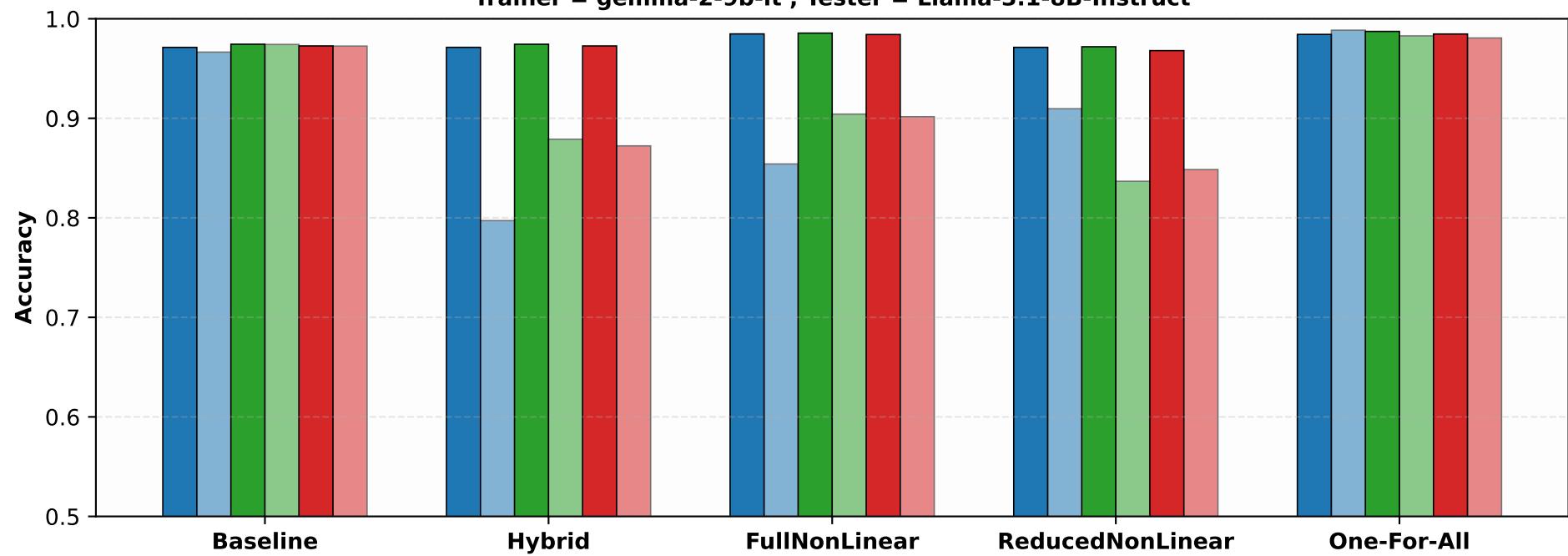
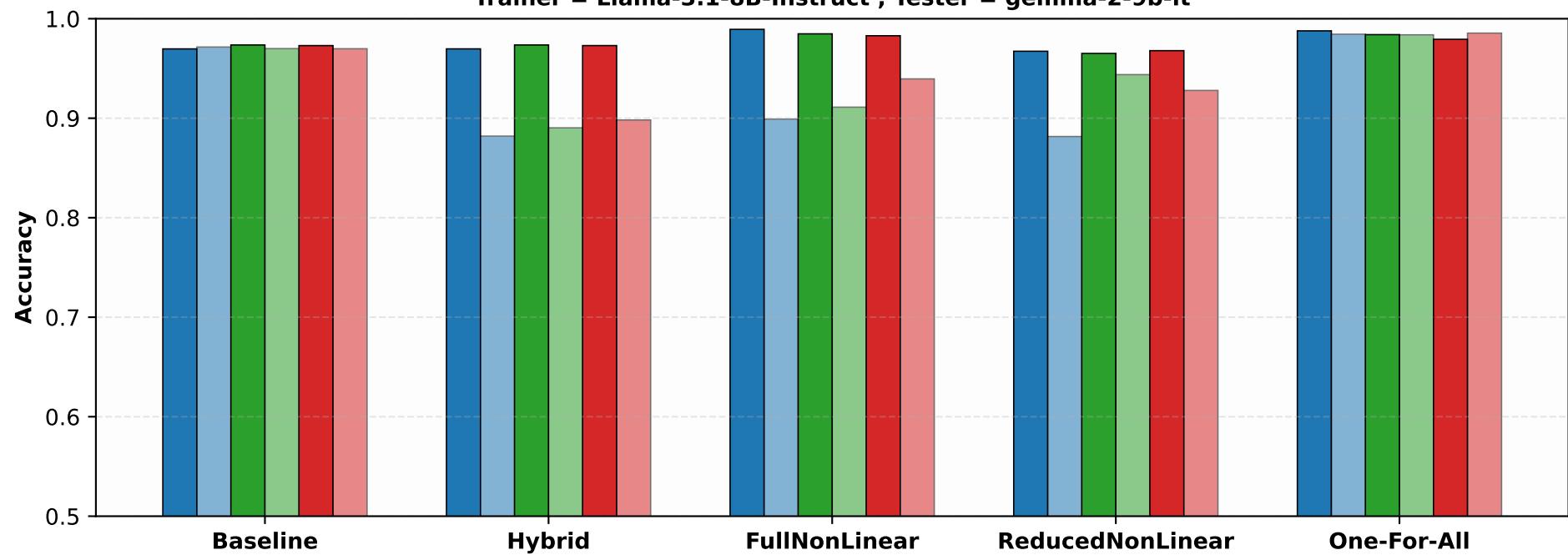


Trainer = gemma-2-9b-it , Tester = Llama-3.1-8B-Instruct



Trainer = Llama-3.1-8B-Instruct , Tester = gemma-2-9b-it



Layer (Model)

[Blue Box] attn (Trainer) [Light Blue Box] attn (Tester) [Green Box] hidden (Trainer) [Light Green Box] hidden (Tester) [Red Box] mlp (Trainer) [Light Red Box] mlp (Tester)