

Transformer test loss vs. learning rate and depth, at  $n = 512$

