Monosemanticity spectrum (decoder.layer3) 1500 2000 2500 Neuron index