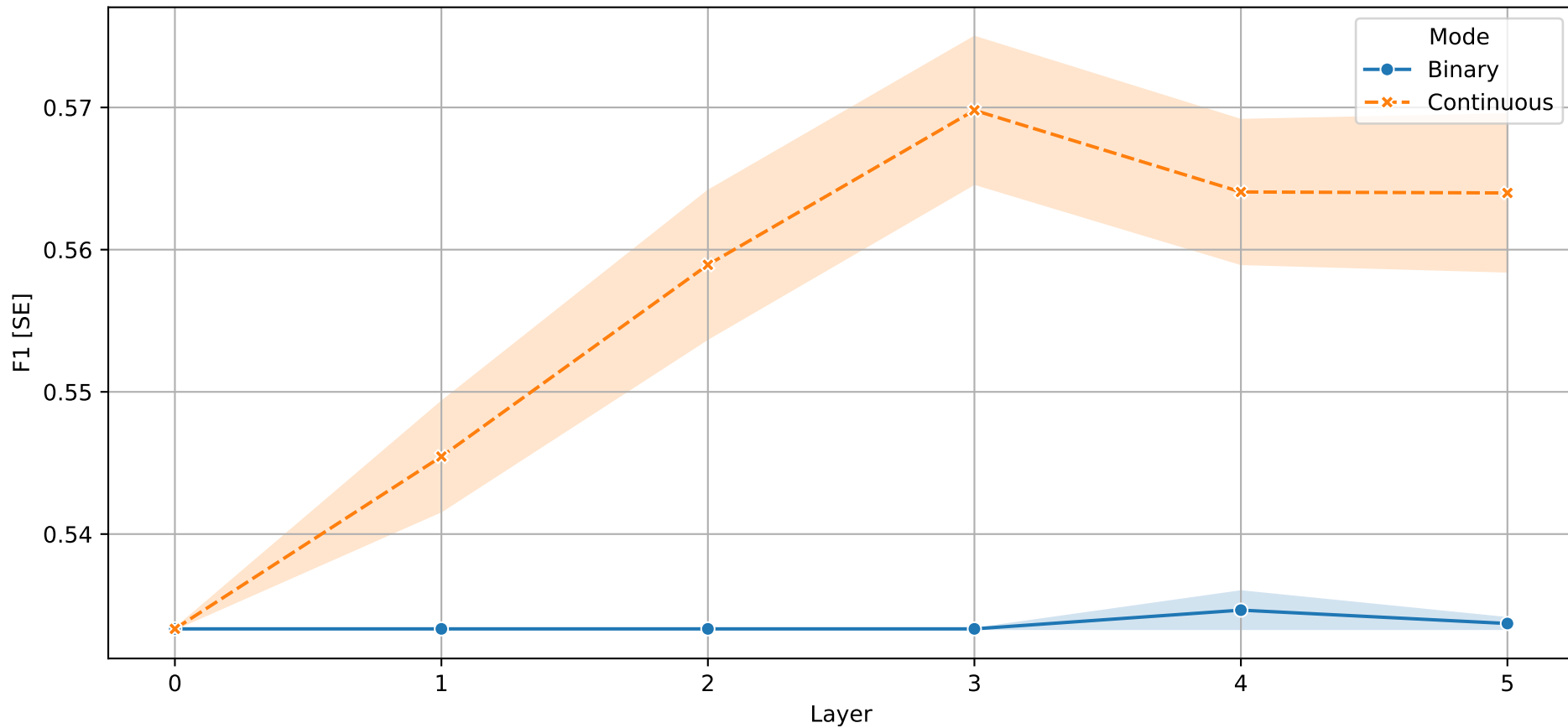
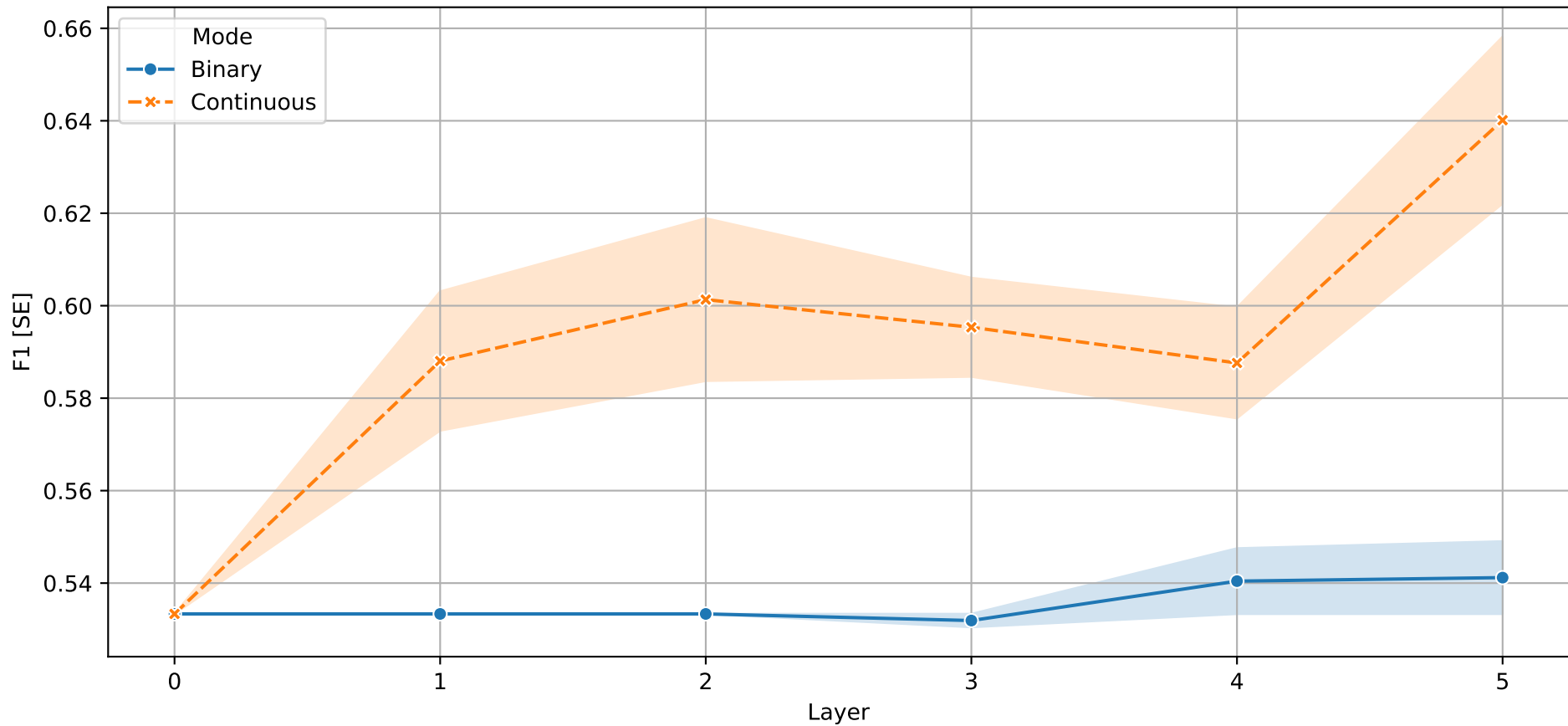


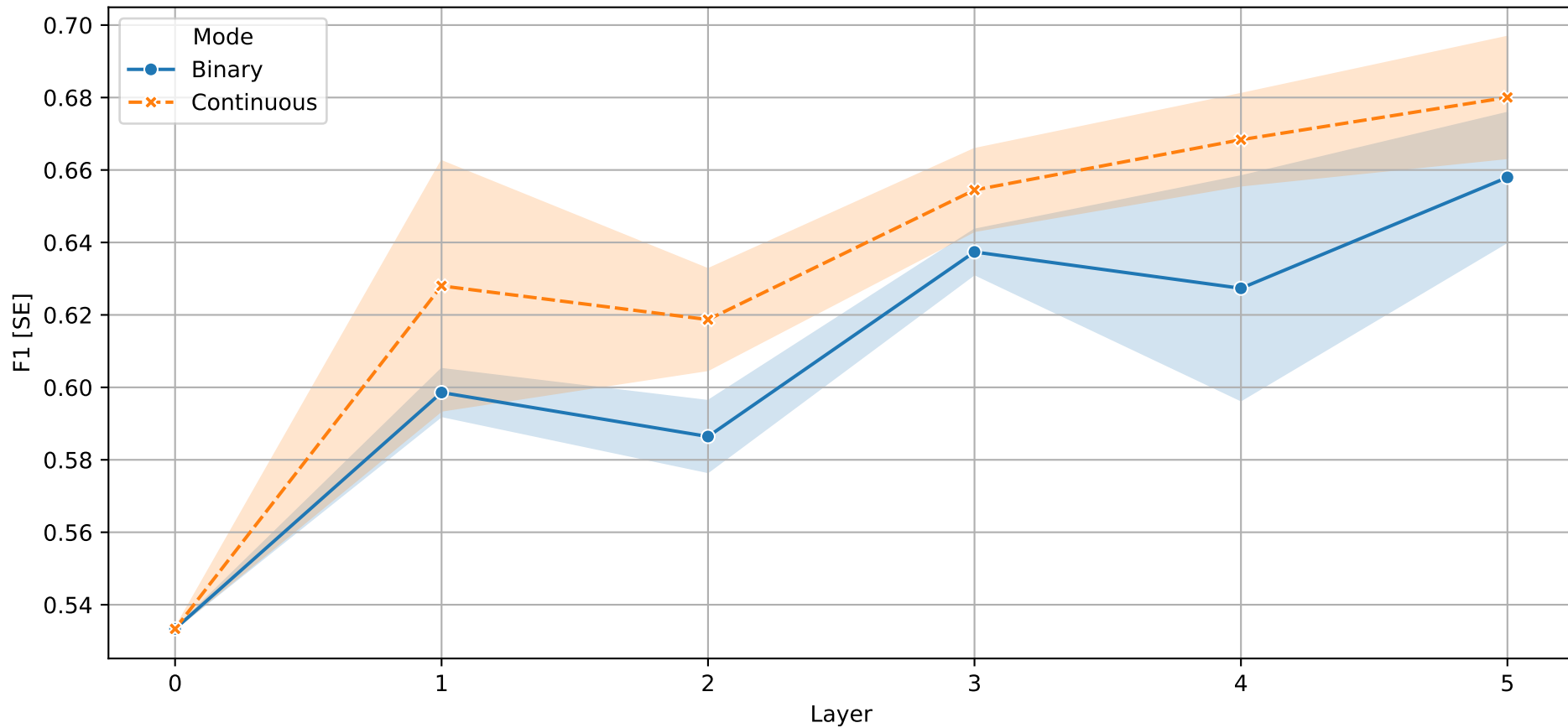
F1 per Layer - Single Neuron Probing



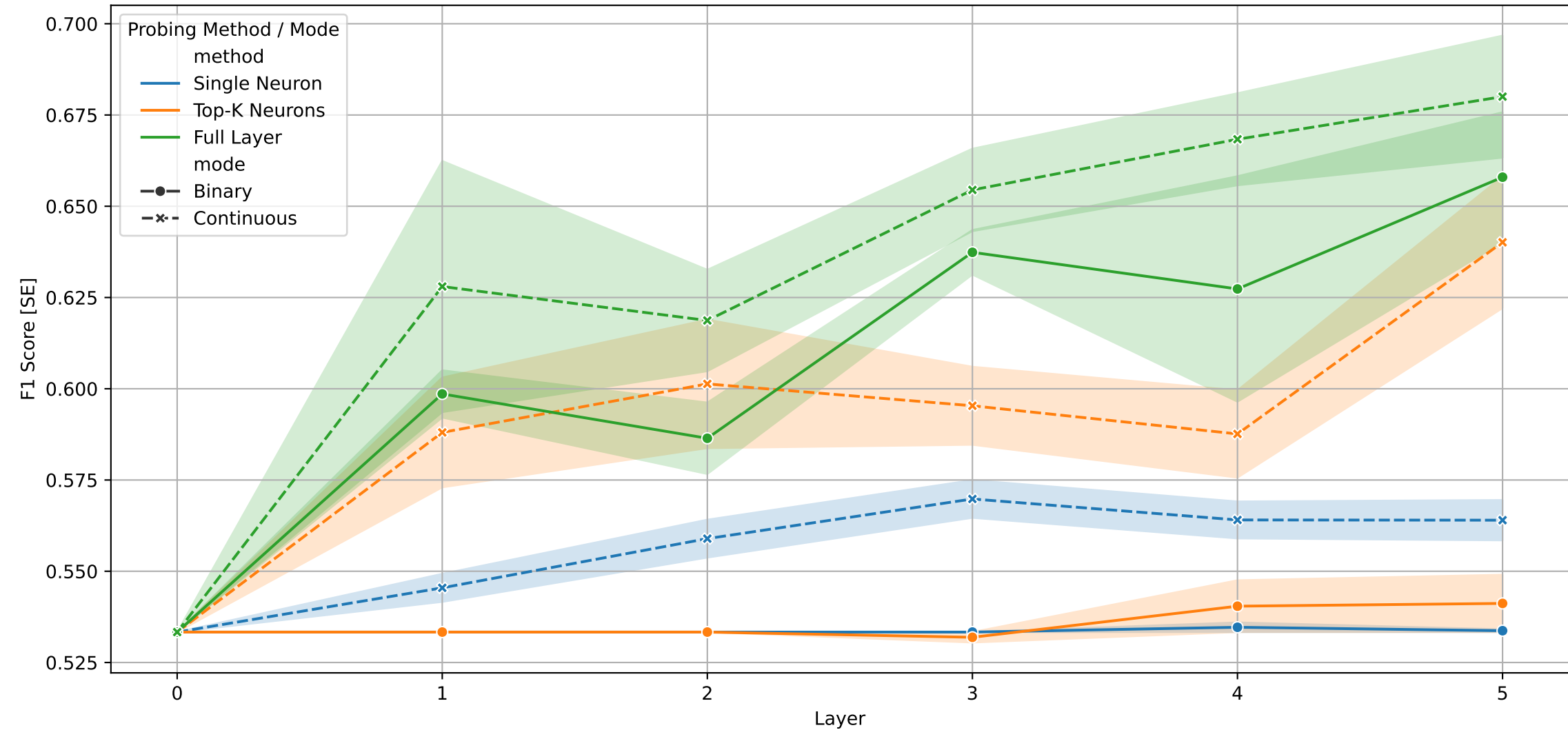
F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



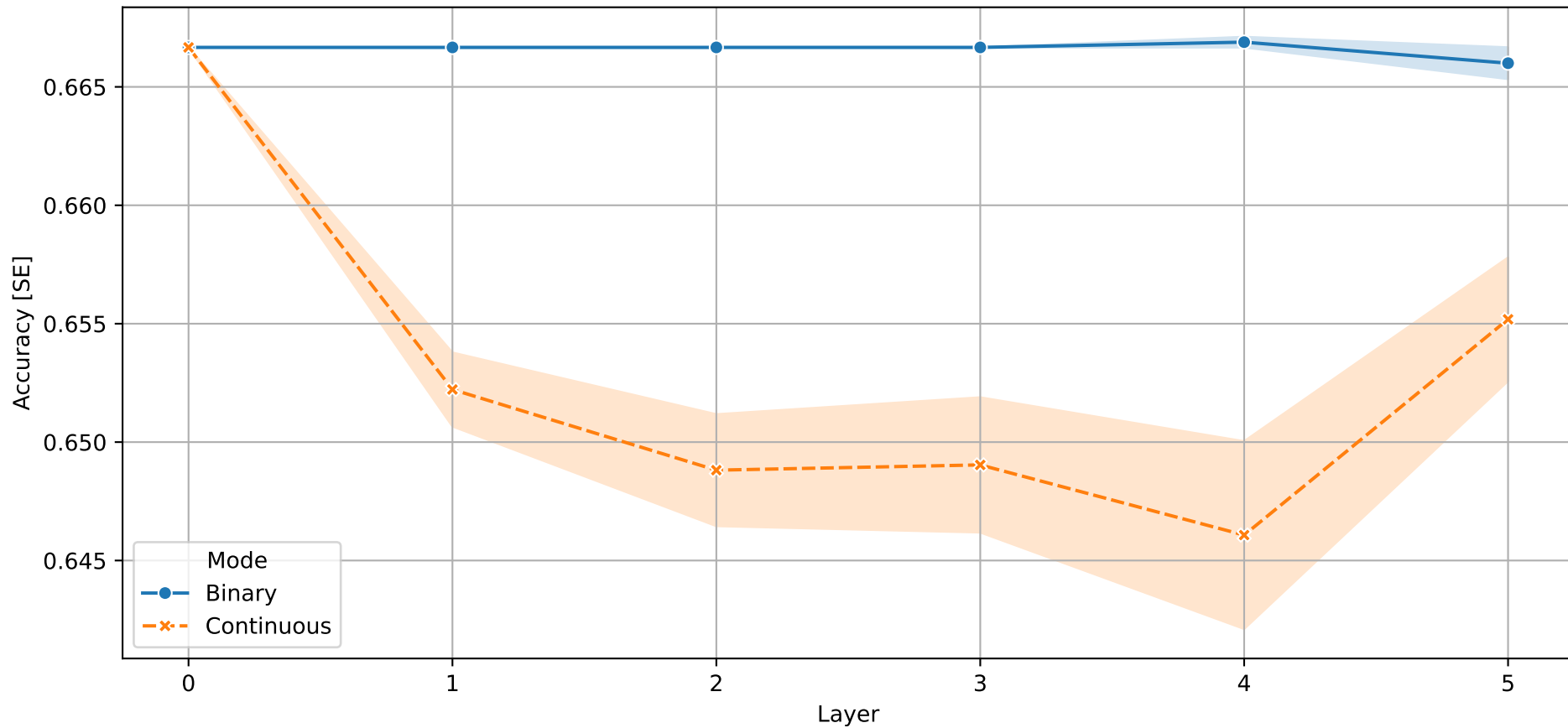
Overall F1 per Layer - All Methods



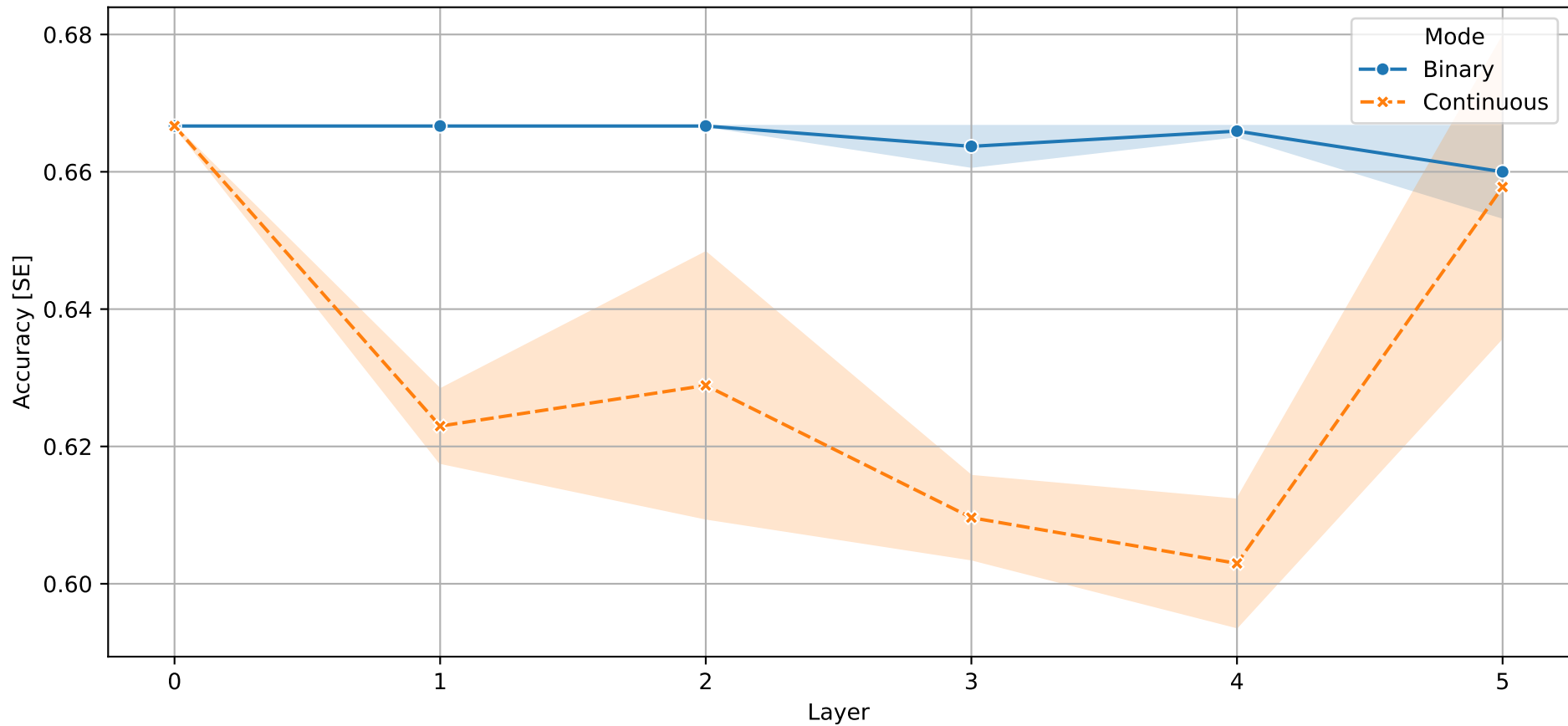
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	5.0	5.0
Full Layer	f1_max	0.6935	0.7071
Full Layer	f1_mean	0.6068	0.6305
Full Layer	f1_std	0.0474	0.0563
Single Neuron	f1_best_layer	4.0	3.0
Single Neuron	f1_max	0.5729	0.6662
Single Neuron	f1_mean	0.5336	0.5559
Single Neuron	f1_std	0.0031	0.0277
Top-K Neurons	f1_best_layer	5.0	5.0
Top-K Neurons	f1_max	0.5569	0.6718
Top-K Neurons	f1_mean	0.5356	0.591
Top-K Neurons	f1_std	0.0074	0.0379

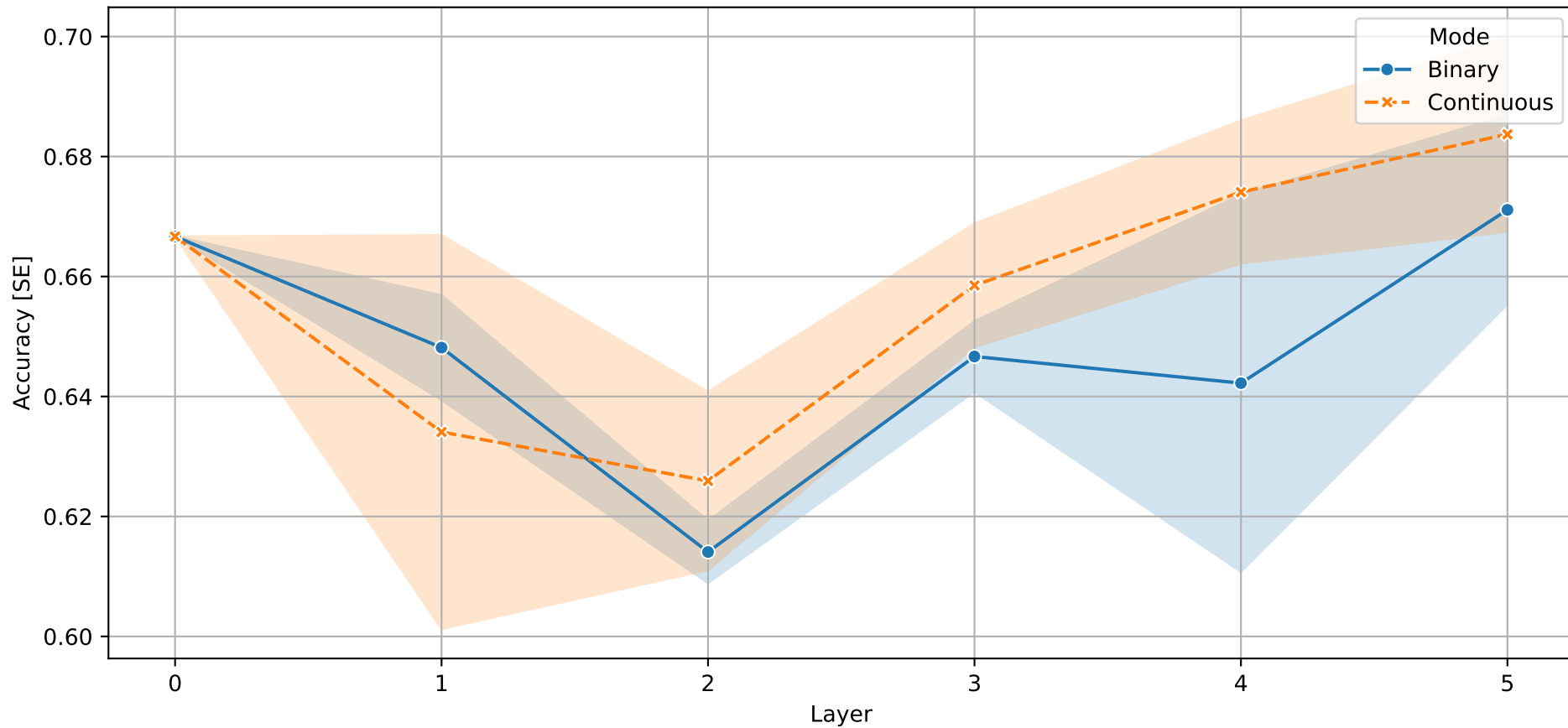
Accuracy per Layer - Single Neuron Probing



Accuracy per Layer - Top-K Neurons Probing

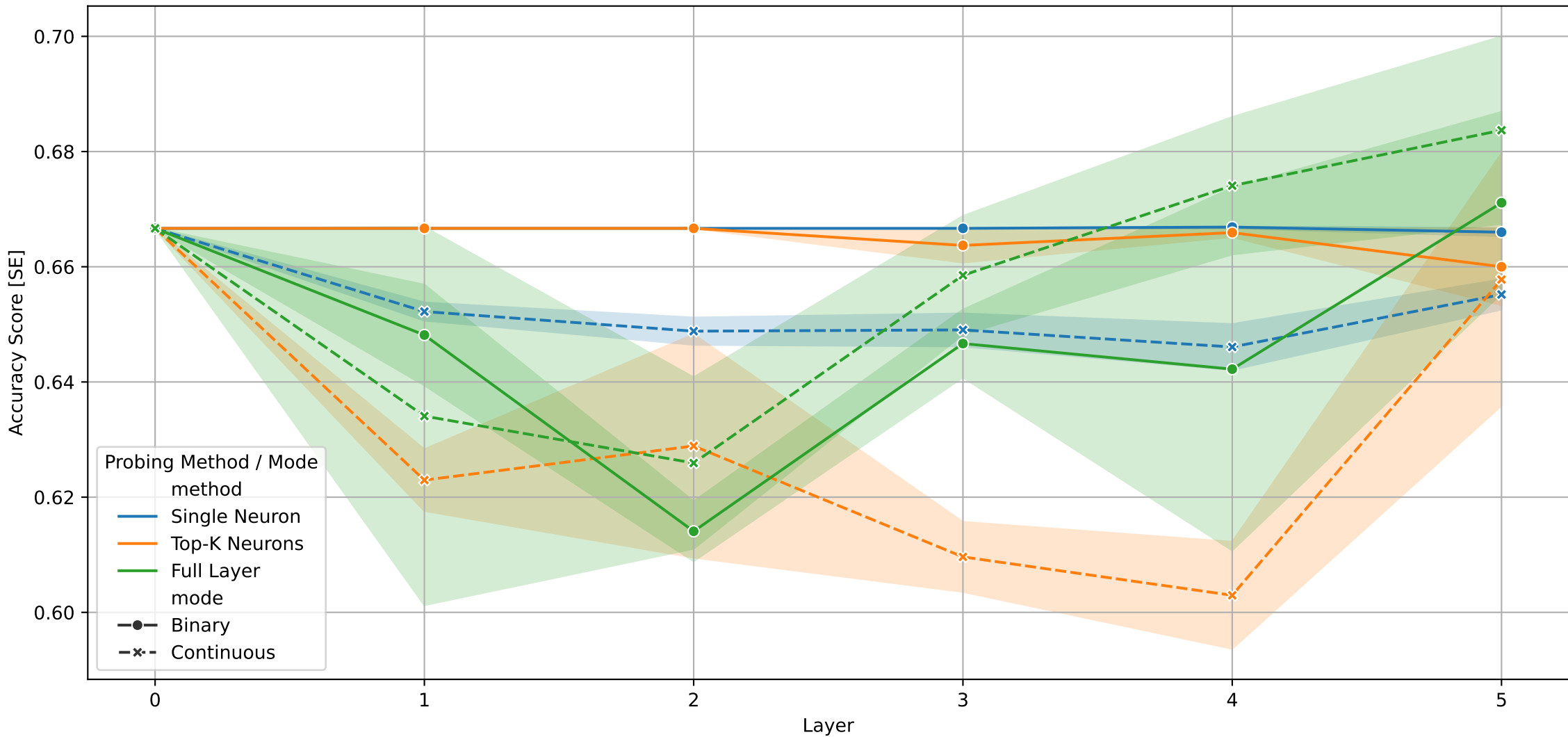


Accuracy per Layer - Full Layer Probing





Overall Accuracy per Layer - All Methods



Accuracy Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	accuracy_best_layer	5.0	5.0
Full Layer	accuracy_max	0.7022	0.7089
Full Layer	accuracy_mean	0.6481	0.6572
Full Layer	accuracy_std	0.0292	0.0331
Single Neuron	accuracy_best_layer	4.0	0.0
Single Neuron	accuracy_max	0.6733	0.6978
Single Neuron	accuracy_mean	0.6666	0.653
Single Neuron	accuracy_std	0.0016	0.0153
Top-K Neurons	accuracy_best_layer	0.0	0.0
Top-K Neurons	accuracy_max	0.6667	0.6889
Top-K Neurons	accuracy_mean	0.6649	0.6315
Top-K Neurons	accuracy_std	0.005	0.0306