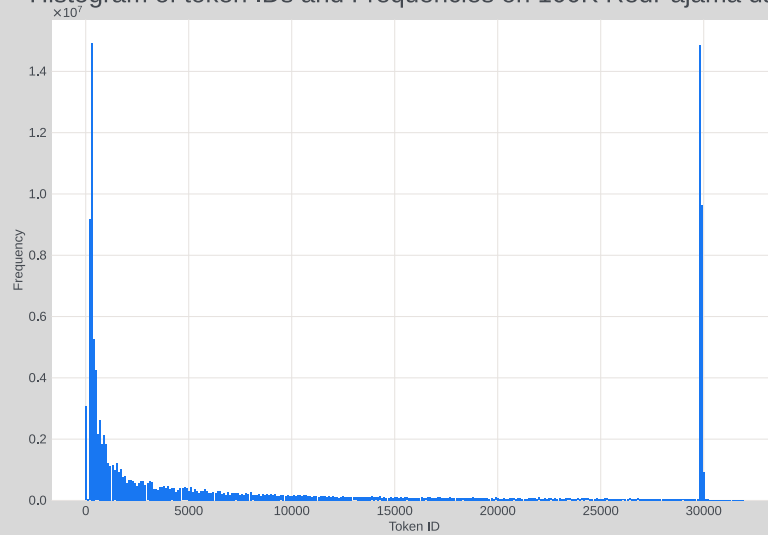
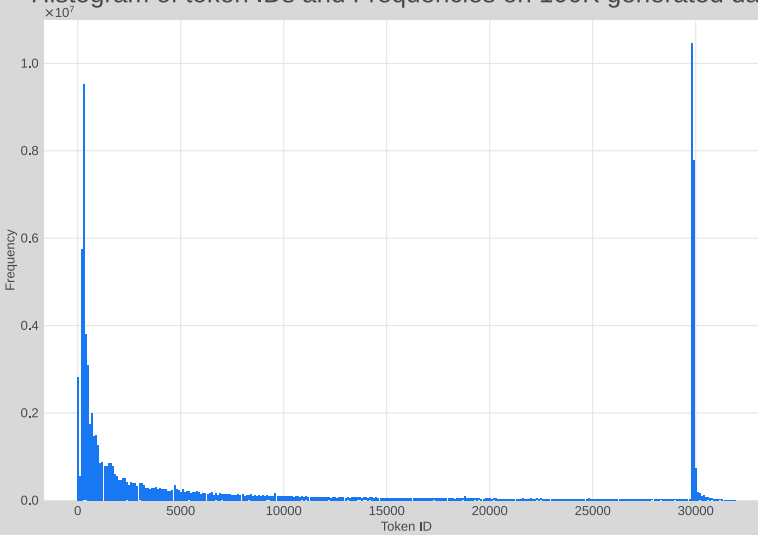


Histogram of token IDs and Frequencies on 100K RedPajama data



Histogram of token IDs and Frequencies on 100K generated data



Histogram of token IDs and Frequencies on 200K generated data

