

Transformer train loss vs. learning rate and number of samples, batch size 1024

