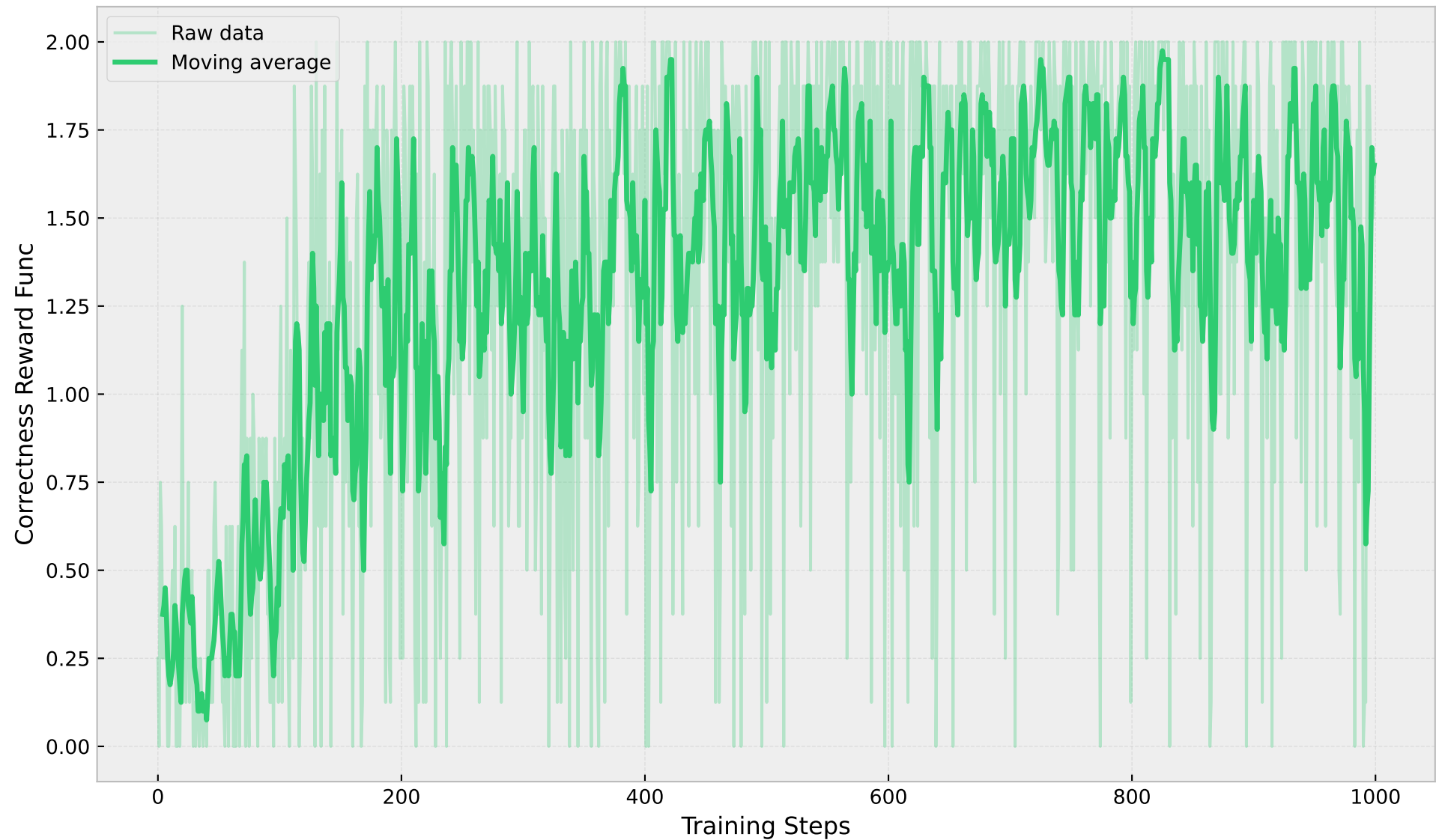
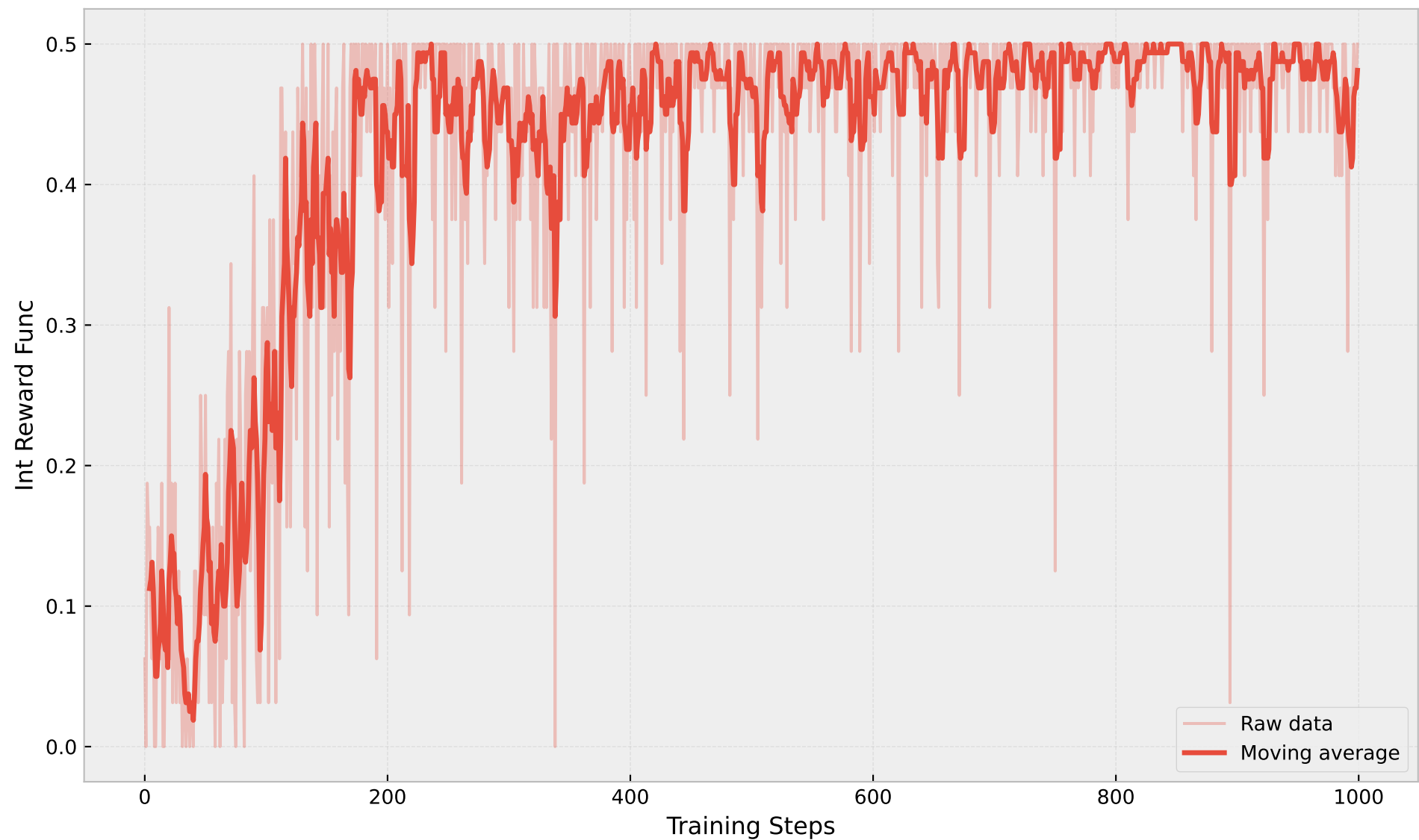


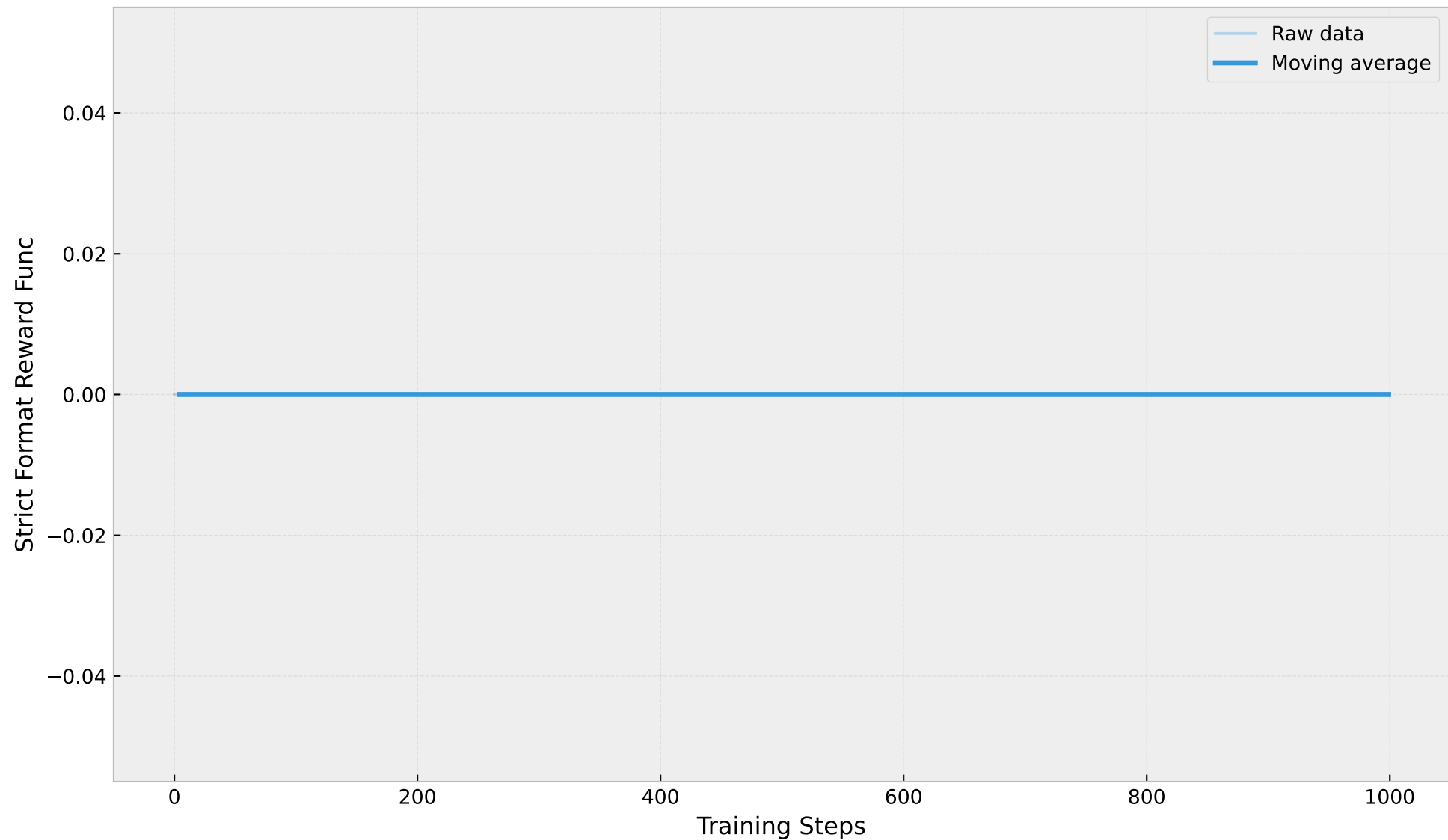
Correctness Reward Func



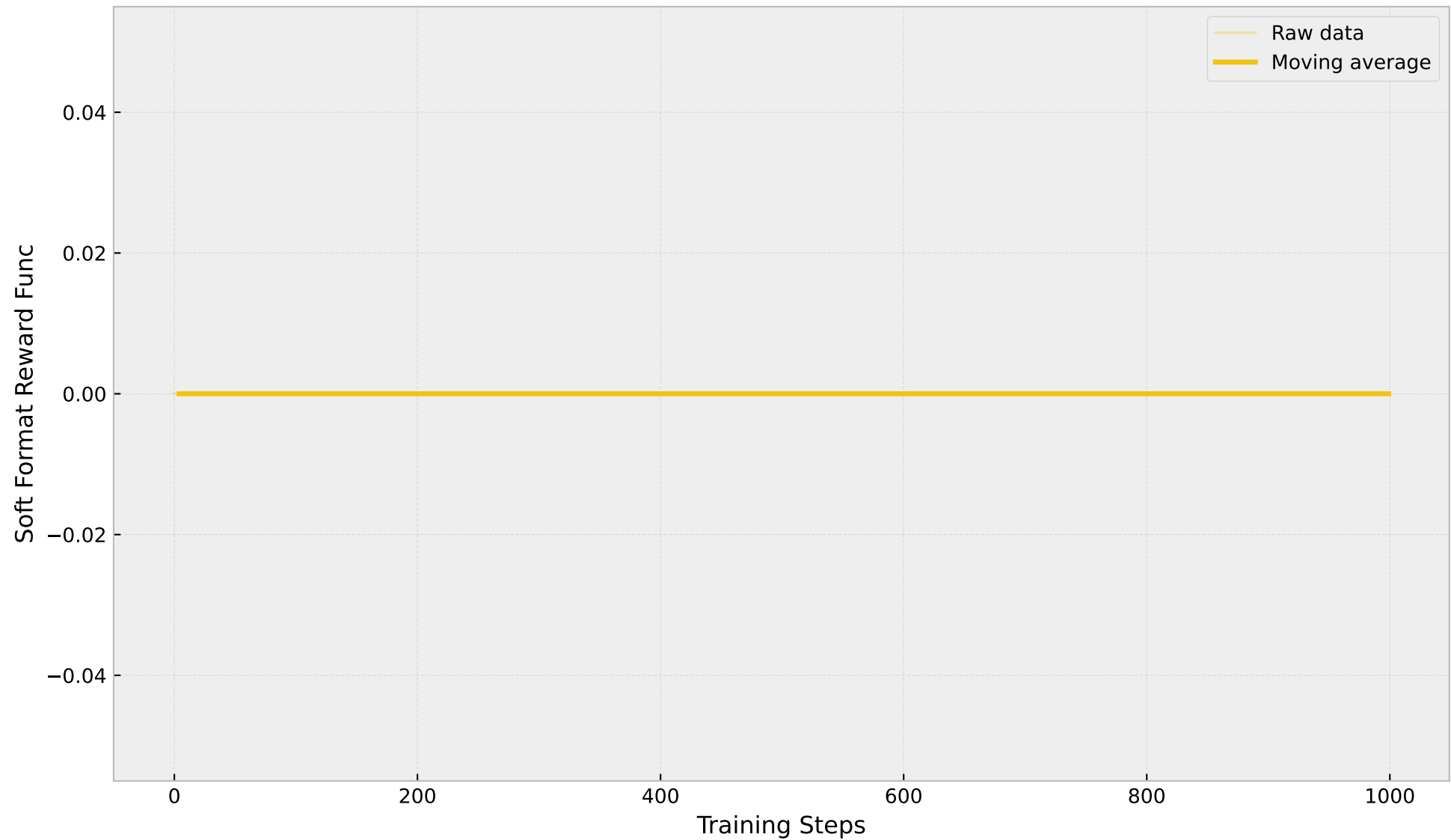
# Int Reward Func



Strict Format Reward Func



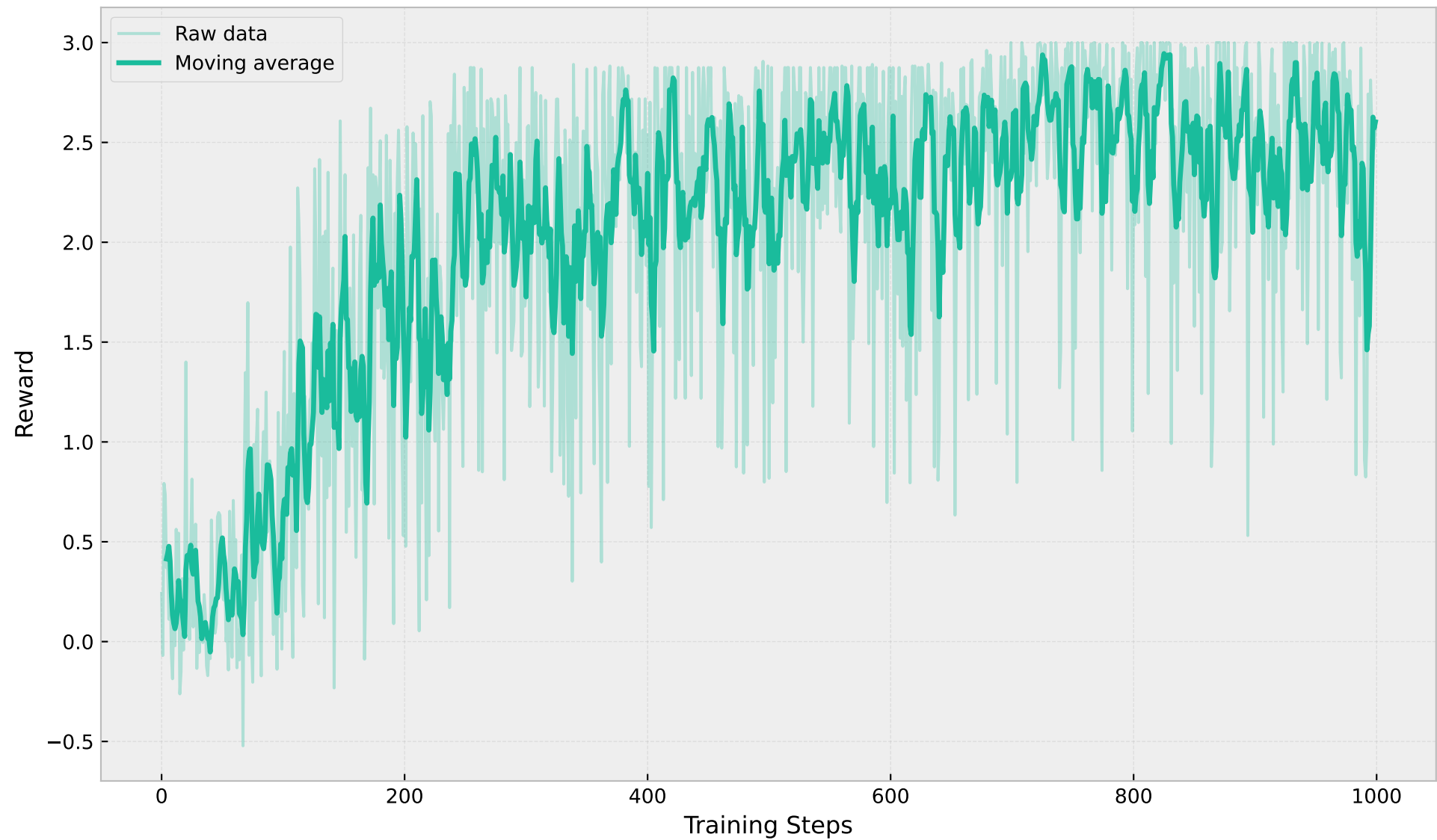
Soft Format Reward Func



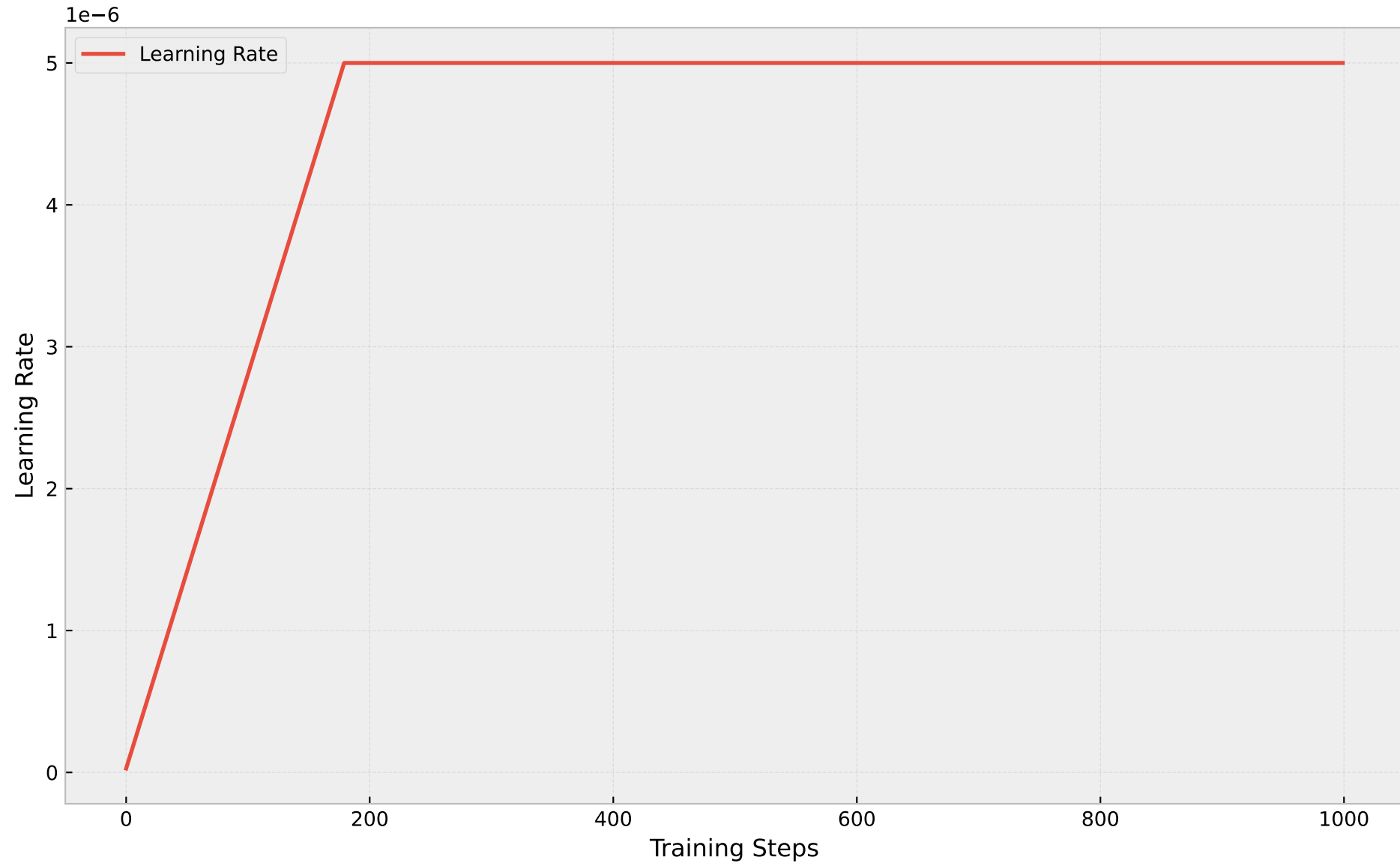
# Xmlcount Reward Func



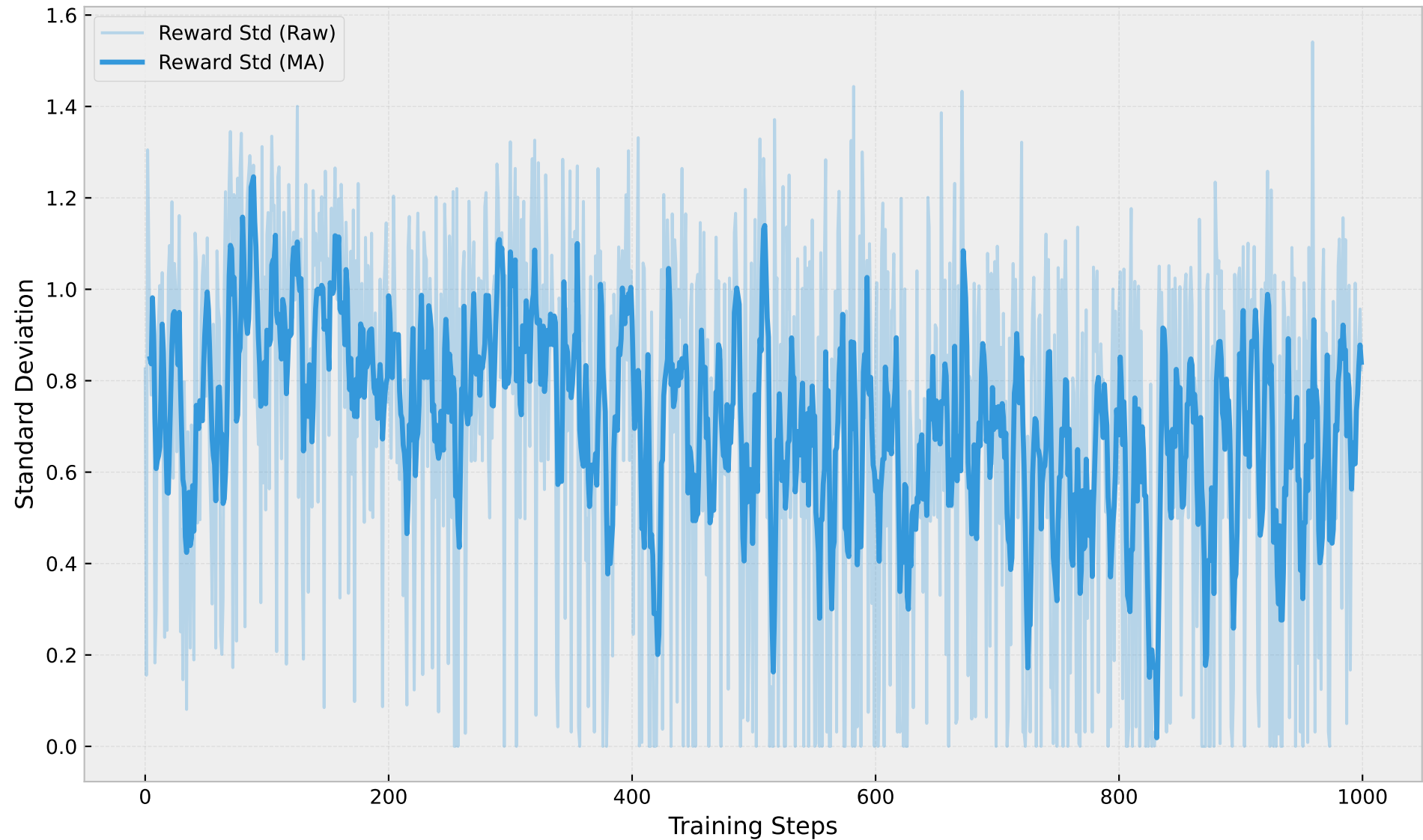
# Reward



# Learning Rate Schedule



# Reward Standard Deviation

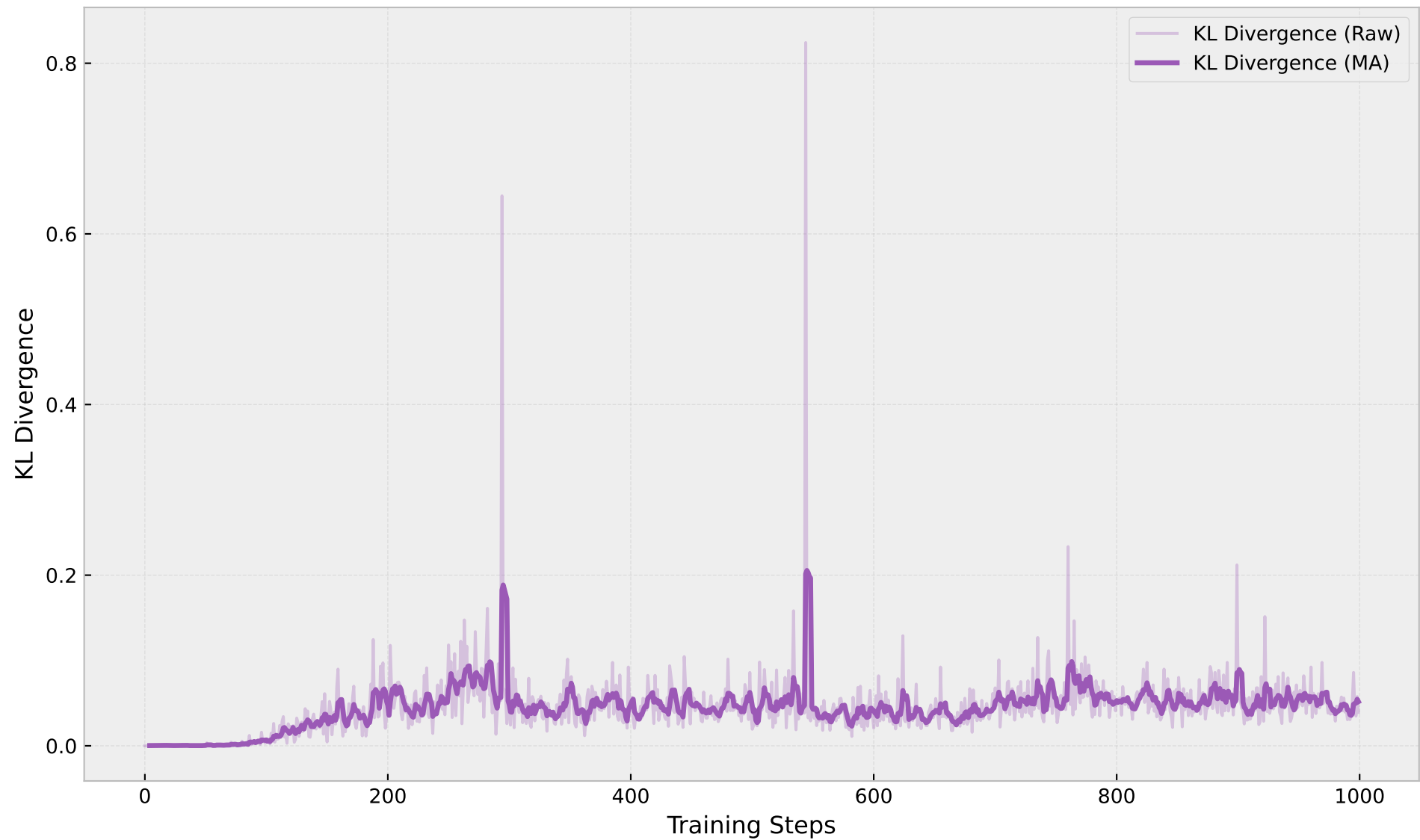




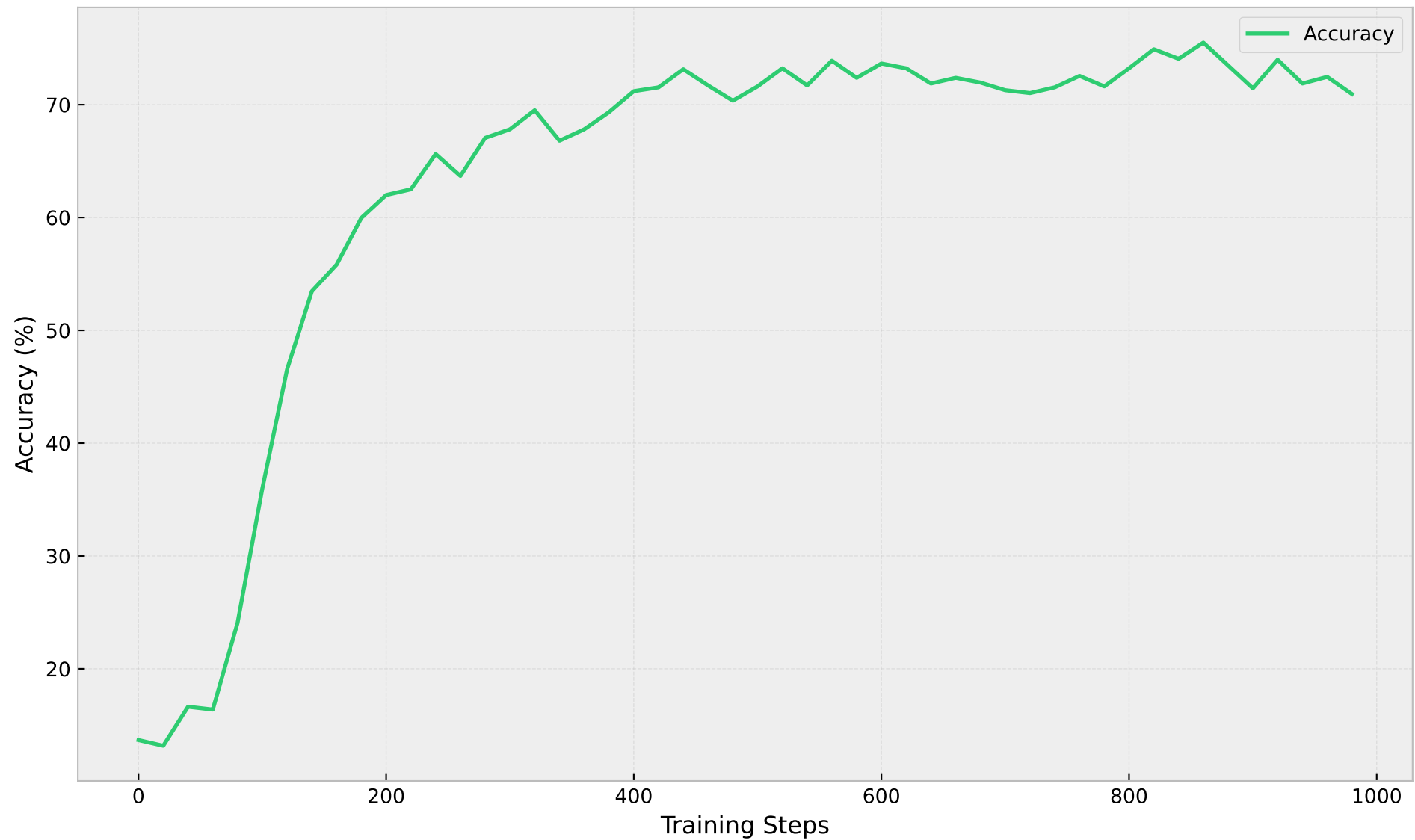
# Training Loss



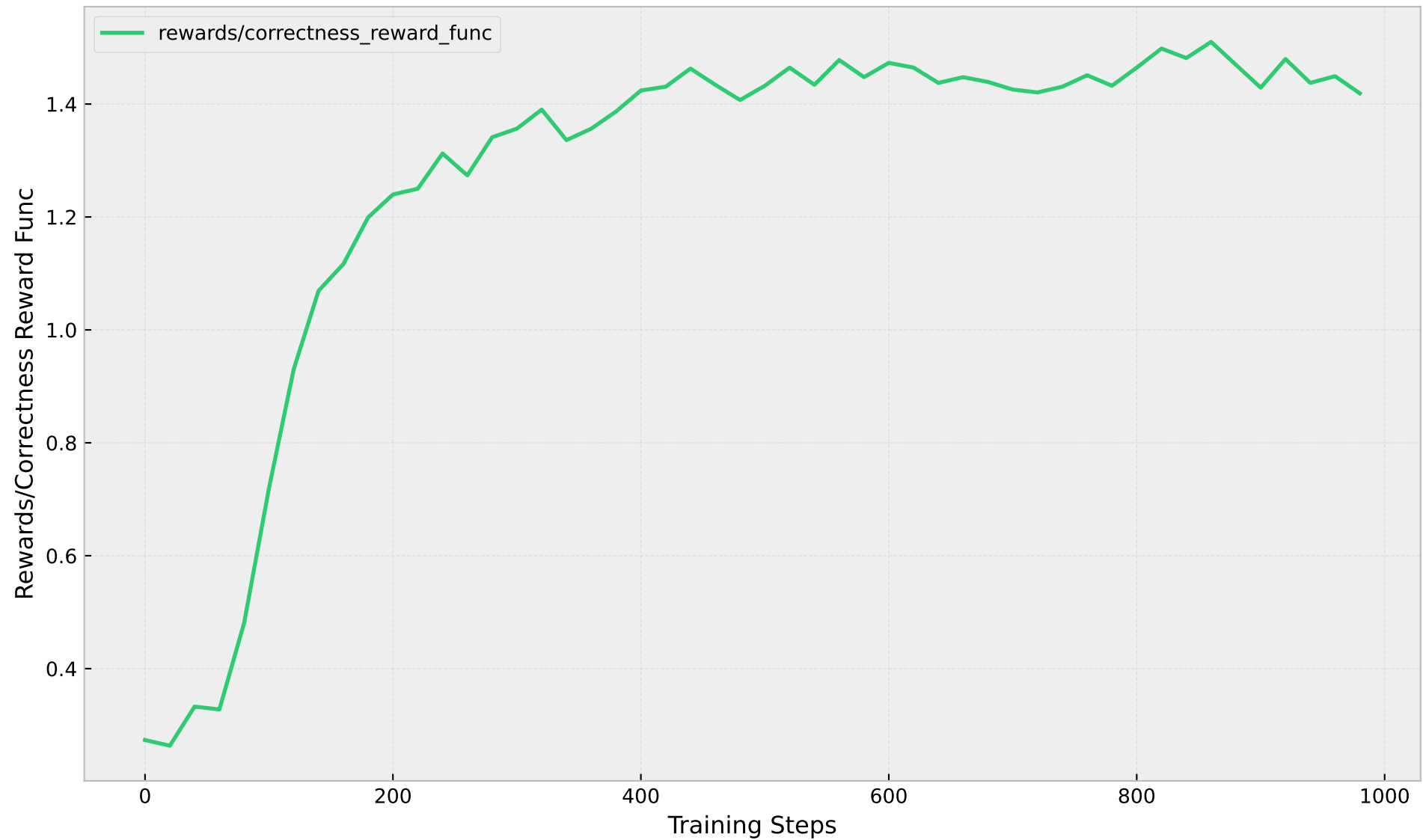
# KL Divergence



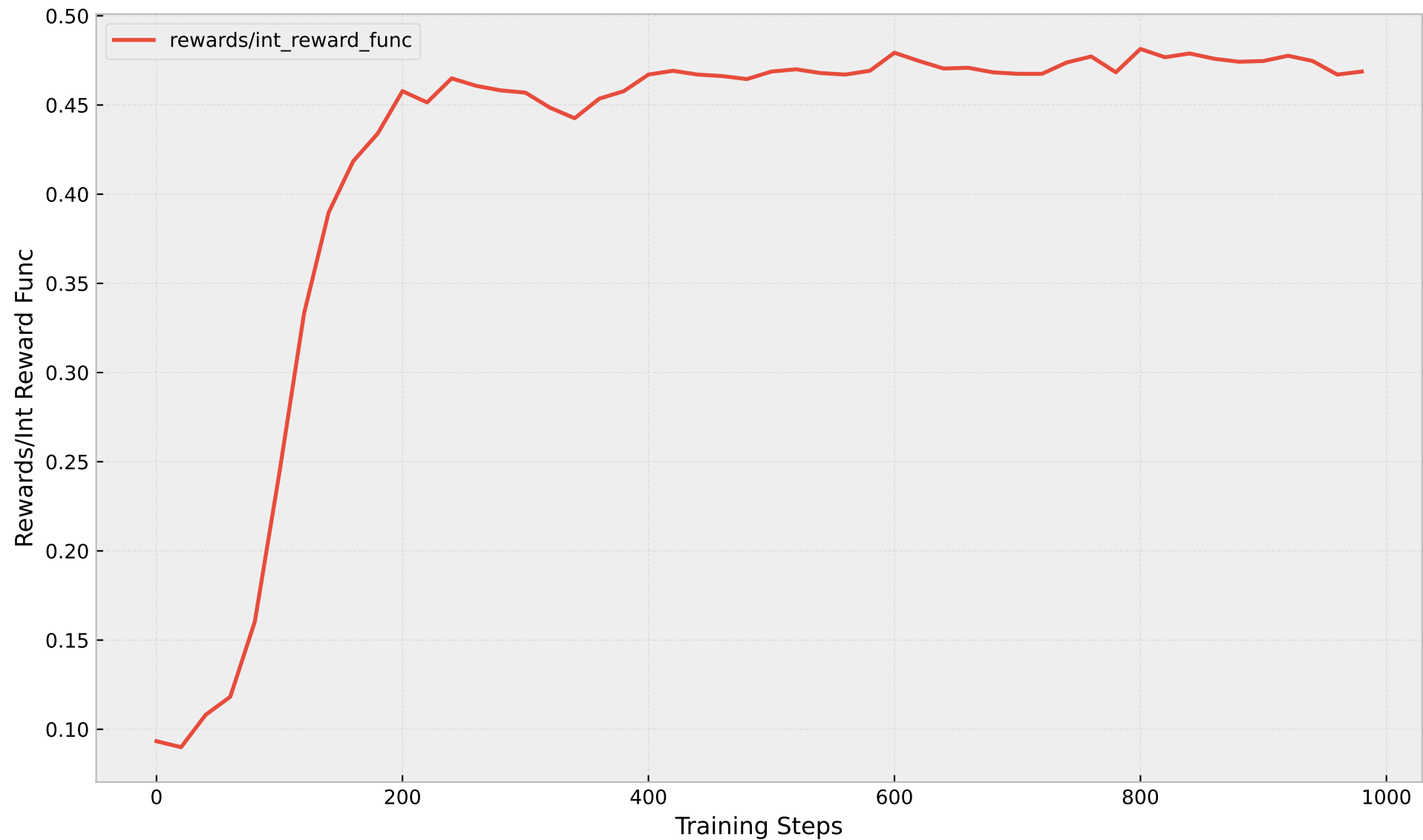
# Evaluation Accuracy



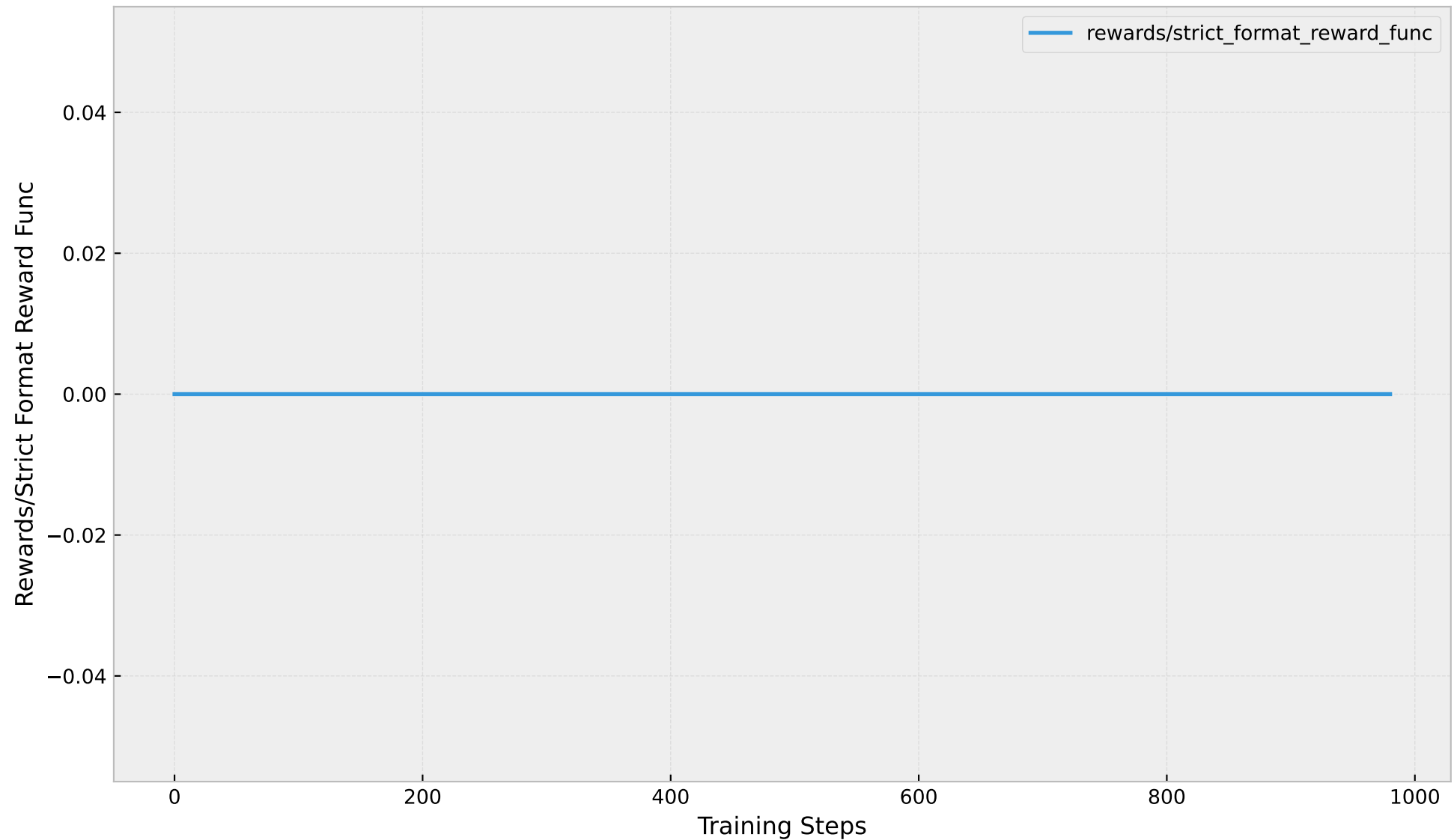
Evaluation Rewards/Correctness Reward Func



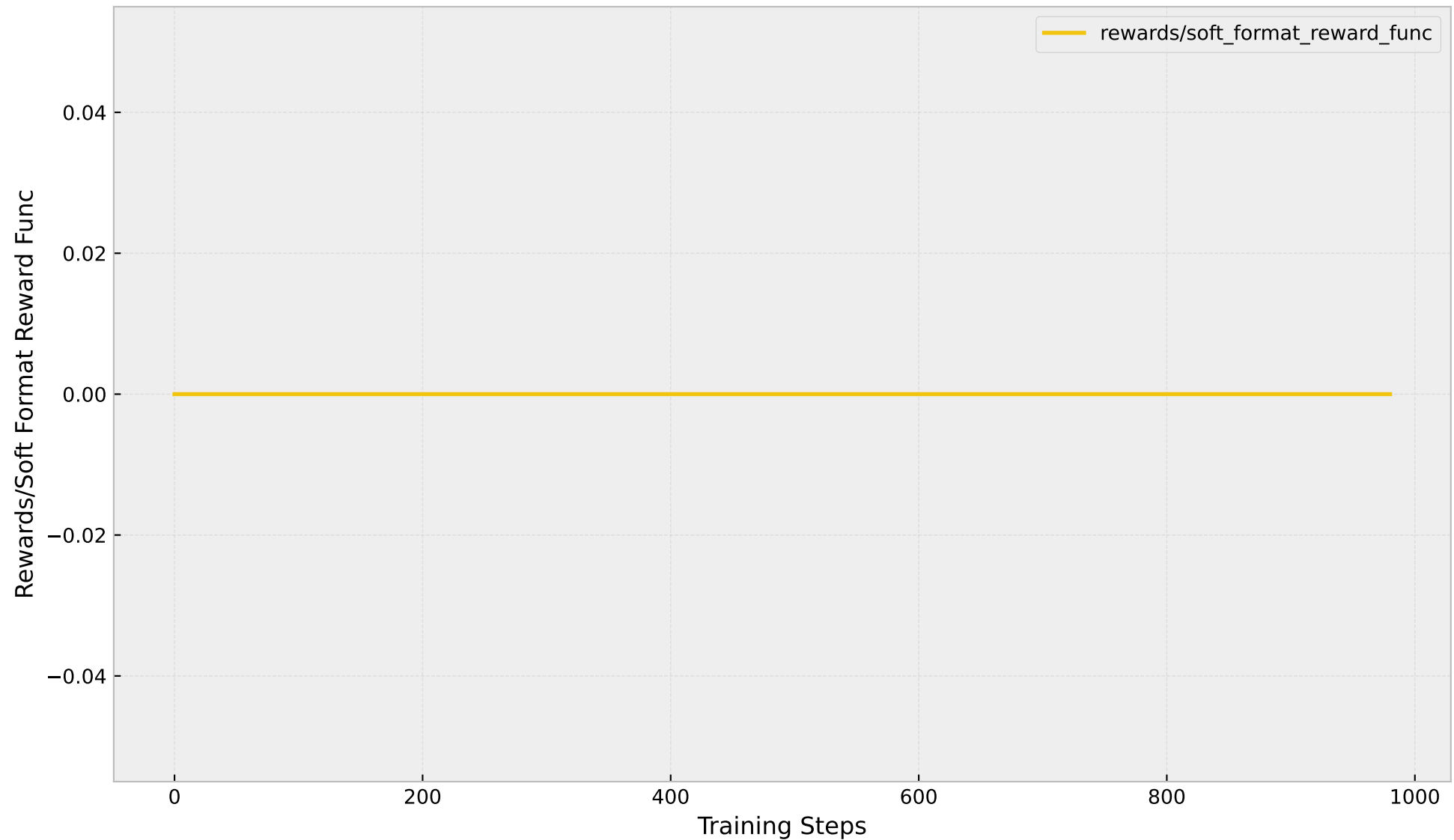
Evaluation Rewards/Int Reward Func



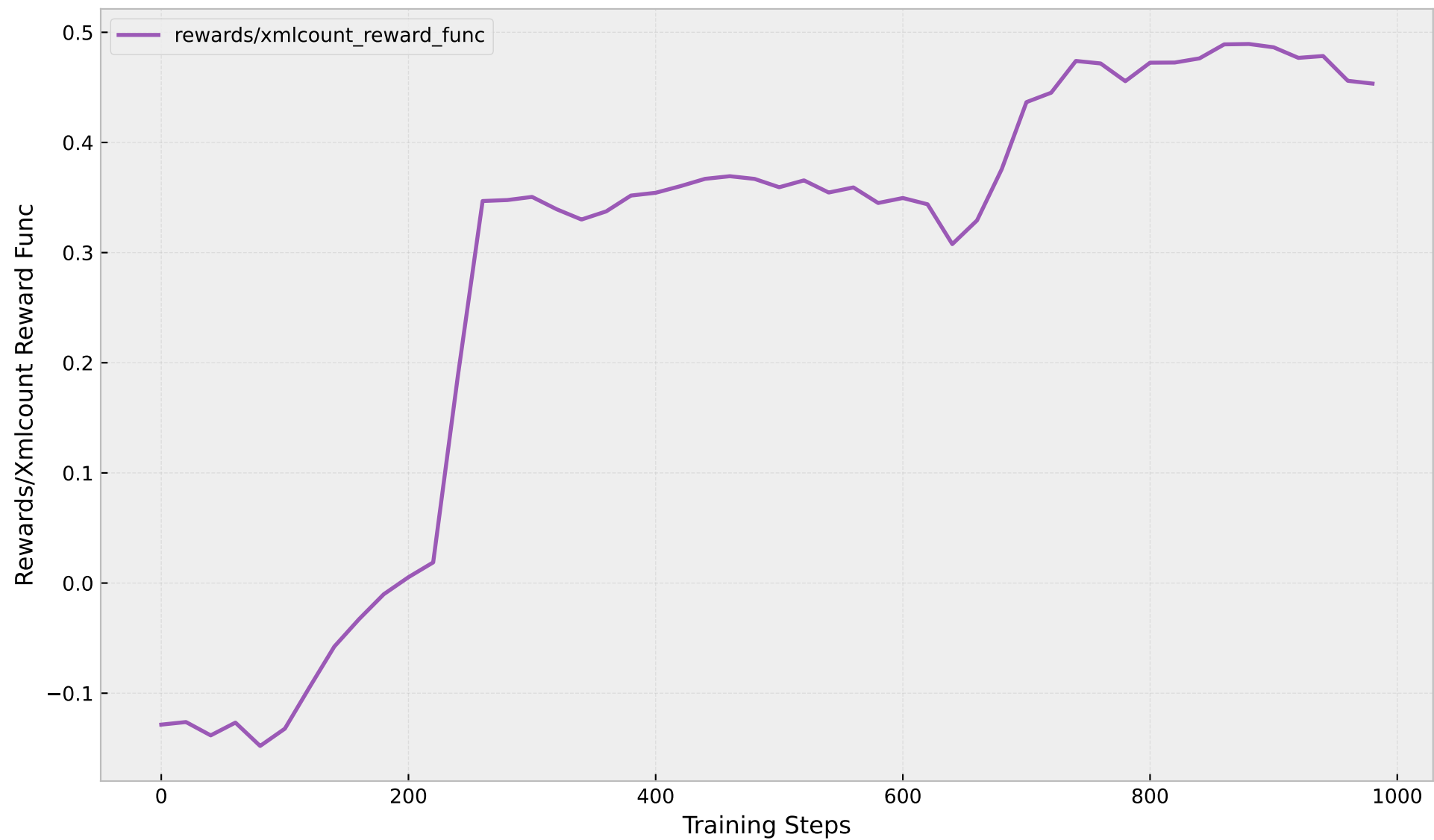
Evaluation Rewards/Strict Format Reward Func



Evaluation Rewards/Soft Format Reward Func

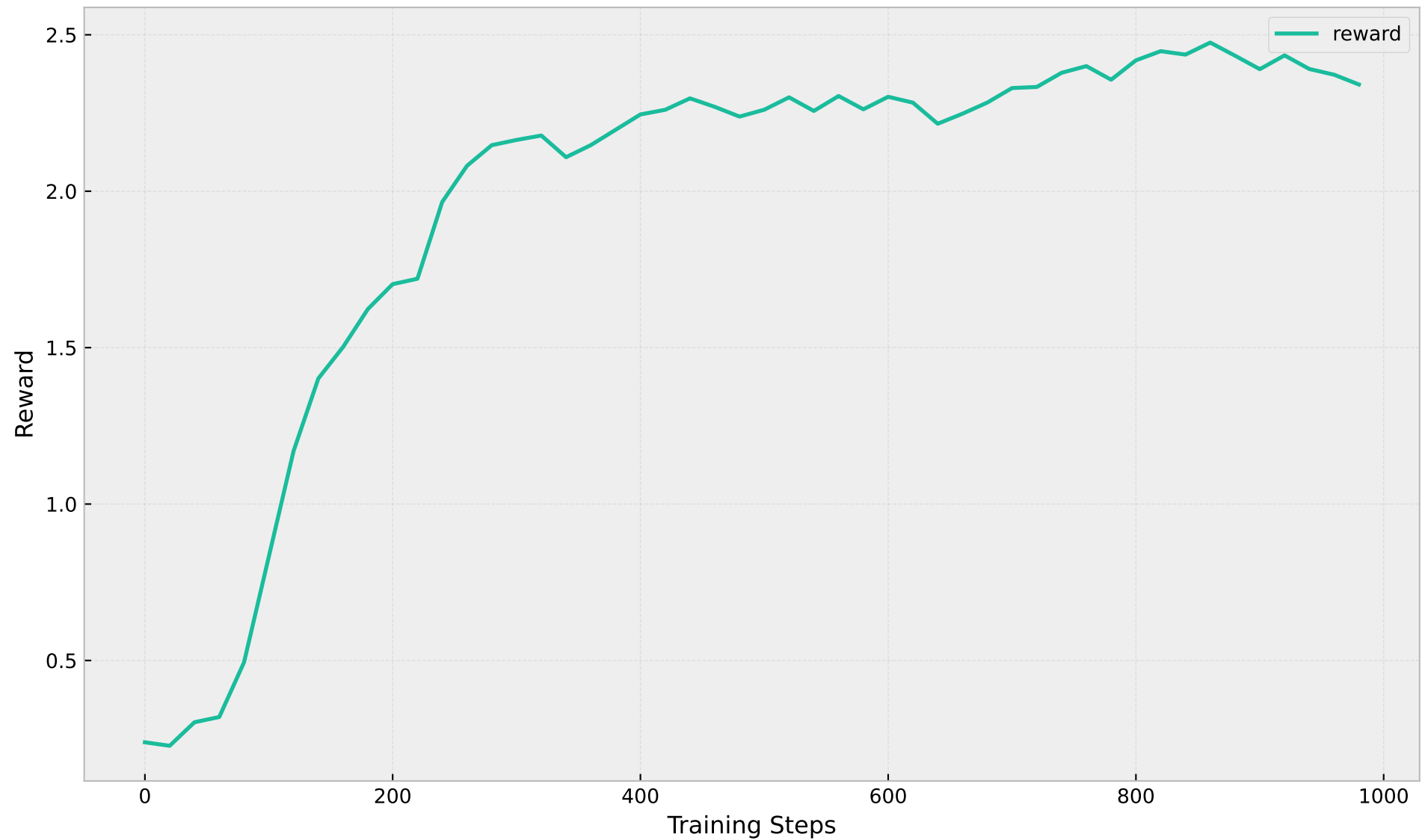


Evaluation Rewards/XMLcount Reward Func





# Evaluation Reward



# Evaluation Accuracy

