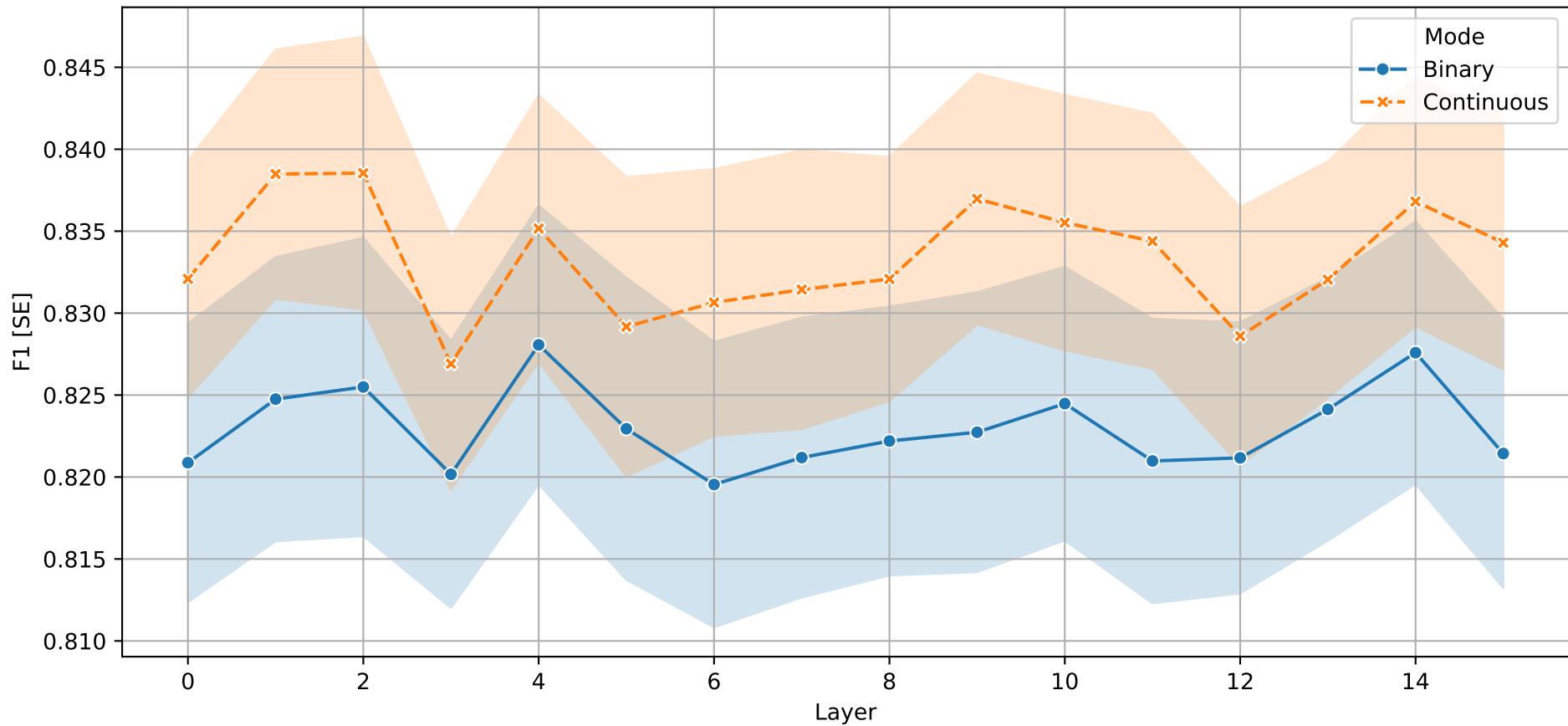
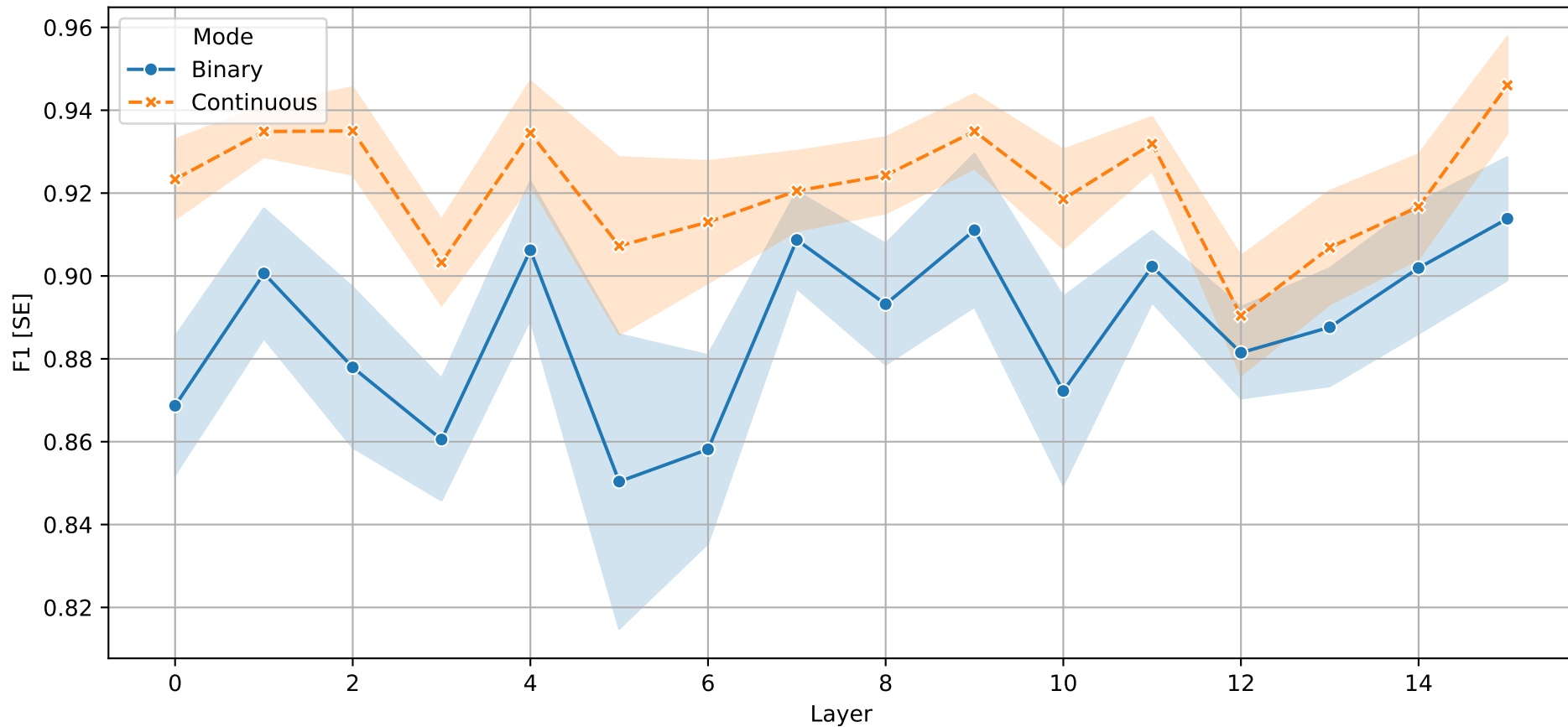


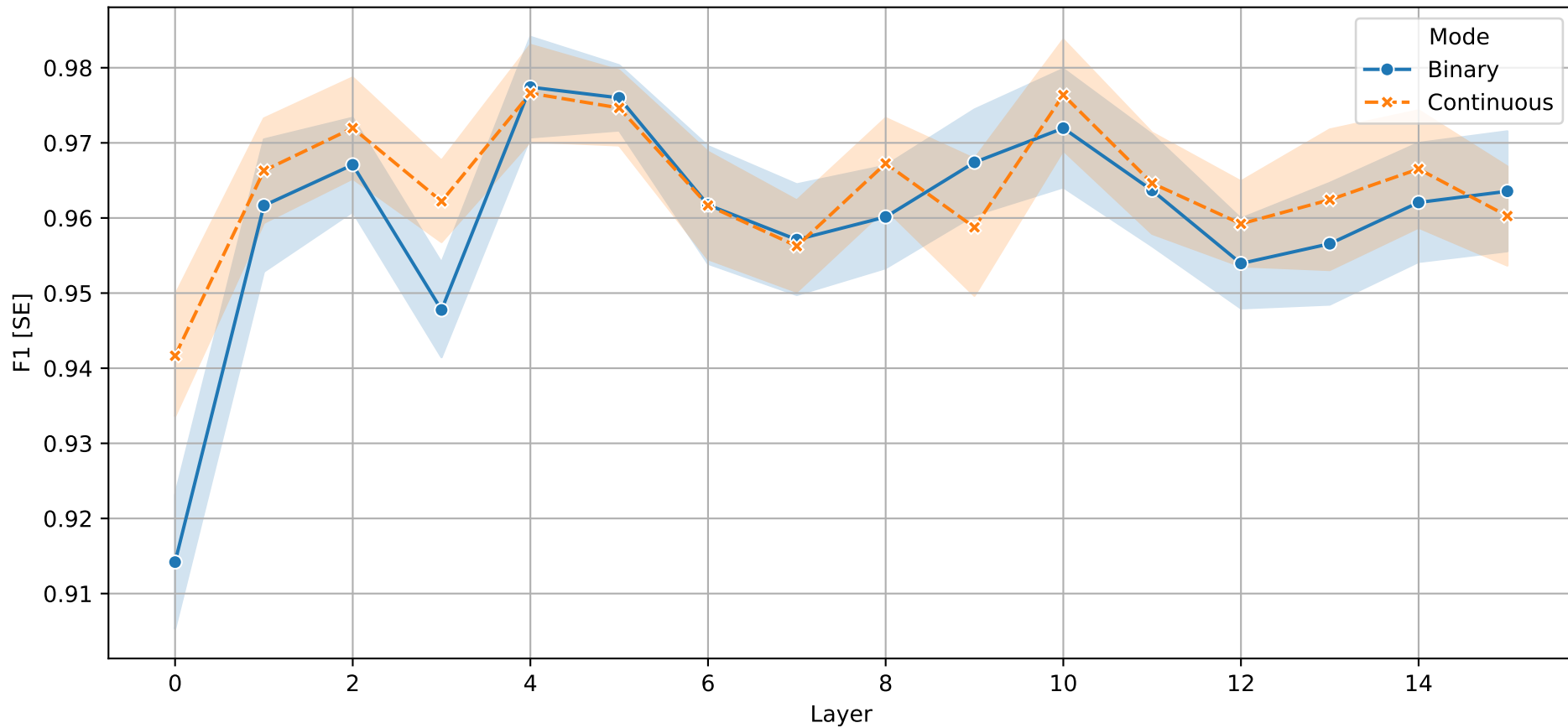
F1 per Layer - Single Neuron Probing



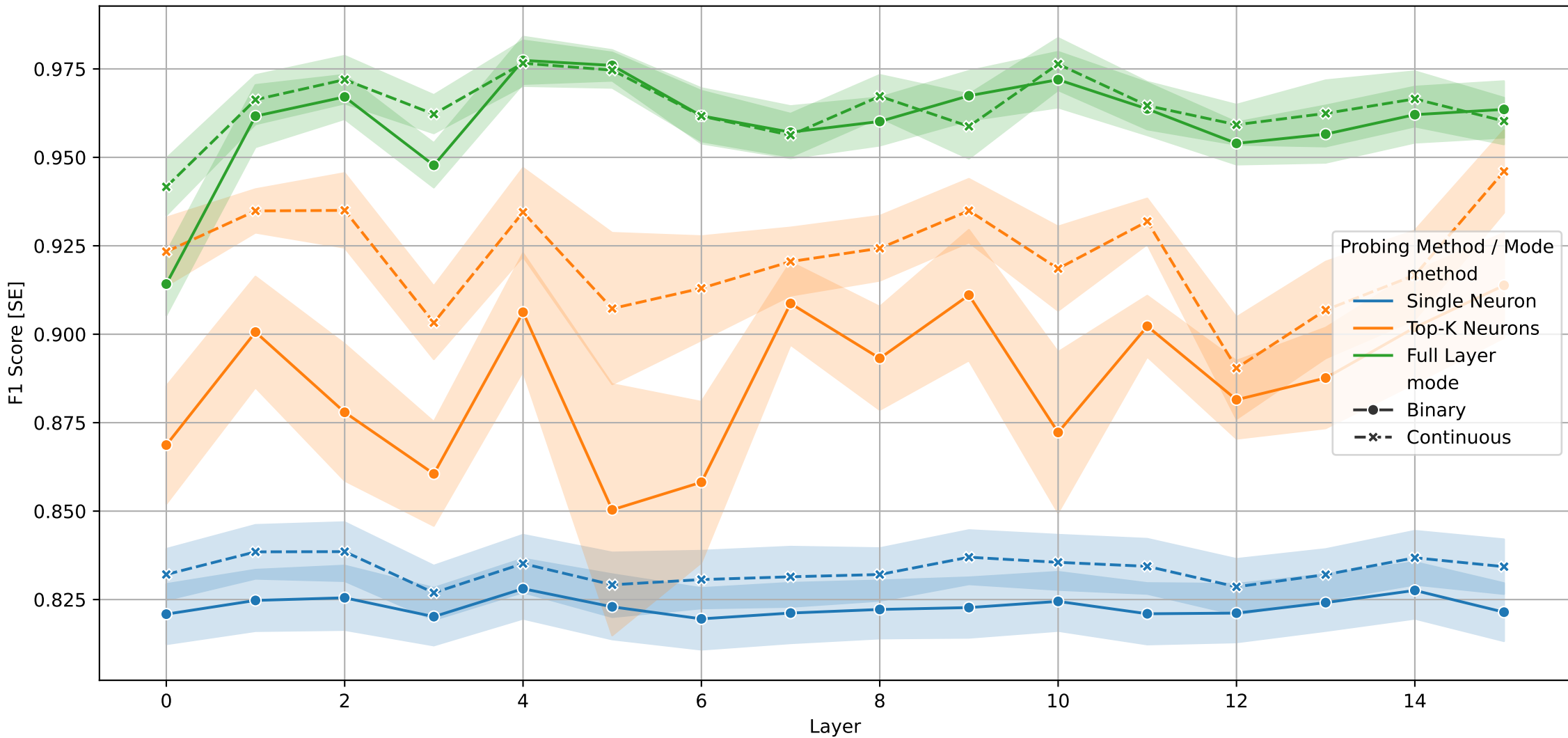
F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



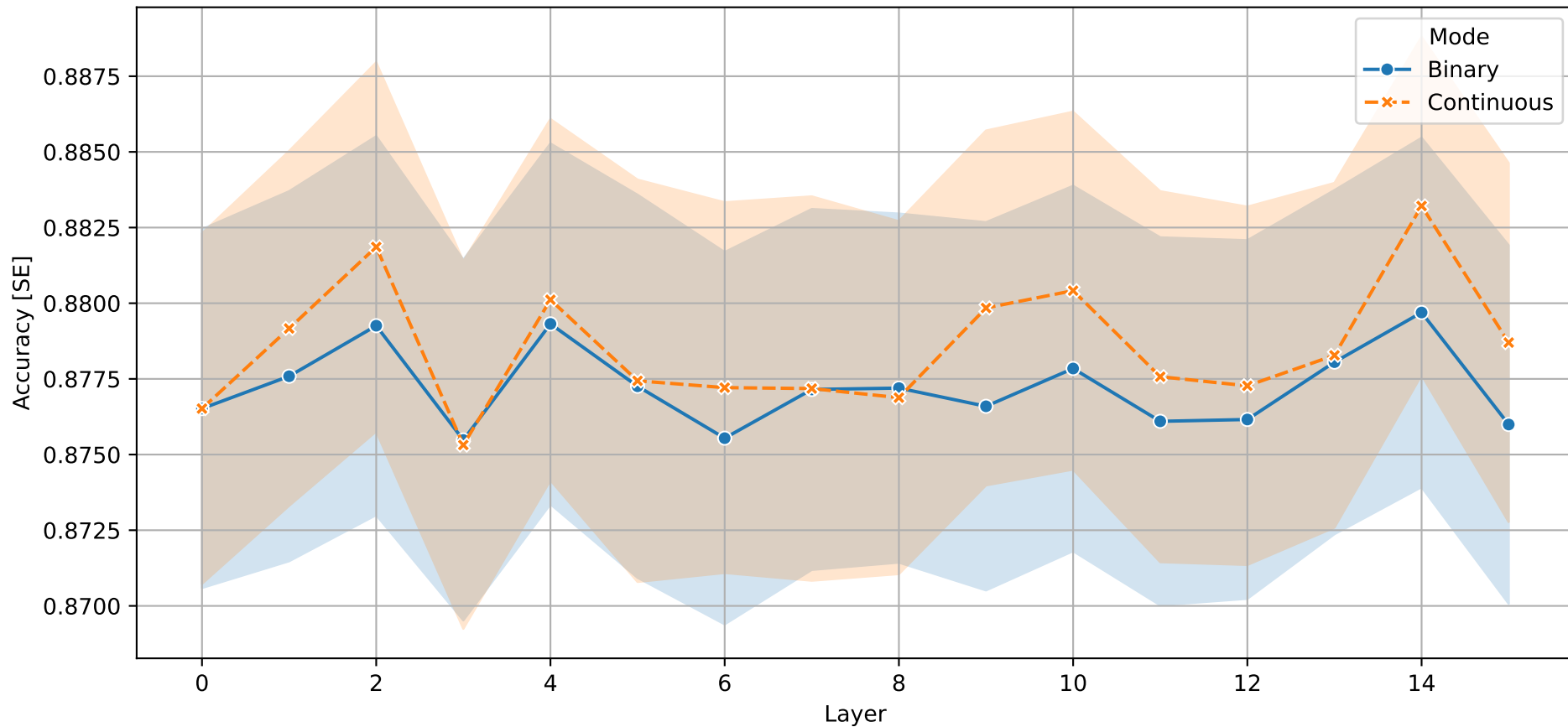
Overall F1 per Layer - All Methods



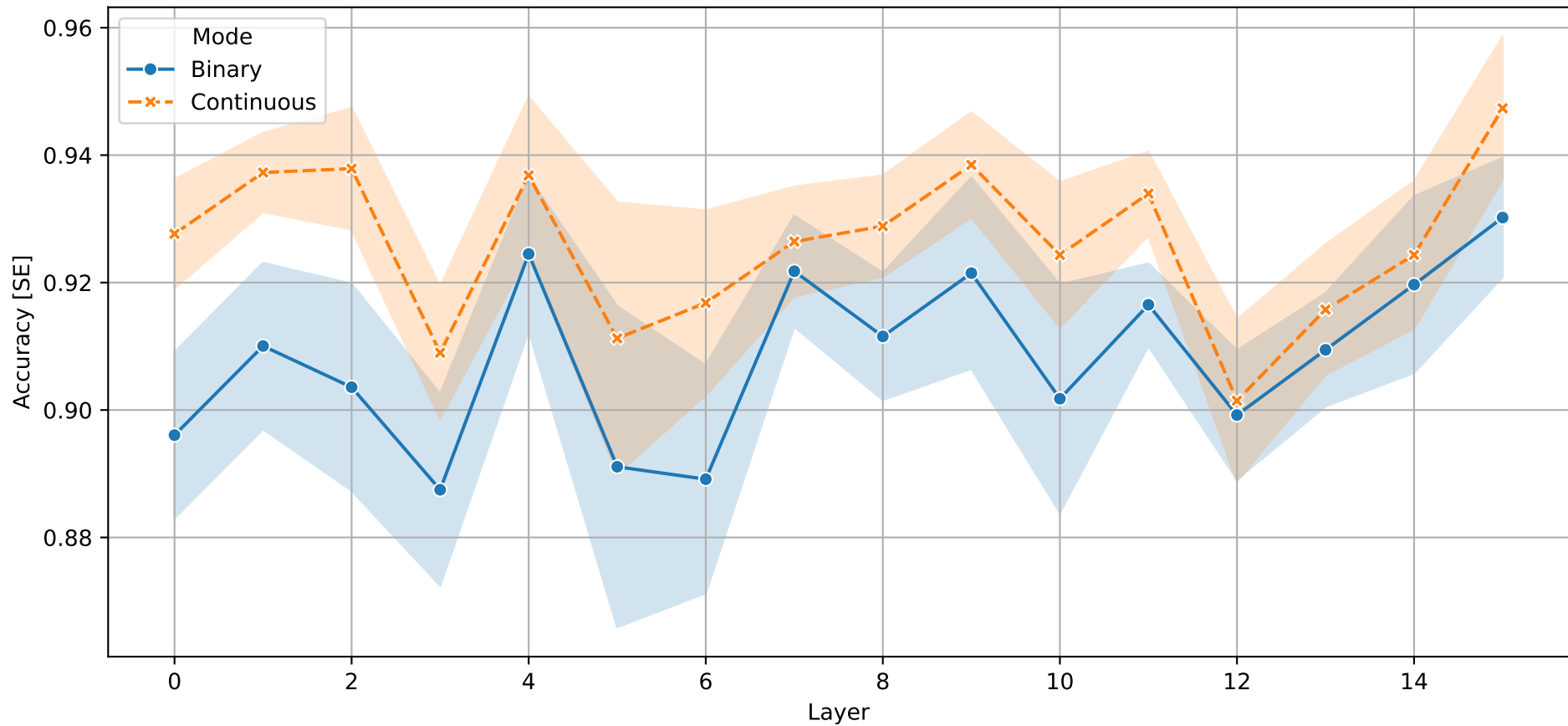
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	4.0	4.0
Full Layer	f1_max	0.9988	0.9988
Full Layer	f1_mean	0.9601	0.9642
Full Layer	f1_std	0.024	0.0204
Single Neuron	f1_best_layer	4.0	2.0
Single Neuron	f1_max	0.9976	0.9976
Single Neuron	f1_mean	0.823	0.8333
Single Neuron	f1_std	0.0754	0.0701
Top-K Neurons	f1_best_layer	15.0	15.0
Top-K Neurons	f1_max	0.9976	0.9988
Top-K Neurons	f1_mean	0.8872	0.9213
Top-K Neurons	f1_std	0.0521	0.035

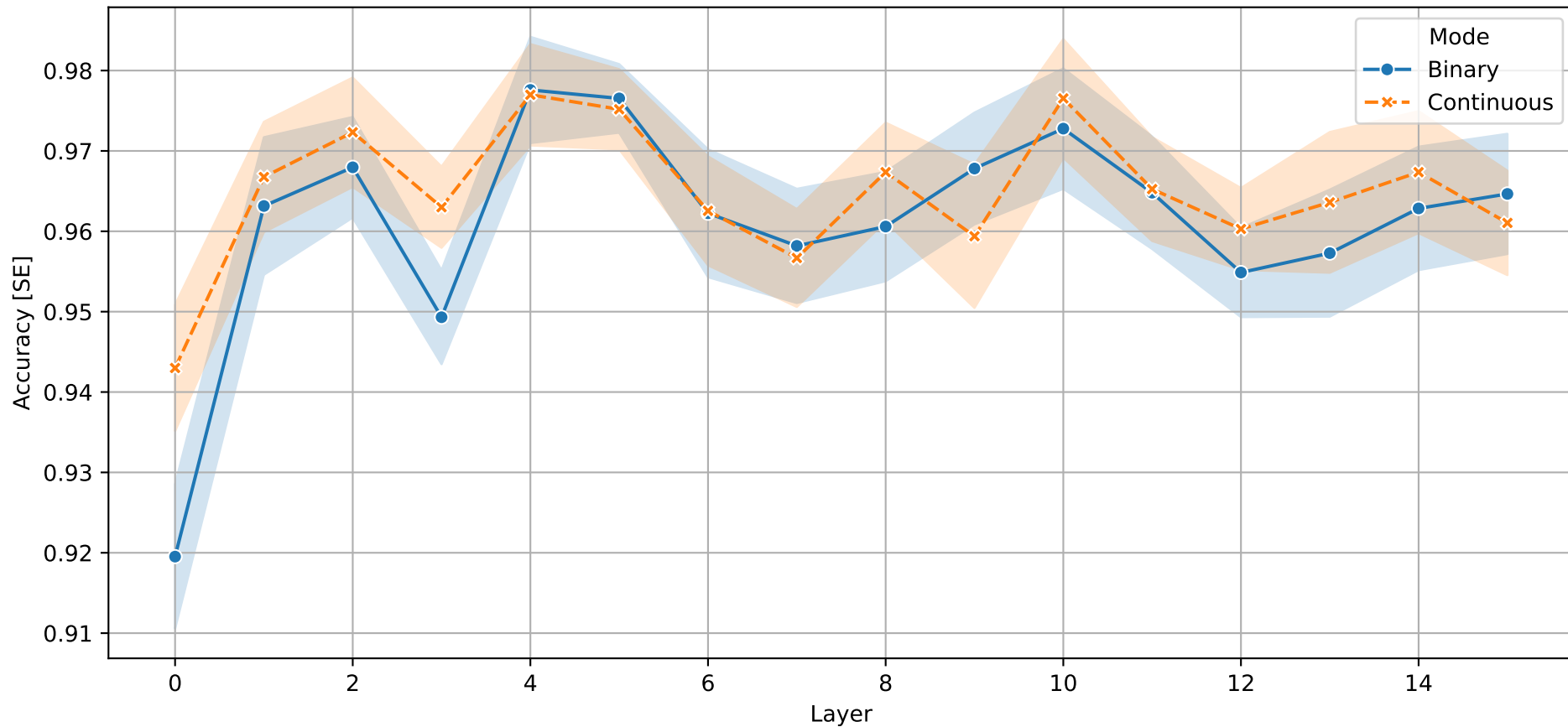
Accuracy per Layer – Single Neuron Probing



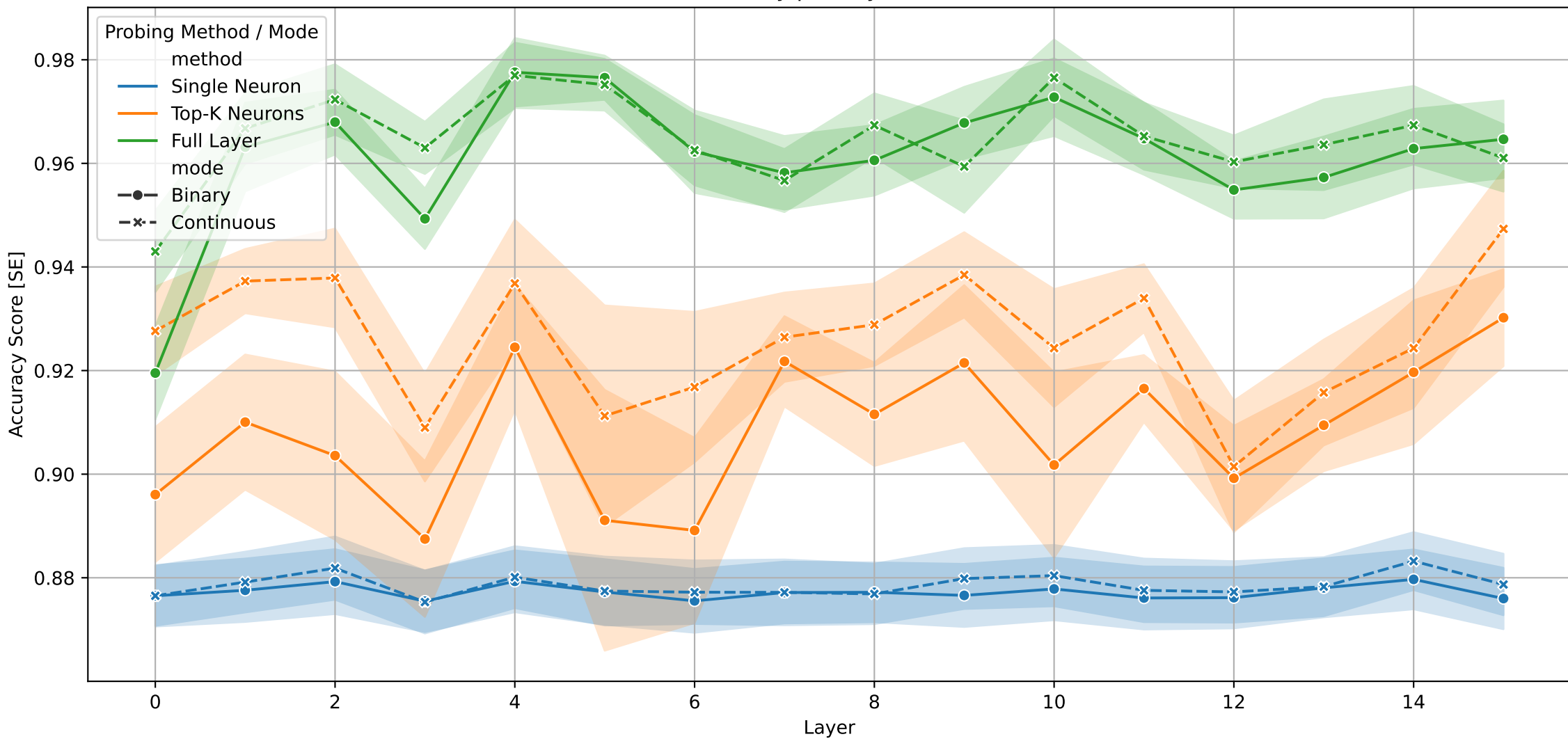
Accuracy per Layer - Top-K Neurons Probing



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	4.0	4.0
Full Layer	accuracy_max	0.9988	0.9988
Full Layer	accuracy_mean	0.9613	0.9648
Full Layer	accuracy_std	0.023	0.0199
Single Neuron	accuracy_best_layer	14.0	14.0
Single Neuron	accuracy_max	0.9976	0.9976
Single Neuron	accuracy_mean	0.8772	0.8786
Single Neuron	accuracy_std	0.0533	0.0533
Top-K Neurons	accuracy_best_layer	15.0	15.0
Top-K Neurons	accuracy_max	0.9976	0.9988
Top-K Neurons	accuracy_mean	0.9083	0.9261
Top-K Neurons	accuracy_std	0.0395	0.0323