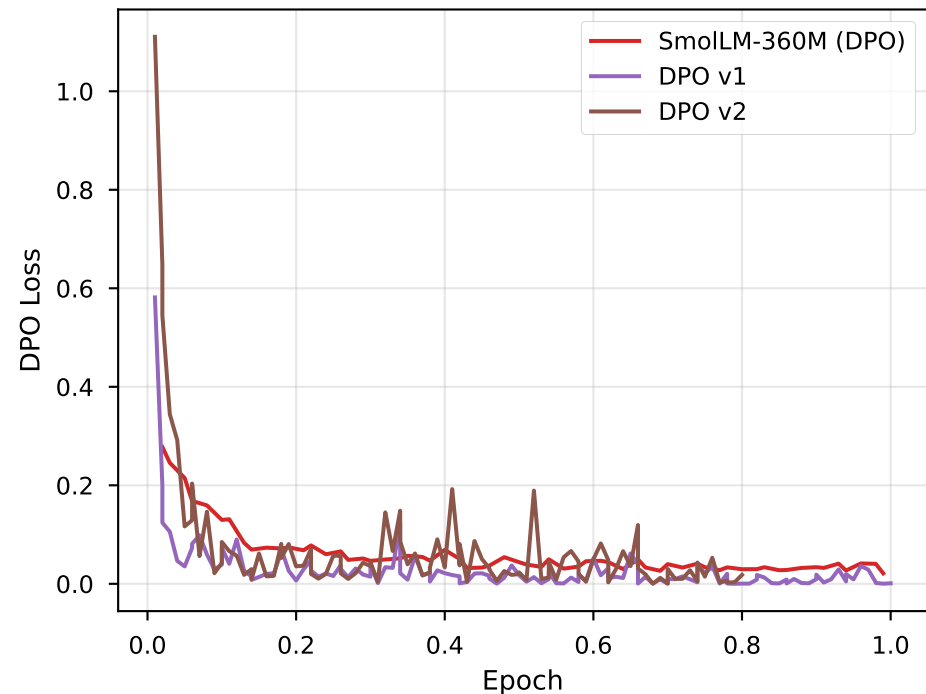
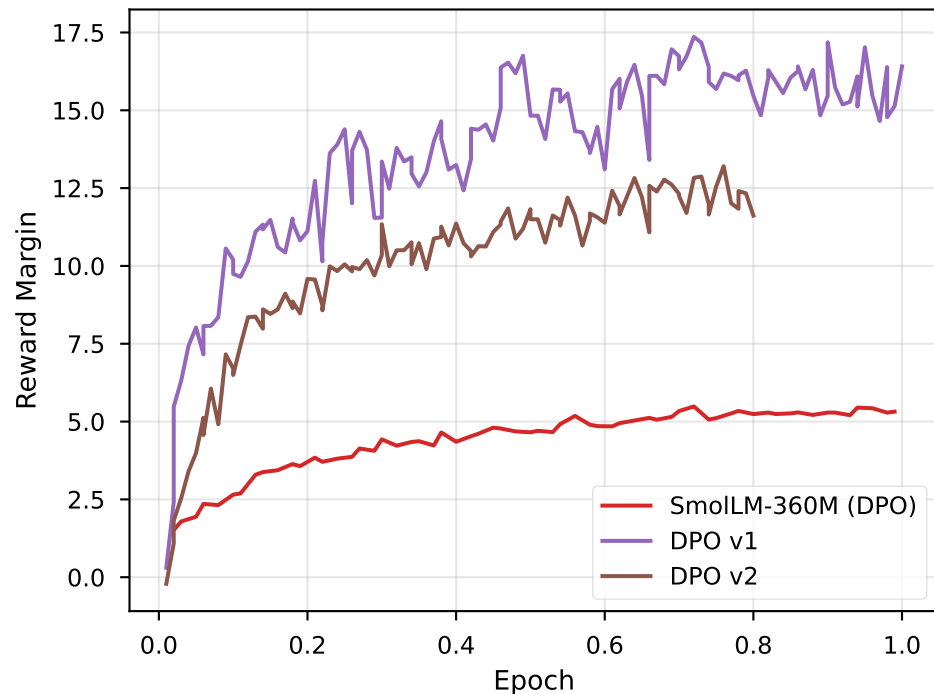


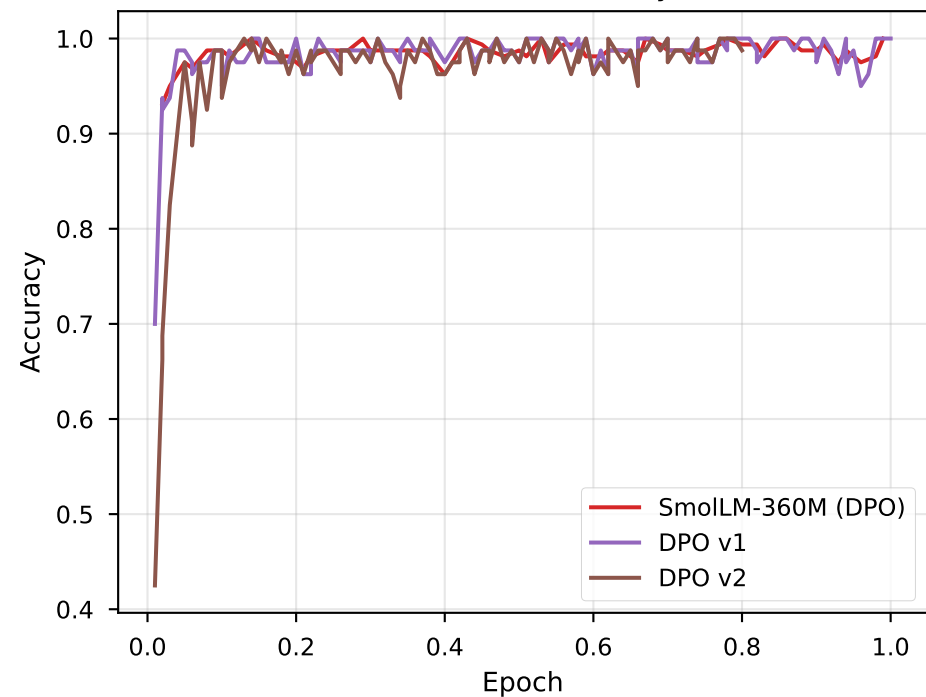
DPO Training Loss



Reward Margins (Chosen - Rejected)



Preference Accuracy



Chosen vs Rejected Rewards

