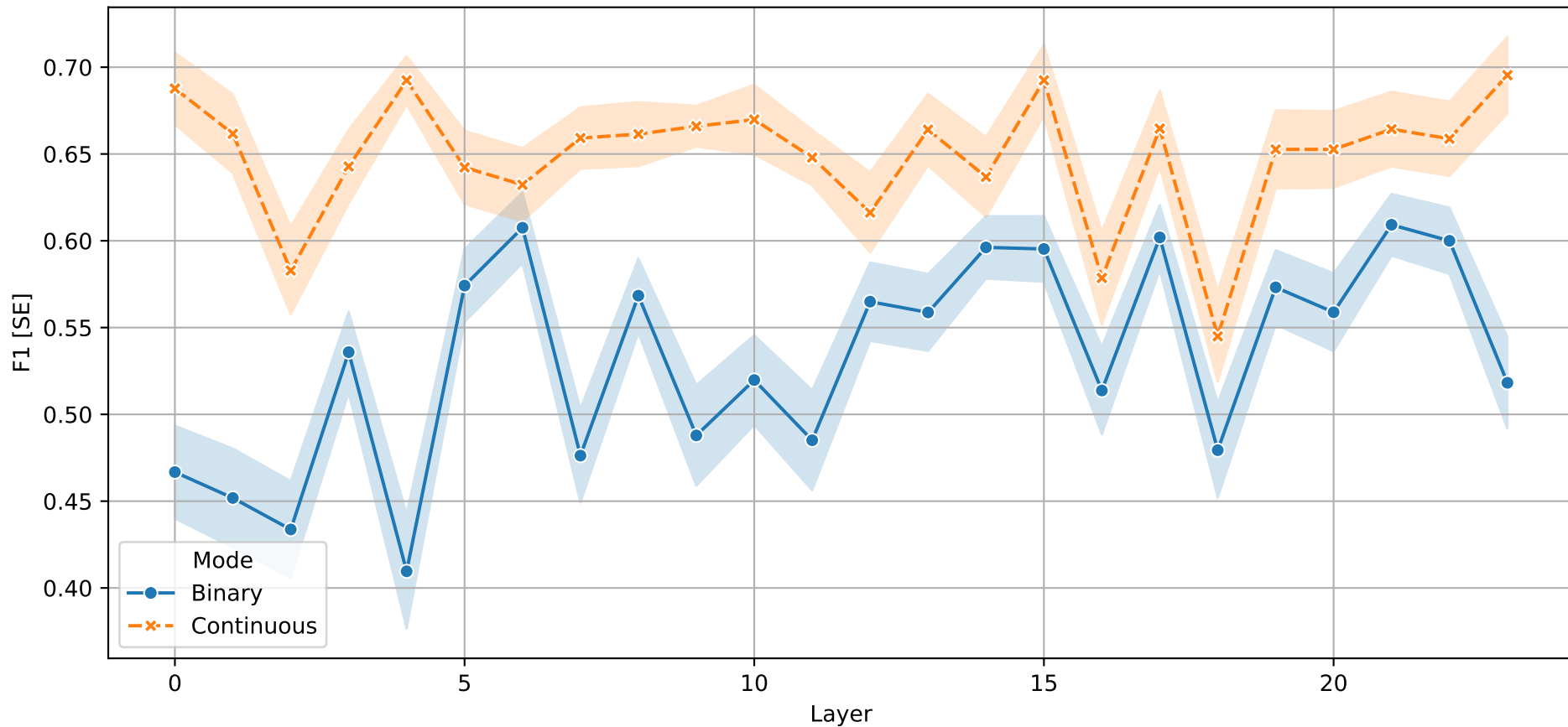
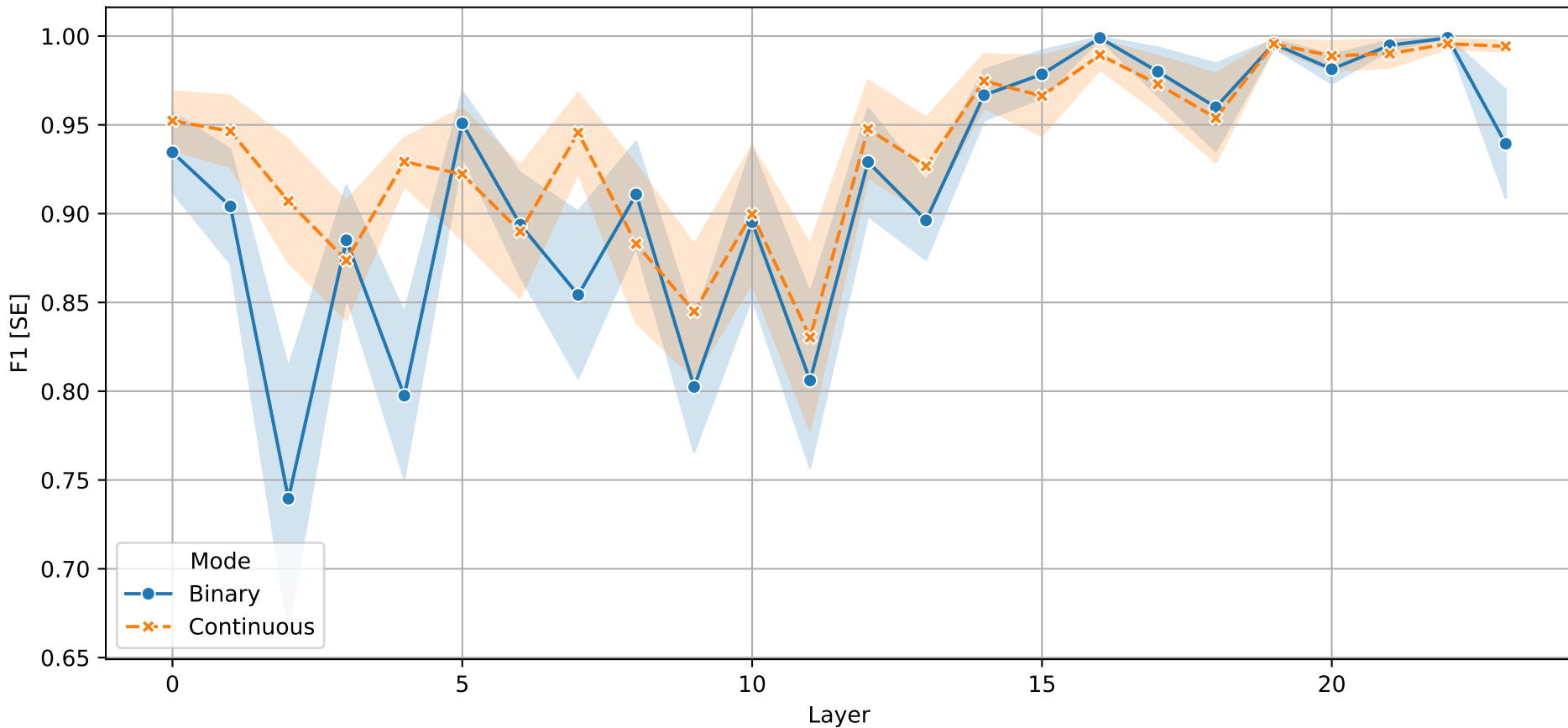


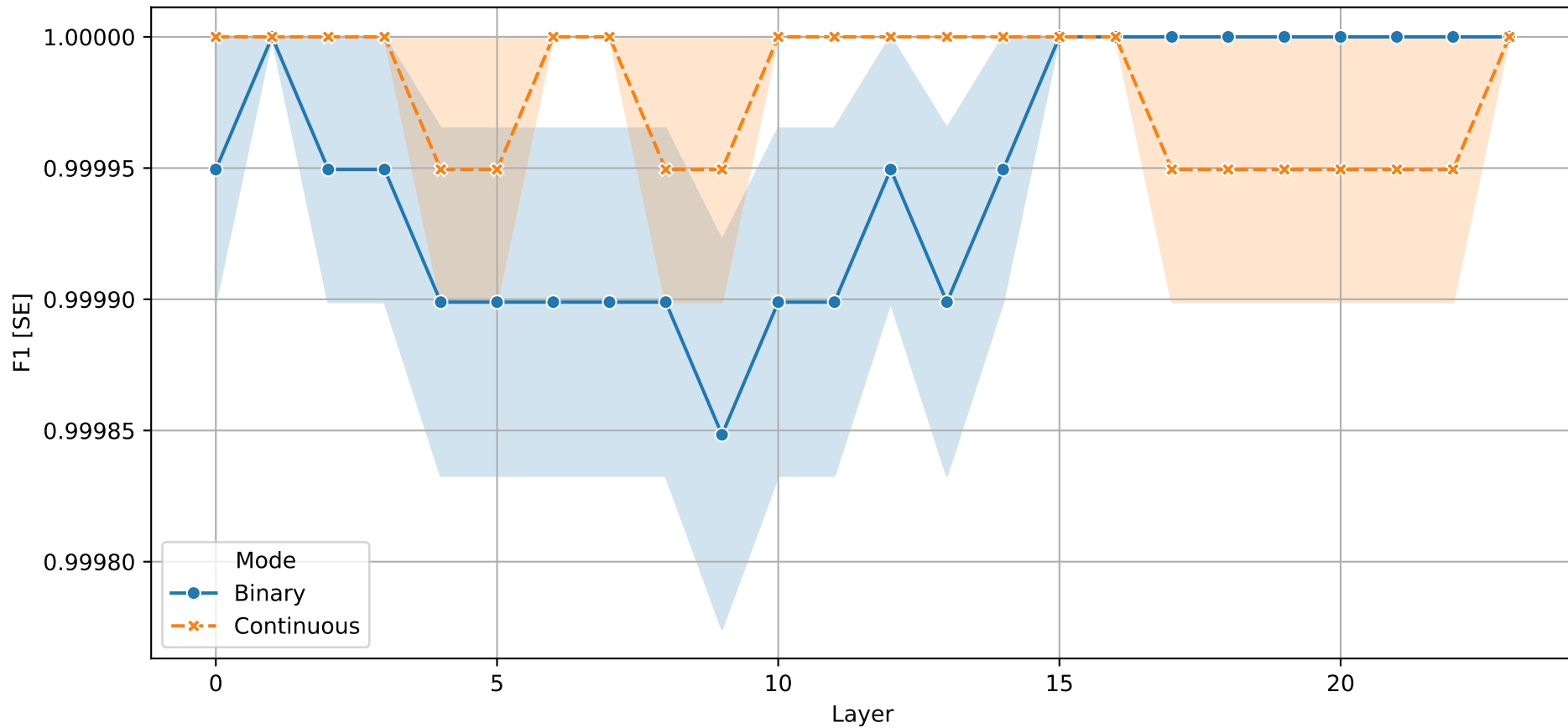
F1 per Layer - Single Neuron Probing



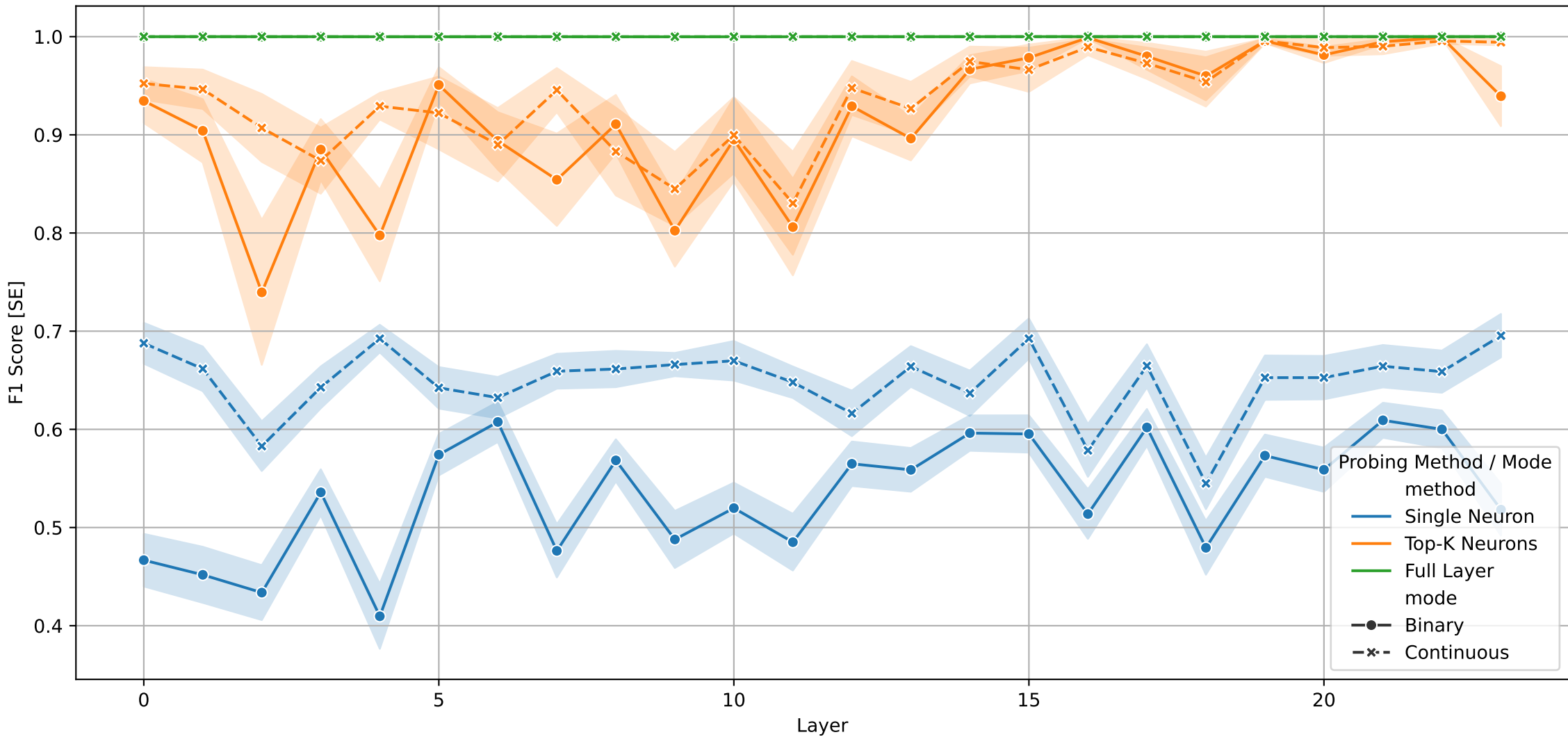
F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



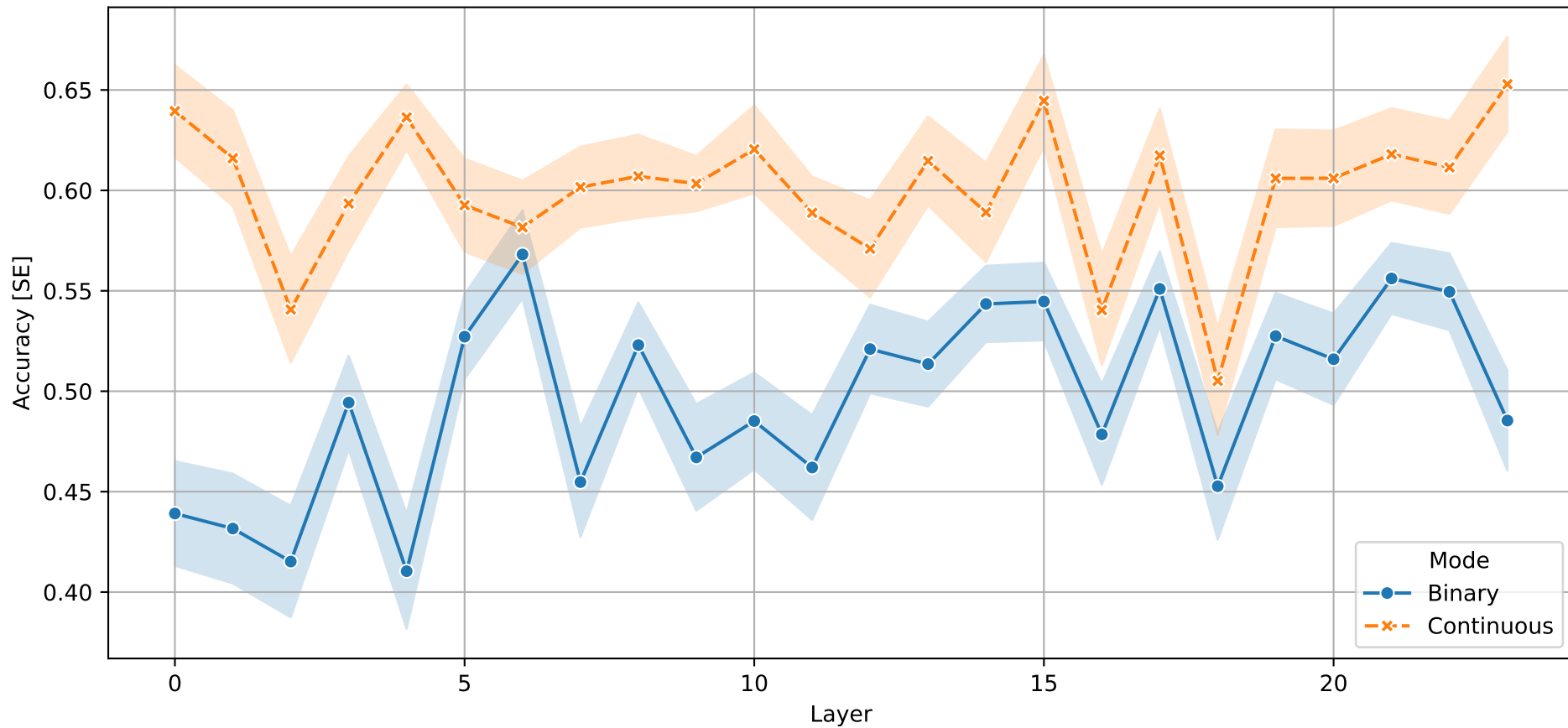
Overall F1 per Layer - All Methods



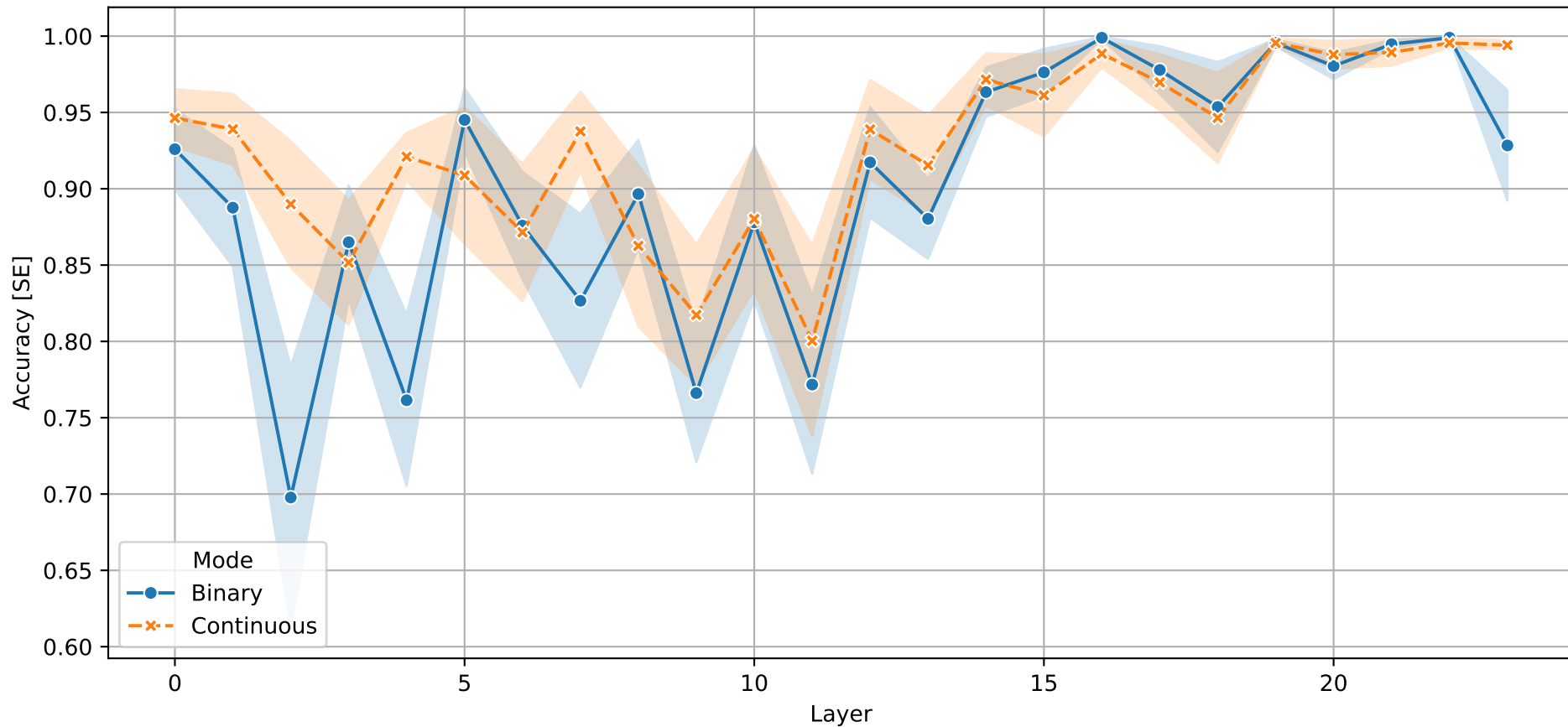
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	1.0	0.0
Full Layer	f1_max	1.0	1.0
Full Layer	f1_mean	0.9999	1.0
Full Layer	f1_std	0.0001	0.0001
Single Neuron	f1_best_layer	21.0	23.0
Single Neuron	f1_max	1.0	1.0
Single Neuron	f1_mean	0.5328	0.6487
Single Neuron	f1_std	0.2233	0.1918
Top-K Neurons	f1_best_layer	22.0	22.0
Top-K Neurons	f1_max	1.0	1.0
Top-K Neurons	f1_mean	0.9164	0.9383
Top-K Neurons	f1_std	0.1095	0.0859

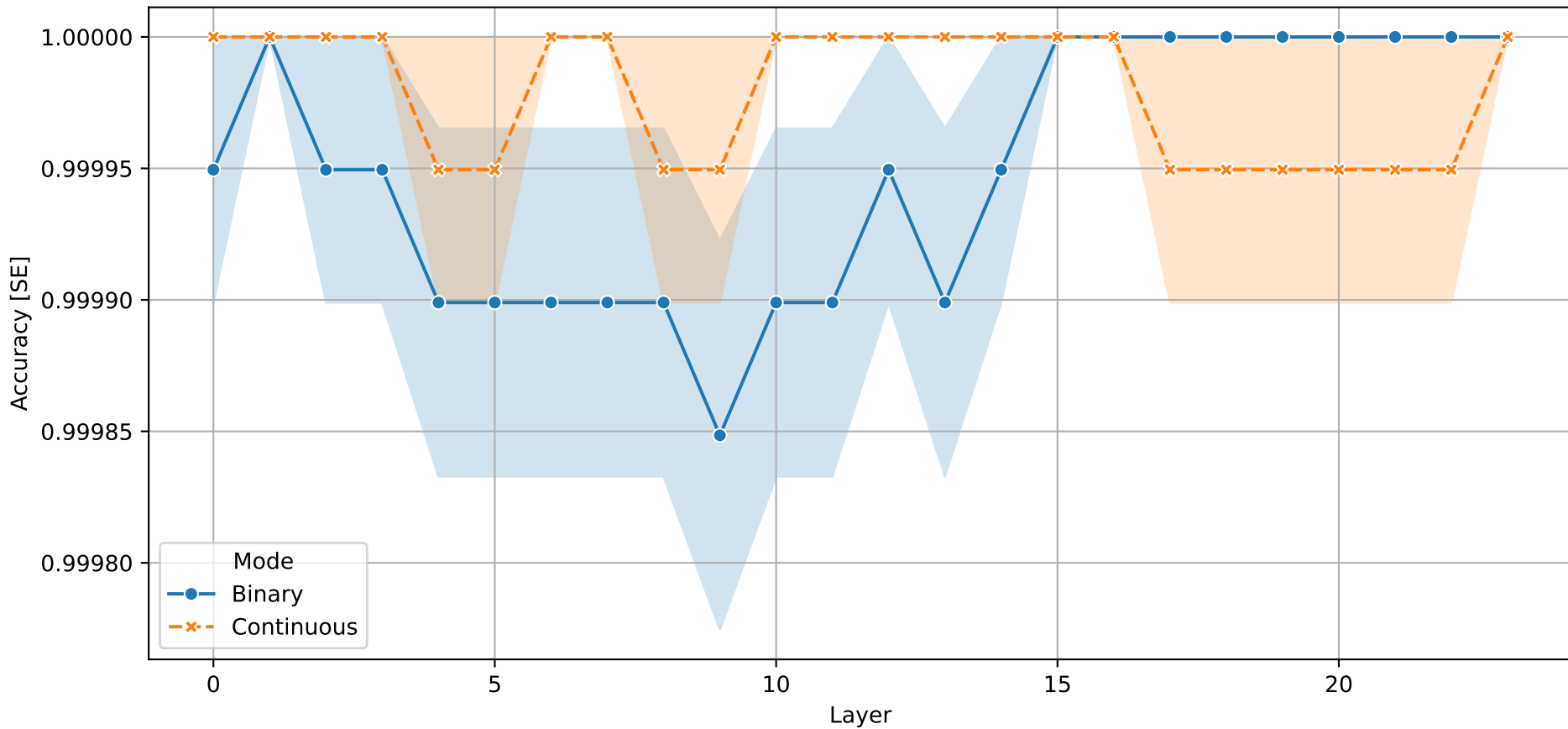
Accuracy per Layer - Single Neuron Probing



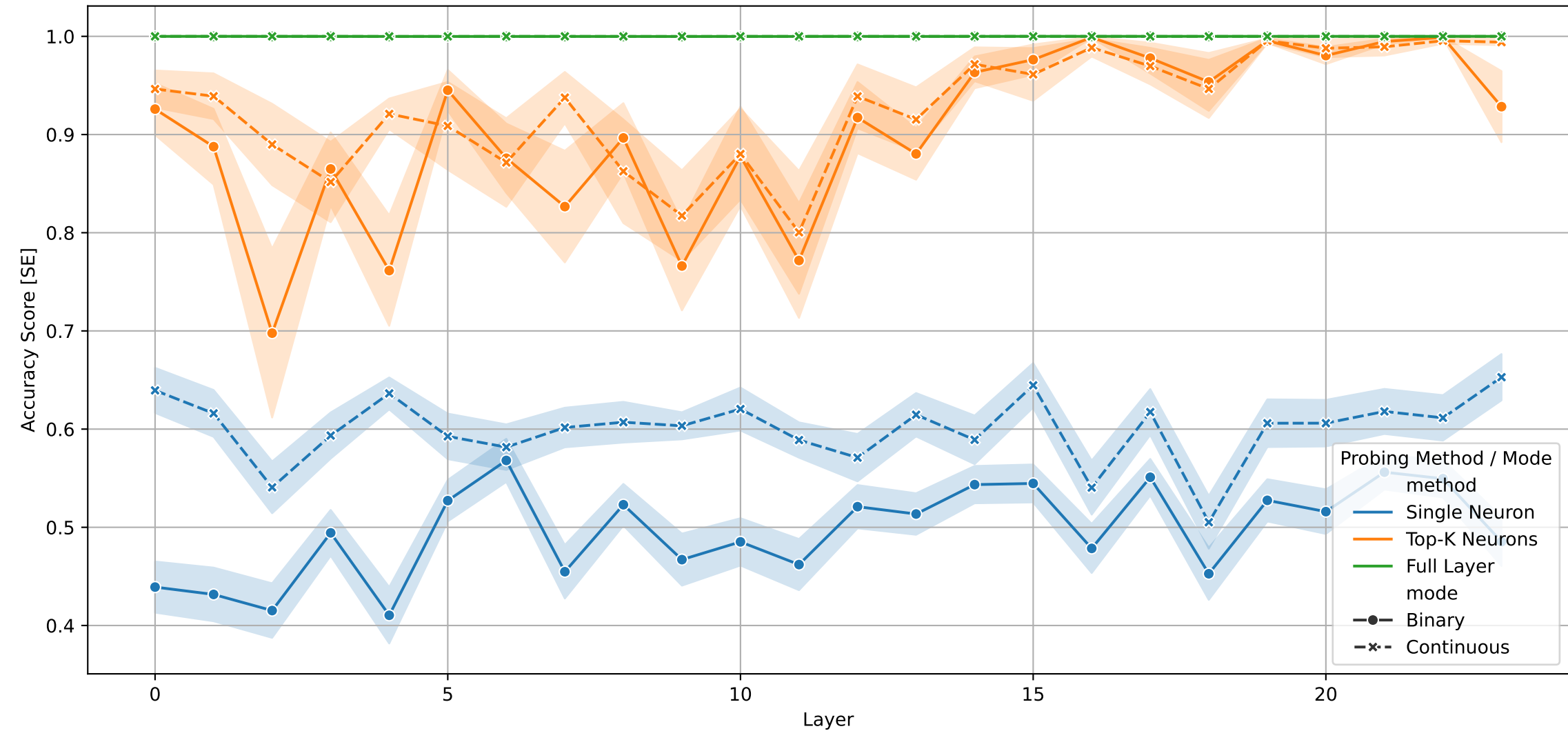
Accuracy per Layer - Top-K Neurons Probing



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	1.0	0.0
Full Layer	accuracy_max	1.0	1.0
Full Layer	accuracy_mean	0.9999	1.0
Full Layer	accuracy_std	0.0001	0.0001
Single Neuron	accuracy_best_layer	6.0	23.0
Single Neuron	accuracy_max	1.0	1.0
Single Neuron	accuracy_mean	0.4966	0.5999
Single Neuron	accuracy_std	0.2138	0.2055
Top-K Neurons	accuracy_best_layer	22.0	19.0
Top-K Neurons	accuracy_max	1.0	1.0
Top-K Neurons	accuracy_mean	0.9026	0.9284
Top-K Neurons	accuracy_std	0.1298	0.1023