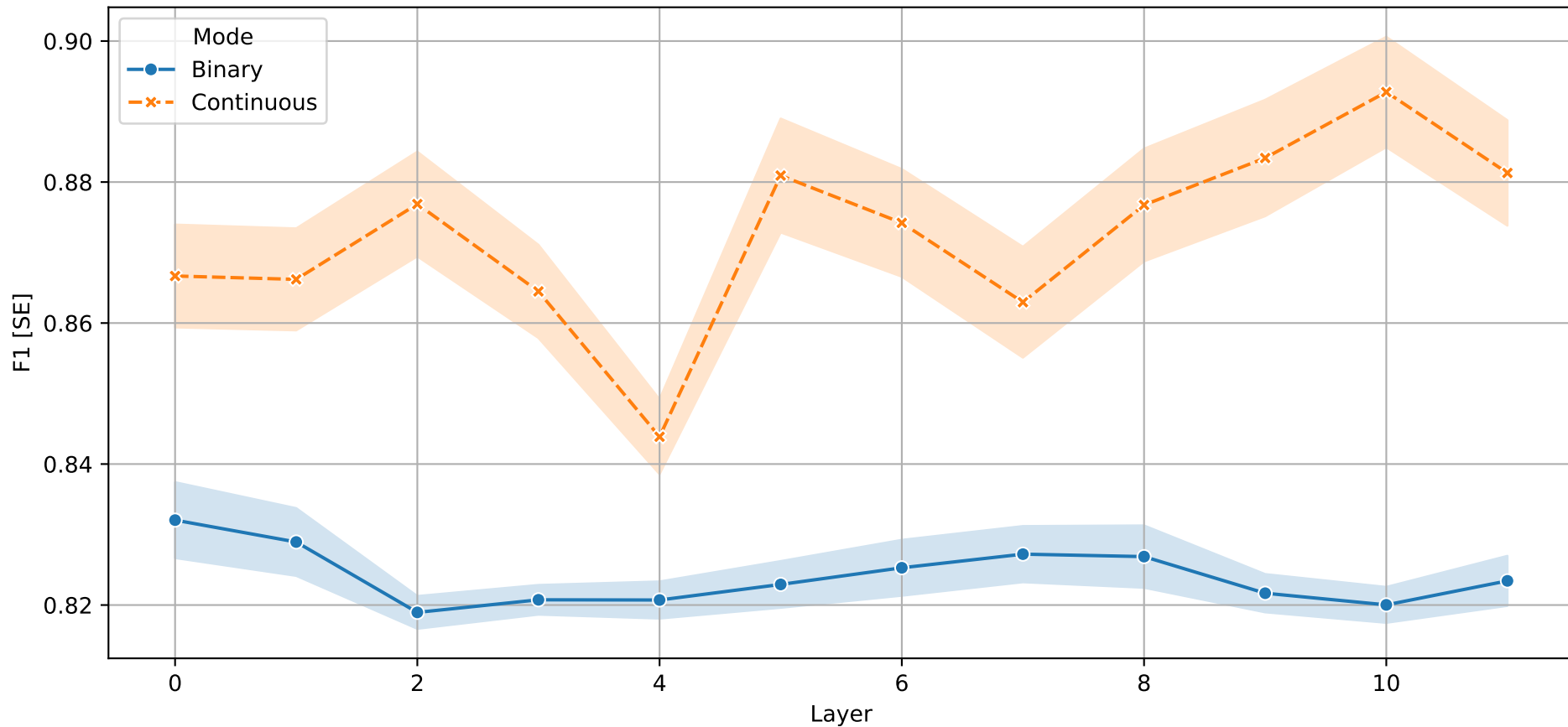
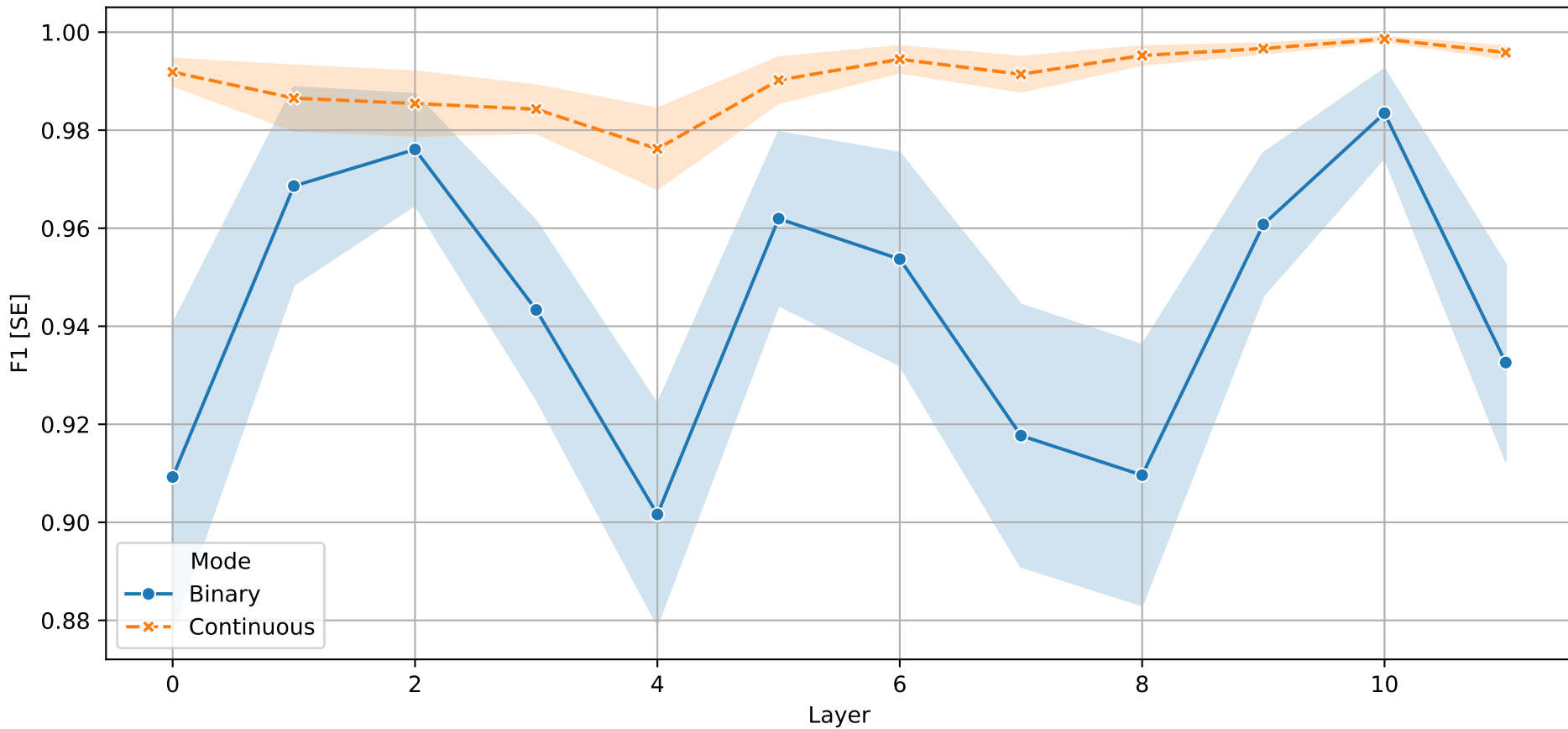


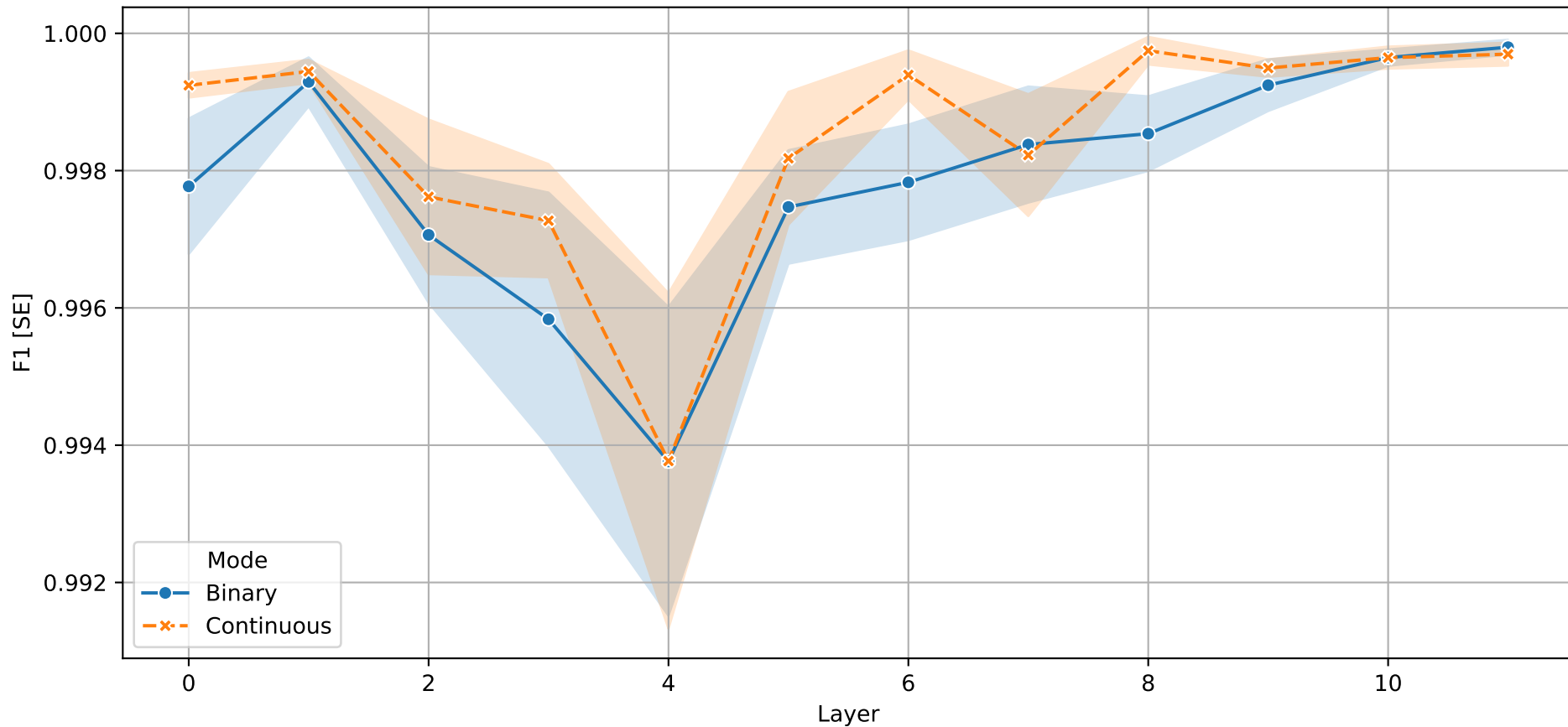
F1 per Layer - Single Neuron Probing



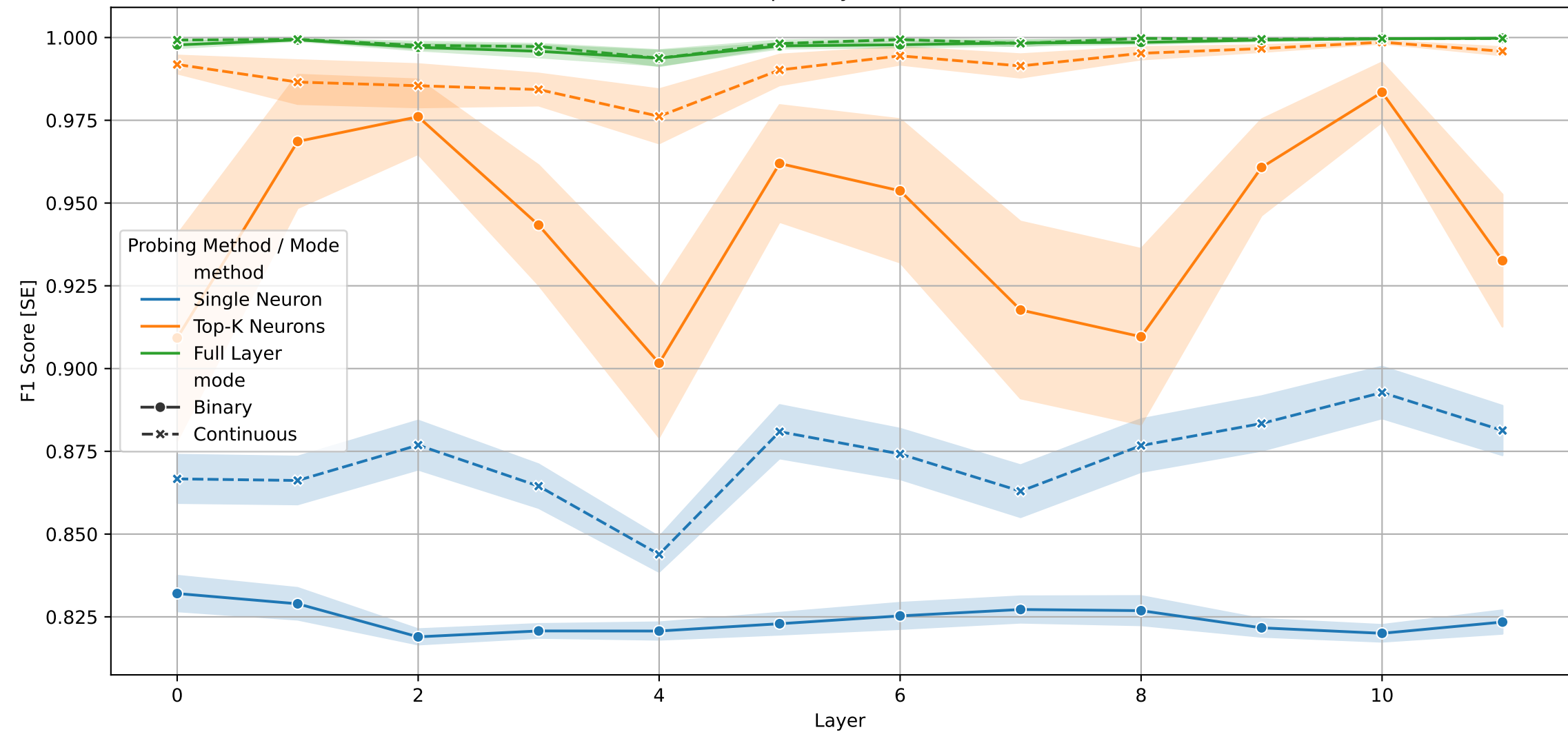
# F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



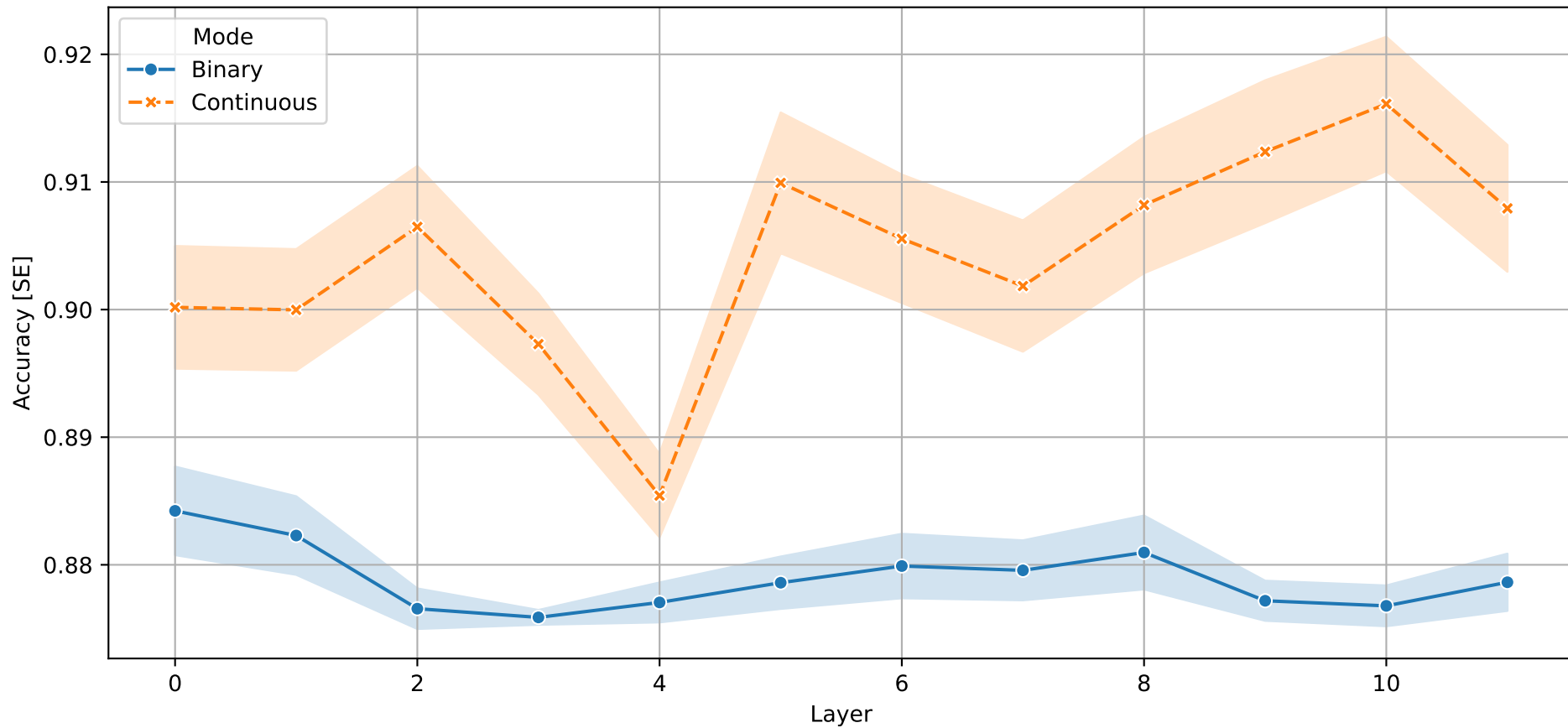
Overall F1 per Layer - All Methods



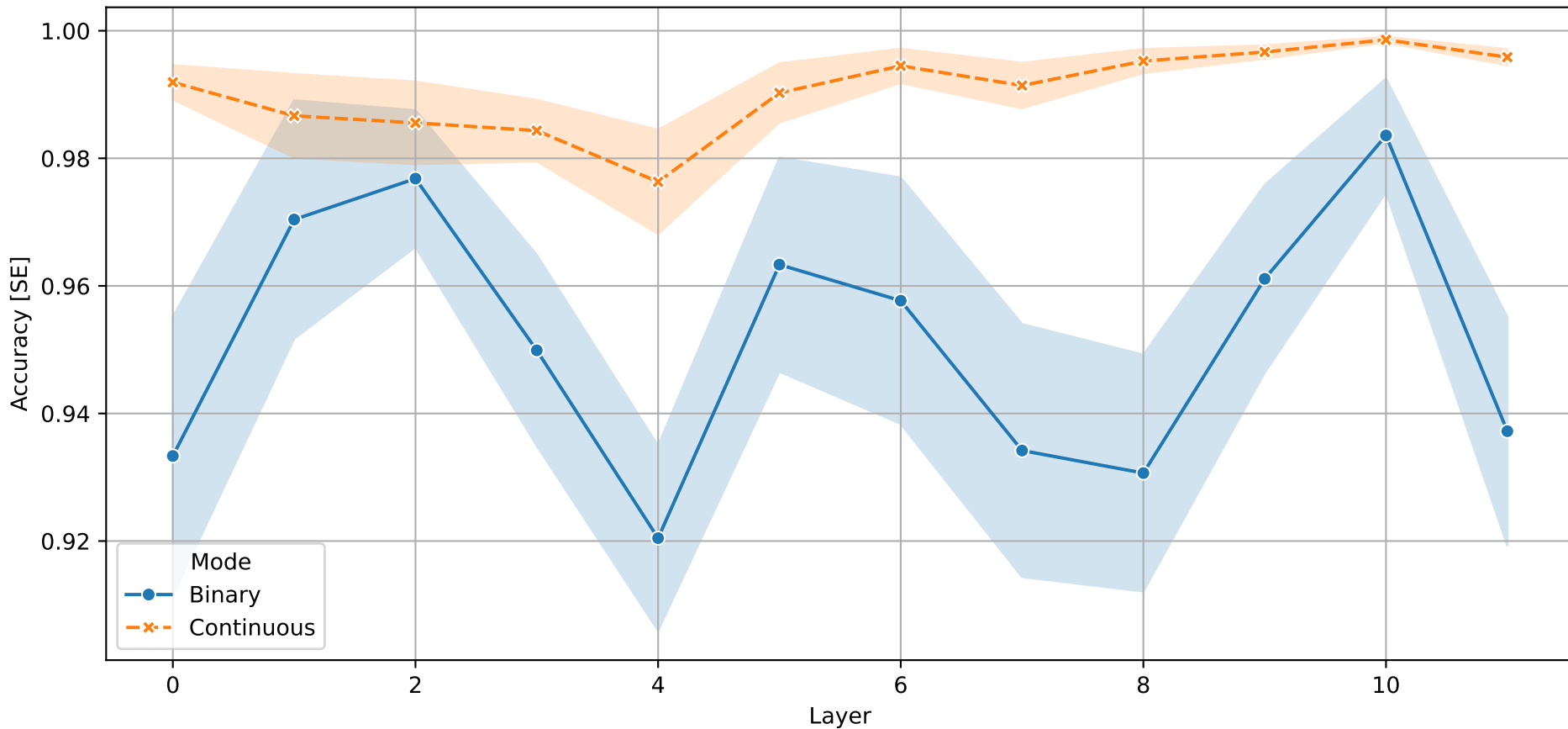
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	11.0	8.0
Full Layer	f1_max	1.0	1.0
Full Layer	f1_mean	0.9979	0.9985
Full Layer	f1_std	0.0032	0.0029
Single Neuron	f1_best_layer	0.0	10.0
Single Neuron	f1_max	1.0	1.0
Single Neuron	f1_mean	0.8241	0.8725
Single Neuron	f1_std	0.0324	0.0673
Top-K Neurons	f1_best_layer	10.0	10.0
Top-K Neurons	f1_max	1.0	1.0
Top-K Neurons	f1_mean	0.9432	0.9906
Top-K Neurons	f1_std	0.0619	0.0132

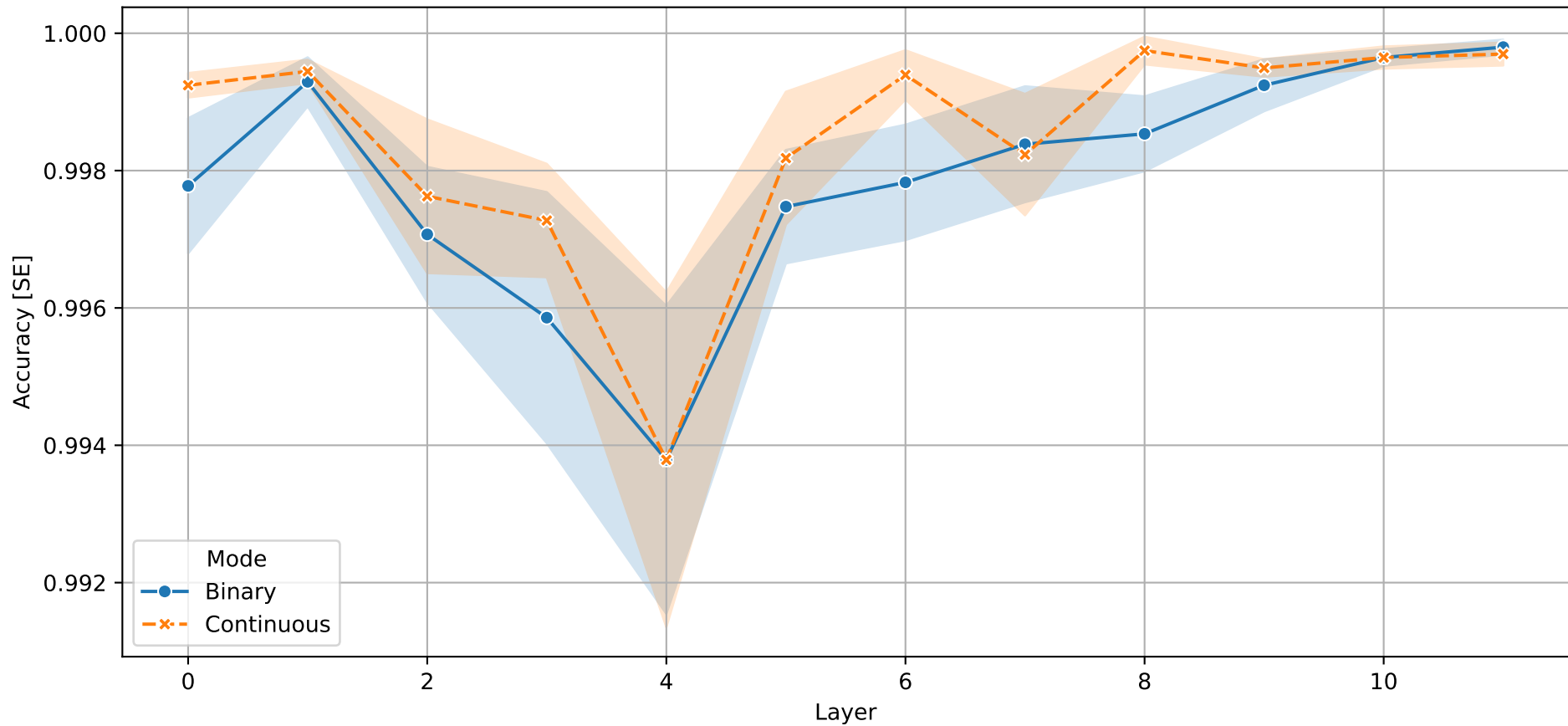
Accuracy per Layer - Single Neuron Probing



Accuracy per Layer - Top-K Neurons Probing

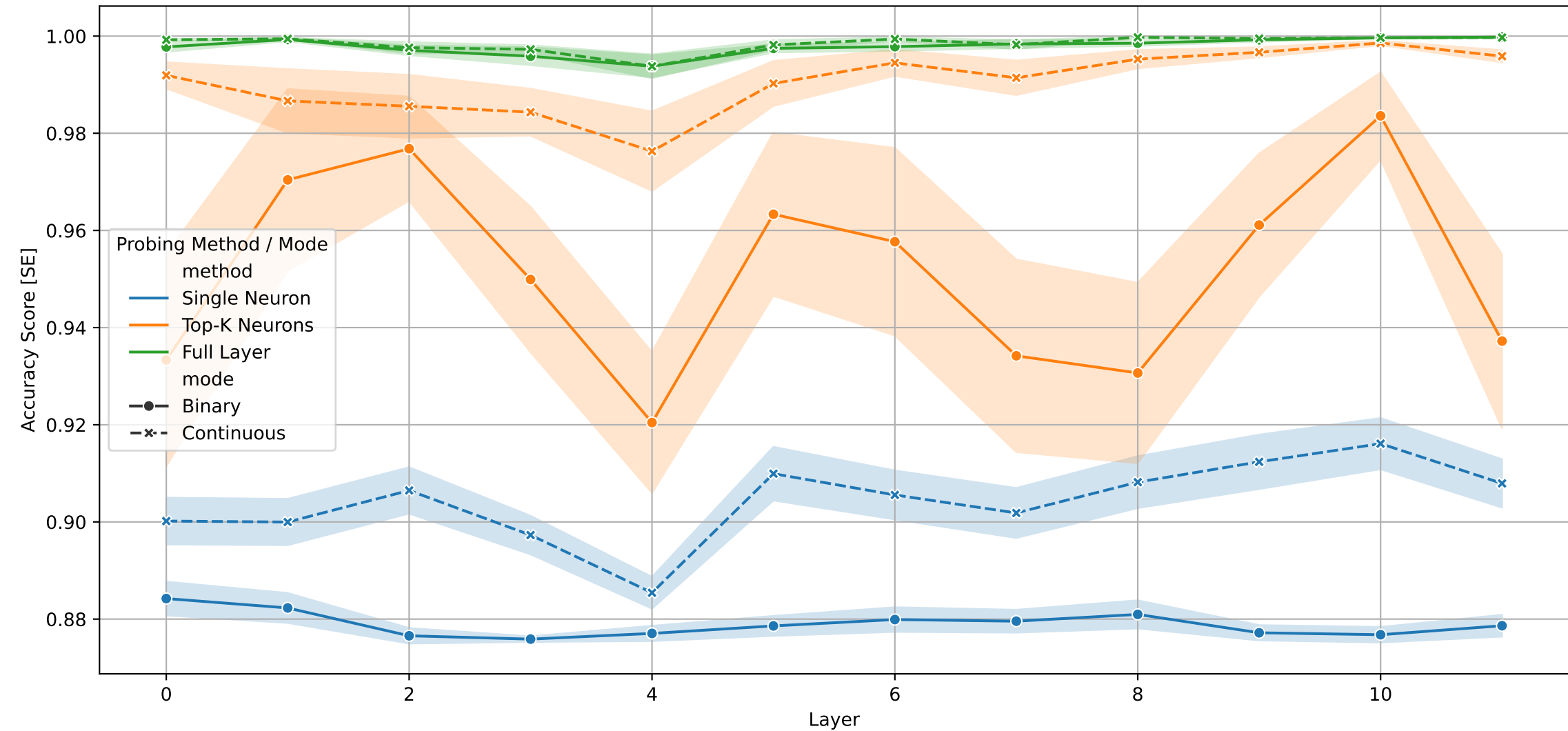


Accuracy per Layer - Full Layer Probing





Overall Accuracy per Layer - All Methods



Accuracy Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	accuracy_best_layer	11.0	8.0
Full Layer	accuracy_max	1.0	1.0
Full Layer	accuracy_mean	0.9979	0.9985
Full Layer	accuracy_std	0.0032	0.0029
Single Neuron	accuracy_best_layer	0.0	10.0
Single Neuron	accuracy_max	1.0	1.0
Single Neuron	accuracy_mean	0.879	0.9043
Single Neuron	accuracy_std	0.0202	0.0443
Top-K Neurons	accuracy_best_layer	10.0	10.0
Top-K Neurons	accuracy_max	1.0	1.0
Top-K Neurons	accuracy_mean	0.9516	0.9906
Top-K Neurons	accuracy_std	0.0487	0.0131