Monosemanticity spectrum (decoder.layer4)

1.0

0.8

0.6

