## Normalized confusion matrix (gpt-4o-mini) action -0.60 0.02 0.00 0.01 0.04 0.01 0.11 0.10 0.01 0.03 0.03 0.02 0.01 delivery -0.16 0.01 0.02 0.40 0.03 0.00 0.17 0.03 0.10 0.05 0.00 0.01 0.00 entrance -0.02 0.00 0.01 0.00 0.83 0.04 0.04 0.01 0.01 0.03 0.00 0.02 0.00 exit -0.01 0.00 0.00 0.00 0.18 0.68 0.05 0.07 0.00 0.00 0.00 0.00 True label interaction -0.10 0.11 0.00 0.00 0.01 0.00 <mark>0.75</mark> 0.02 0.00 0.00 0.01 0.00 0.00 movement -0.08 0.00 0.00 0.00 0.08 0.14 0.08 0.60 0.00 0.02 0.00 0.00 0.01 music -0.01 0.00 0.02 0.01 0.01 0.00 0.00 0.01 0.86 0.04 0.01 0.03 0.01 narration -0.09 0.05 0.00 0.01 0.03 0.02 0.17 0.02 0.03 0.53 0.01 0.04 0.00 object -0.40 0.03 0.00 0.09 0.01 0.00 0.13 0.01 0.01 0.01 0.28 0.01 0.00 setting -0.02 0.01 0.01 0.00 0.02 0.02 0.01 0.09 0.00 0.04 0.01 0.79 0.00 toward -0.00 0.00 0.00 0.00 0.03 0.00 0.06 0.01 0.01 0.00 0.00 0.14 0.75 actionession interaction narration MUSIC Predicted label