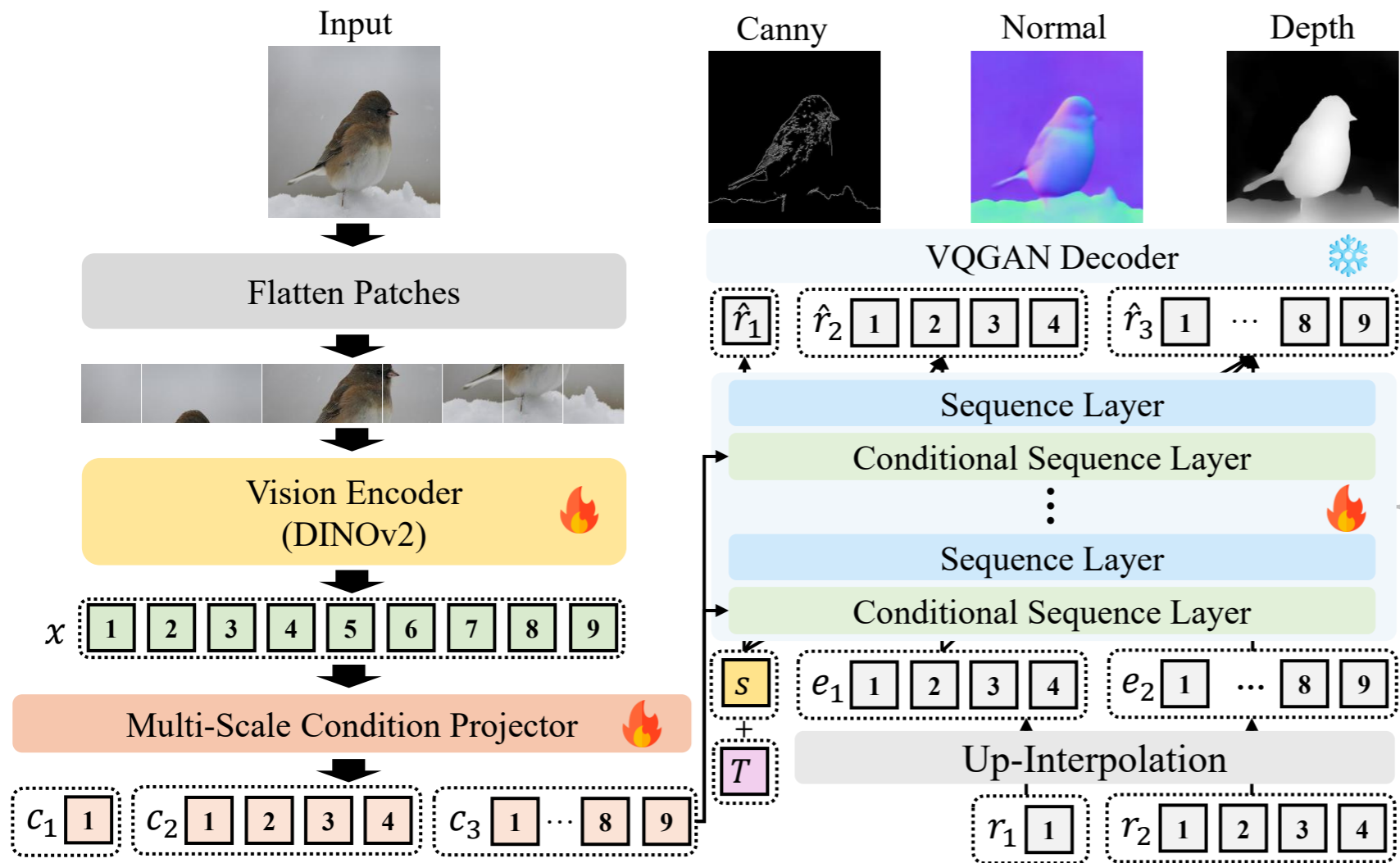


(a) Supervised Finetuning



(b) Online DPO

