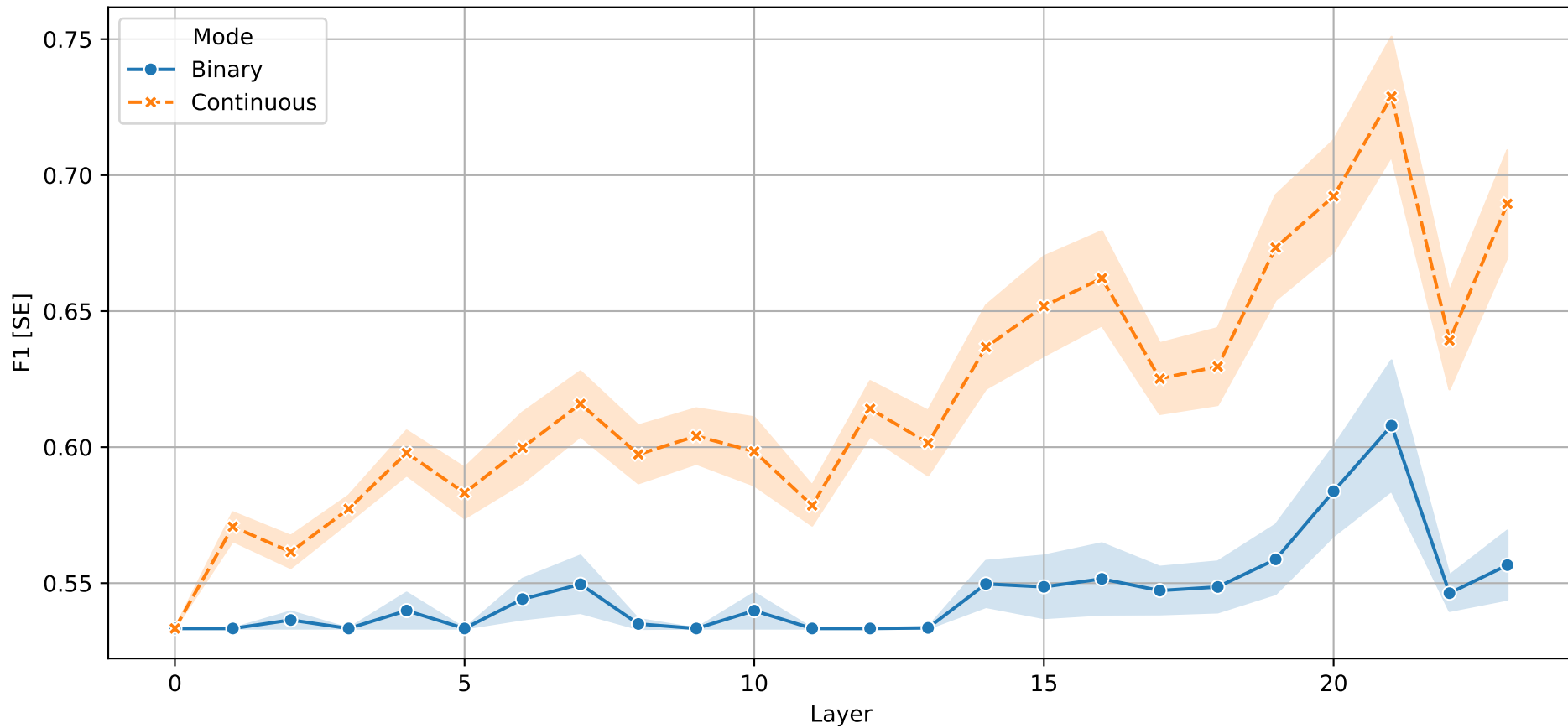
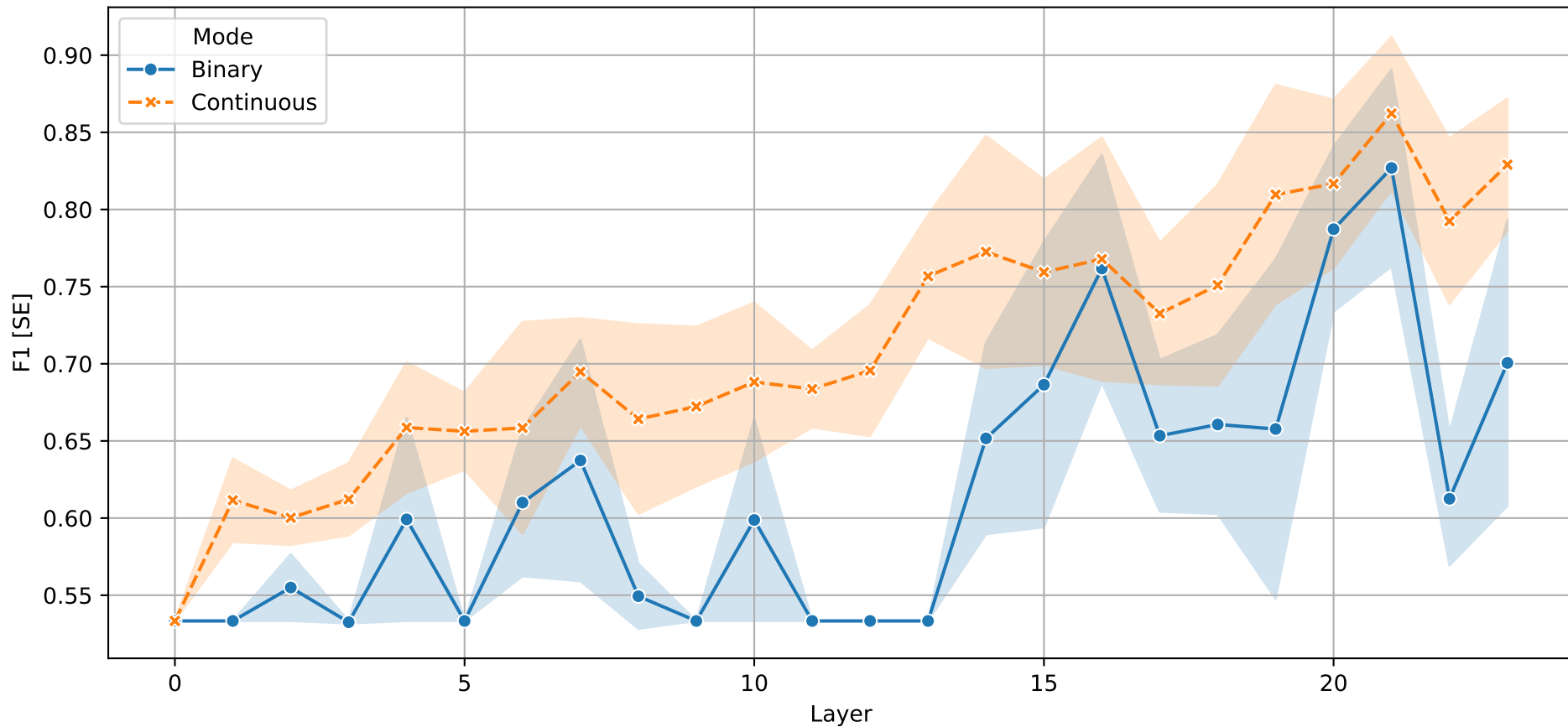


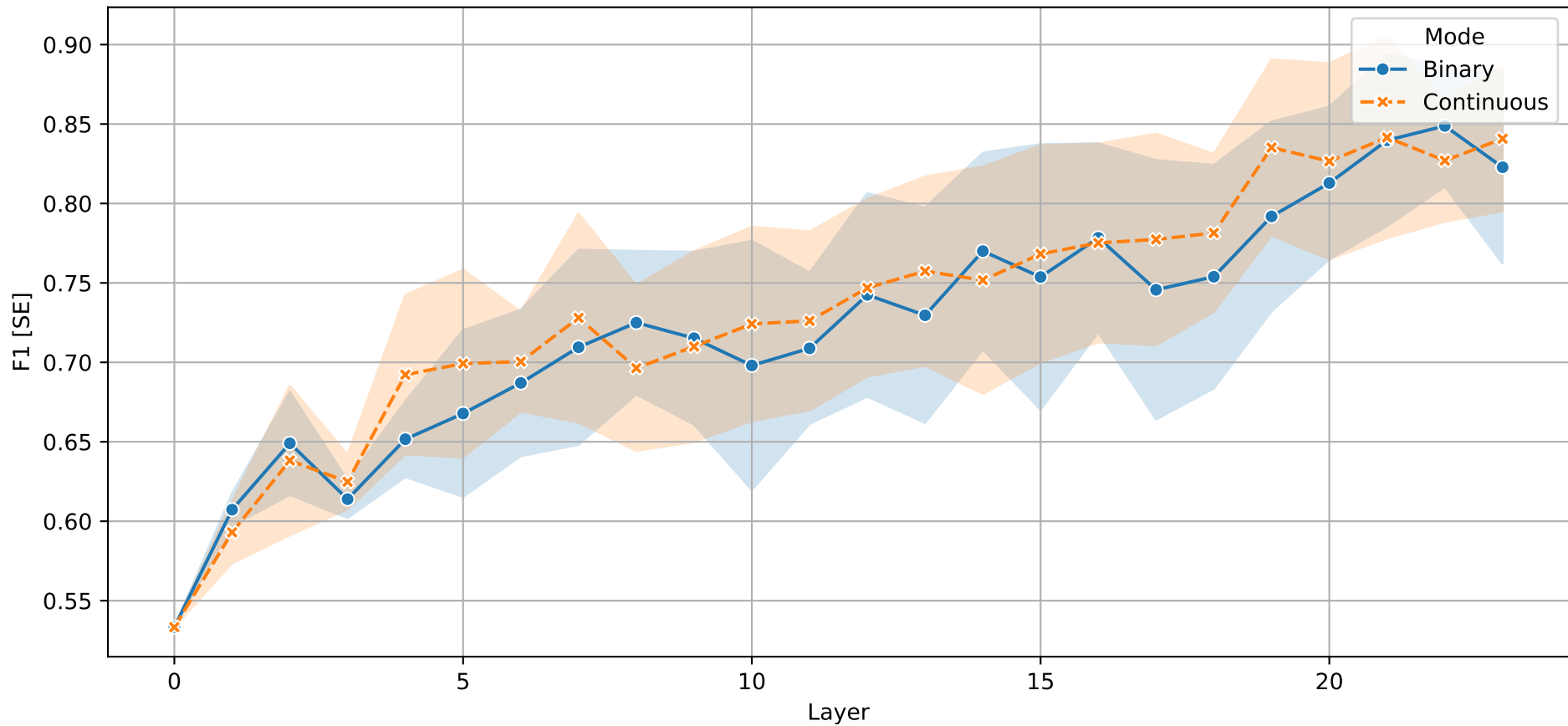
F1 per Layer - Single Neuron Probing



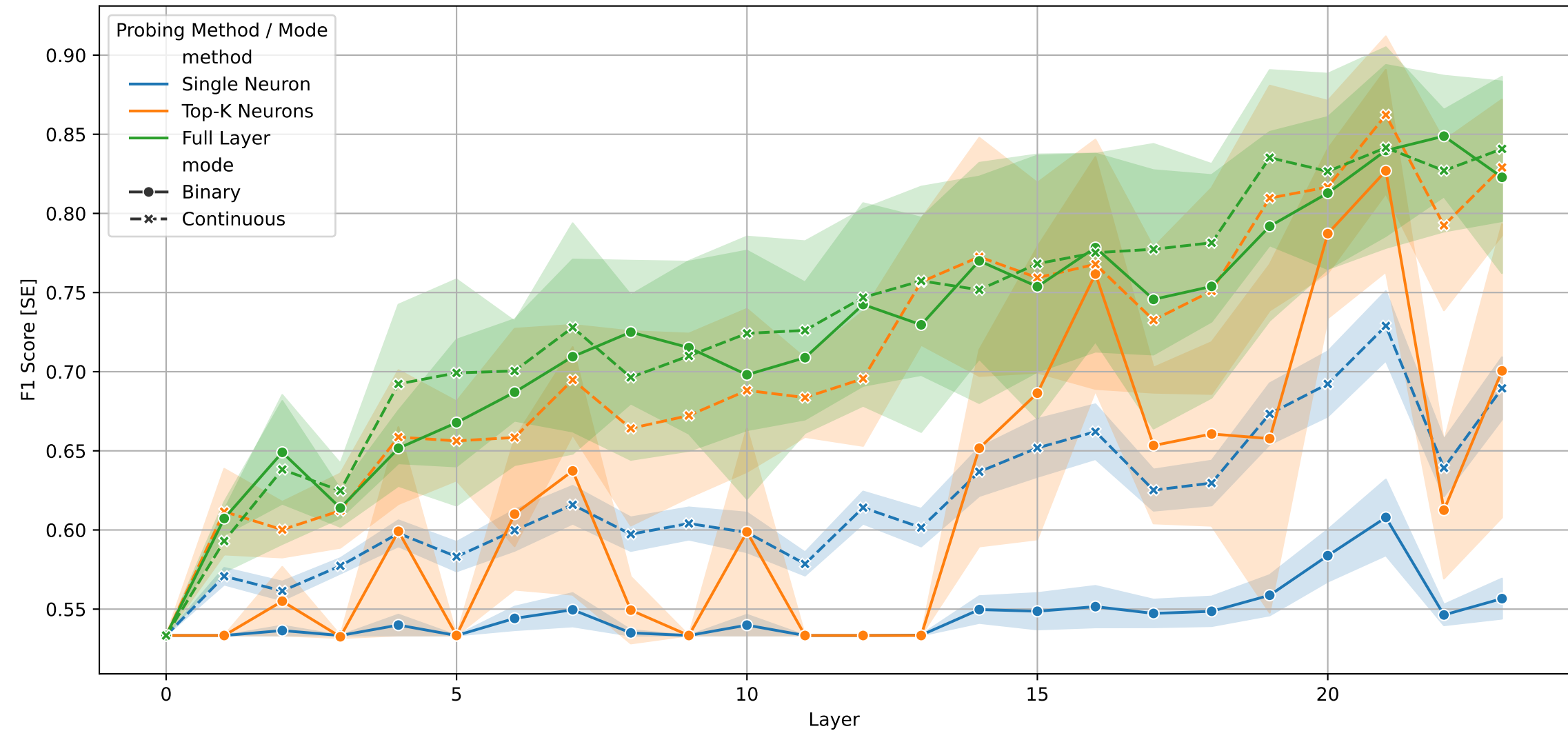
F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



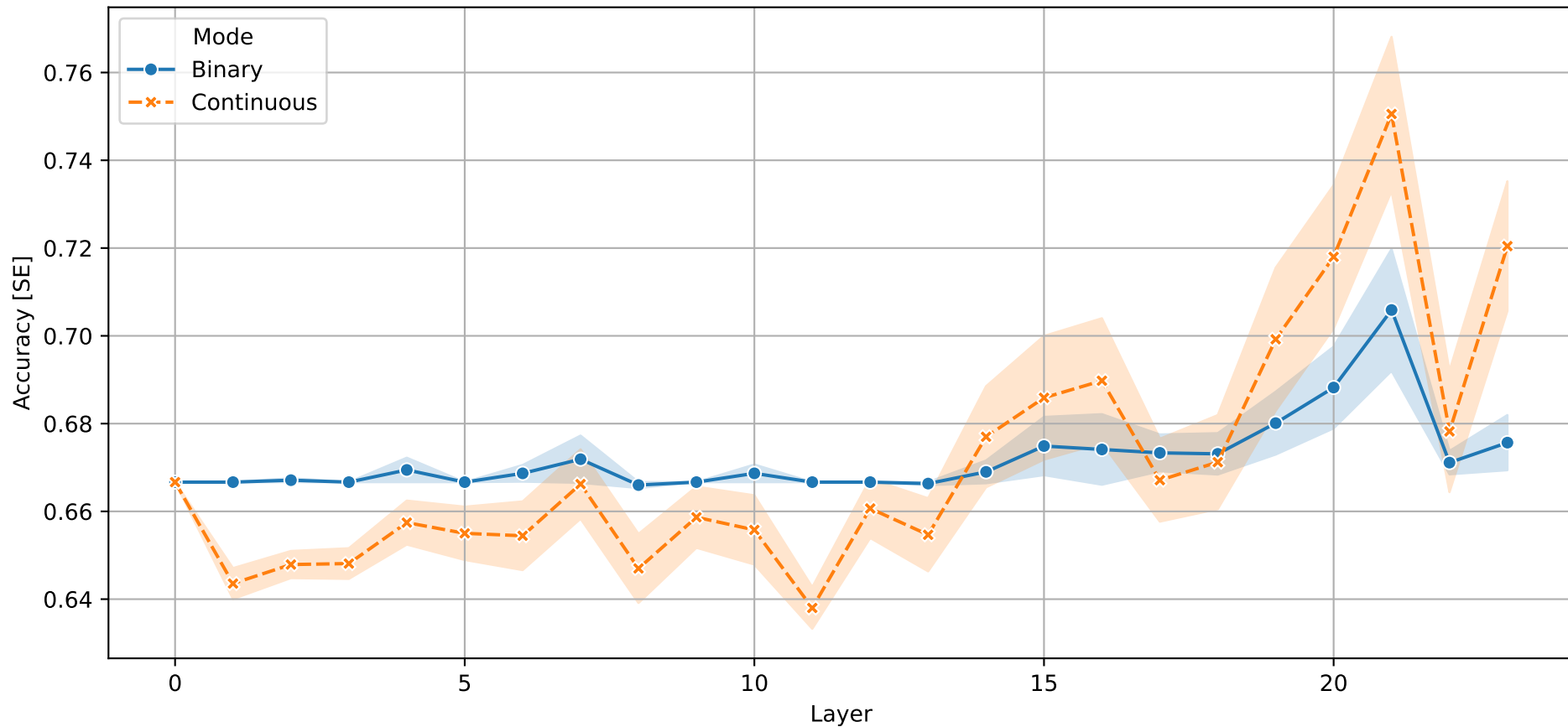
Overall F1 per Layer - All Methods



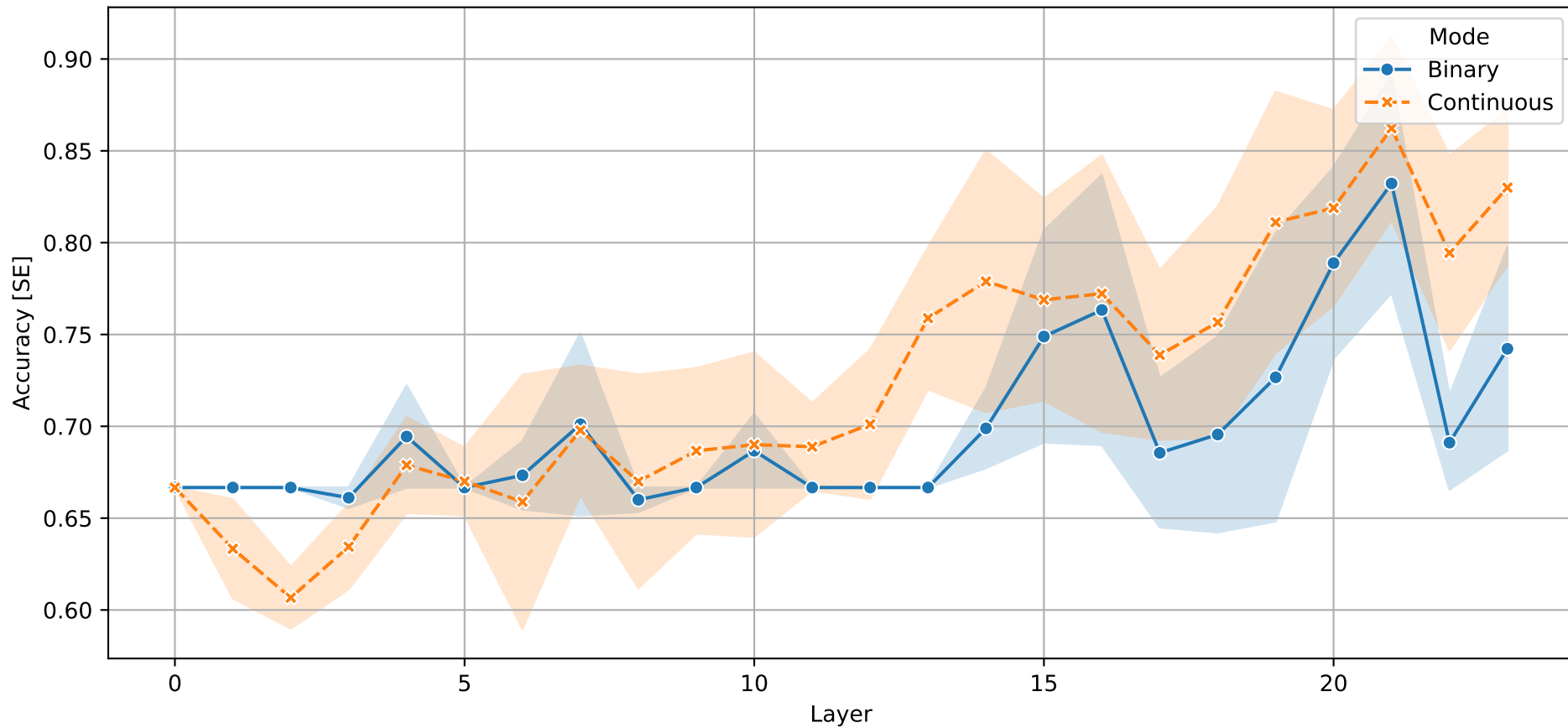
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	22.0	21.0
Full Layer	f1_max	0.9465	0.9599
Full Layer	f1_mean	0.7232	0.7331
Full Layer	f1_std	0.1093	0.1096
Single Neuron	f1_best_layer	21.0	21.0
Single Neuron	f1_max	0.9454	0.9599
Single Neuron	f1_mean	0.5463	0.6193
Single Neuron	f1_std	0.0523	0.0862
Top-K Neurons	f1_best_layer	21.0	21.0
Top-K Neurons	f1_max	0.9454	0.9599
Top-K Neurons	f1_mean	0.6173	0.7116
Top-K Neurons	f1_std	0.1162	0.1079

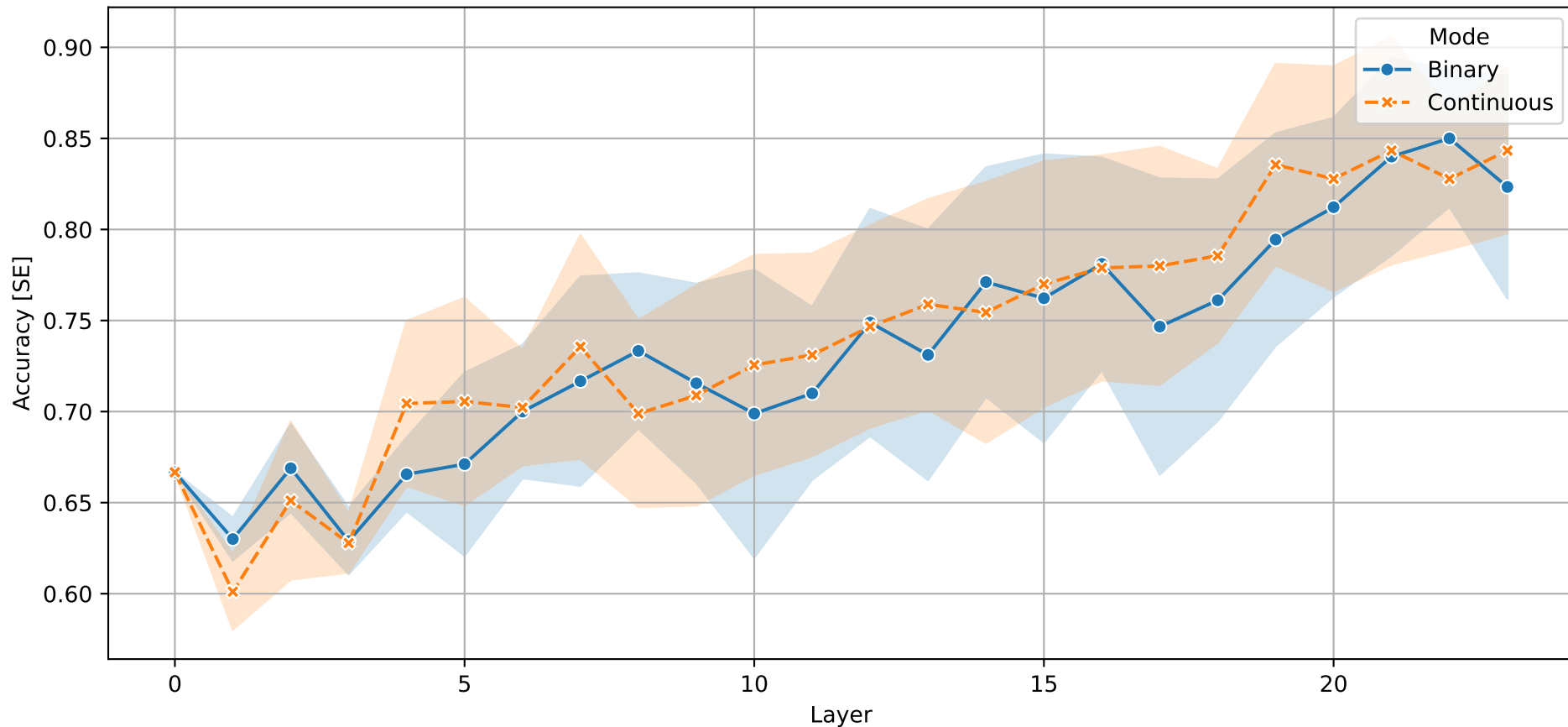
Accuracy per Layer - Single Neuron Probing



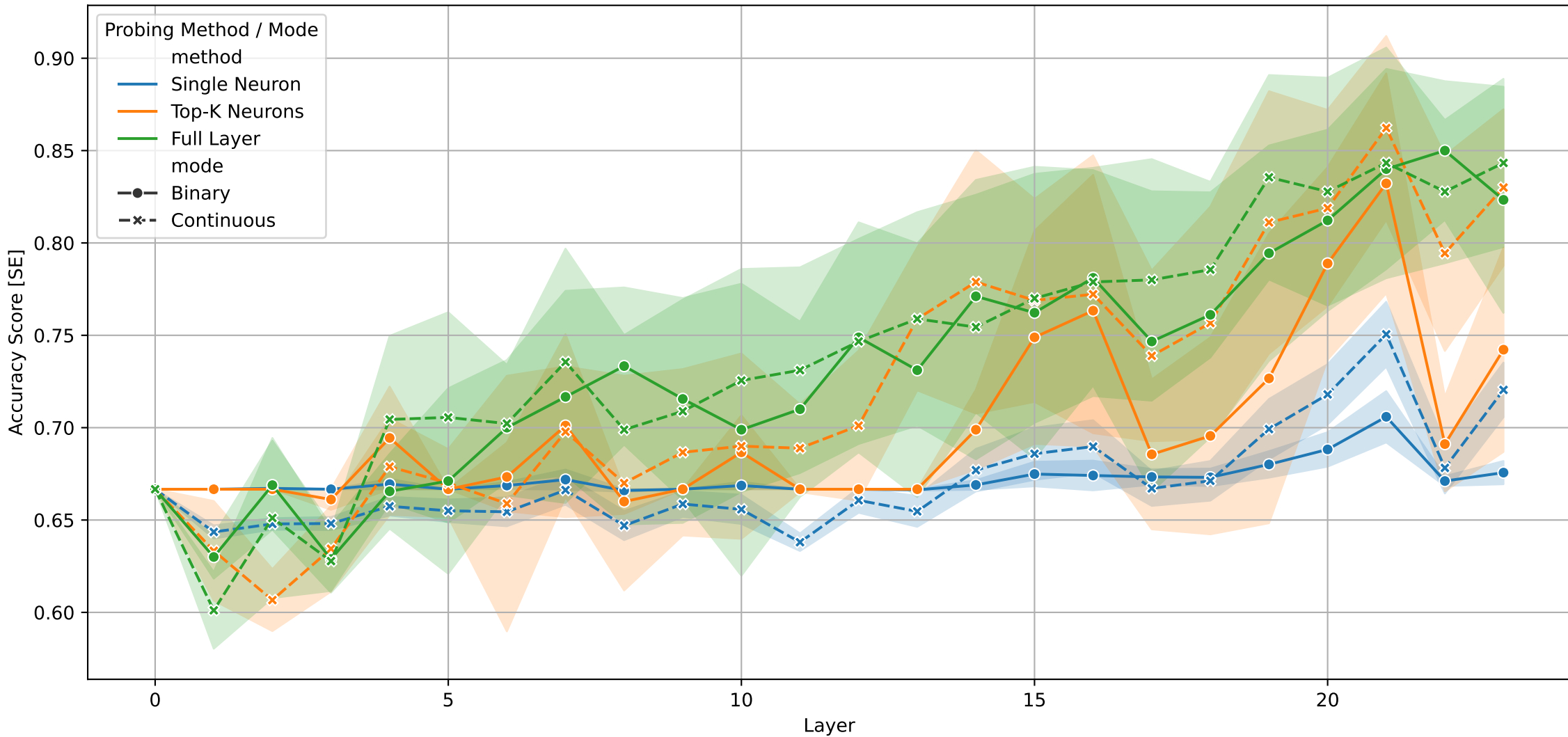
Accuracy per Layer - Top-K Neurons Probing



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	22.0	21.0
Full Layer	accuracy_max	0.9467	0.96
Full Layer	accuracy_mean	0.7345	0.7421
Full Layer	accuracy_std	0.0983	0.1001
Single Neuron	accuracy_best_layer	21.0	21.0
Single Neuron	accuracy_max	0.9467	0.96
Single Neuron	accuracy_mean	0.6721	0.6713
Single Neuron	accuracy_std	0.0278	0.0614
Top-K Neurons	accuracy_best_layer	21.0	21.0
Top-K Neurons	accuracy_max	0.9467	0.96
Top-K Neurons	accuracy_mean	0.6993	0.7239
Top-K Neurons	accuracy_std	0.0698	0.0969