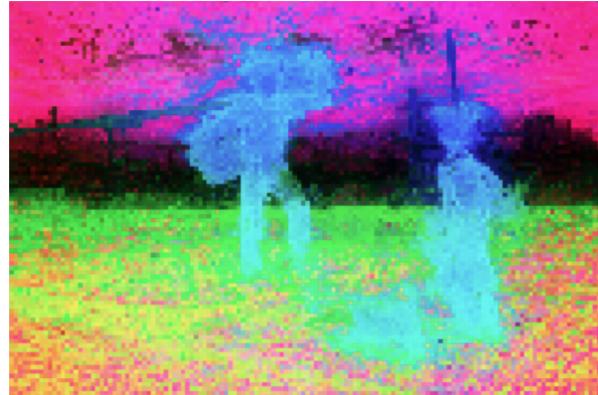


Early-fusion Grounding VLM

PCA maps: SigLIP2



DINOv3



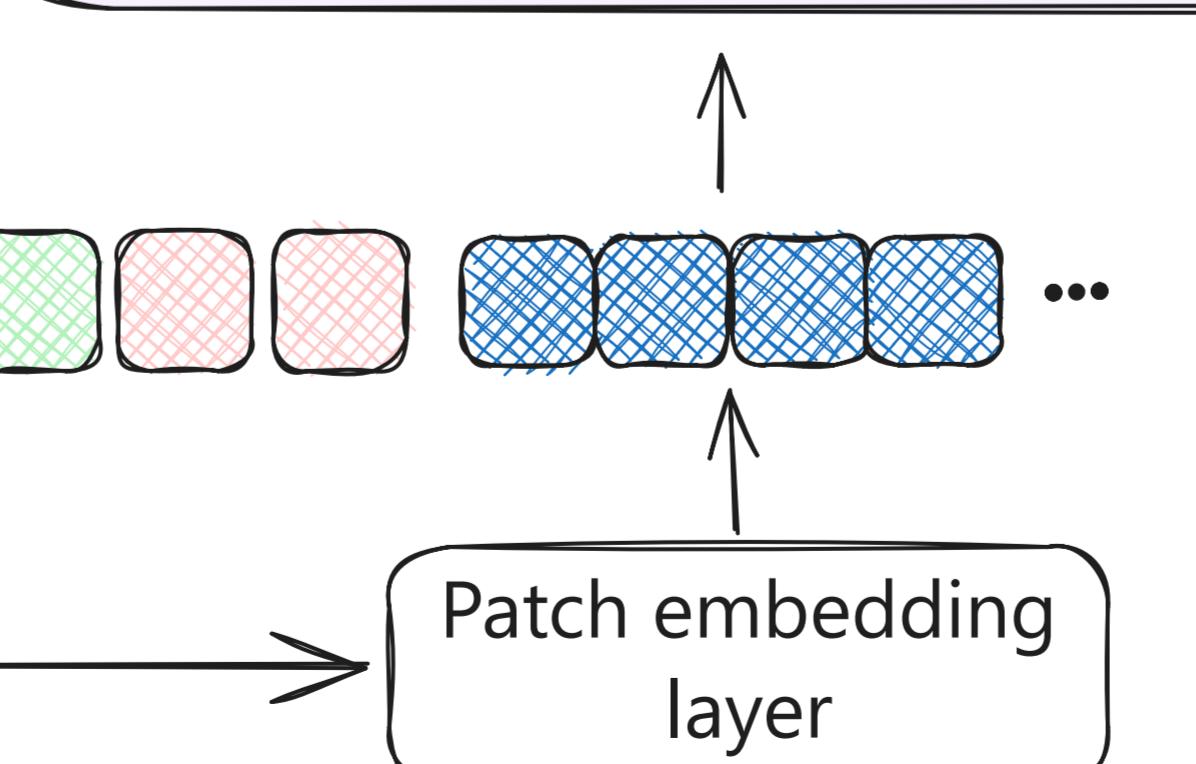
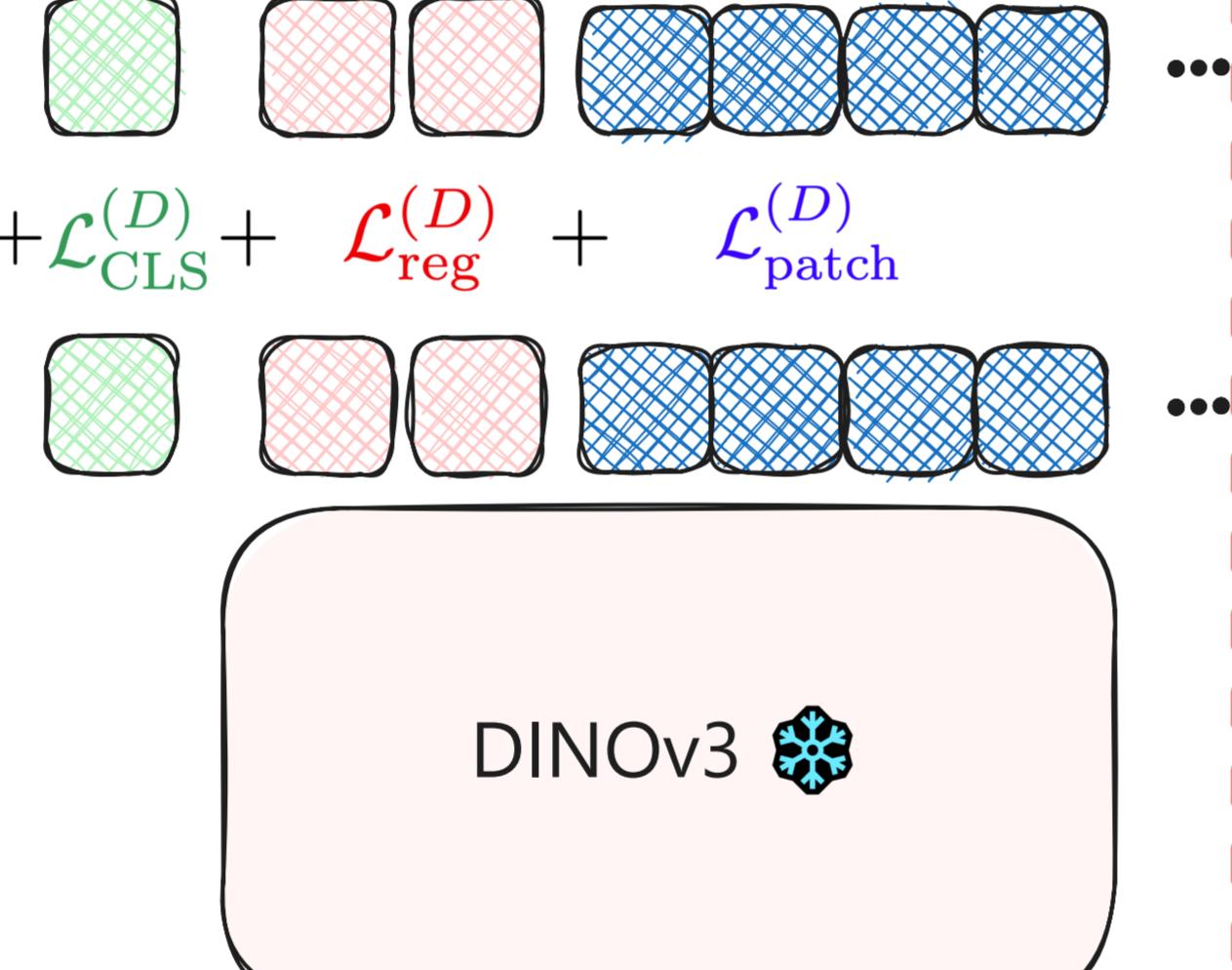
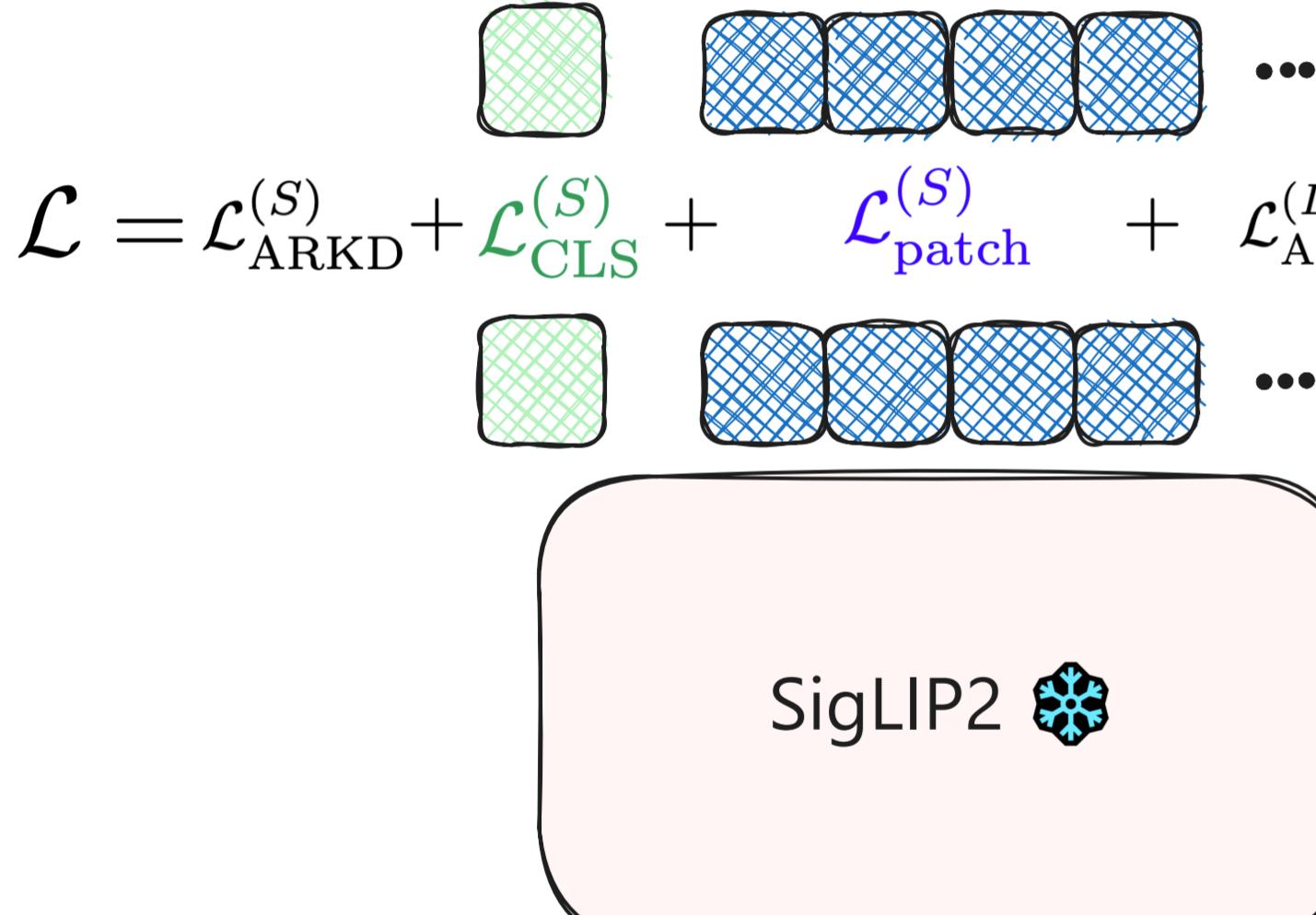
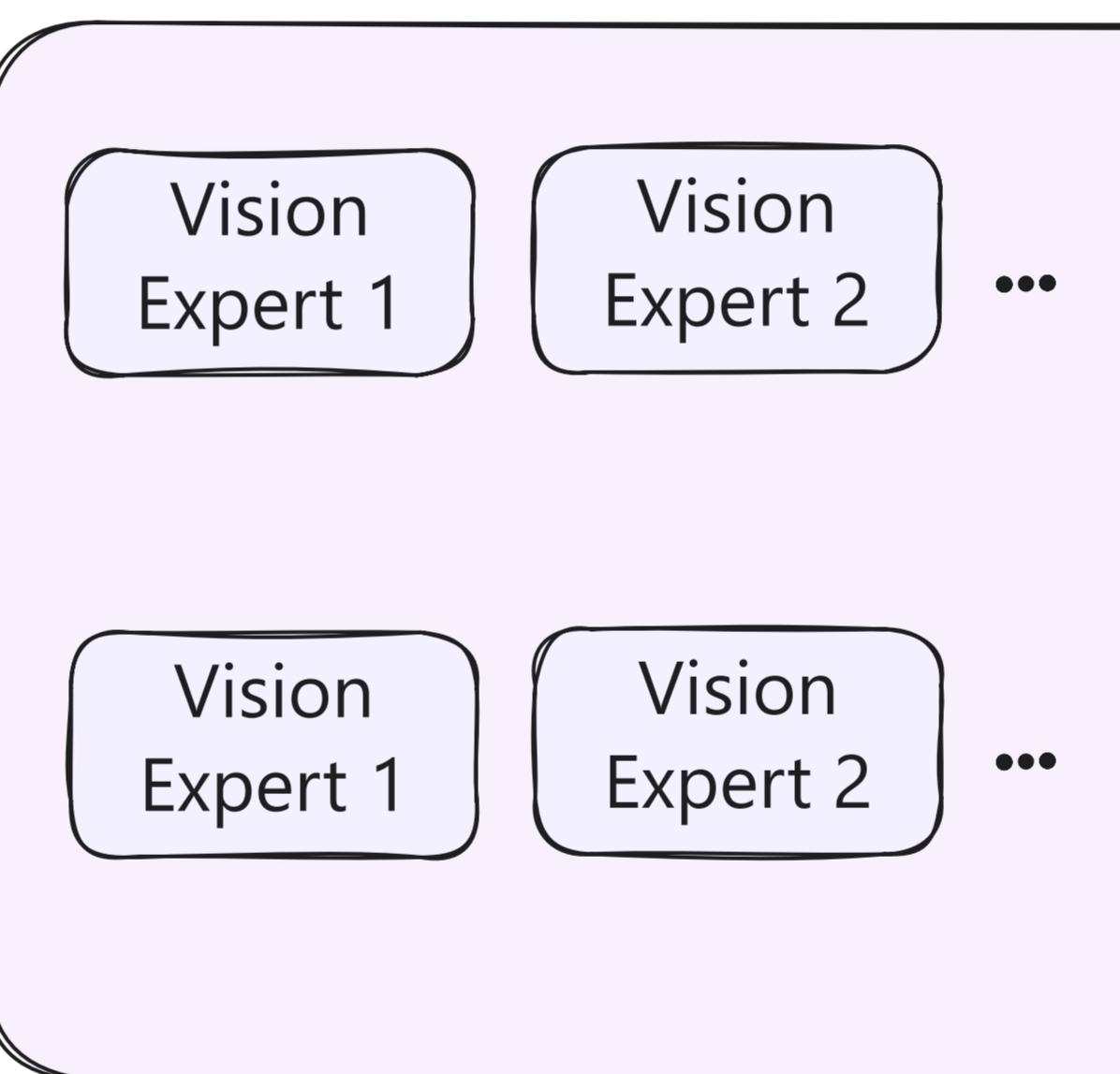
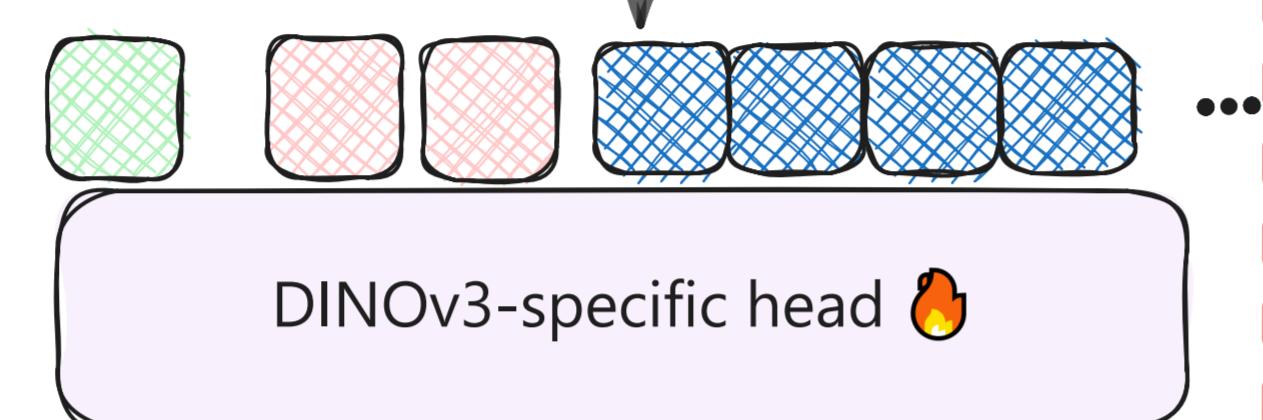
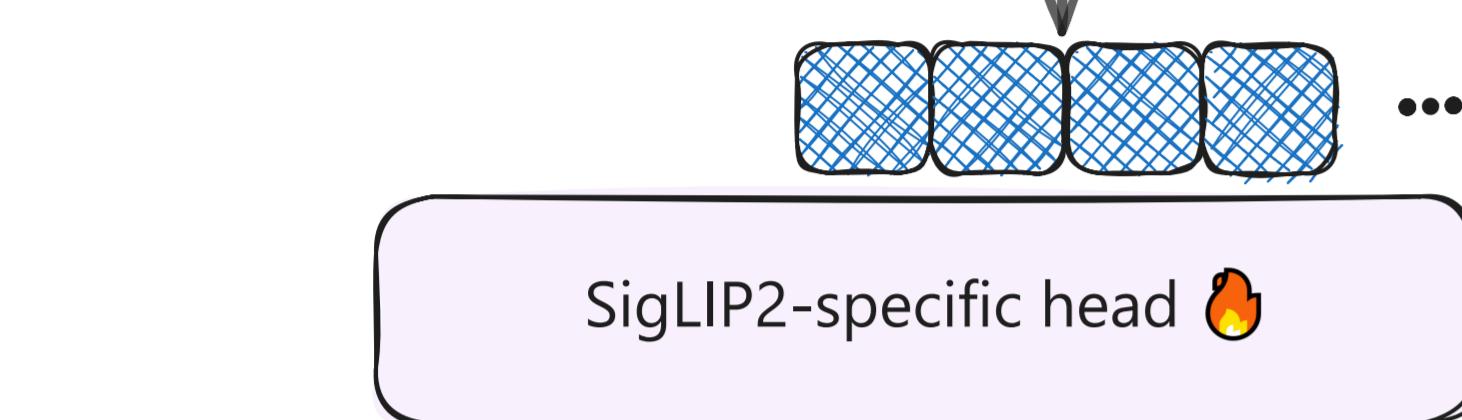
Ours



Bounding box



Segmentation mask



- [white square] Text token
- [pink square] Register token
- [green square] Global representation token
- [blue square] Patch token

Stage 0: MT-Distillation



The statue of the person holding two baskets