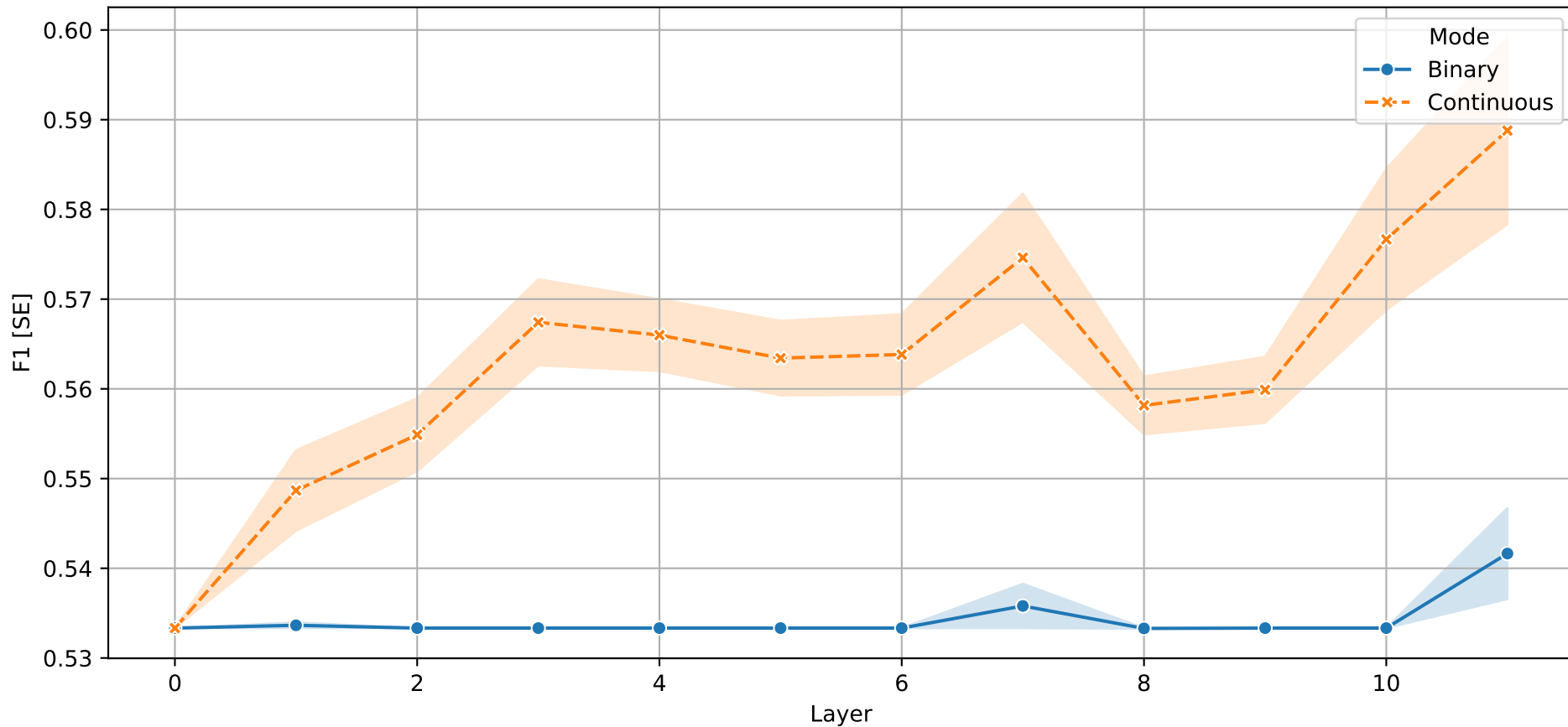
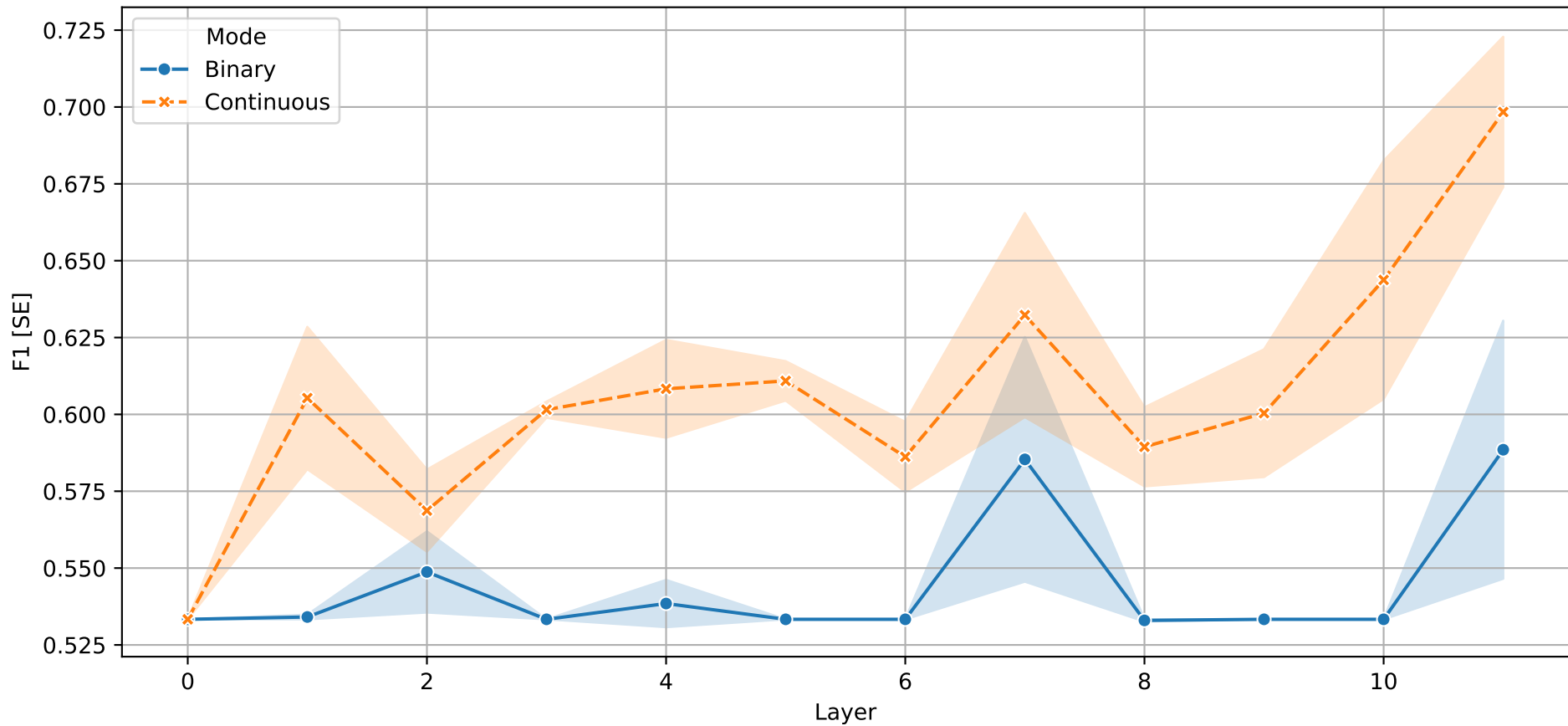


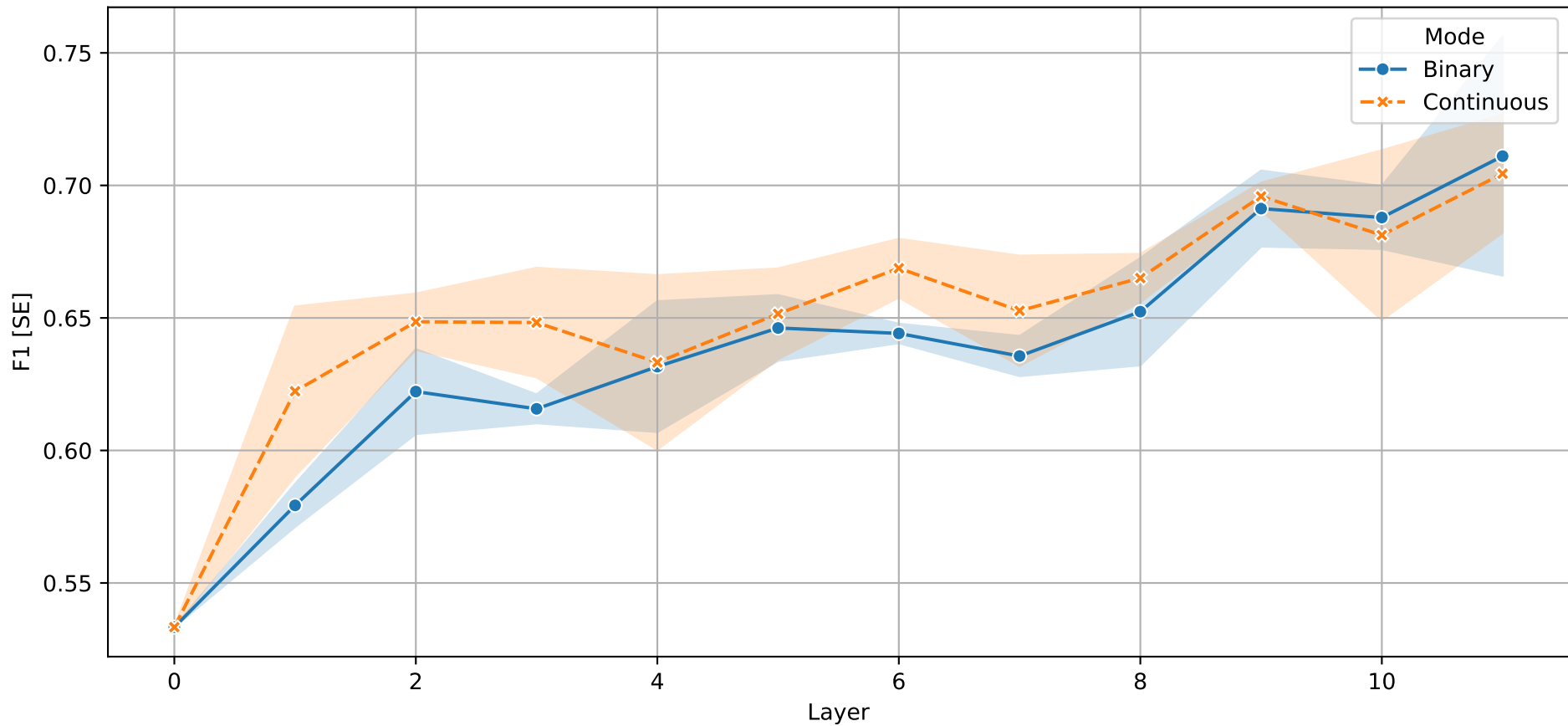
F1 per Layer - Single Neuron Probing



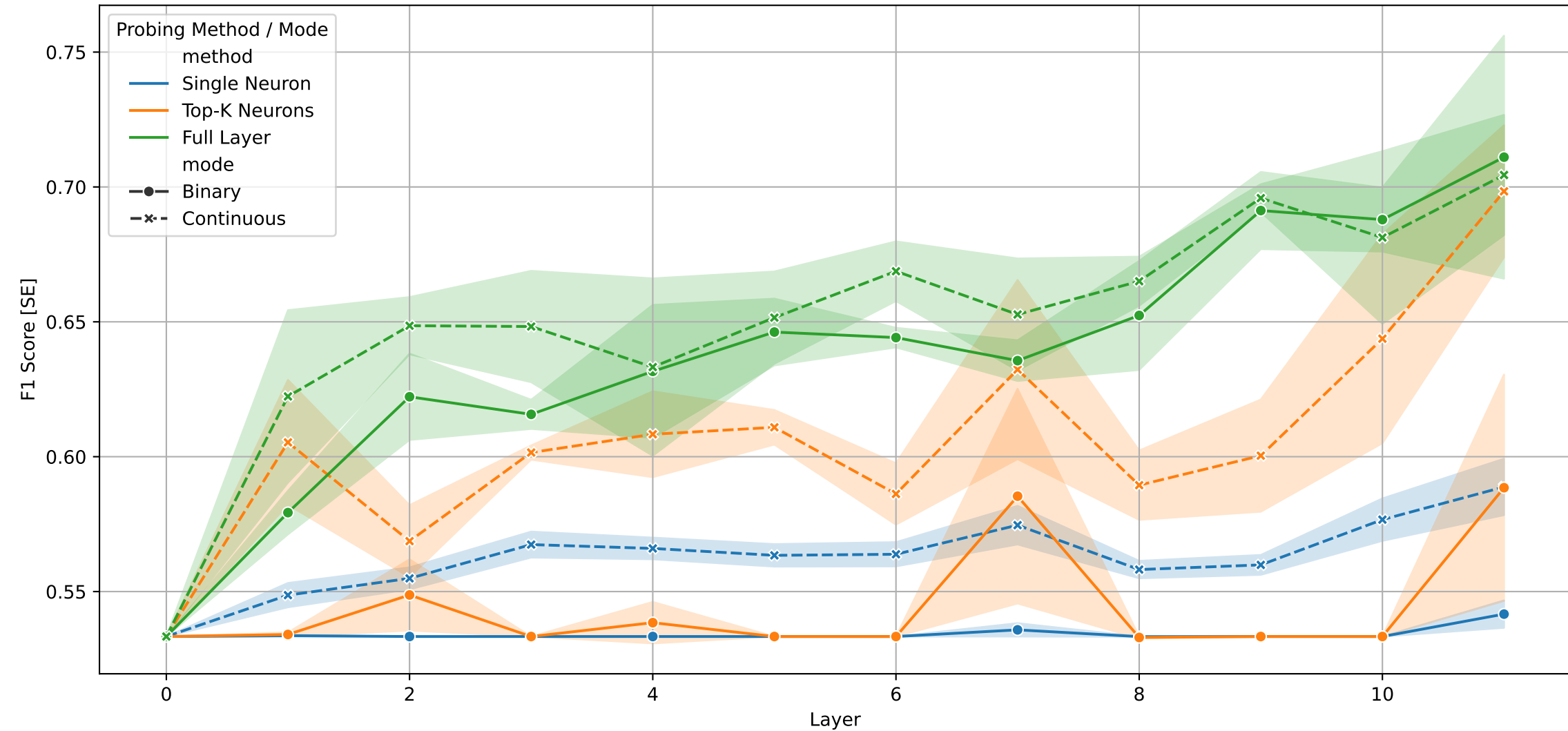
F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



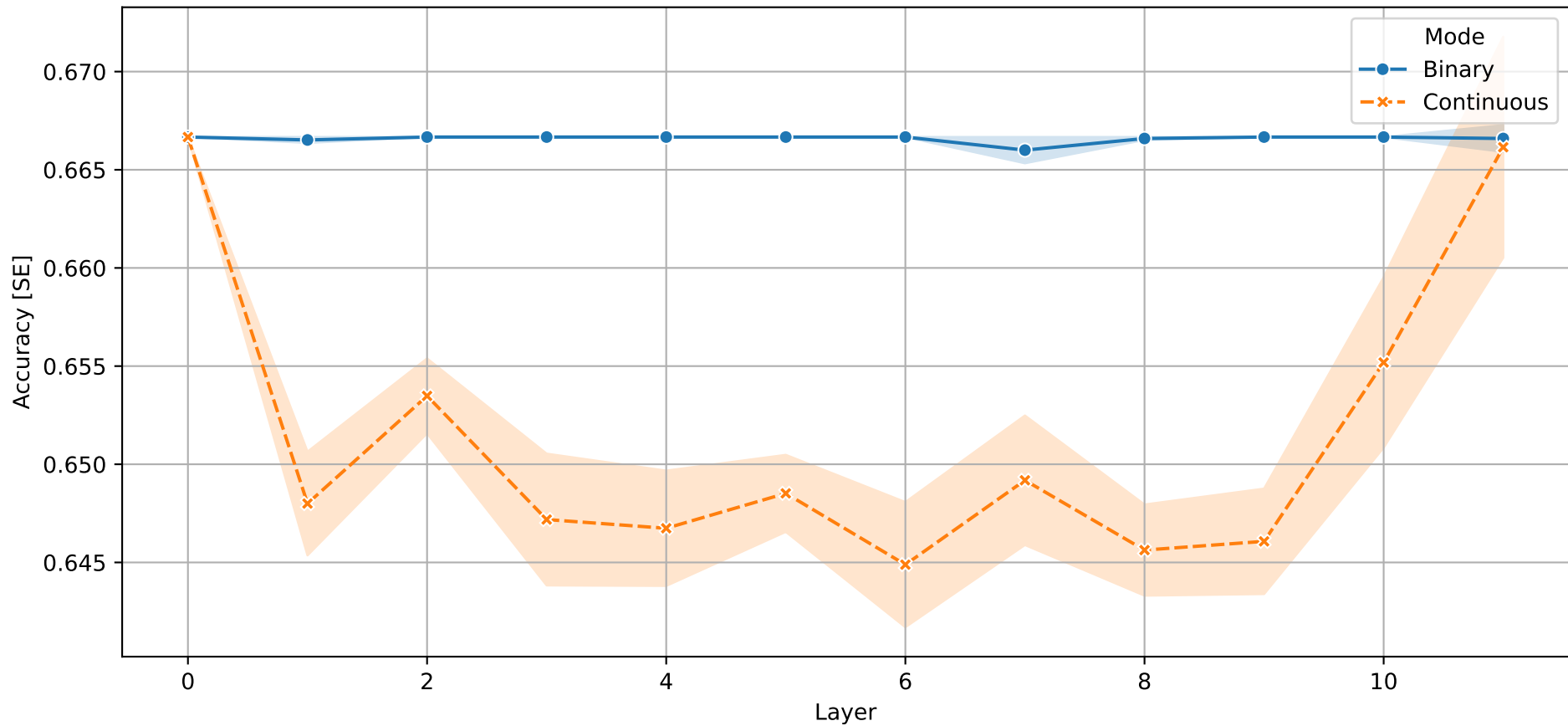
Overall F1 per Layer - All Methods



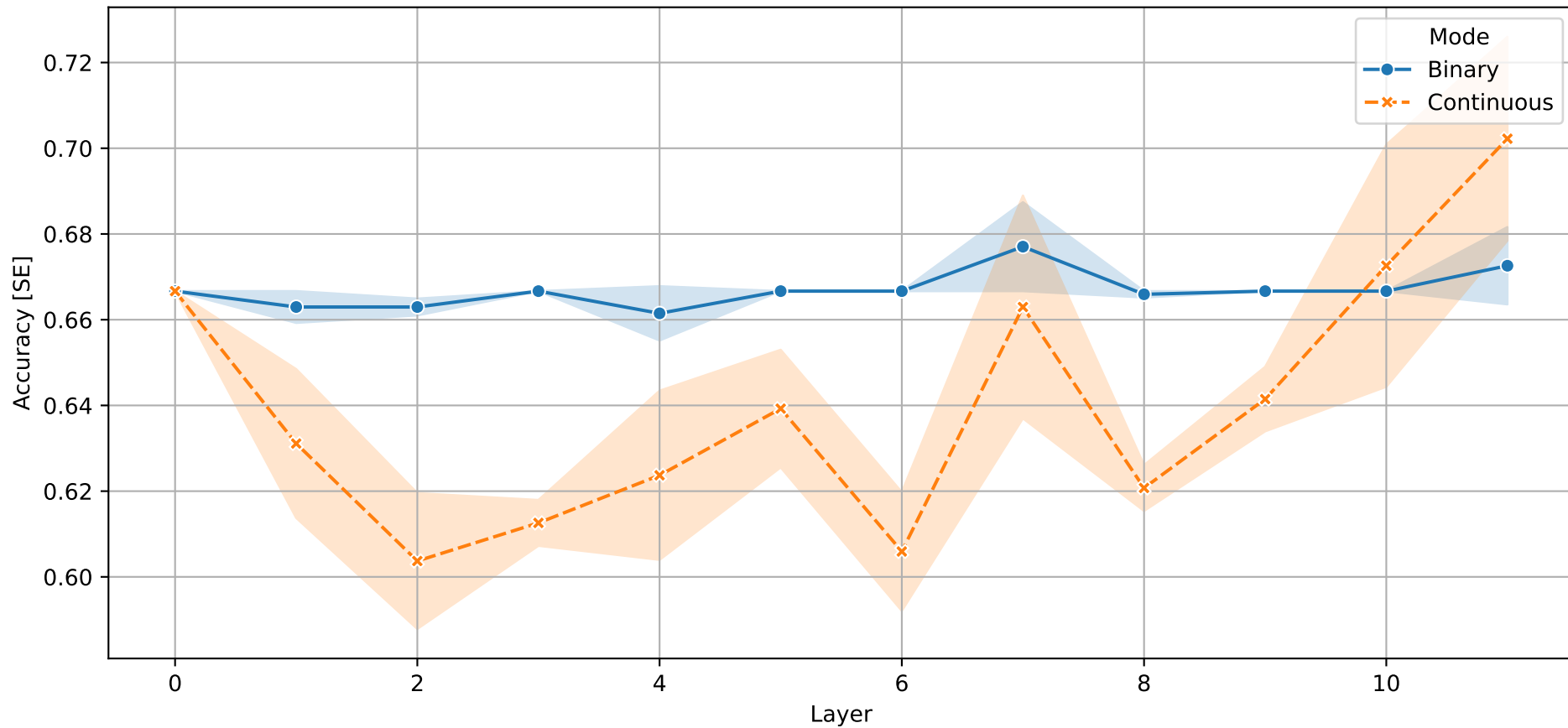
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	11.0	11.0
Full Layer	f1_max	0.7812	0.749
Full Layer	f1_mean	0.6376	0.6504
Full Layer	f1_std	0.0541	0.052
Single Neuron	f1_best_layer	11.0	11.0
Single Neuron	f1_max	0.6725	0.749
Single Neuron	f1_mean	0.5343	0.563
Single Neuron	f1_std	0.0091	0.0325
Top-K Neurons	f1_best_layer	11.0	11.0
Top-K Neurons	f1_max	0.6708	0.7429
Top-K Neurons	f1_mean	0.544	0.6065
Top-K Neurons	f1_std	0.0318	0.049

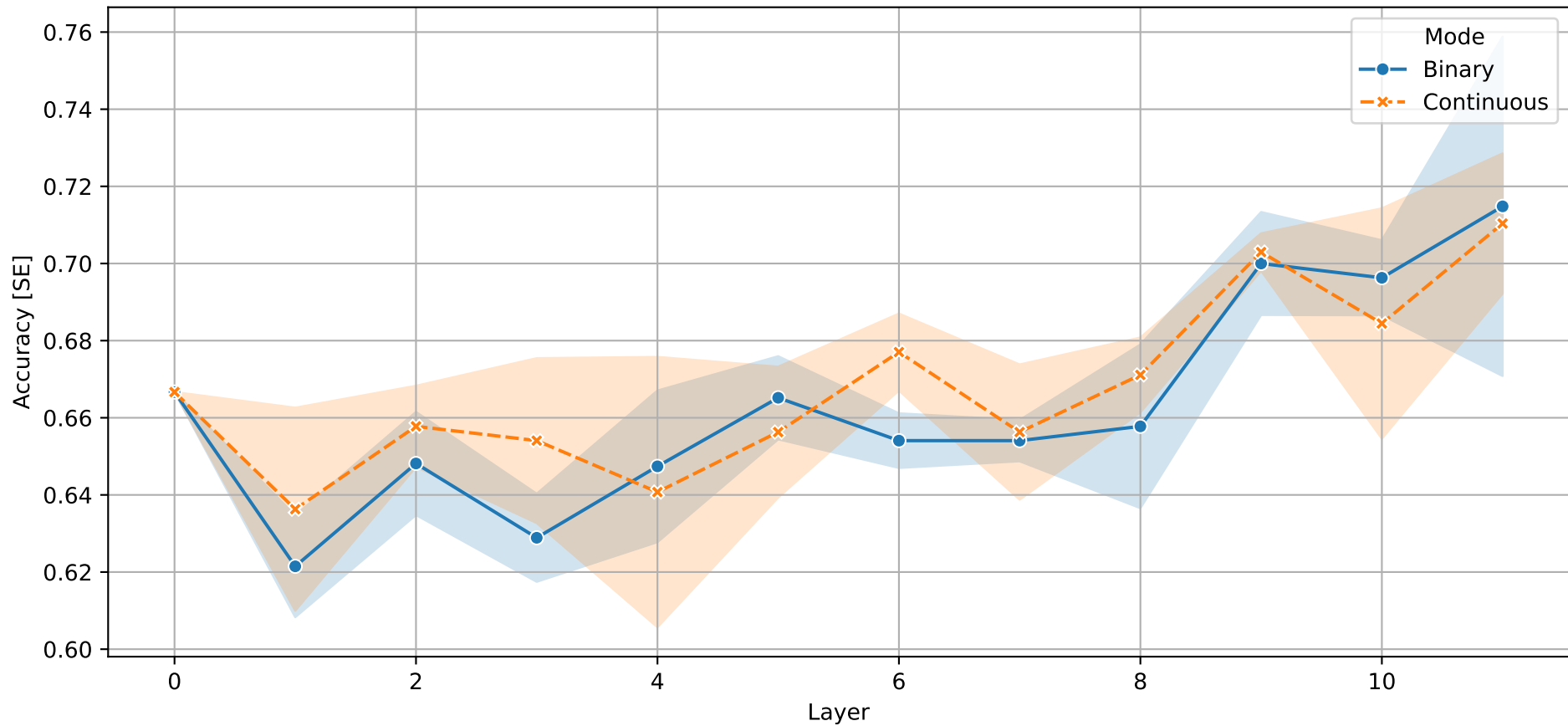
Accuracy per Layer - Single Neuron Probing



Accuracy per Layer - Top-K Neurons Probing

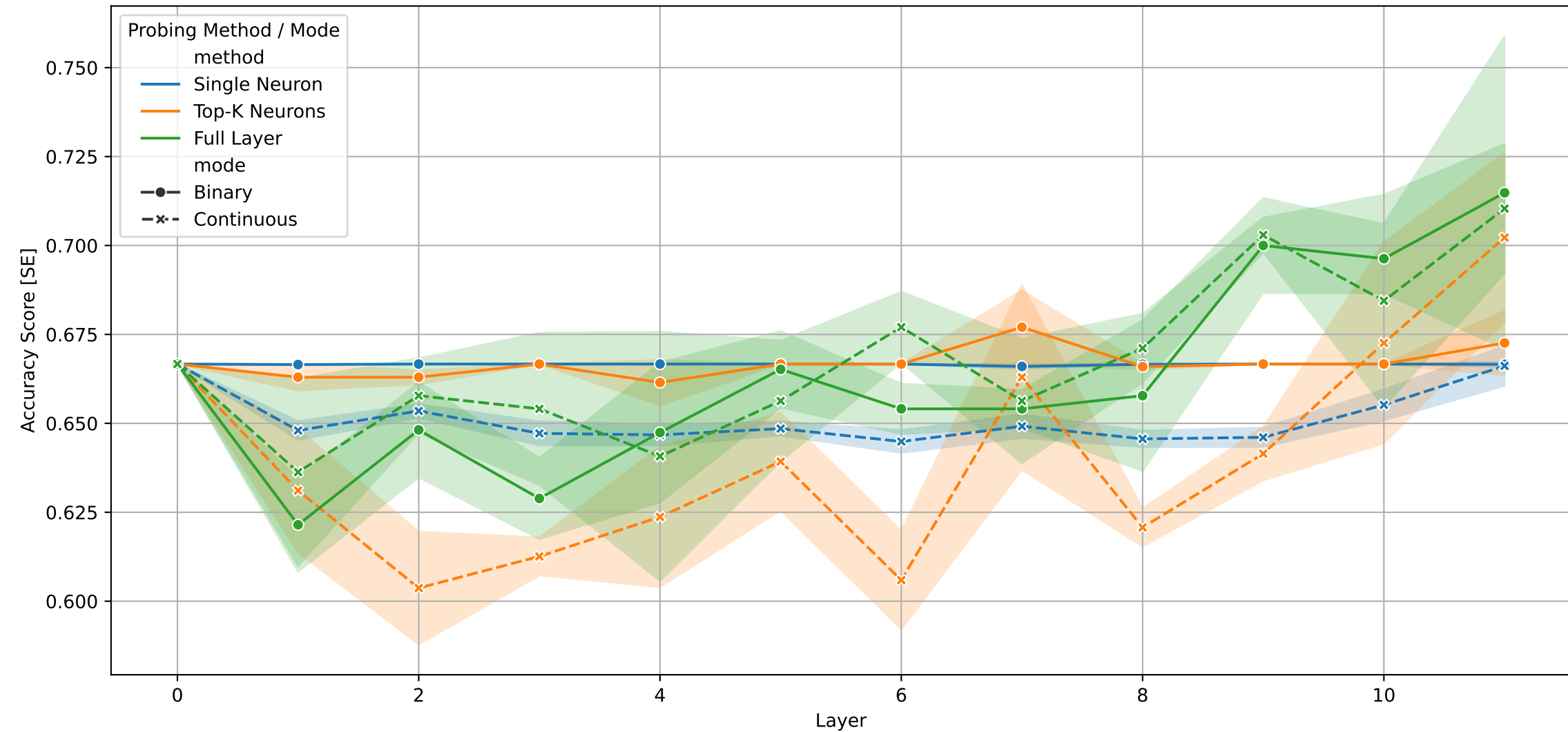


Accuracy per Layer - Full Layer Probing





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	0.78	0.7467
Full Layer	accuracy_mean	0.6629	0.6678
Full Layer	accuracy_std	0.0372	0.0355
Single Neuron	accuracy_best_layer	0.0	0.0
Single Neuron	accuracy_max	0.6822	0.7556
Single Neuron	accuracy_mean	0.6666	0.6515
Single Neuron	accuracy_std	0.0015	0.0184
Top-K Neurons	accuracy_best_layer	7.0	11.0
Top-K Neurons	accuracy_max	0.6978	0.7444
Top-K Neurons	accuracy_mean	0.6669	0.6402
Top-K Neurons	accuracy_std	0.0077	0.0382