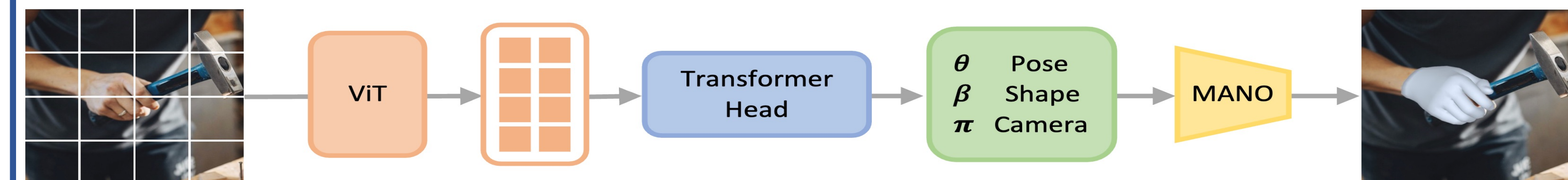




Introduction

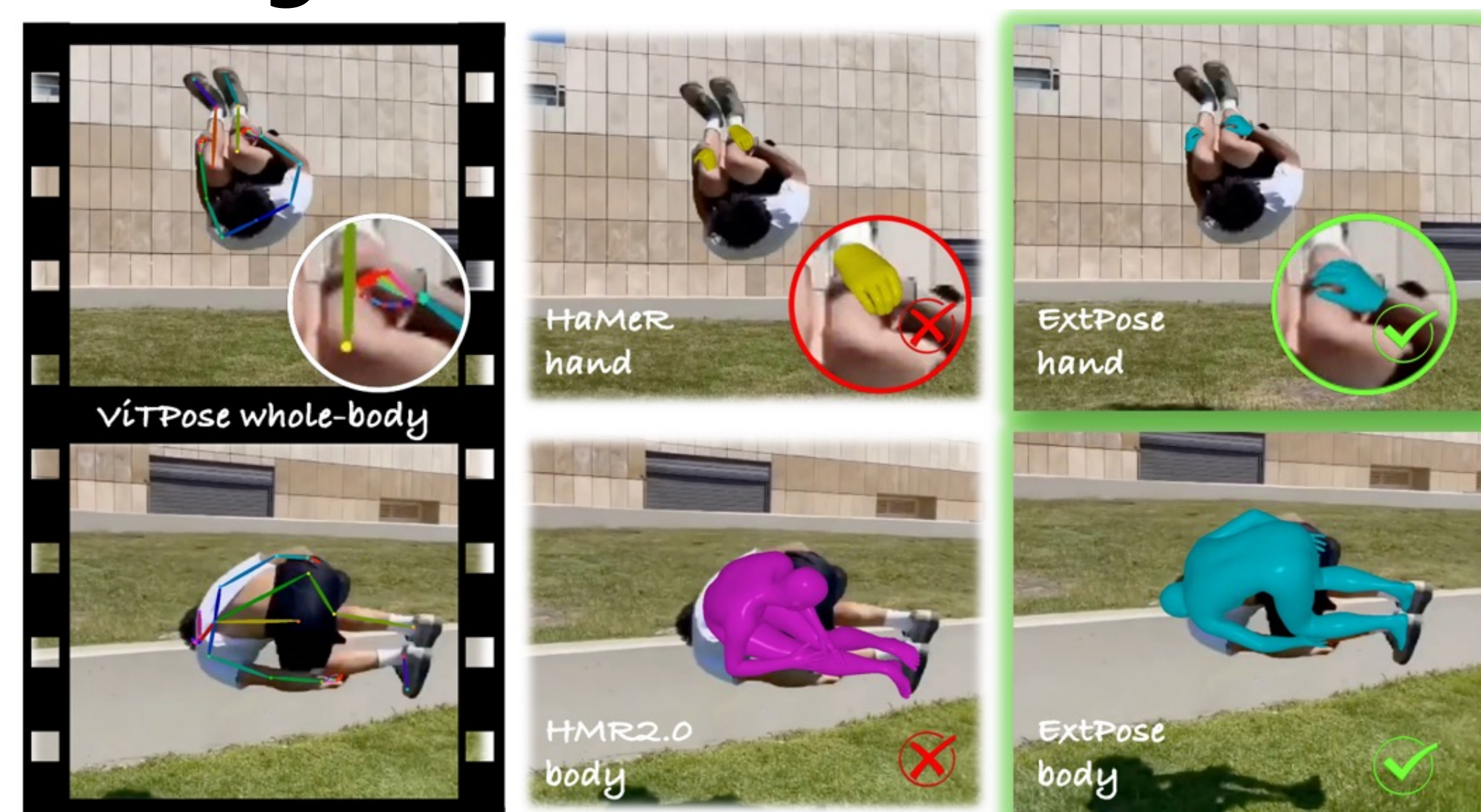
Task: 3D Human/Hand Pose Estimation (HPE)



•**Output:** **3D pose parameters** of the Human & Hand Model. Projected to 2D for visualization

•**Architecture:** **Vision Transformers (ViTs)** working on **image patches** of size 16 x 16

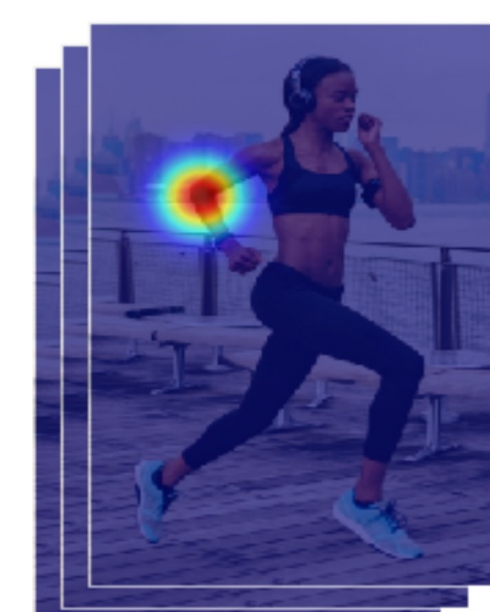
Shortcomings



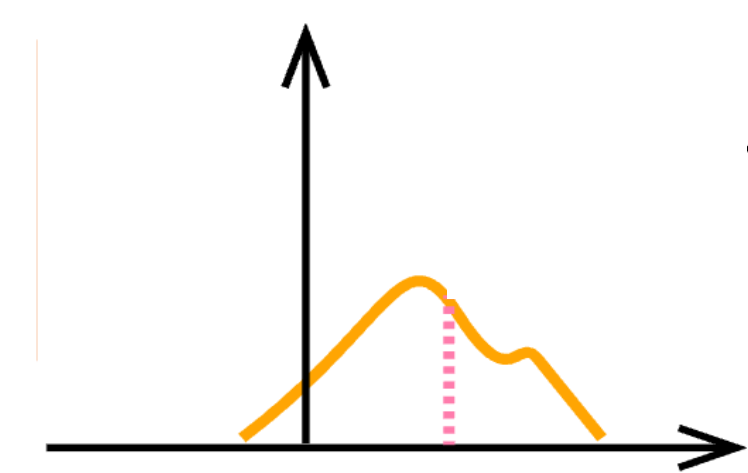
✗**Robustness:** misalignment between 3D poses& images, e.g. wrong orient. for complex motion in Col. 2 & occlusions

✗**Coherence:** the ViT itself does not consider the **temporal info** for videos, requiring an additional temporal module on top of frame features to alleviate jitter

Insights



Template Matching VS. Regression



•**Image alignment:** 2D HPE based on **template matching** is better than **regression**-based 3D HPE

•**Attention freely collects info on any relationship**, e.g. those **1.** between hands & bodies (spatial), **2.** between image & 2D pose modalities, & **3.** across frames (temporal)

Takeaways

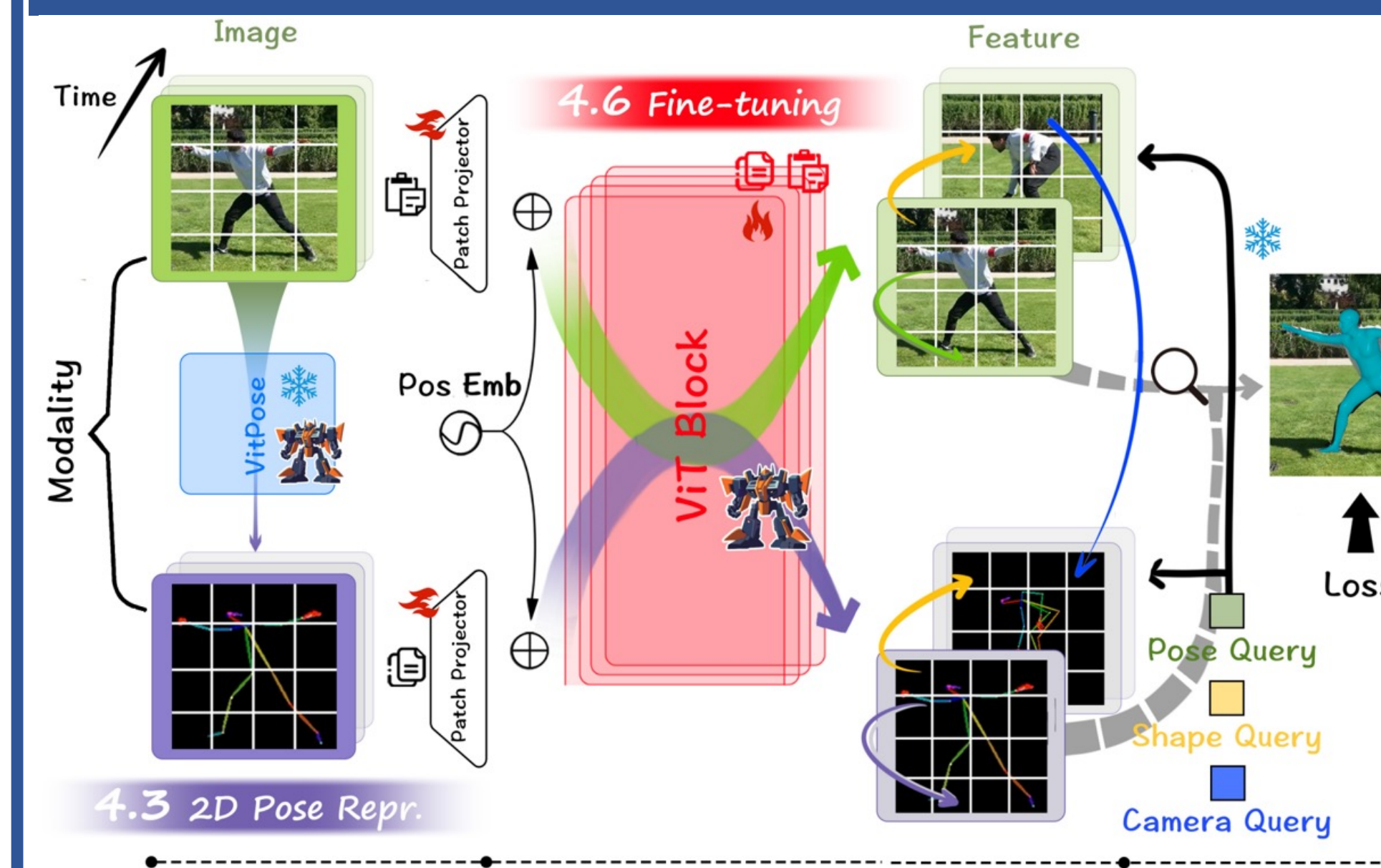
Unified HPE Framework with Modular Attention

•We extend the pose ViT into the first Video ViT pose estimator by re-programming the attention, which enhances robust & coherent features & can incorporate the info from multiple modalities, frames, & views, etc.

•**Logo:** complex Pose, attention (the cube) Extension

•**Quotes:** *You can enjoy a grander sight, By climbing to a greater height* (Tang Peoms). I.e. leverage all available info incl. that from other parts, modalities, & frames

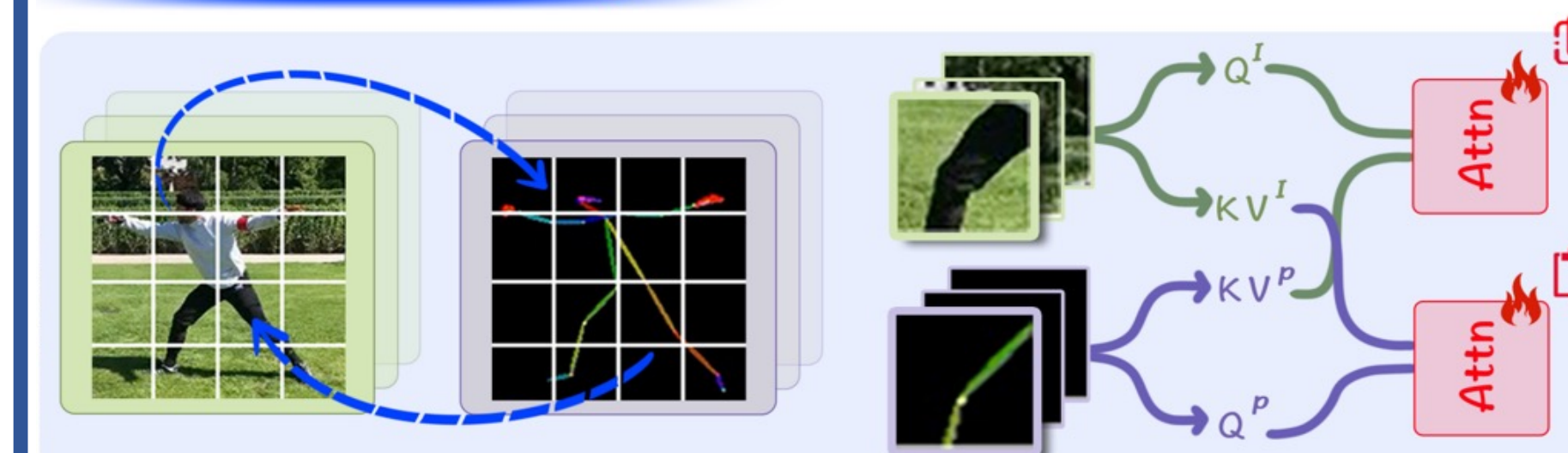
Extending Pose Vision Transformers



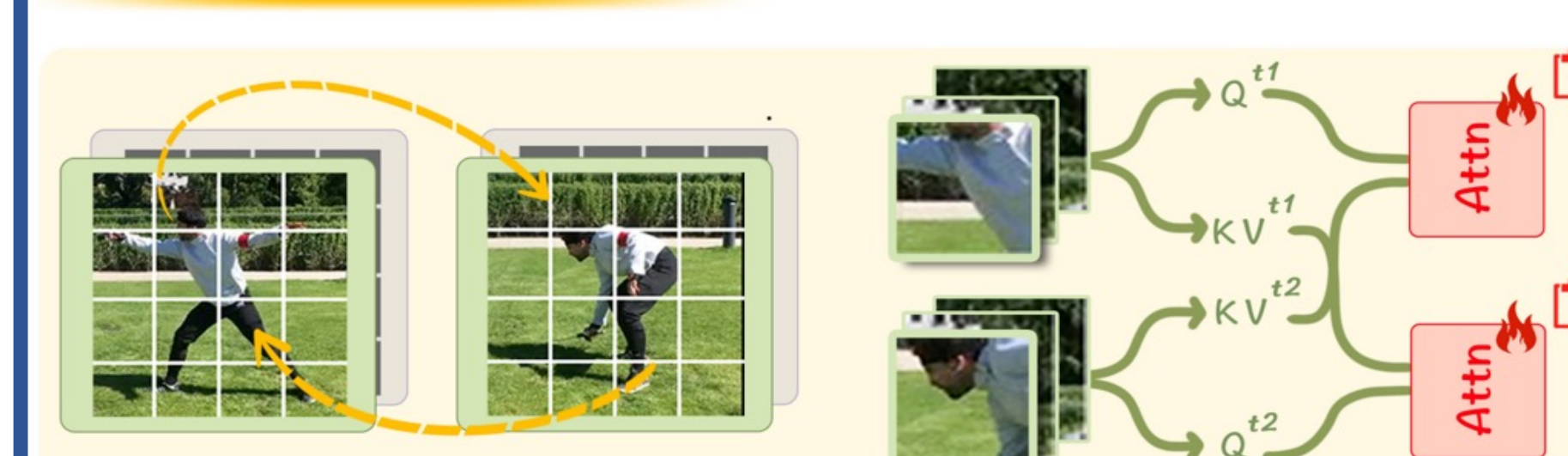
+ Multi-Modal Pose ViTs

2D pose images & RGB images can be well processed & seamlessly fused by **one shared** vanilla ViTs with the **Multi-Modal Attention**, exploiting the layout but not being misled by 2D pose errors like Concat & ControlNet

4.4 Cross Modality



4.5 Cross Frame



2D ViT Pose Form

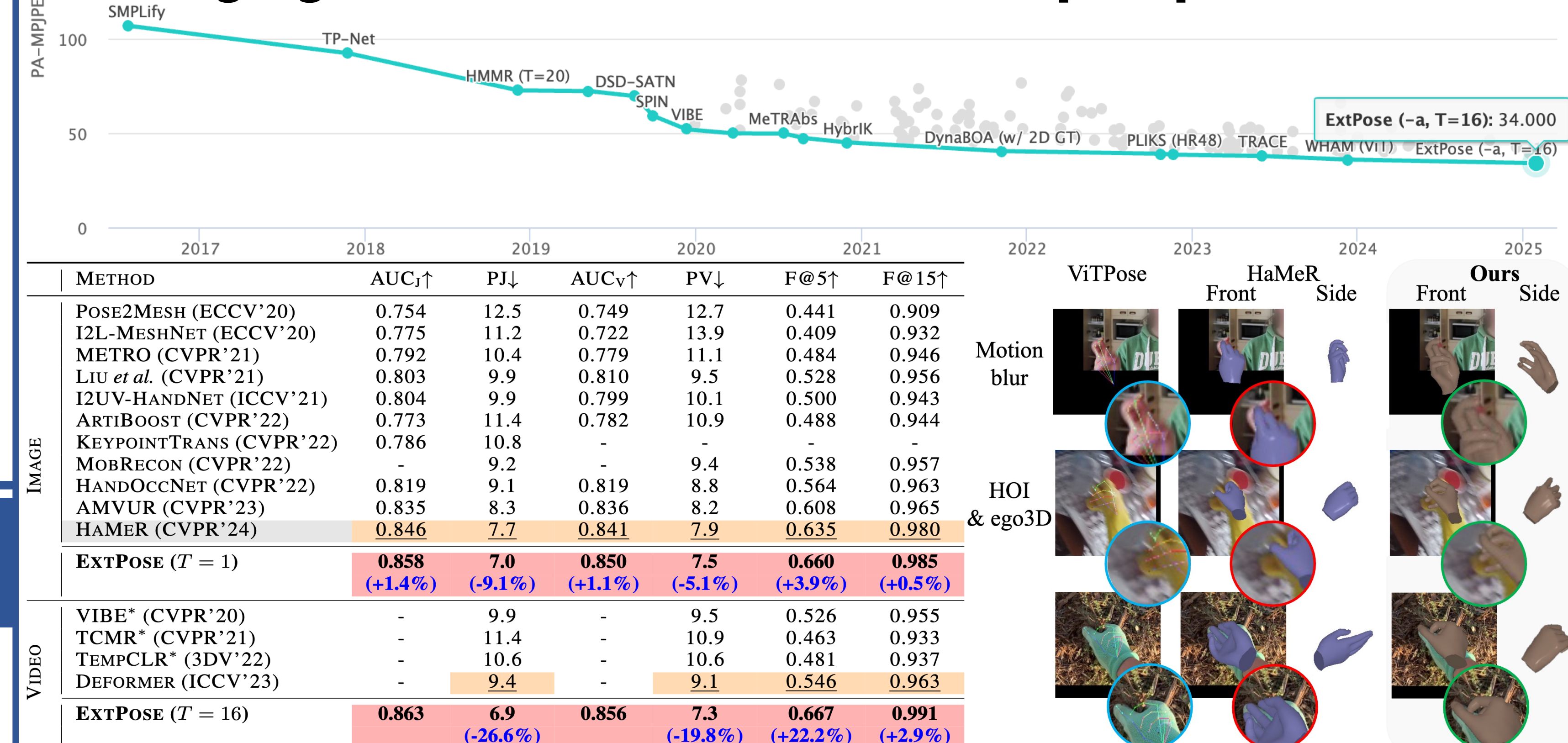
Pose images (Col. 1) instead of 1D arrays are chosen for the same spatial layout as **images** & depicting joints & human configs

Video Pose ViTs

Attending & fusing features from multiple frames at **each layer** is more effective than just fine-tuning a **temporal head** on frame features

Extensive Experiments

🏆 **SOTA on 5 human & hand datasets: 23%** accuracy (PA-MPJPE) improvement on the 3DPW, ✓robust & coherent in challenging **motion blur, occlusion, & perspective**

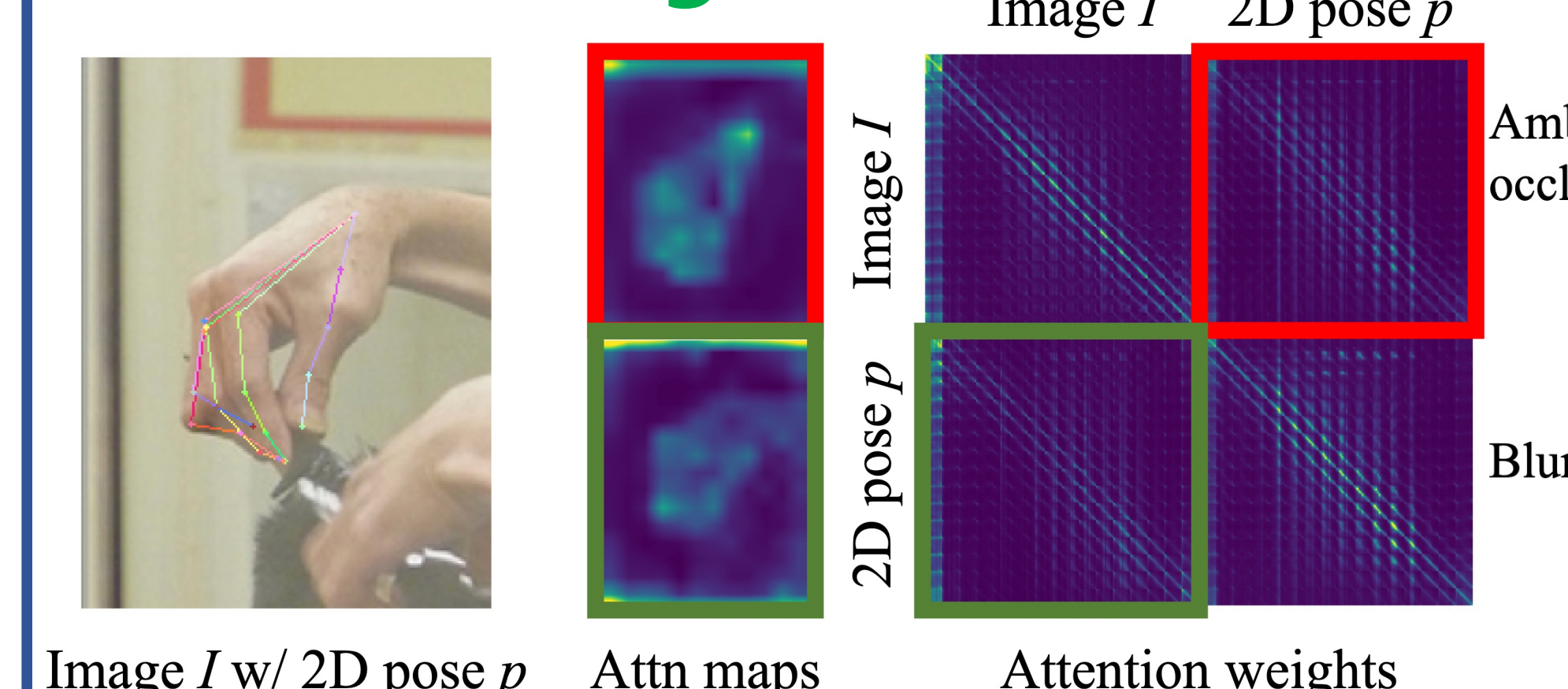


•Ablations of 2D pose forms, modality fusions, & strategies

IMG	2D POSE	PA-MPJPE↓	PA-MPVPE↓	F@5↑	F@15↑
1D	HEATMAP	6.5	6.6	0.724	0.983
	SKEL. IMAGE	6.3	6.3	0.747	0.984
	SKEL. IMAGE	6.2	6.3	0.742	0.985
✓	SKEL. IMAGE	6.0	5.7	0.783	0.991
✓	SKEL. IMAGE	4.9	5.1	0.823	0.993

METHOD	NEW DAYS	VISOR
	@0.05 @0.1 @0.15 @0.05 @0.1 @0.15	
HAMER	48.0 78.0 88.8 43.0 76.9 89.3	
FUSION		
LATE FUSION	50.5 82.4 92.5 52.5 87.1 95.6	
CHANNEL CONCAT*	56.3 83.6 92.2 55.9 87.3 95.3	
CONTROLNET*	55.6 83.5 92.3 57.7 87.5 95.5	
TRAINING		
FROM ViTPose	49.9 82.2 92.2 46.4 85.3 95.2	
ONLY Q, K	50.0 81.9 92.3 49.1 85.3 95.1	
1 st HALF	50.8 82.2 92.3 50.2 85.8 95.2	
EXTPOSE	59.6 84.8 92.7 61.1 88.5 95.6	

⚠️ **Multi-Modal Attention:** image features focus on the **keypoint** of 2D poses, while 2D pose features also examine the depth info in the **RGB background**



➡️ **Thus, robust to 2D pose errors; yet, both branches fail in extreme cases where we may seek more cues**

