**Average Sequence Perplexity Heatmaps for Style Backdoor Trigger**