## Normalized confusion matrix (gpt-4o-mini) action -0.63 0.02 0.00 0.02 0.03 0.01 0.15 0.09 0.02 0.01 0.02 0.01 0.00 delivery -0.18 0.02 0.01 0.41 0.02 0.00 0.15 0.02 0.11 0.05 0.00 0.01 0.00 entrance -0.01 0.00 0.00 0.00 0.87 0.04 0.02 0.02 0.01 0.02 0.00 0.02 0.00 exit -0.01 0.01 0.00 0.00 0.17 0.76 0.02 0.03 0.00 0.00 0.00 0.00 0.00 True label interaction -0.09 0.07 0.00 0.00 0.02 0.00 <mark>0.79</mark> 0.03 0.00 0.00 0.00 0.00 0.00 movement -0.06 0.00 0.00 0.00 0.03 0.17 0.15 0.58 0.00 0.00 0.00 0.01 0.01 music -0.01 0.00 0.09 0.04 0.01 0.00 0.02 0.01 0.72 0.02 0.03 0.03 0.01 narration -0.11 0.05 0.00 0.01 0.03 0.01 0.18 0.03 0.06 0.48 0.01 0.04 0.00 object -0.31 0.04 0.00 0.20 0.01 0.00 0.14 0.01 0.01 0.01 0.25 0.01 0.00 toward -0.01 0.00 0.02 0.01 0.03 0.00 0.06 0.01 0.00 0.00 0.00 0.06 0.79 action interaction narration MUSIC Predicted label