## Normalized confusion matrix (gpt-4o-mini) action -0.38 0.02 0.00 0.03 0.03 0.01 0.20 0.11 0.01 0.03 0.12 0.02 0.03 delivery -0.23 0.04 0.02 0.40 0.02 0.00 0.07 0.03 0.05 0.11 0.00 0.00 0.04 entrance -0.05 0.00 0.00 0.01 0.77 0.06 0.00 0.01 0.02 0.05 0.00 0.02 0.00 exit -0.02 0.01 0.00 0.00 0.14 0.77 0.00 0.02 0.00 0.02 0.01 0.00 0.00 True label interaction -0.14 0.09 0.00 0.02 0.04 0.00 0.61 0.03 0.00 0.02 0.05 0.00 0.01 movement -0.21 0.01 0.00 0.00 0.07 0.13 0.08 0.43 0.01 0.02 0.00 0.01 0.03 music -0.02 0.00 0.00 0.05 0.01 0.00 0.00 0.00 0.82 0.06 0.02 0.01 0.00 narration -0.12 0.07 0.01 0.00 0.02 0.02 0.04 0.02 0.01 0.64 0.02 0.04 0.00 object -0.12 0.03 0.00 0.02 0.01 0.00 0.08 0.04 0.00 0.01 0.65 0.01 0.00 setting -0.04 0.01 0.01 0.01 0.02 0.02 0.02 0.06 0.01 0.04 0.00 0.77 0.01 toward -0.01 0.00 0.00 0.00 0.04 0.00 0.03 0.00 0.01 0.00 0.01 0.03 0.86 action movement interaction narration MUSIC

Predicted label