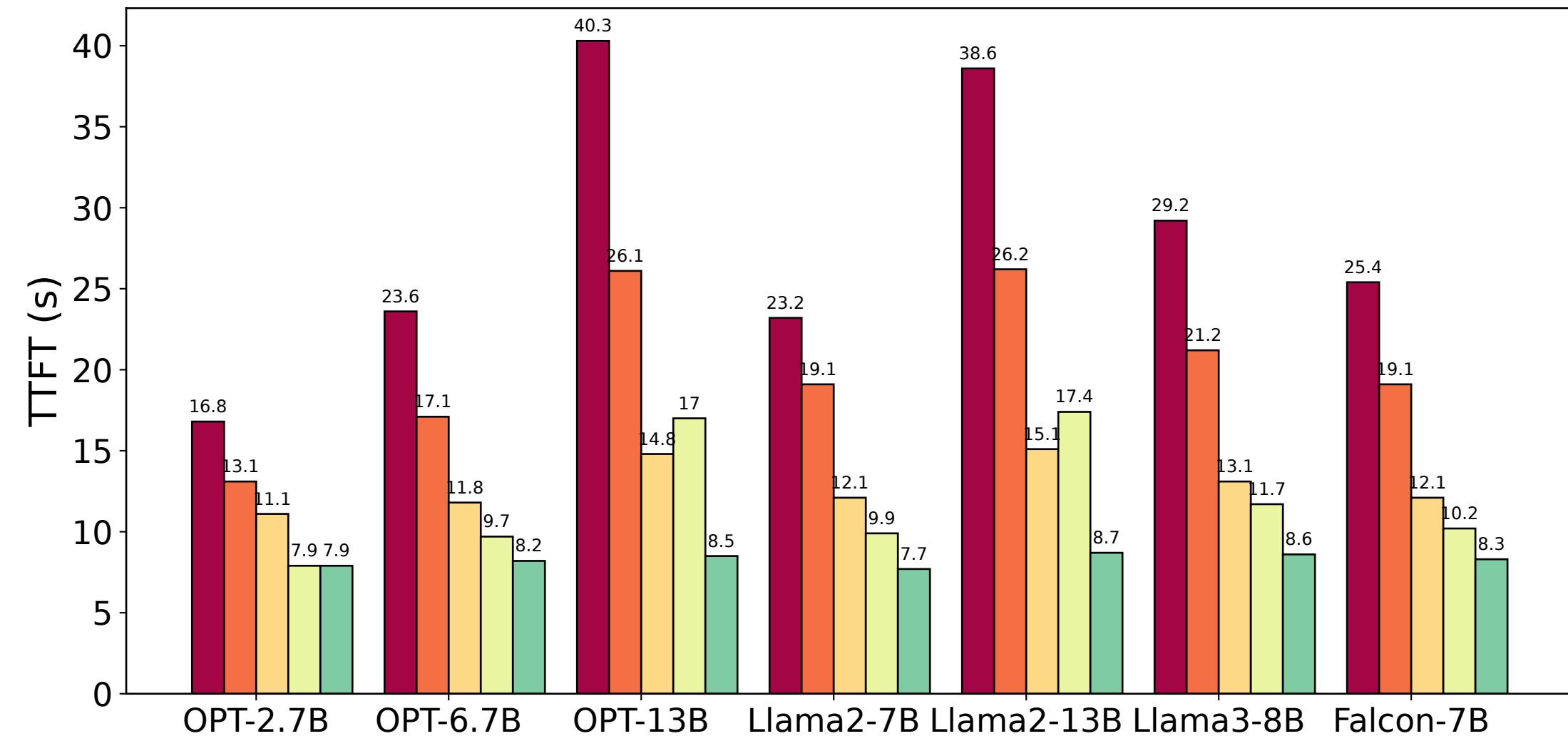
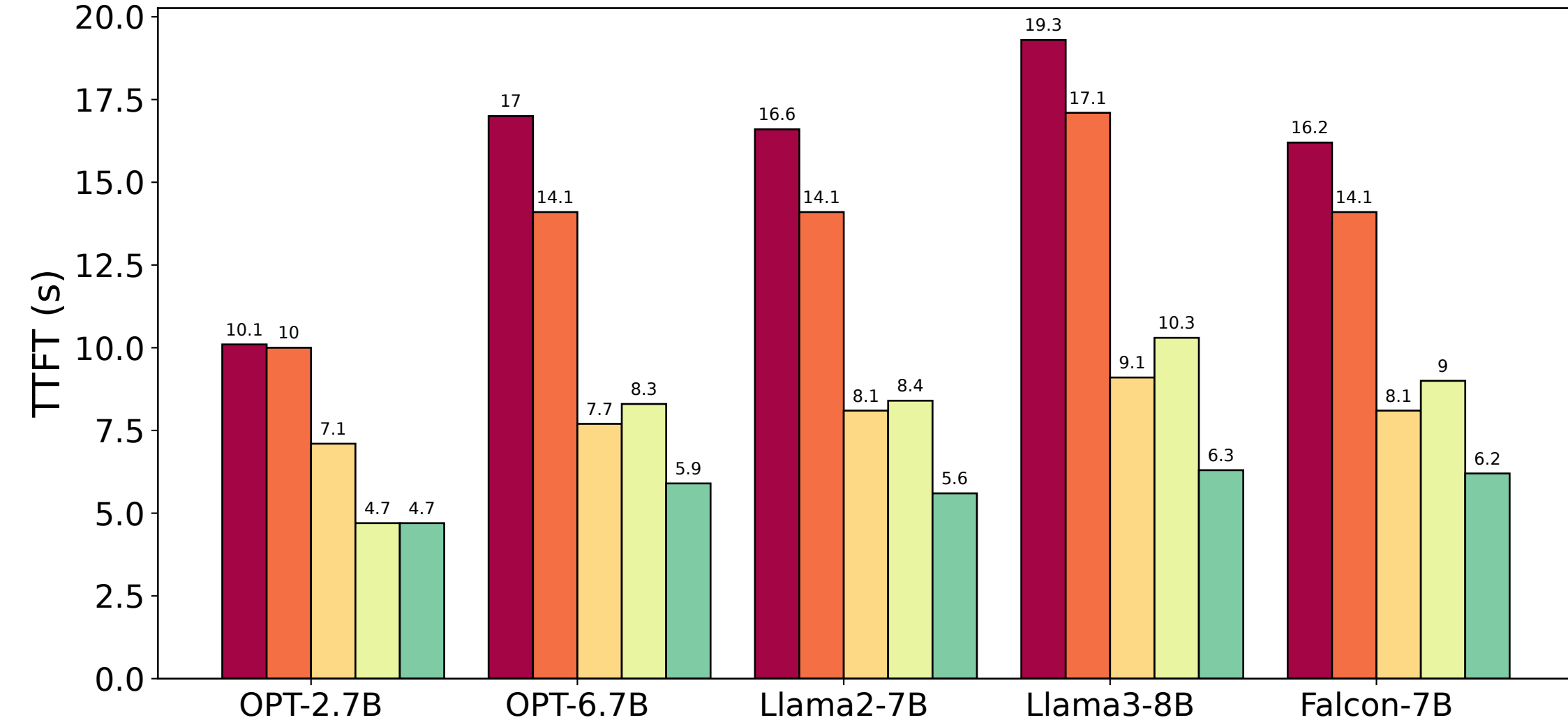


■ Serverless vLLM    
 ■ ServerlessLLM    
 ■ ServerlessLLM with cached model    
 ■ HydraServe with single worker    
 ■ HydraServe



(a) Models on V100



(b) Models on A10