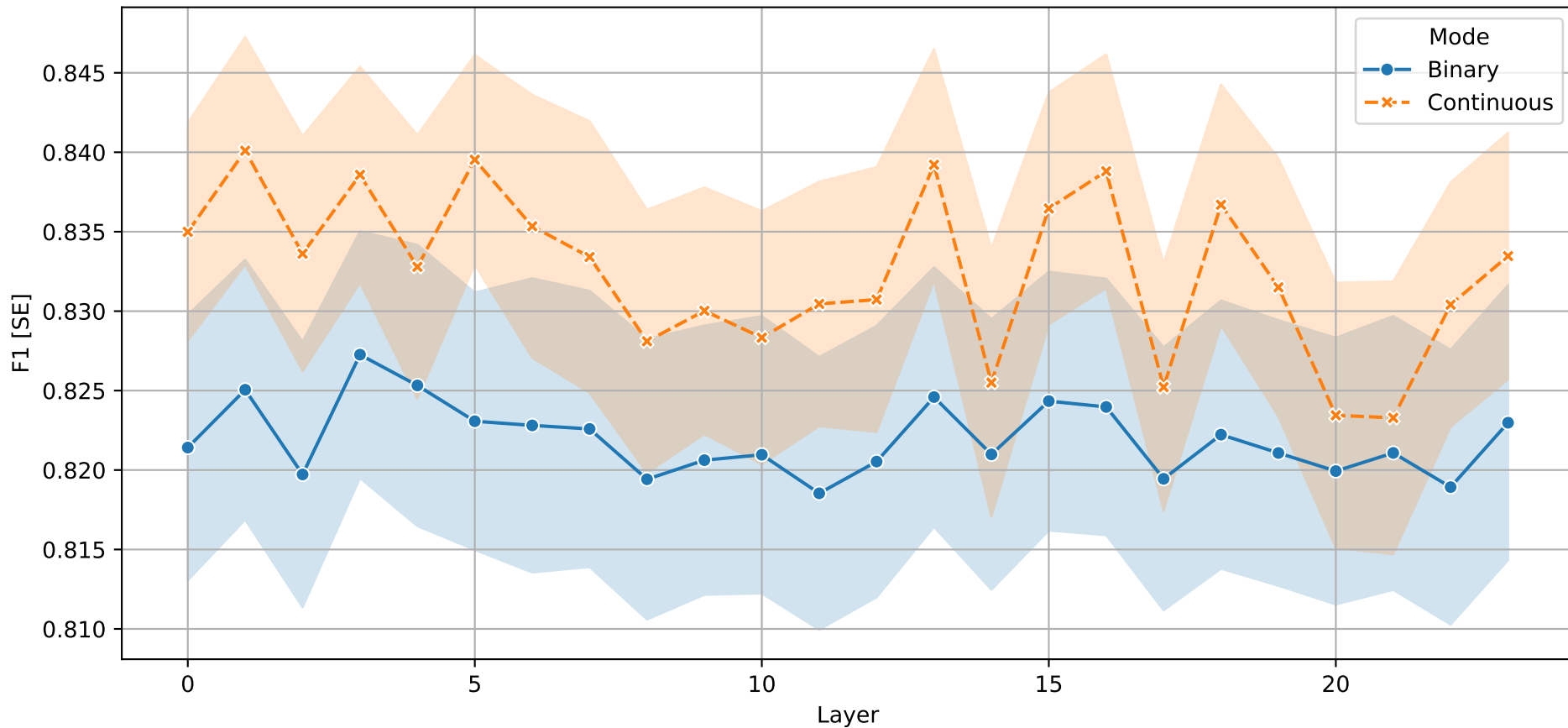
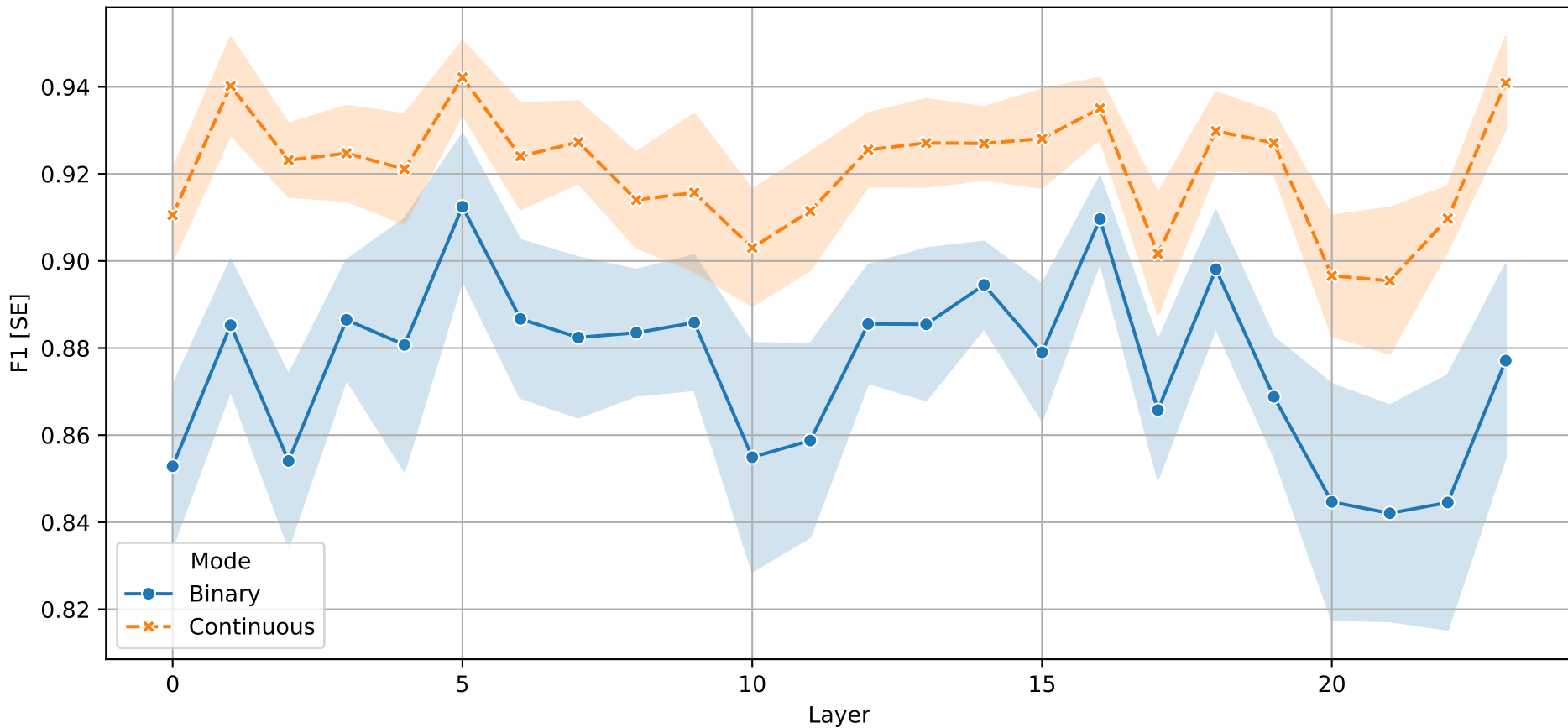


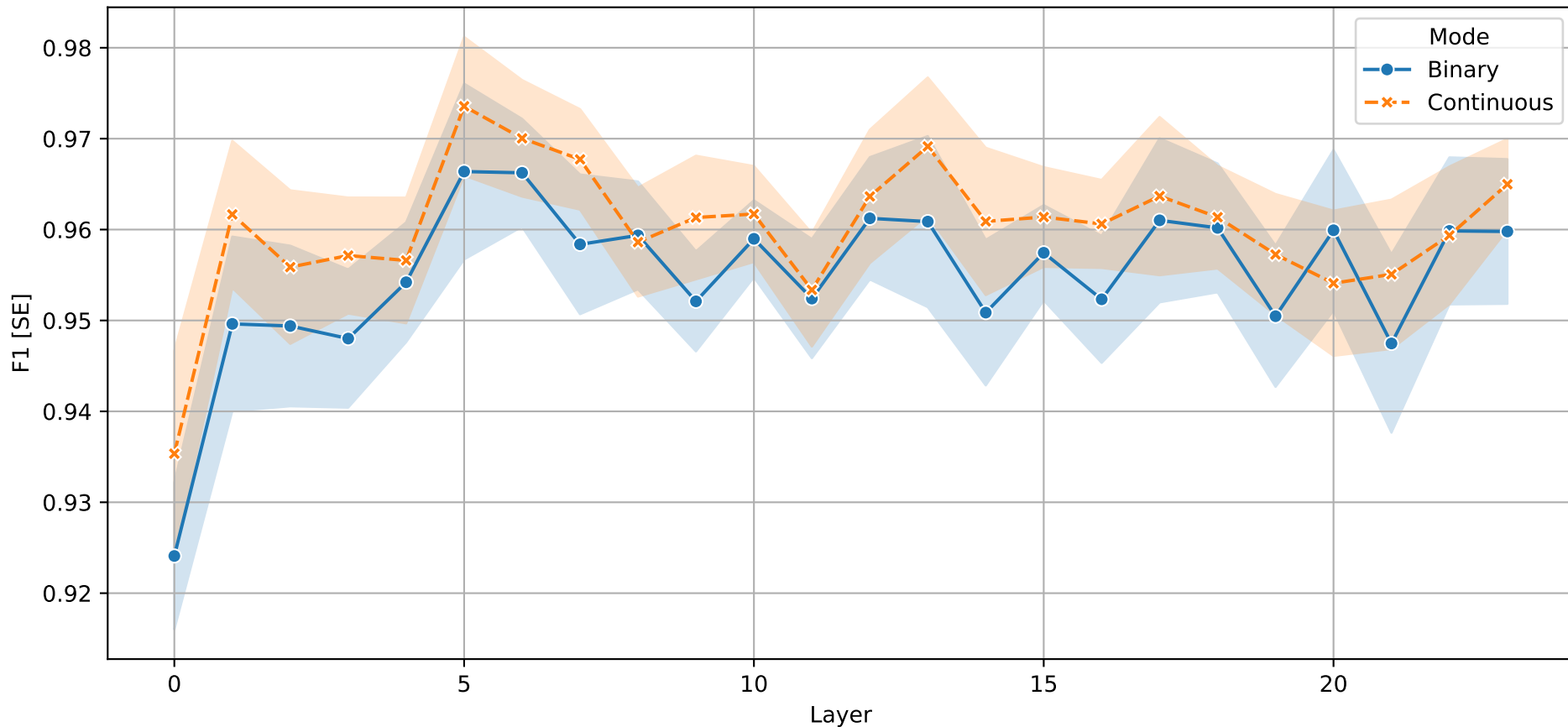
F1 per Layer - Single Neuron Probing



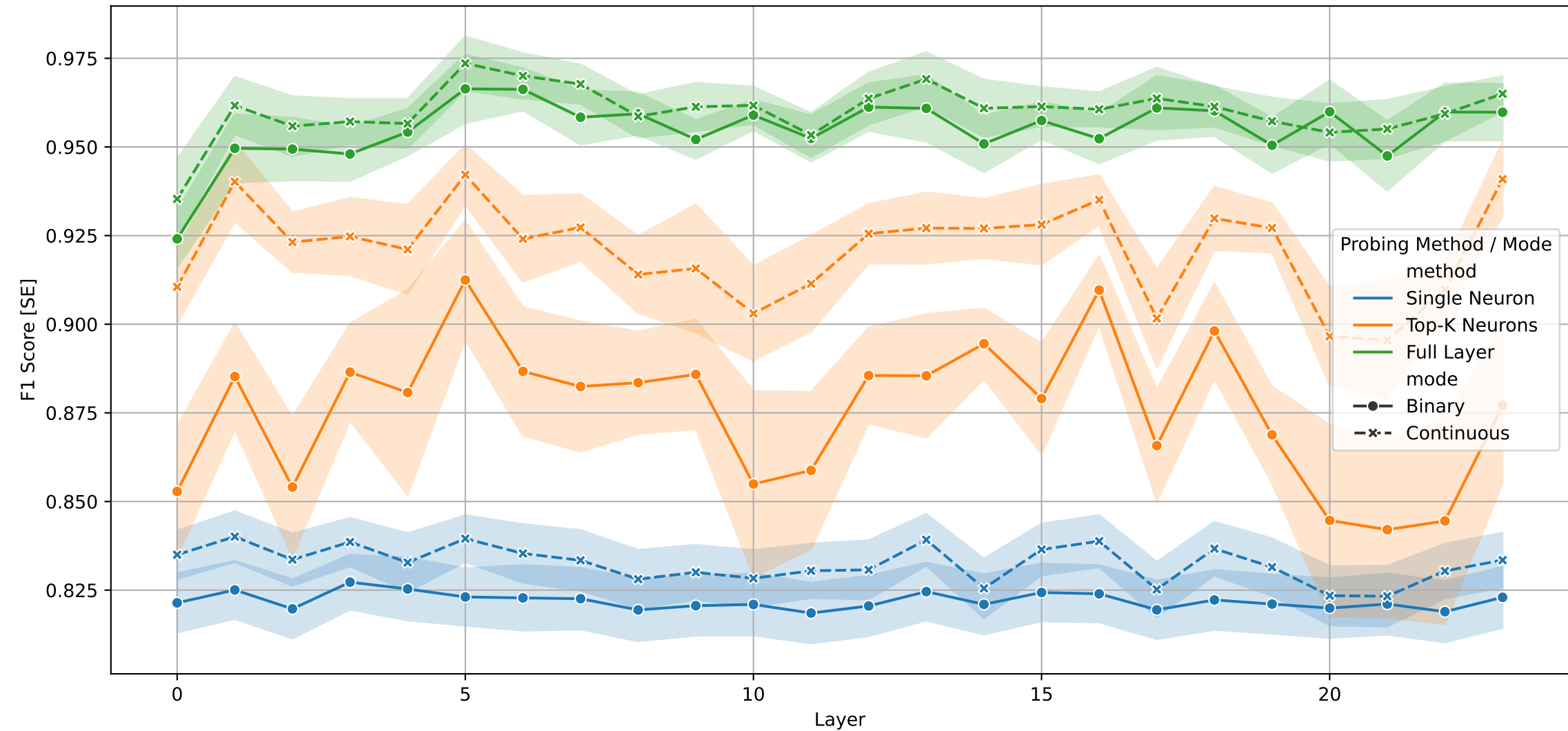
F1 per Layer - Top-K Neurons Probing



F1 per Layer - Full Layer Probing



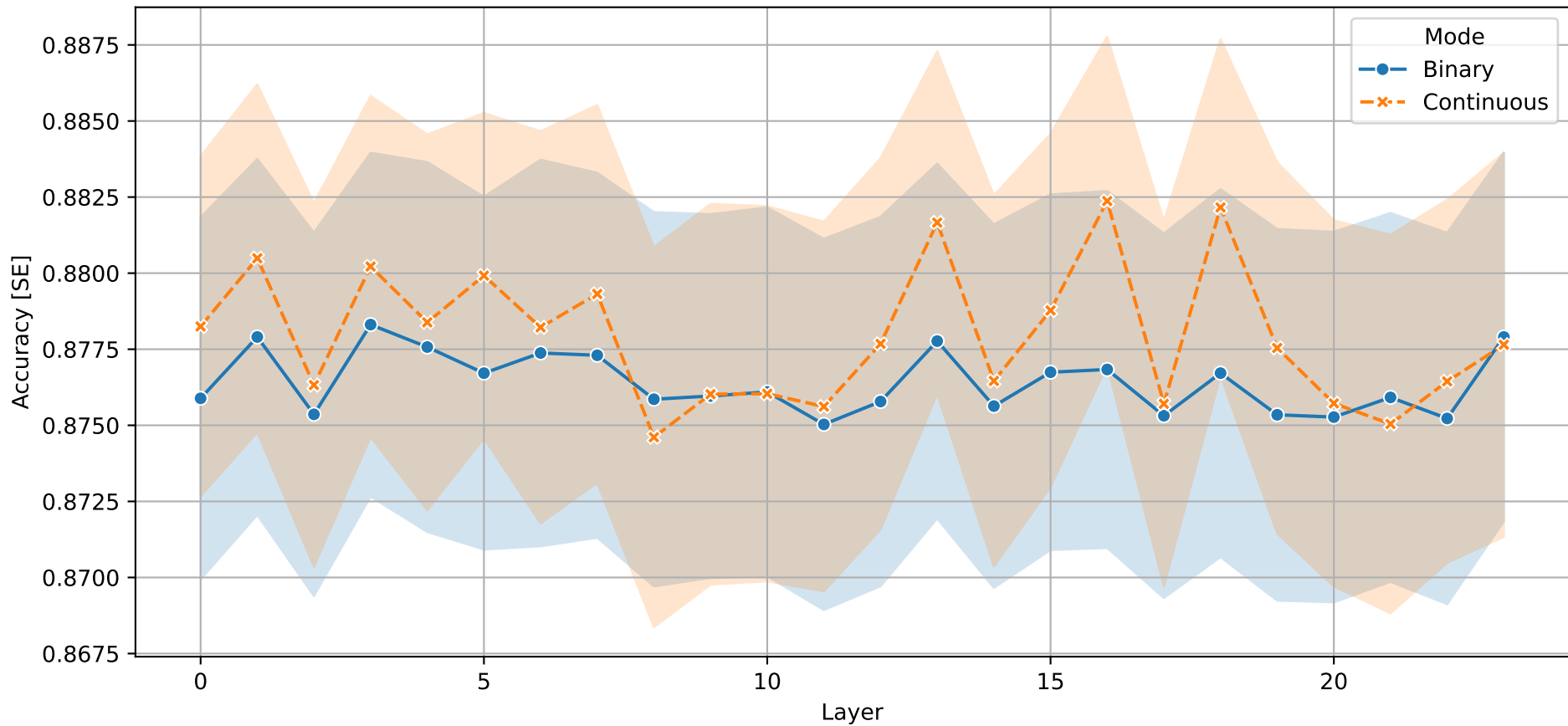
Overall F1 per Layer - All Methods



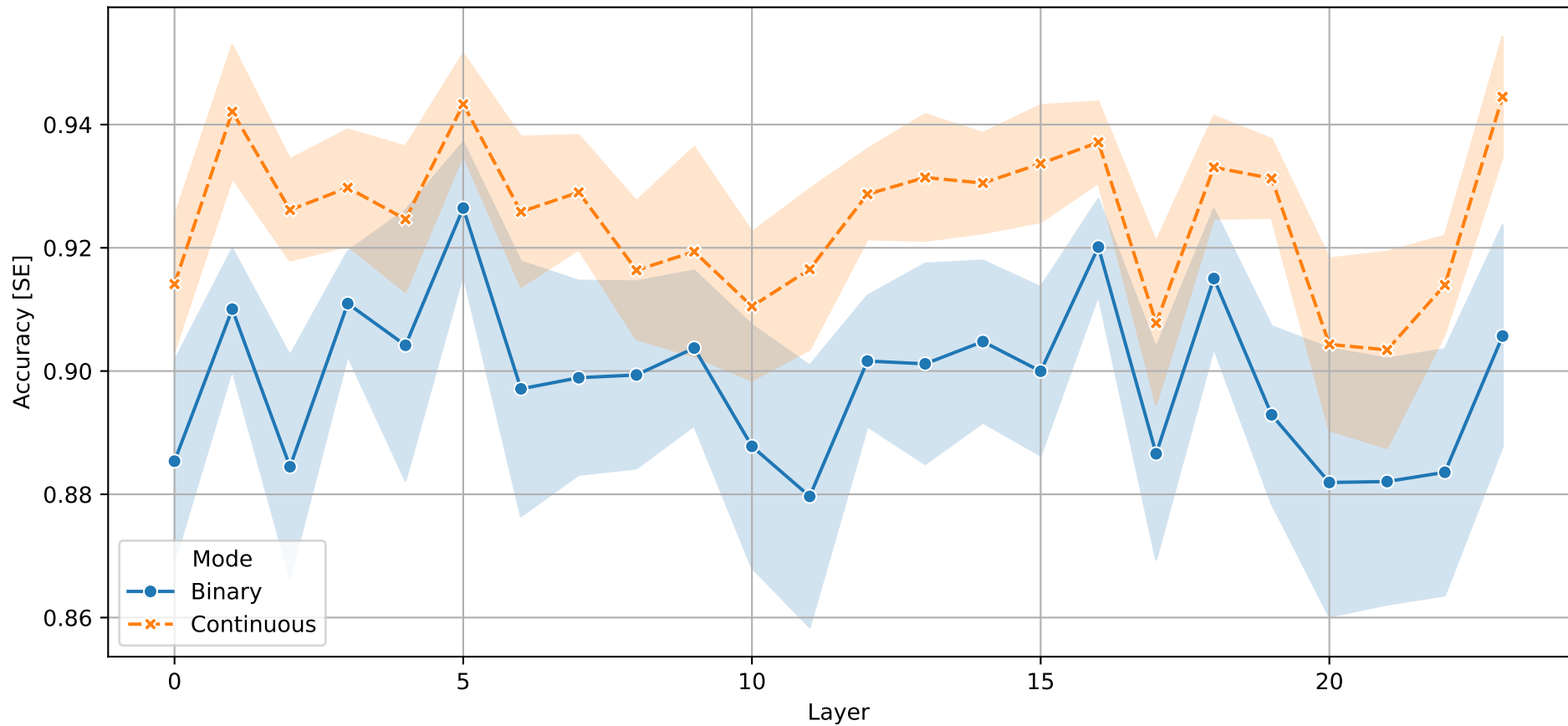
F1 Score Summary by Probing Method

Method	Metric	Binary	Continuous
Full Layer	f1_best_layer	5.0	5.0
Full Layer	f1_max	0.9988	0.9988
Full Layer	f1_mean	0.955	0.9602
Full Layer	f1_std	0.0221	0.0203
Single Neuron	f1_best_layer	3.0	1.0
Single Neuron	f1_max	0.9833	0.9876
Single Neuron	f1_mean	0.822	0.8325
Single Neuron	f1_std	0.0753	0.0693
Top-K Neurons	f1_best_layer	5.0	5.0
Top-K Neurons	f1_max	0.9857	0.9903
Top-K Neurons	f1_mean	0.8758	0.9209
Top-K Neurons	f1_std	0.0543	0.0327

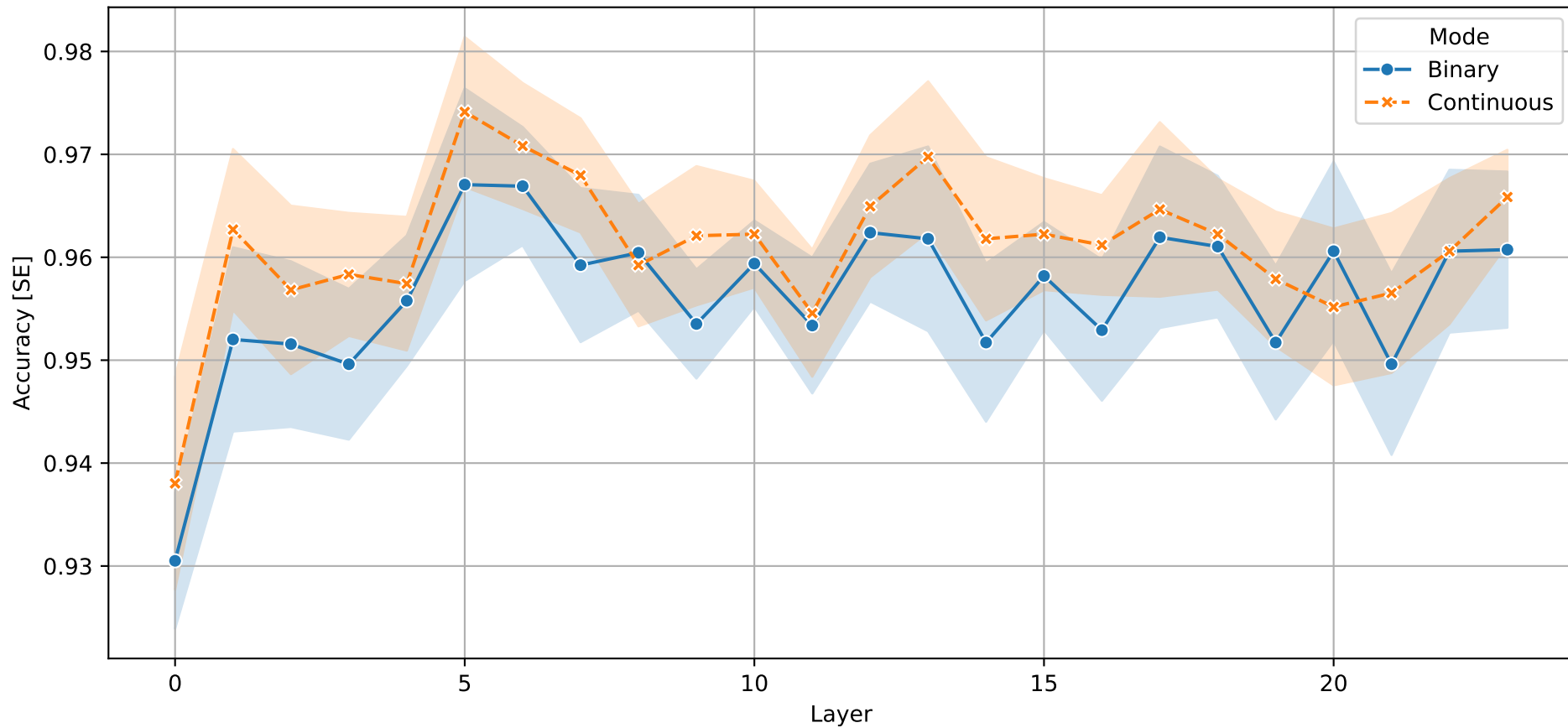
Accuracy per Layer - Single Neuron Probing



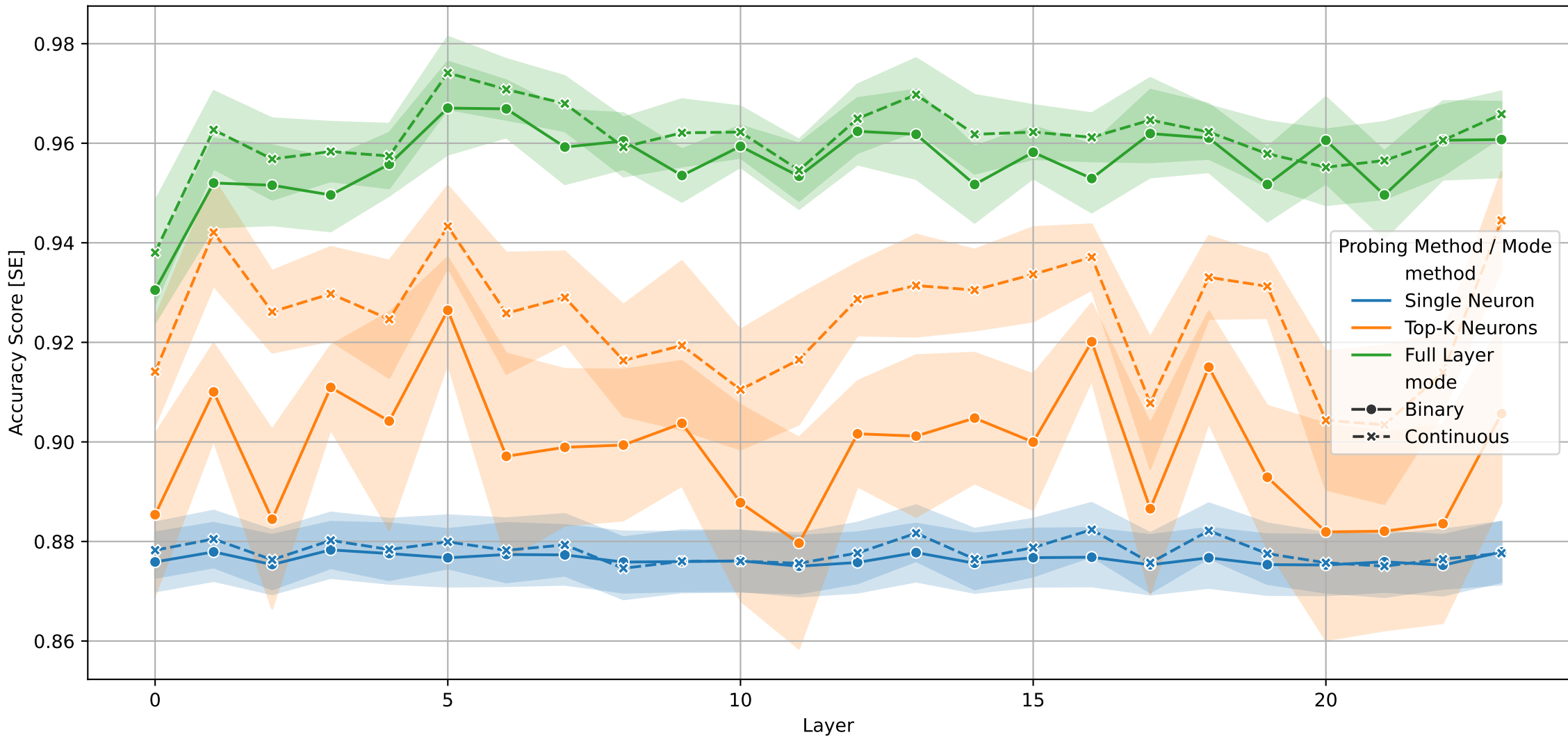
Accuracy per Layer - Top-K Neurons Probing



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	0.9988	0.9988
Full Layer	accuracy_mean	0.9564	0.9611
Full Layer	accuracy_std	0.0208	0.0195
Single Neuron	accuracy_best_layer	3.0	16.0
Single Neuron	accuracy_max	0.9832	0.988
Single Neuron	accuracy_mean	0.8764	0.8779
Single Neuron	accuracy_std	0.0533	0.053
Top-K Neurons	accuracy_best_layer	5.0	23.0
Top-K Neurons	accuracy_max	0.9856	0.9904
Top-K Neurons	accuracy_mean	0.8985	0.9249
Top-K Neurons	accuracy_std	0.0447	0.031