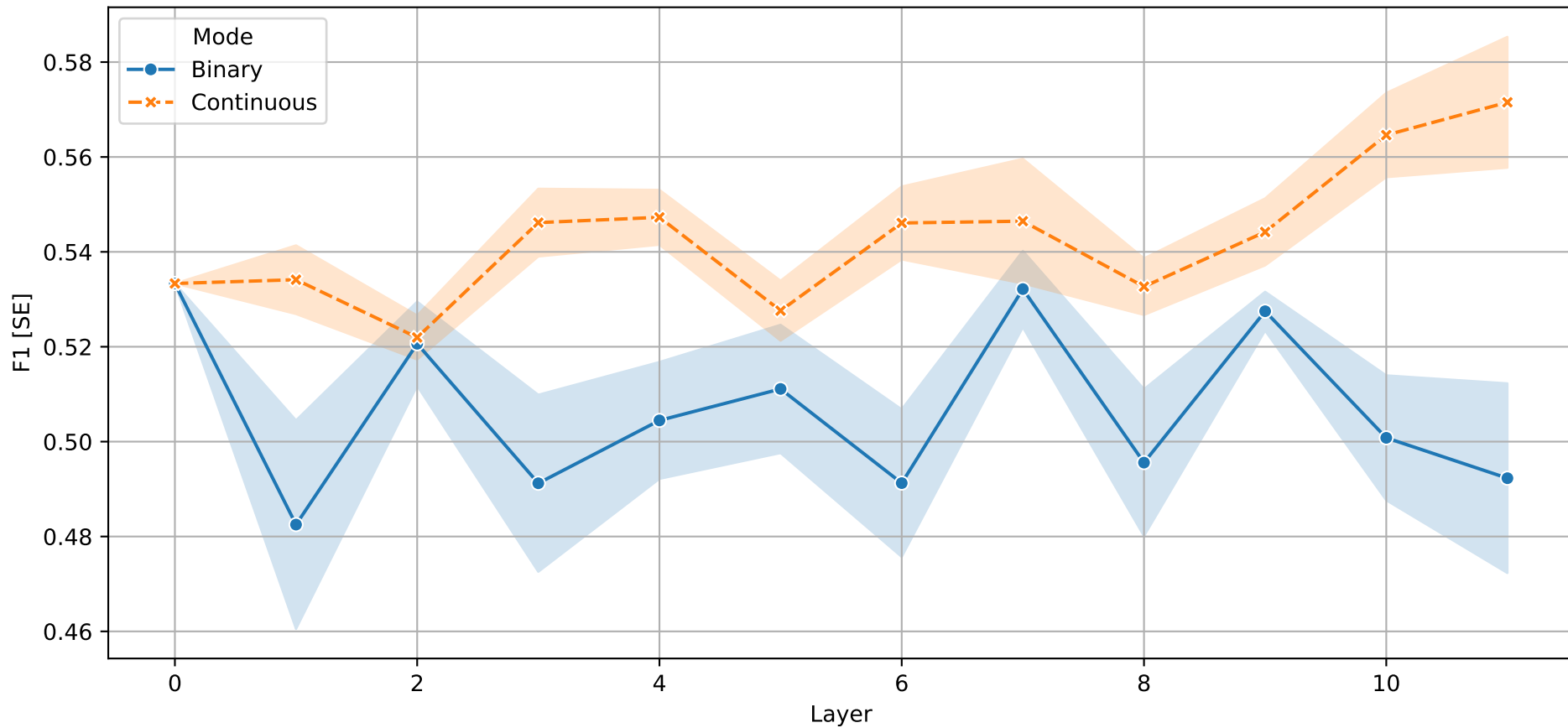
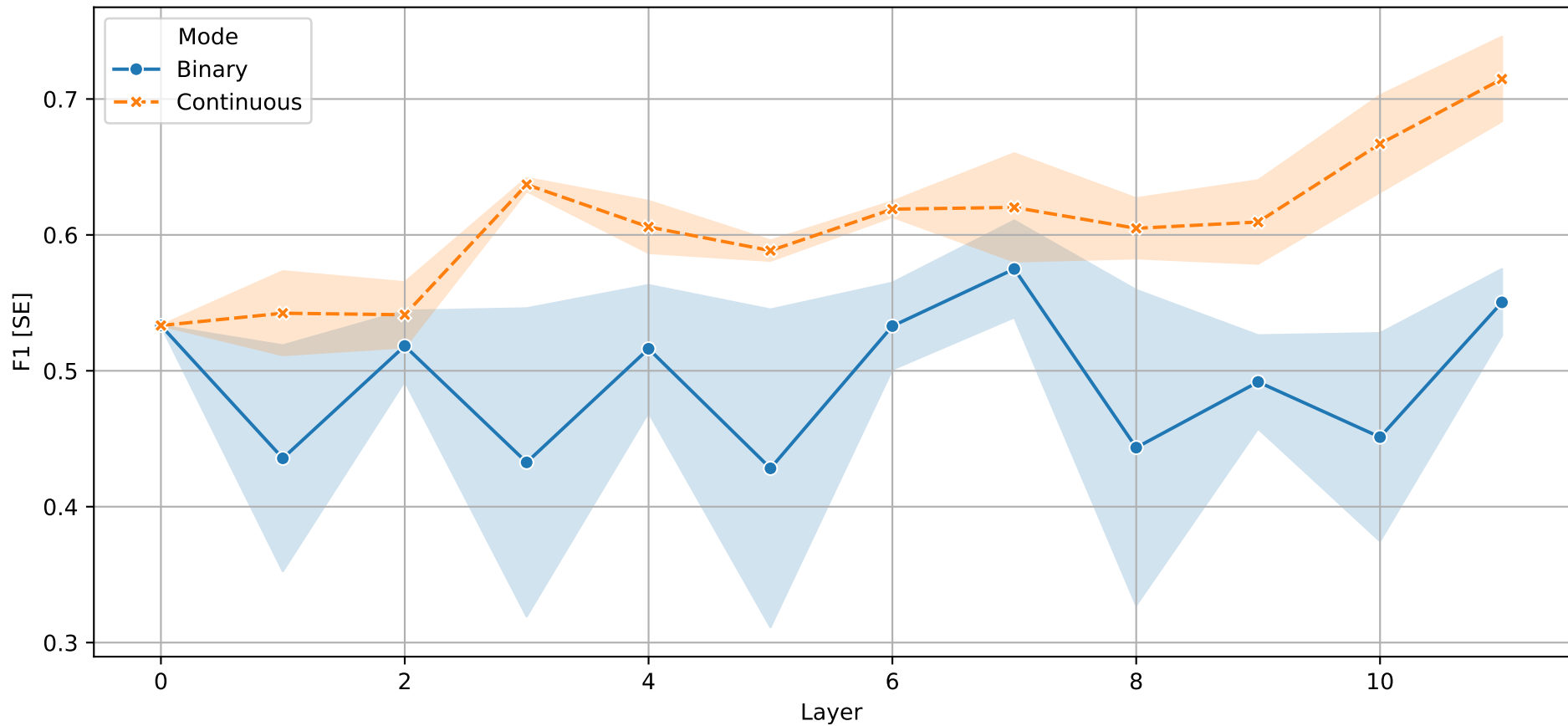


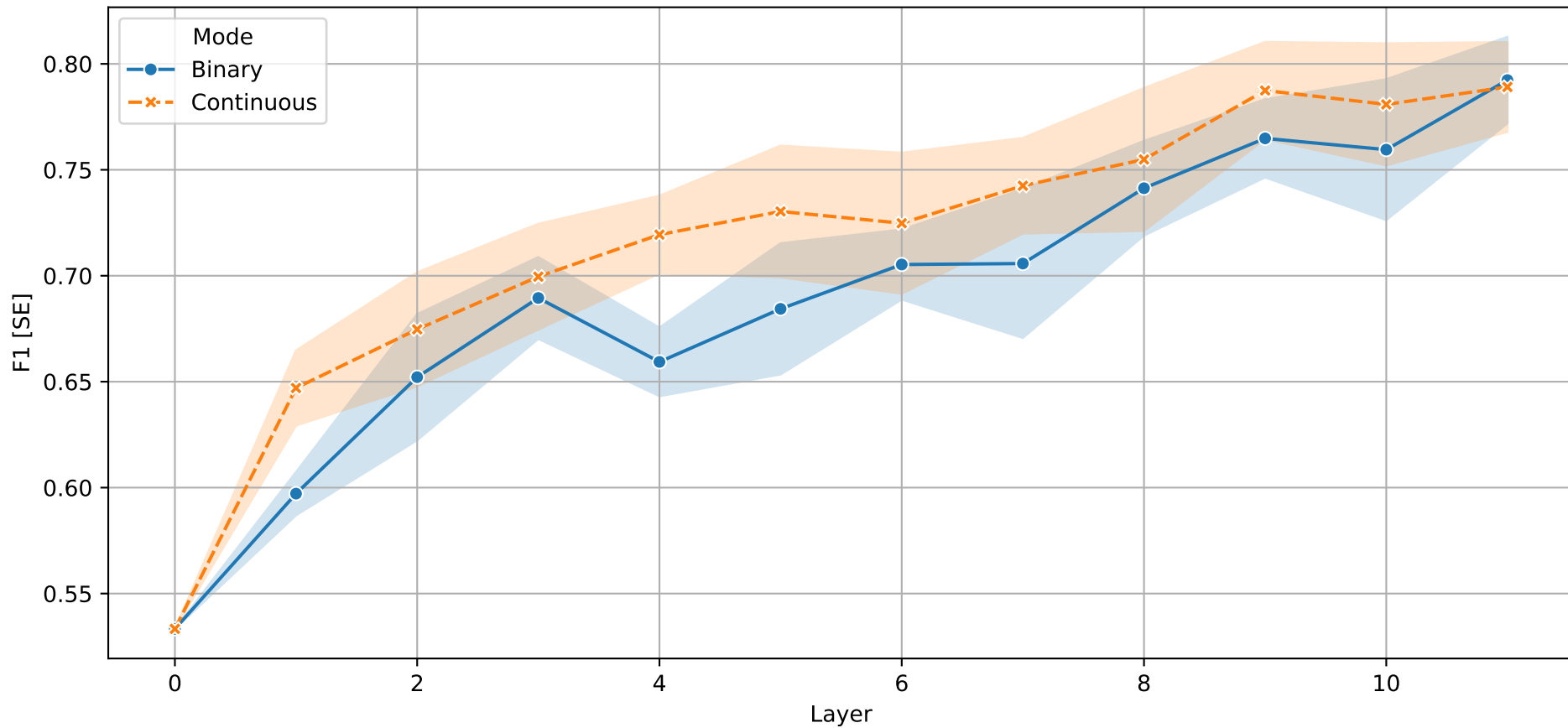
F1 per Layer - Single Neuron Probing



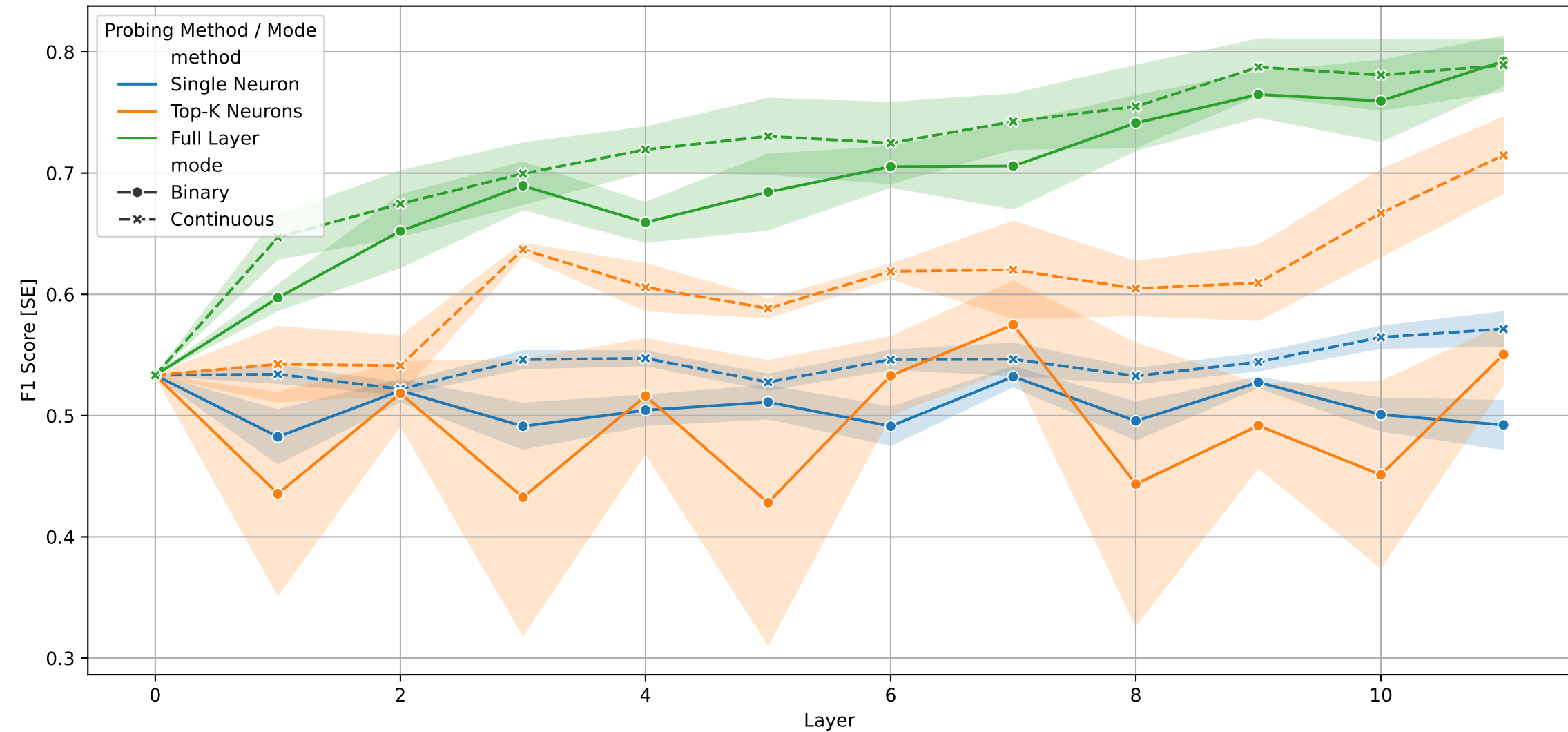
F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



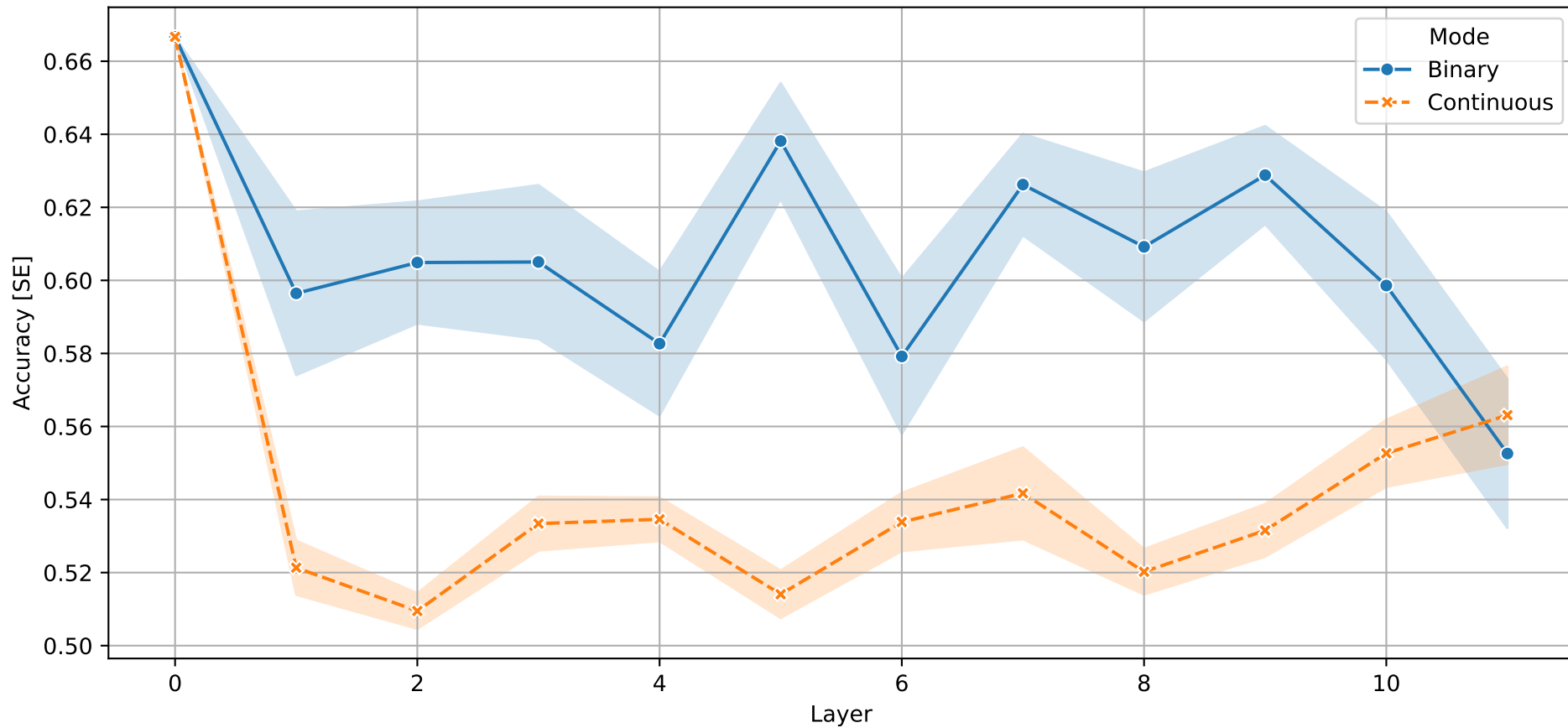
Overall F1 per Layer - All Methods



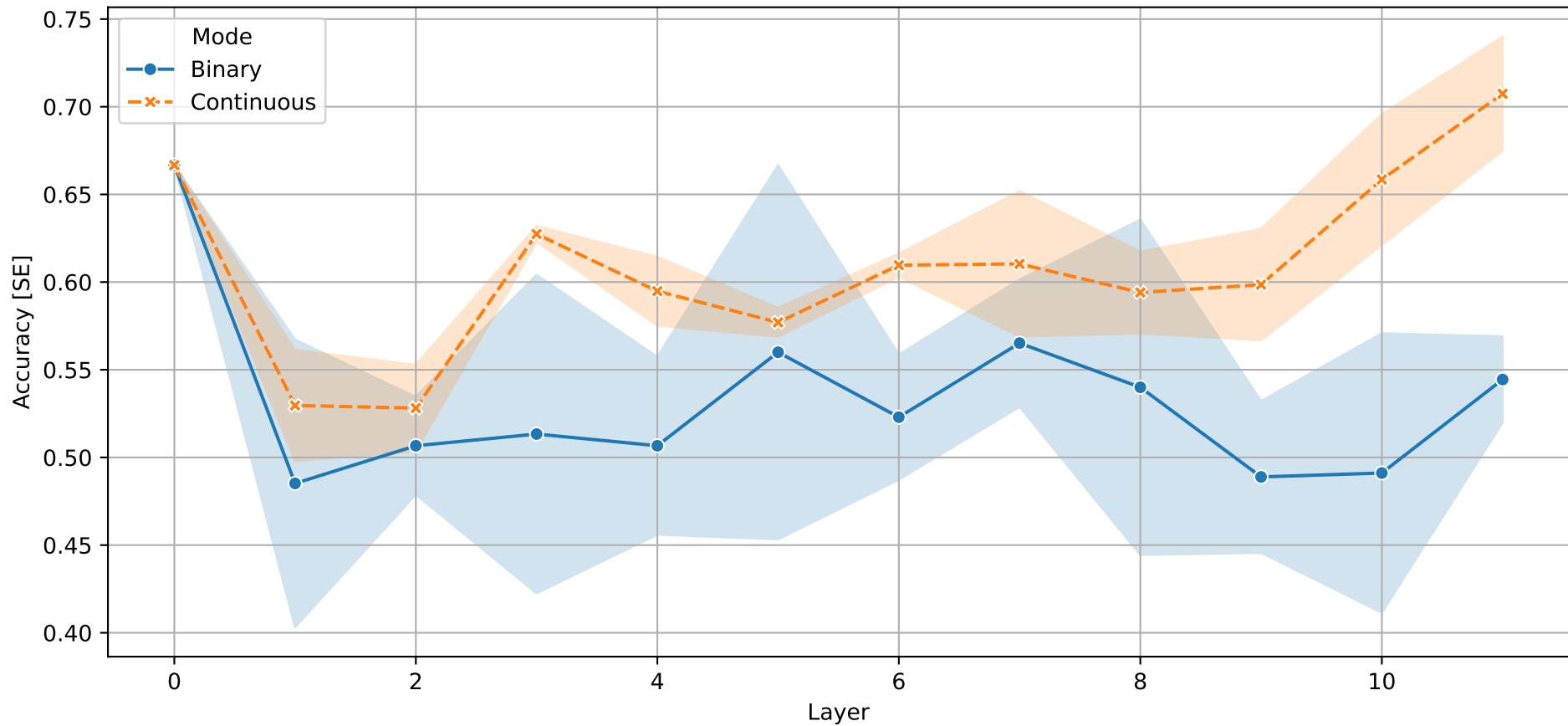
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	11.0	11.0
Full Layer	f1_max	0.8323	0.8384
Full Layer	f1_mean	0.6904	0.7153
Full Layer	f1_std	0.0788	0.0787
Single Neuron	f1_best_layer	0.0	11.0
Single Neuron	f1_max	0.6725	0.7557
Single Neuron	f1_mean	0.5069	0.543
Single Neuron	f1_std	0.0778	0.0457
Top-K Neurons	f1_best_layer	7.0	11.0
Top-K Neurons	f1_max	0.646	0.7756
Top-K Neurons	f1_mean	0.4924	0.607
Top-K Neurons	f1_std	0.113	0.0622

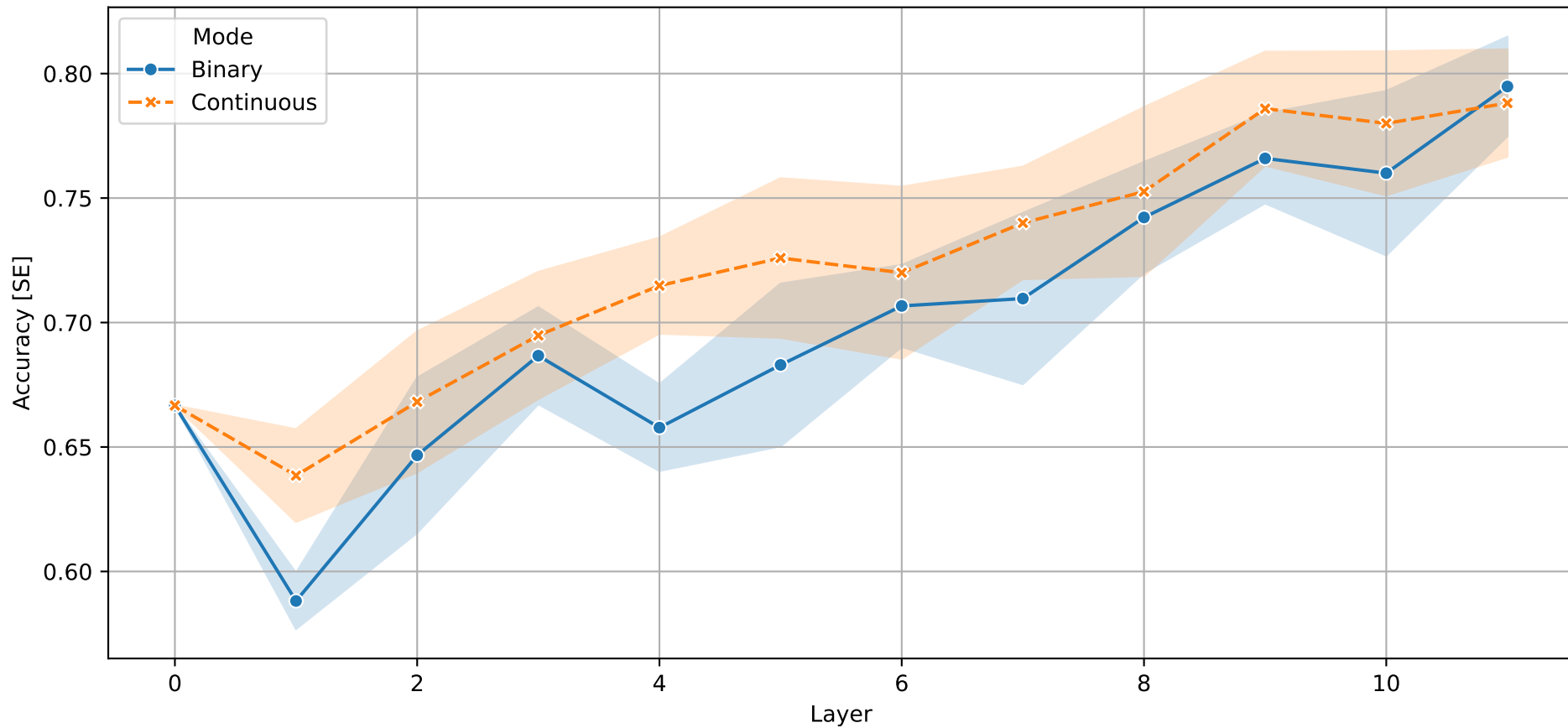
Accuracy per Layer - Single Neuron Probing



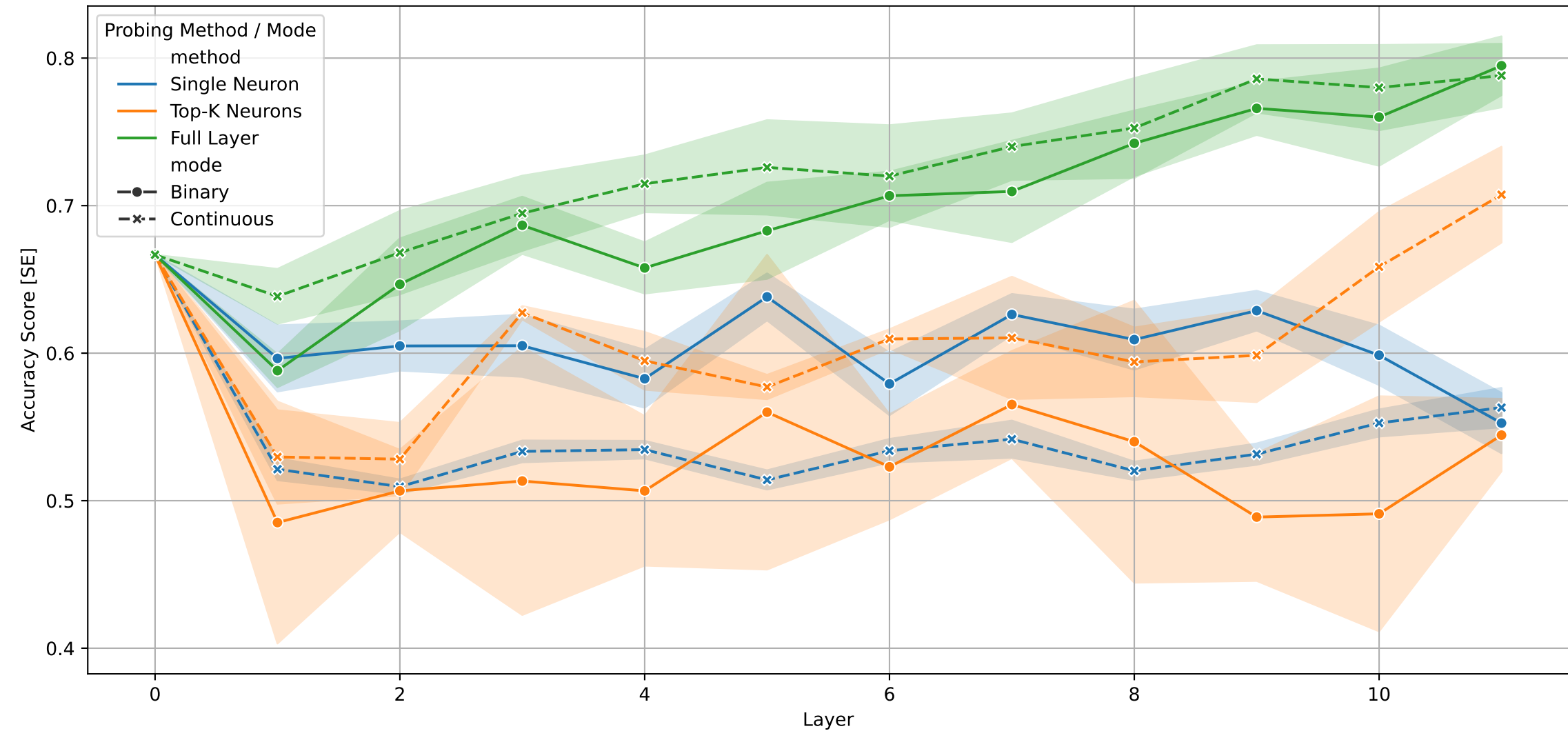
Accuracy per Layer - Top-K Neurons Probing



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	11.0	11.0
Full Layer	accuracy_max	0.8333	0.8378
Full Layer	accuracy_mean	0.7007	0.723
Full Layer	accuracy_std	0.0659	0.0604
Single Neuron	accuracy_best_layer	0.0	0.0
Single Neuron	accuracy_max	0.6822	0.7533
Single Neuron	accuracy_mean	0.6074	0.5436
Single Neuron	accuracy_std	0.1017	0.0592
Top-K Neurons	accuracy_best_layer	0.0	11.0
Top-K Neurons	accuracy_max	0.6689	0.7711
Top-K Neurons	accuracy_mean	0.5326	0.6085
Top-K Neurons	accuracy_std	0.1048	0.0627