Cumulative Token Count Percentage vs. Number of Tokens 100 Train Dataset Test Dataset 90 Unseen Dataset (%) 80 Cumulative Percentage 70 60 50 40 30 20 10 0 0 20 40 60 80 20 20 20 20 20 20 20 20 20 30 30 30 30 **Number of Tokens**