



ICCV 2019
Seoul, Korea

Generative Multi-View Human Action Recognition

Lichen Wang¹, Zhengming Ding², Zhiqiang Tao¹, Yunyu Liu¹, and Yun Fu¹

wanglichenxj@gmail.com

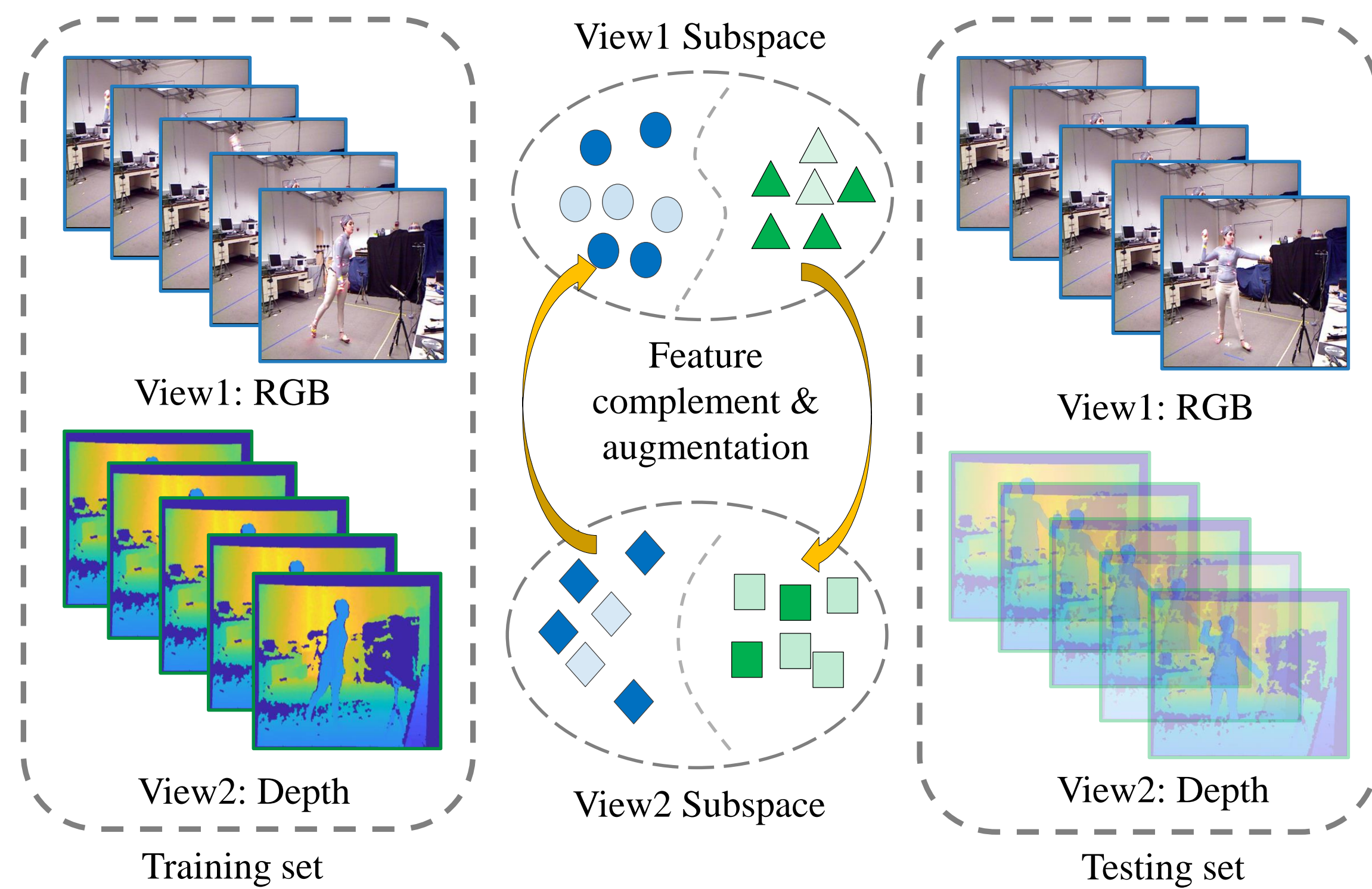
¹Northeastern University, ²Indiana University-Purdue University Indianapolis



Introduction

❖ Multi-view Action Recognition

- **Input:** Multi-view action sequence (e.g., RGB + Depth)
- **Output:** Action prediction results



Concept of multi-view human action recognition

❖ Challenges

- Heterogeneous (significantly different) multi-view feature domains (e.g., RGB and depth, RGB and electronic signal)
- Incomplete/missing view sequences (e.g., only one view is available)
- Inconsistent view-specific prediction (e.g., RGB and depth have different prediction results)

❖ Motivations

- Obtain more distinctive feature representations
- Explore cross-view feature relations
- Explore the multi-view prediction results in high-level label space

Our model

❖ Three modules are proposed for solving the challenges

- **View-specific encoders**
 - Seek distinctive action representations in subspaces
 - Label information + triplet loss objective
- **Cross-view Adversarial Generation**
 - Increase cross-view representation diversity
 - Enhance model robustness
 - Address missing/incomplete view sequences
- **View Correlation Discovery Network**
 - View-specific initial classification is firstly obtained
 - Pair-wise label correlation matrix is generated
 - VCDN fully explore the latent high-level label correlation for higher performance

$$L_{Em} = \sum_{i=1}^M \max \left(\left\| E_m(X_{tr_i}^a) - E_m(X_{tr_i}^p) \right\|_2^2 - \left\| E_m(X_{tr_i}^a) - E_m(X_{tr_i}^n) \right\|_2^2 + \alpha \right)$$

$$L_{G_1d} = -E_{z \sim p_z(z)} \log \left(1 - D_1(G_1(z|E_1(X_{tr}^1))) \right)$$

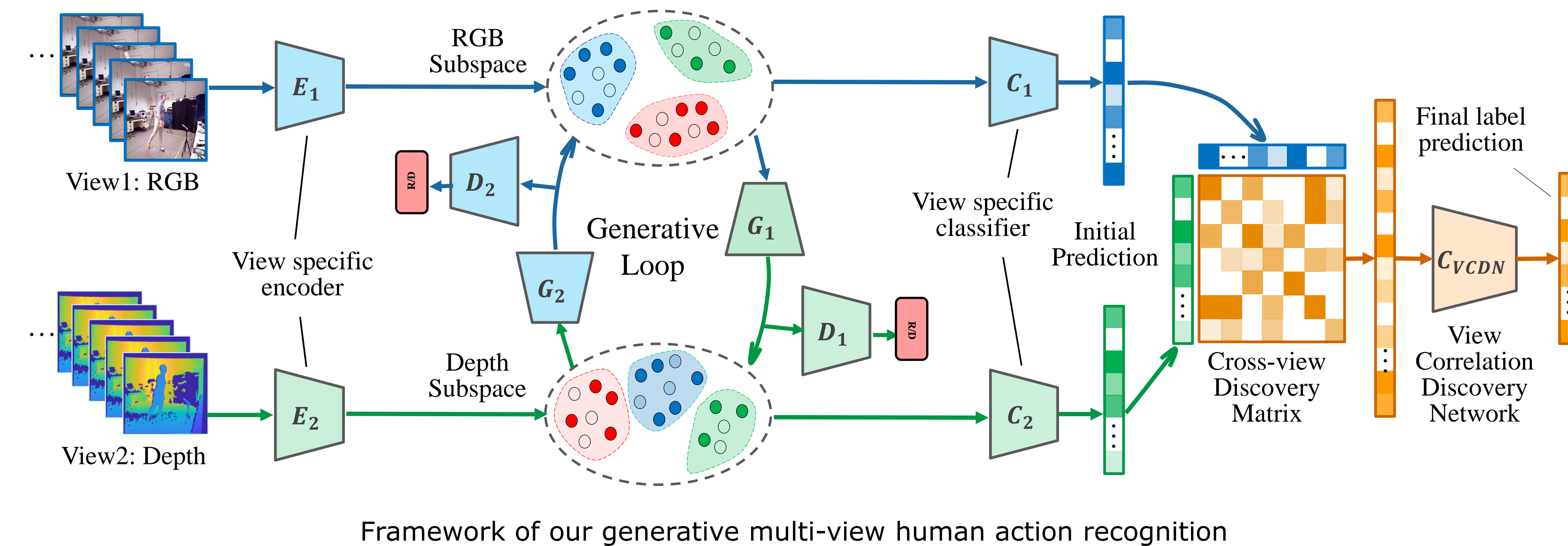
$$L_{G_1s} = E_{z \sim p_z(z)} \left(\left\| G_1(z|E_1(X_{tr}^1)) - E_2(X_{tr}^2) \right\|_F^2 \right)$$

$$L_{D_1} = E_{X \sim p_X(X)} \log D_1(E_2(X_{tr}^2)) + E_{z \sim p_z(z)} \log \left(1 - D_1(G_1(z|E_1(X_{tr}^1))) \right)$$

$$L_{C_1g} = \|Y_{tr} - C_1(G_2(z|E_2(X_{tr}^2)))\|_F^2$$

$$L_{C_2g} = \|Y_{tr} - C_2(G_1(z|E_1(X_{tr}^1)))\|_F^2$$

$$L_{C_{VCDN}} = \sum_{i=1}^{n_{tr}} \|y_i - C_{VCDN}(y_{tr_i}^2 \cdot y_{tr_i}^{1\top})\|_2^2$$



Framework of our generative multi-view human action recognition

❖ Single-view & multi-view action classification

Method	RGB	R→D	Depth	D→R	R+D	Method	RGB	R→D	Depth	D→R	R+D	Method	RGB	R→D	Depth	D→R	R+D
LSR	67.59	69.17	45.45	37.73	68.77	LSR	96.46	97.17	47.63	42.51	97.17	LSR	65.02	65.43	82.30	48.56	77.36
SVM [36]	69.44	68.53	34.92	34.33	72.72	SVM [36]	96.09	96.80	45.39	45.13	96.80	SVM [36]	66.11	70.24	78.92	78.18	83.47
VLAD [14]	71.54	-	-	-	-	VLAD [14]	97.17	-	-	-	-	VLAD [14]	67.13	-	-	-	-
TSN [51]	71.01	-	-	-	-	TSN [51]	97.31	-	-	-	-	TSN [51]	67.85	-	-	-	-
WDMM [1]	-	-	46.58	-	-	WDMM [1]	-	-	66.41	-	-	WDMM [1]	-	-	81.05	-	-
AMGL [30]	69.17	71.54	39.92	35.96	68.53	AMGL [30]	96.46	97.11	30.03	29.96	94.70	AMGL [30]	64.61	59.05	72.84	67.33	74.89
MLAN [29]	67.19	67.19	33.28	33.61	66.64	MLAN [29]	96.05	96.10	41.48	41.25	96.46	MLAN [29]	67.91	67.91	72.96	72.83	76.13
PM-GANs [49]	-	71.36	-	49.01	-	PM-GANs [49]	-	96.76	-	66.84	-	PM-GANs [49]	-	68.72	-	76.02	-
Ours	-	73.53	-	50.35	76.28	Ours	-	98.23	-	68.32	98.94	Ours	-	69.72	-	83.48	88.72

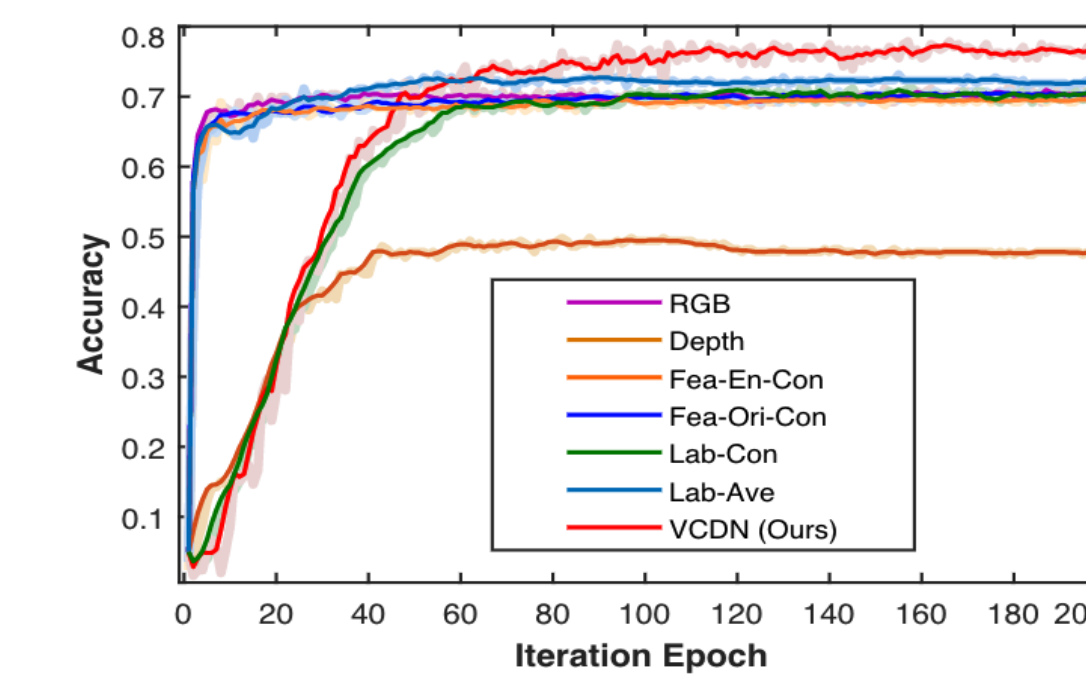
UWA dataset

MHAD dataset

DHA dataset

❖ Ablation study

- Intentionally remove or modify the View Correlation Discovery Network



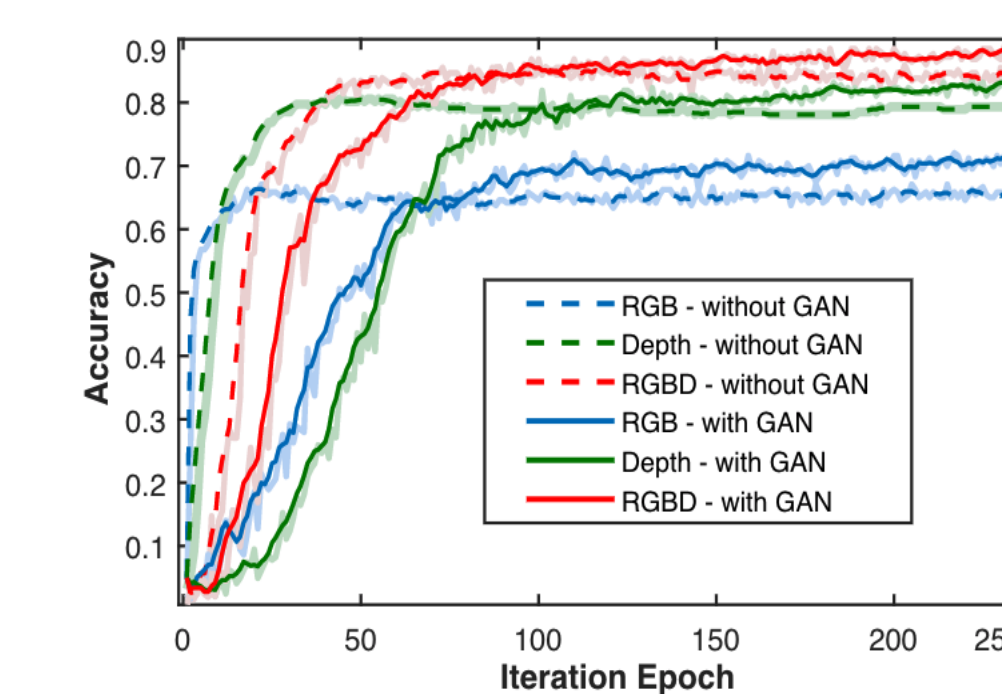
Performance with different label fusion strategies

Setting	UWA	MHAD	DHA
RGB- C_1	69.18	96.42	68.15
Depth- C_2	45.28	63.05	79.79
RGBD-Fea-En-Con	68.78	96.82	70.85
RGBD-Fea-Ori-Con	69.22	97.32	70.83
RGBD-Lab-Con	70.38	96.28	80.95
RGBD-Lab-Ave	71.84	97.56	83.28
RGBD-Lab-Wei	71.15	97.17	83.95
RGBD-VCDN (Ours)	74.07	98.06	84.32

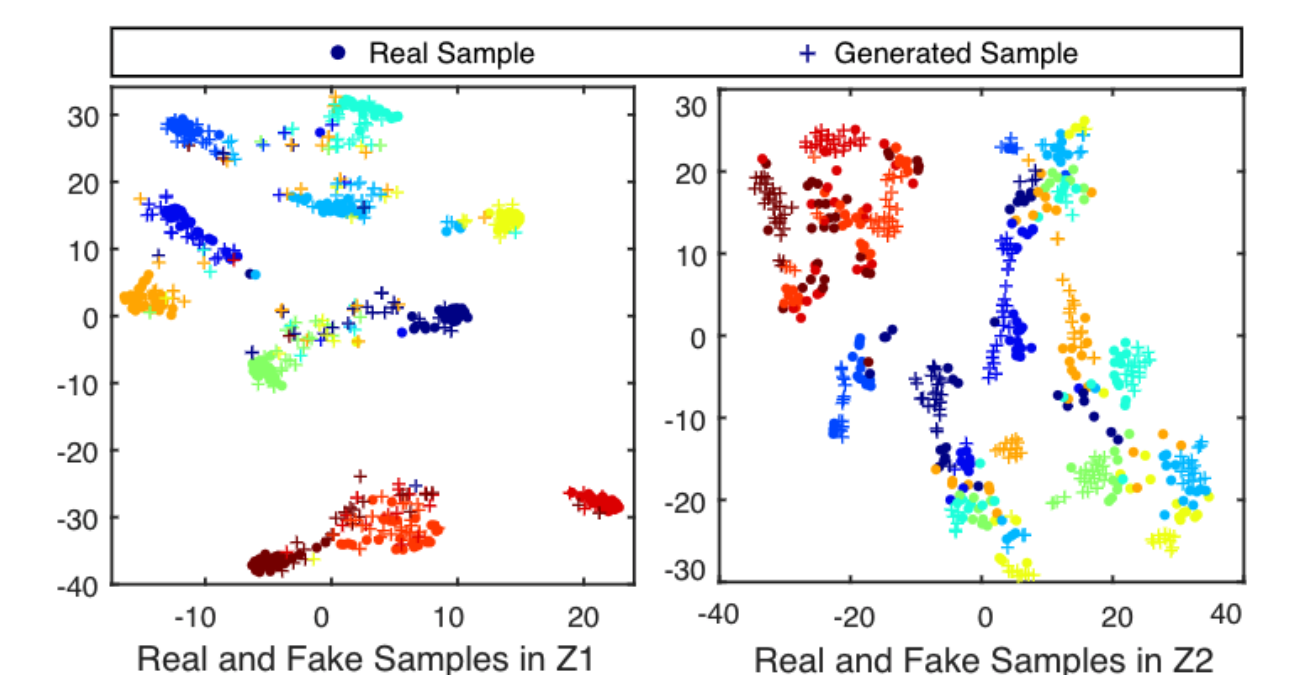
Performance with different label fusion strategies

Performance with different label fusion strategies

- The performance with and without the generative model



Performance with and without the GAN module



t-SNE visualization of real and generated samples

Experiments

- Conventional multi-view action recognition setting.
- Single-view action recognition setting, where another view is considered as missing view which is generated by cross-view generation strategy.
- Ablation study for cross-view generation module and view correlation discovery module

Conclusion

We proposed three modules to address the challenges of multi-view human action recognition. View-specific encoder learns distinctive action representations. Cross-view generation extend the representation distributions. View correlation discovery network explore the high-level correlations in label space. All modules are trained simultaneously to achieve the best performance.