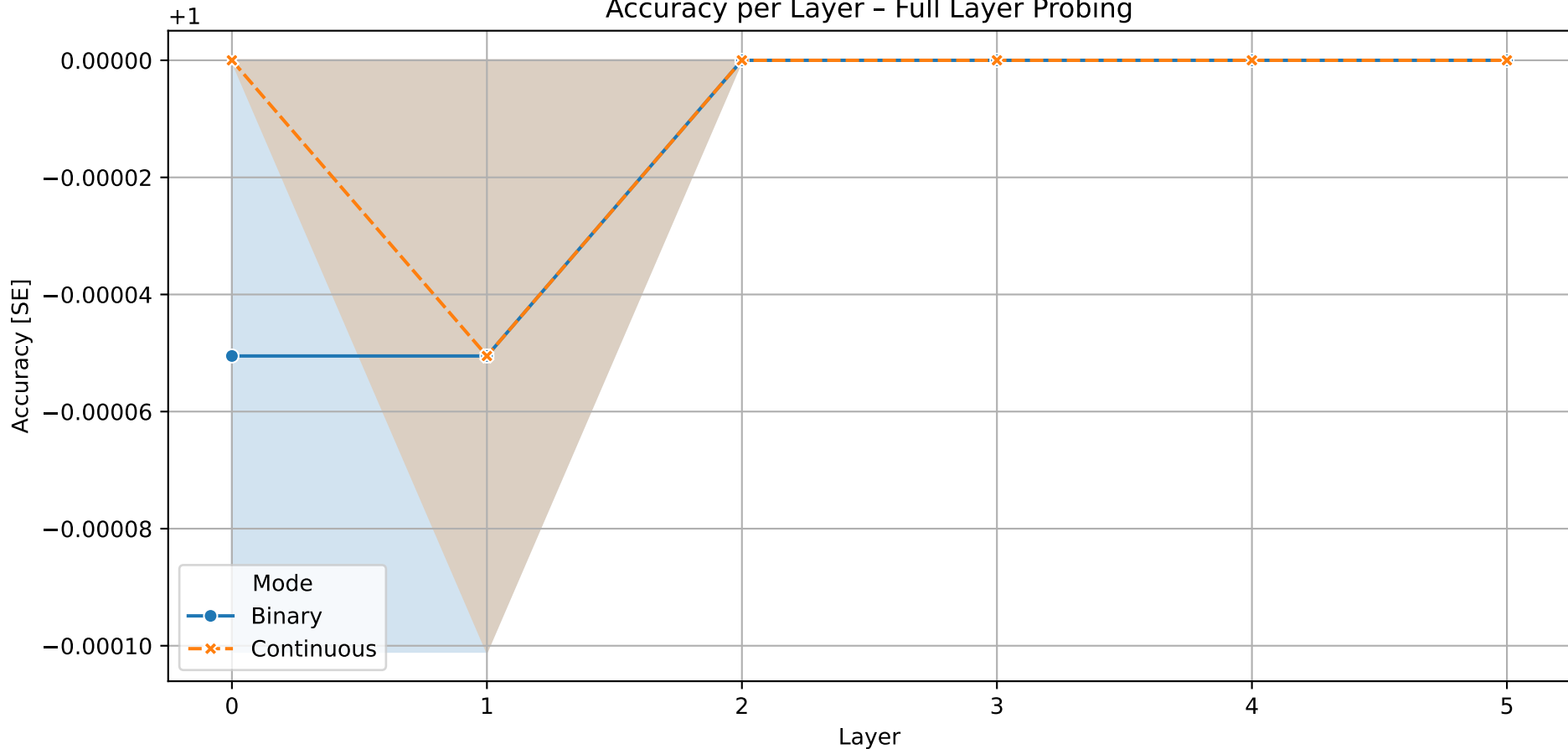
Accuracy per Layer - Single Neuron Probing



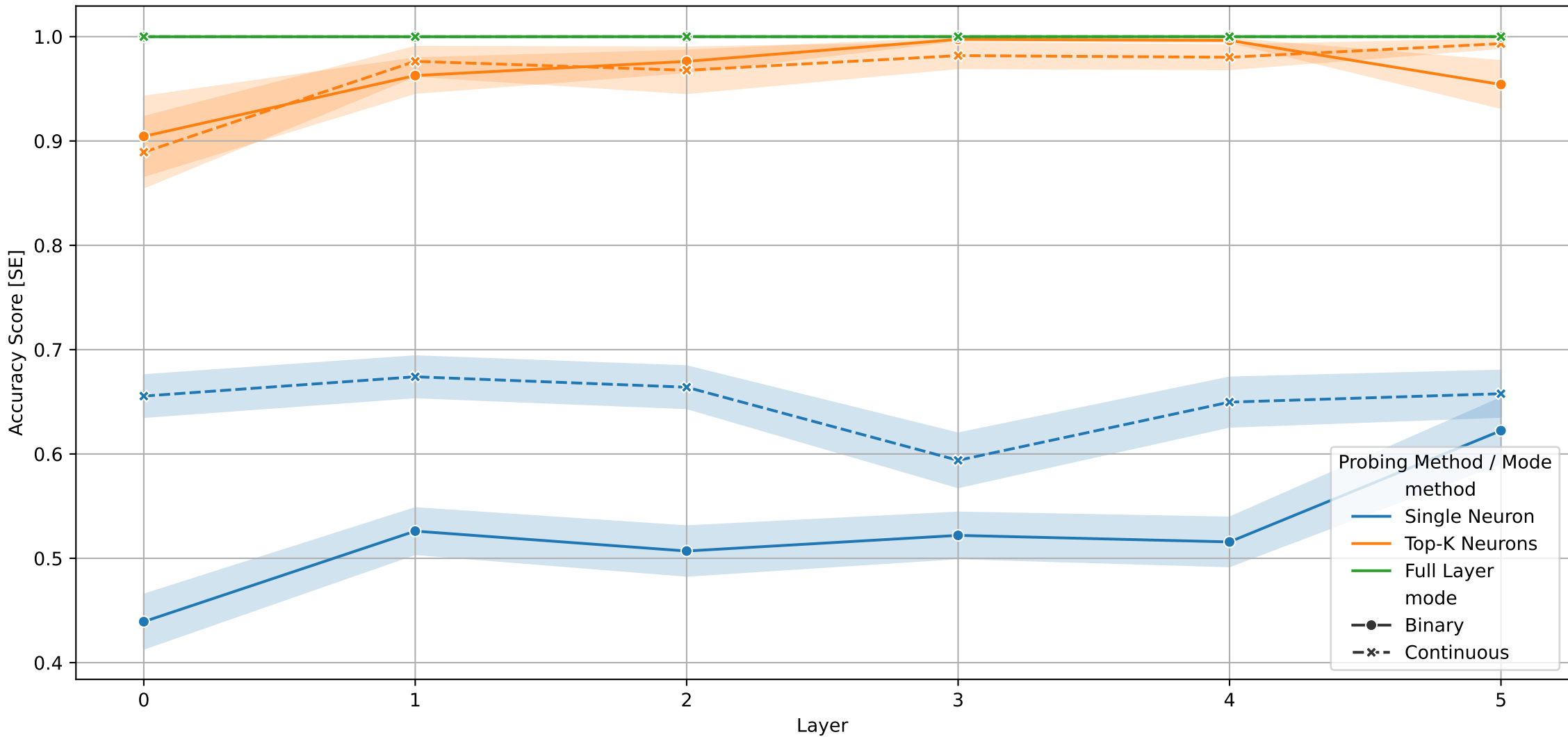
Accuracy per Layer - Top-K Neurons Probing



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	2.0	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	5.0	1.0
Single Neuron	accuracy_max	1.0	1.0
Single Neuron	accuracy_mean	0.5221	0.6491
Single Neuron	accuracy_std	0.2283	0.1977
Top-K Neurons	accuracy_best_layer	3.0	5.0
Top-K Neurons	accuracy_max	1.0	1.0
Top-K Neurons	accuracy_mean	0.9652	0.9648
Top-K Neurons	accuracy_std	0.0617	0.0613