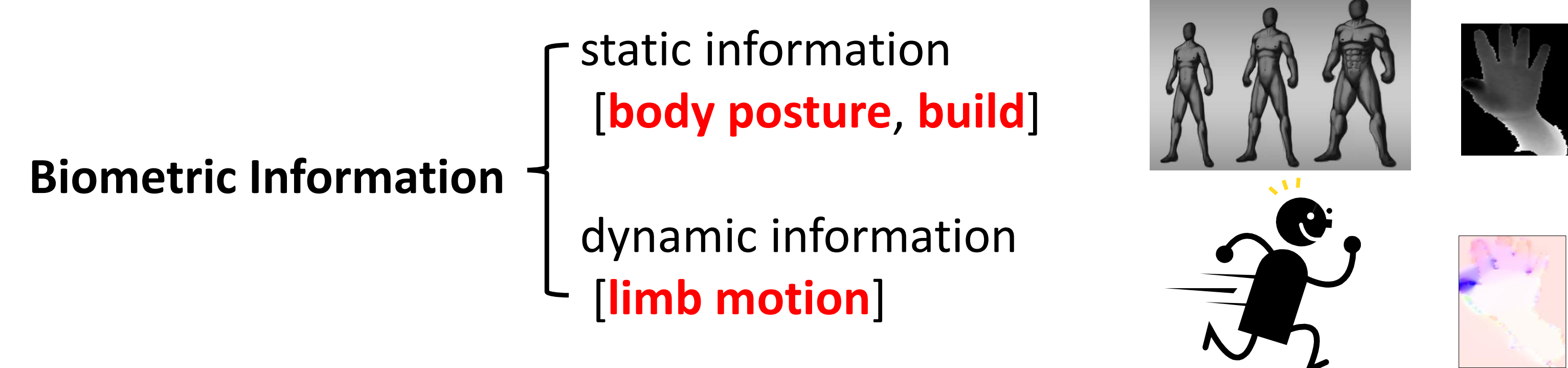


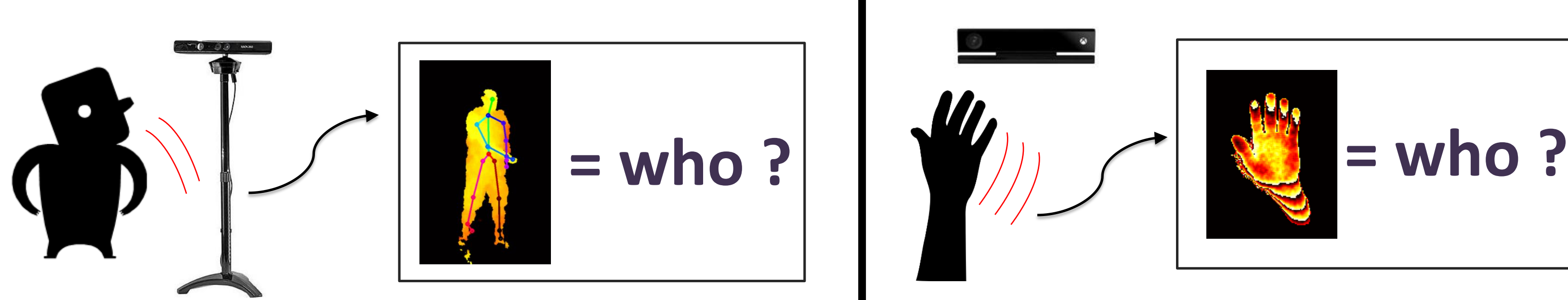


Motivation

- Use **body** or **hand** gestures to recognize users



- Can obtain gesture depth-maps with **time-of-flight** cameras (Kinect V1 and V2)



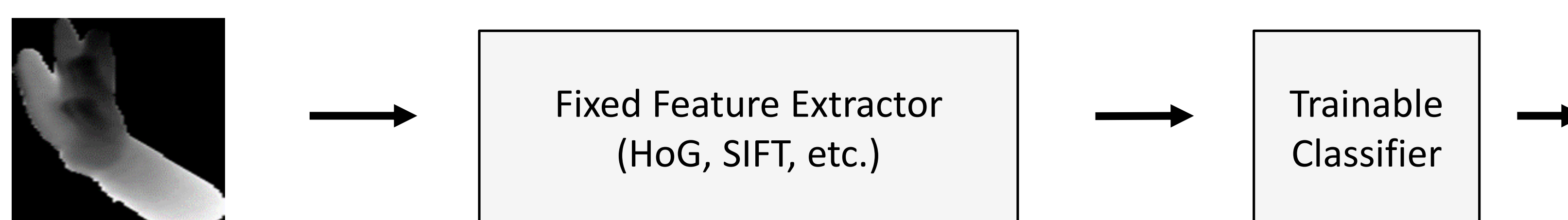
- Previous works:** compared gesture representations, assessed the value of multiple views, studied spoof attacks

- Focus of this study:**

- Use two-stream convolutional networks to recognize users
- Evaluate gesture generalization performance: **learn** user style
- Visualize deep gesture features with t-SNE

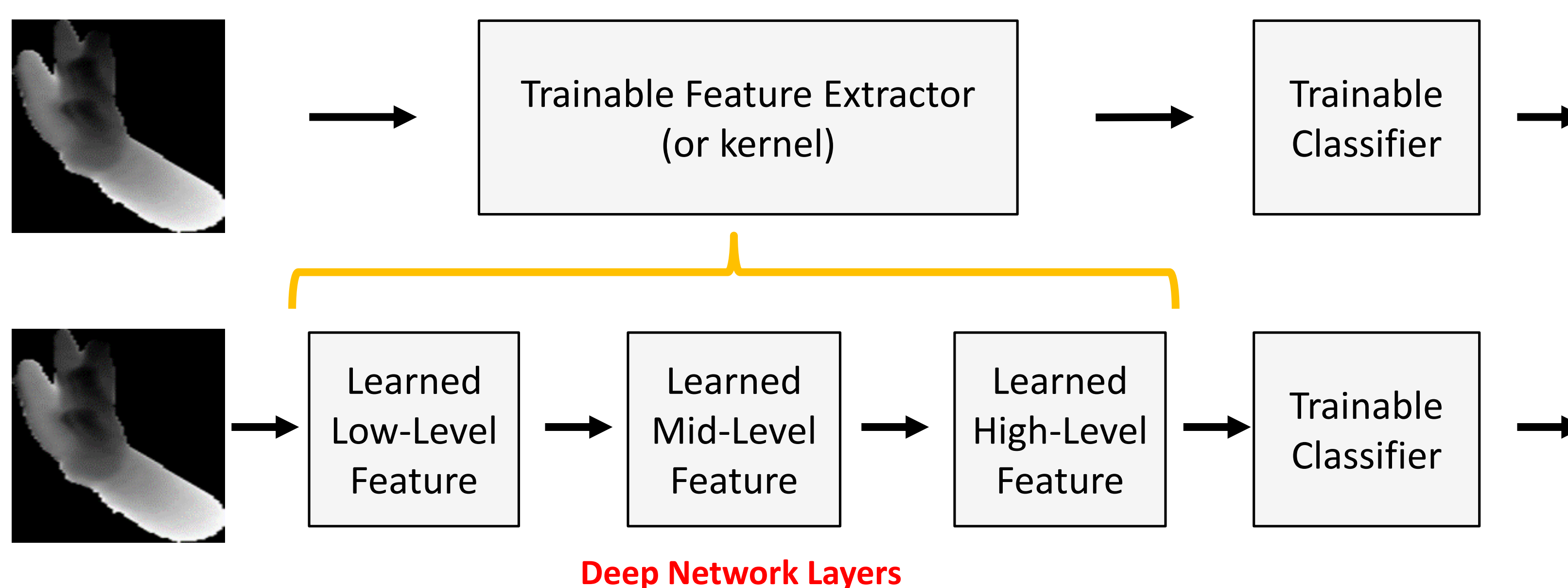
Deep Learning Recap

Traditional Learning Pipeline



Deep Learning Pipeline

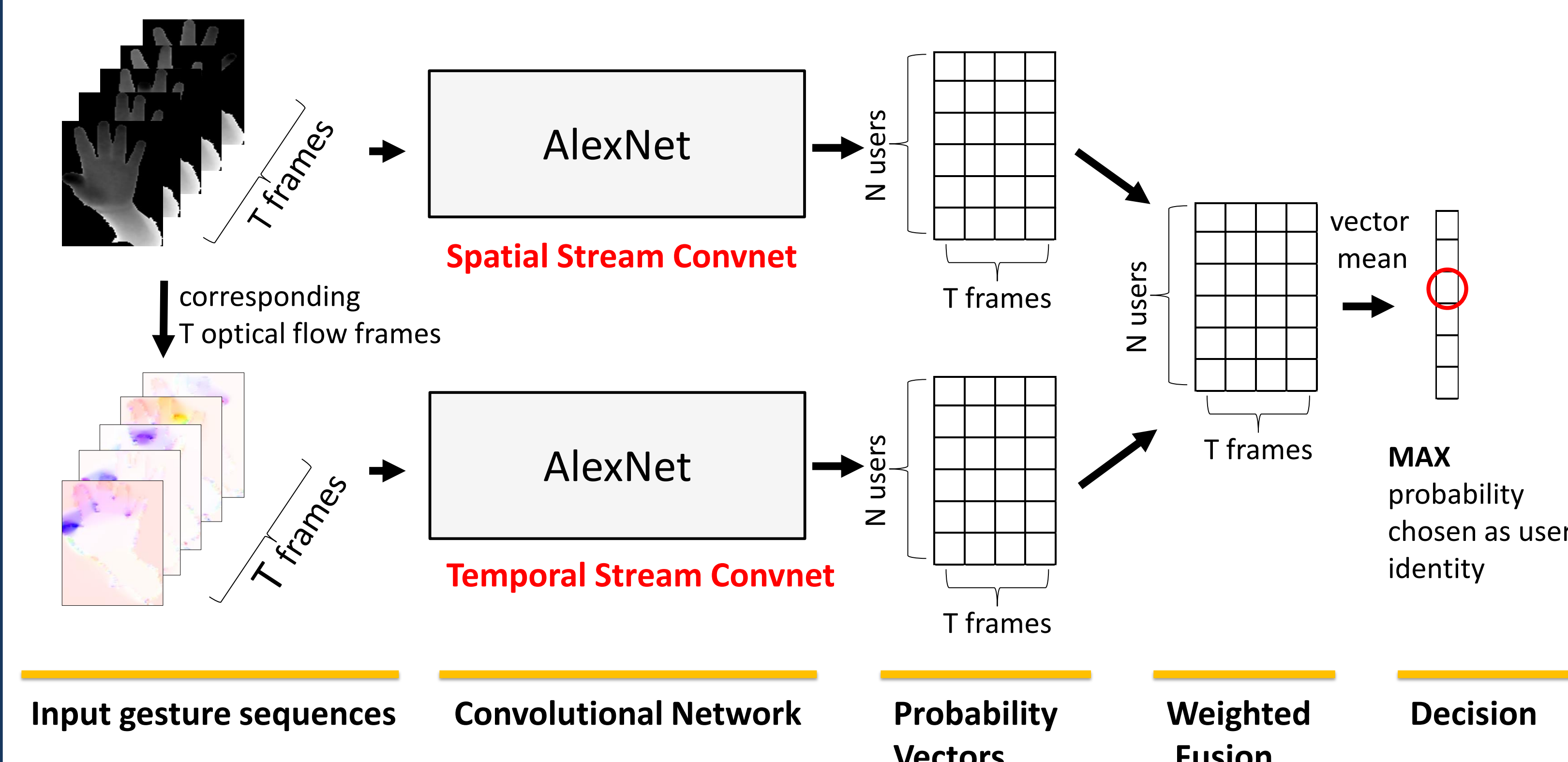
- Learn feature representation directly from image (end-to-end learning)
- Hidden "weight" layers are a composition of non-linear transformations



Two-Stream Convolutional Networks

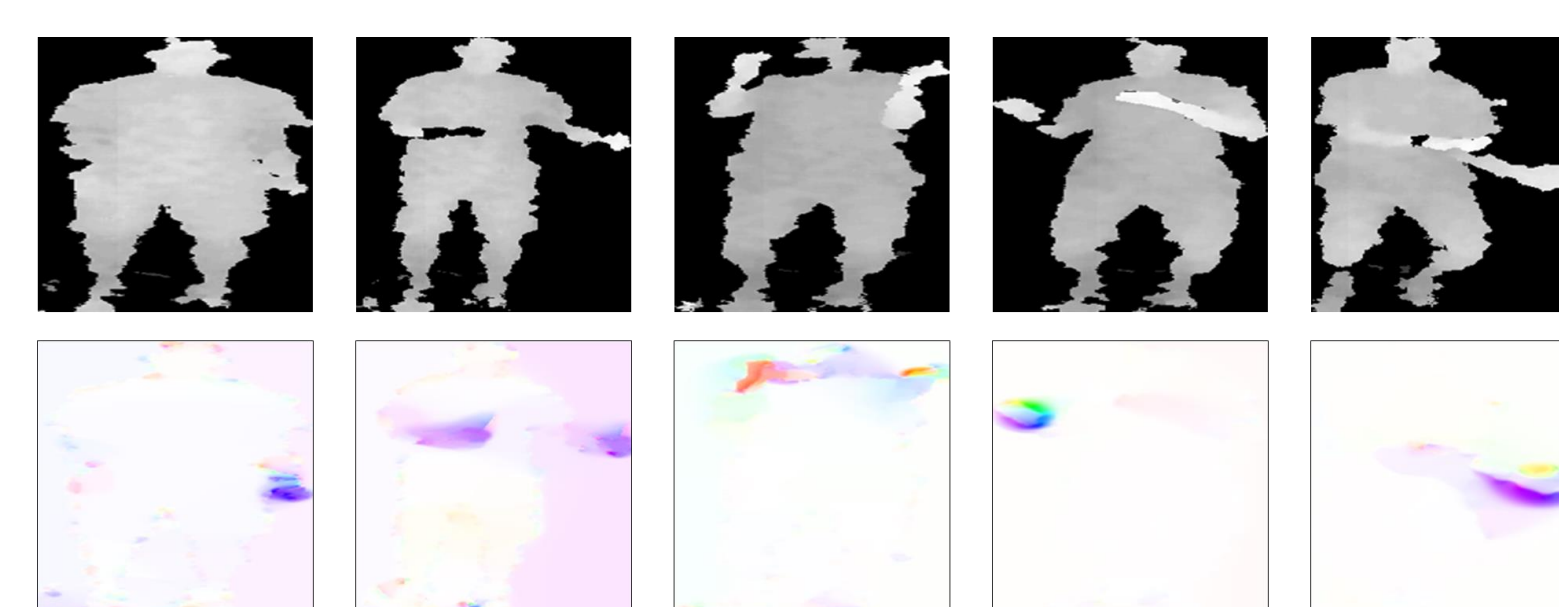
Method: Adapt a two-stream convolutional network [1] for user **identification**

- Leverages static and dynamic information of a gesture
- Learns two separate image-based convolutional networks
- AlexNet [2] used as network of choice (5 conv layers, 3 fc layers)
- Pre-trained from ImageNet [3], then fine-tuned

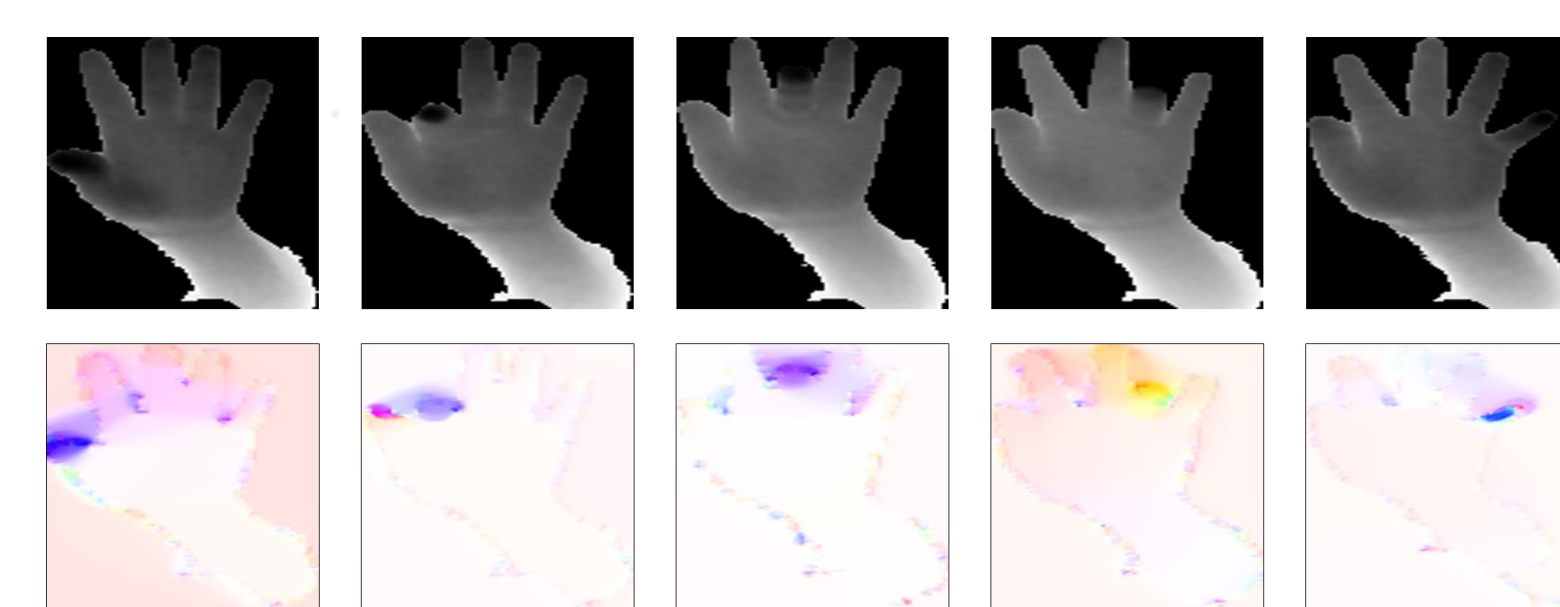


Public Gesture Datasets

Body Gesture Dataset (BodyLogin):
40 users, 5 gestures (1 user-defined)

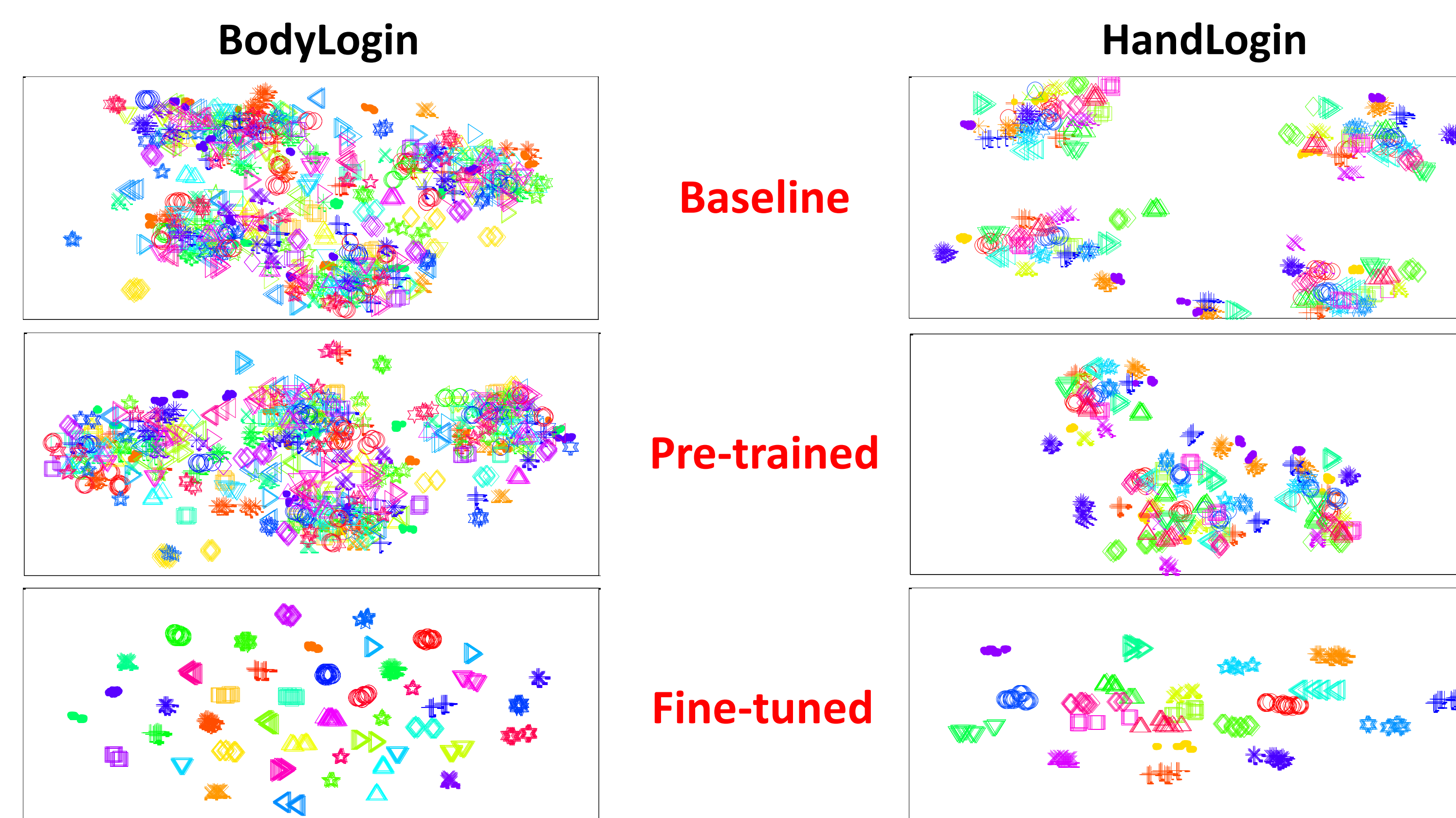


Hand Gesture Dataset (HandLogin):
21 users, 4 gestures



Feature Visualization

- Marker Point ↔ Gesture, Marker Color and Shape ↔ user identity
- t-SNE [4] shows strong user separation after fine-tuning



User Identification Experiments

Evaluate **Correct Classification Error (CCE = 100% - CCR)** for various scenarios:

- Training and testing with all gestures

Dataset	Spatial		Temporal		Baseline*
	(1,0)	(0.66,0.33)	(0.5,0.5)	(0.33,0.66)	
HandLogin	0.24%	0.24%	0.24%	0.71%	6.43%
BodyLogin	0.05%	0.05%	0.05%	0.05%	1.15%

- Testing with gestures unseen in training (left-out), evaluates generalization performance

Generalizing Gesture	Spatial		Temporal		Baseline*
	(1,0)	(0.66,0.33)	(0.5,0.5)	(0.33,0.66)	
HandLogin Compass	2.38%	2.86%	4.76%	8.57%	82.38%
HandLogin Piano	1.91%	0.48%	1.43%	1.91%	68.10%
HandLogin Push	44.29%	49.05%	54.29%	67.62%	79.52%
HandLogin Fist	16.67%	15.71%	17.14%	20.00%	72.38%
BodyLogin S motion	0.75%	1.00%	1.25%	1.75%	75.75%
BodyLogin Left-Right	0.88%	1.25%	1.50%	1.88%	80.88%
BodyLogin 2-Hand Arch	0.13%	0.13%	0.13%	0.38%	74.50%
BodyLogin Balancing	9.26%	10.01%	13.27%	19.52%	77.97%
BodyLogin User Defined	5.28%	5.53%	6.16%	8.54%	71.61%

- Suppression of dynamics (use first few frames only in training and testing)

Dataset	Scenario	Spatial Stream	Baseline
HandLogin	All frames	0.24%	6.43%
	No dynamics	1.90% ↑	9.29% ↑
BodyLogin	All frames	0.05%	1.15%
	No dynamics	1.00% ↑	32.60% ↑

*Baseline: temporal hierarchy of depth-aware silhouette tunnels [5]

Conclusions

- Quite feasible to learn a user's gesture style from a bank of gesture types
- Possible to generalize user style to similar gestures with only slight degradation in performance
- Convolutional networks offer drastic improvements over state-of-the-art
- Temporal/dynamic information always valuable

- Check out paper:** additional experiments and analysis
- Data and trained models **available** online [see QR-code] (<http://vip.bu.edu/projects/hcis/deep-login>)