FrozenLake-Hard: Direct-SFT vs Textual-CoT Performance
(Test Set: 500 Problems)