

Transformer train loss vs. learning rate and number of samples

