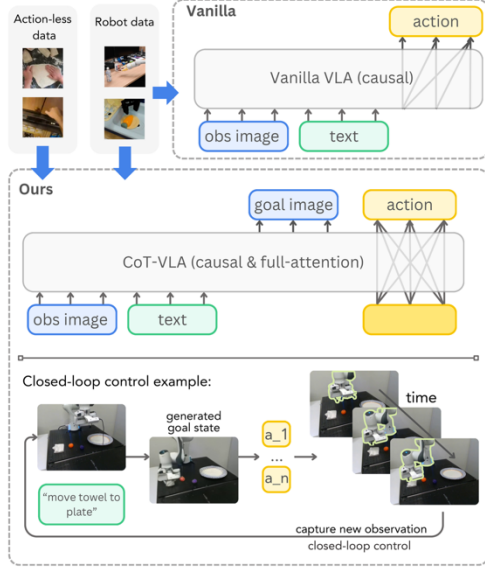


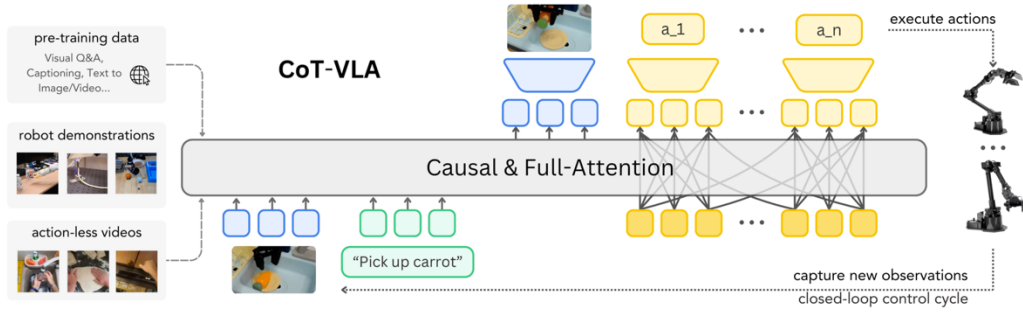
# CoT-VLA: Visual Chain-of-Thought Reasoning for Vision-Language-Action Models

Qingqing Zhao<sup>1,2,\*</sup> Yao Lu<sup>1</sup> Moo Jin Kim<sup>2</sup> Zipeng Fu<sup>2</sup>  
 Zhuoyang Zhang<sup>3</sup> Yecheng Wu<sup>1,3</sup> Zhaoshuo Li<sup>1</sup> Qianli Ma<sup>1</sup> Song Han<sup>1,3</sup> Chelsea Finn<sup>2</sup>  
 Ankur Handa<sup>1</sup> Ming-Yu Liu<sup>1</sup> Donglai Xiang<sup>1†</sup> Gordon Wetzstein<sup>2†</sup> Tsung-Yi Lin<sup>1†</sup>  
<sup>1</sup>NVIDIA <sup>2</sup>Stanford University <sup>3</sup>MIT

25



**Figure 1. Comparison between vanilla VLA and CoT-VLA frameworks.** Prior VLA models (top) directly predict robot actions from task inputs without explicit reasoning steps and only use action-annotated robot demonstration data for training. Unlike vanilla VLAs, CoT-VLA (bottom) can also leverage action-less datasets like EPIC-KITCHEN-100 [27] to enhance subgoal image generation ability, unlocking the potential of using abundant unlabeled video data to improve VLA’s visual reasoning capability. CoT-VLA first generates a subgoal image as an intermediate reasoning step, and then generate a short action sequence to achieve the subgoal. We outline the robot arm for better visualization.



**Figure 2. Overview of CoT-VLA framework.** We build our model on VILA-U [67], a generative multimodal model pretrained on interleaved text-image data. The base model then trains on robot demonstrations [48] and action-less videos [20, 27]. During deployment, given a visual observation and a text instruction, the model performs visual chain-of-thought reasoning by generating a subgoal image (upper blue) with causal attention. It then generates a short action sequence with full attention ( $a_1 \dots a_n$ ) for robot execution. The system operates in a closed-loop control manner by capturing new observations after executing predicted action sequences.



**Figure 3. Hybrid attention mechanism in CoT-VLA.** We use causal attention for image or text generation and full attention for action generation.  $[x]$ ,  $[\theta]$  and  $[g]$  are special tokens for parallel decoding of actions.