Average Sequence Perplexity heatmaps with Threshold = 1.0