# Unsupervised Video Domain Adaptation: A Disentanglement Perspective

## Motivation & Contribution

#### TL;DR

- We introduce the Transfer Sequential VAE (TransVAE) framework to tackle unsupervised video domain adaptation tasks from a disentanglement view.
- Our key idea is to handle the spatial and temporal domain divergence separately through disentanglement; we pursue an explicit decoupling of the domain-specific info from other info via generative modeling.



○ Source Domain ○ Target Domain ▲ Source-Related Info ■ Target-Related Info ● Semantic-Related Info

Our approach achieves new state-of-the-art performance on the domain adaptation benchmarks of UCF-HMDB, Jester, Epic-Kitchens, and Sprites.

#### **Domain Generation**

We consider the generation process of cross-domain videos from two sets of latent factors: one set consists of a sequence of random variables, which are dynamic and inclined to encode the semantic information; another set is static and introduces some domain-related spatial information.



- The graphical illustrations of the proposed generative (left) and inference (right) models for video domain disentanglement are shown above.
- The blue and red **nodes** denote the **observed** source domain and target domain **videos**  $\mathbf{x}^{S}$  and  $\mathbf{x}^{T}$ , respectively, over *t* timestamps.
- ✤ The static latent variables  $\mathbf{z}_d^S$  and  $\mathbf{z}_d^T$  follow a joint distribution and combining either of them with dynamic latent variables  $\mathbf{z}_t$  constructs one video data of a domain, with which we developed TransVAE.

# **ByteDance**



Pengfei Wei, Lingdong Kong, Xinghua Qu, Yi Ren, Zhiqiang Xu, Jing Jiang, Xiang Yin

## Methodology

#### The TranSVAE Framework

- The input videos are fed into an encoder to extract the visual features, followed by an LSTM to explore the temporal information.
- ✤ Two groups of mean and variance networks are then applied to model the posterior of the latent factors, i.e.,  $q(\mathbf{z}_t^{\mathcal{D}} | \mathbf{x}_{< t}^{\mathcal{D}})$  and  $q(\mathbf{z}_d^{\mathcal{D}} | \mathbf{x}_{< 1:T}^{\mathcal{D}})$ .



✤ The new representations  $\mathbf{z}_1^{\mathcal{D}}$ ,  $\mathbf{z}_2^{\mathcal{D}}$ , ...,  $\mathbf{z}_T^{\mathcal{D}}$  and  $\mathbf{z}_d^{\mathcal{D}}$  are sampled, and then concatenated and passed to a **decoder** for reconstruction. Four constraints are then proposed to regulate the latent factors for adaptation purposes.

#### Domain Disentanglement

We show domain disentanglement and transfer examples using Sprites as follows.

Input Sequence <ul> <li>Left: "Human", slash, x<sup>P1</sup></li> <li>Right: "Alien", walk, x<sup>P2</sup></li> </ul>	X							<b>A</b>
Reconstruction         • Left: "Human", slash, $\tilde{\mathbf{x}}^{\mathbf{P}_1}$ • Right: "Alien", walk, $\tilde{\mathbf{x}}^{\mathbf{P}_2}$								<b>Ö</b>
Domain Disentanglement <ul> <li>Left: "Null", slash</li> <li>Right: "Null", walk</li> </ul>		<u>.</u>	• <u>\$</u> .					
Domain Transfer <ul> <li>Left: "Alien", slash</li> <li>Right: "Human", walk</li> </ul>		<u>.</u>	. 🤼			<b>*</b>		*
Input Sequence <ul> <li>Left: "Human", spellcard, x<sup>P1</sup></li> <li>Right: "Alien", slash, x<sup>P2</sup></li> </ul>	§ A	*	* 😪 🤱	ŵ ¢		â V	¥ ¥	â
Reconstruction         • Left: "Human", spellcard, $\tilde{x}^{P_1}$ • Right: "Alien", slash, $\tilde{x}^{P_2}$	â a	*	* 😤 🏯		) <u>â</u>	â â	À À	â
Domain Disentanglement         • Left: "Null", spellcard         • Right: "Null", slash			* <u>-</u>					
Domain Transfer • Left: "Alien", spellcard		<b>.</b>	* <mark></mark>	8		<u></u>	<u>&amp;</u> &	8



MOHAMED BIN ZAYED UNIVERSITY OF ARTIFICIAL INTELLIGENCE



Co  $\Rightarrow$  V  $a_1$  Tas U = 1 H = 1  $D_1 = 1$   $D_2 = 2$   $D_2 = 2$   $D_3 = 2$   $D_3 = 2$   $D_3 = 2$   $D_3 = 2$ Avera

> ו Ab

 $\mathcal{L}_{\text{svae}}$ 

✓ ✓





#### Experiments & Analyses

#### **Comparative Study**

We conduct extensive experiments on popular unsupervised video domain adaptation benchmarks: UCF-HMDB, Jester, Epic-Kitchens, and Sprites.

ask	$\mathcal{S}_{ ext{only}}$	MM-SADA	STCDA	CMCD	A3R	CleanAdapt	CycDA	MixDANN	CIA	TranSVAE
$\rightarrow \mathbf{H} \\ \rightarrow \mathbf{U}$	$86.1 \\ 92.5$	$\begin{array}{c} 84.2\\91.1\end{array}$	$\begin{array}{c} 83.1\\92.1\end{array}$	$84.7 \\ 92.8$	/	89.8 99.2	$\begin{array}{c} 88.1 \\ 98.0 \end{array}$	$82.2 \\ 92.8$	<b>88.3</b> 94.1	87.8 (+1.7) 99.0 (+6.5)
rage ↑	89.3	87.7	87.6	88.8	/	94.5	93.1	87.5	91.2	93.4 (+4.1)
$\rightarrow \mathbf{D}_2$	43.2	49.5	52.0	50.3	53.2	52.7	/	56.0	52.5	50.5 (+ <b>7</b> .3)
$\rightarrow \mathbf{D}_3$	42.5	44.1	45.5	46.3	<b>52.1</b>	47.0	/	47.3	47.8	50.3 (+7.8)
$\rightarrow \mathbf{D}_1$	43.0	48.2	49.0	49.5	<b>51.9</b>	46.2	/	50.3	49.8	50.3 (+7.3)
$\rightarrow \mathbf{D}_3$	48.0	52.7	52.5	52.0	55.5	52.7	/	52.4	53.2	58.6 (+10.6)
$ ightarrow \mathbf{D}_1$	43.0	50.9	52.6	48.7	51.5	47.8	/	51.0	52.2	48.0 (+5.0)
$ ightarrow \mathbf{D}_2$	55.5	56.1	55.6	56.3	<b>63</b> .2	54.4	/	54.7	57.6	58.0 (+2.5)
rage ↑	45.9	50.3	51.2	51.0	54.1	50.3	/	52.0	52.2	$52.6 \ (+6.7)$

Across all tasks, TranSVAE consistently outperforms all previous methods using single modality inputs; our approach even achieves better average results than seven out of eight multi-modal approaches (RGB + flow).

#### **Ablation Study**

The loss separation study on different tasks demonstrates the effectiveness for each of the four constraints designed in the TranSVAE framework.

$\mathcal{L}_{cls}$	$\mathcal{L}_{adv}$	$\mathcal{L}_{mi}$	$\mathcal{L}_{ctc}$	PL	$\mathbf{U} \to \mathbf{H}$	$\mathbf{H} \to \mathbf{U}$	$\mathbf{J}_\mathcal{S} \to \mathbf{J}_\mathcal{T}$	$  \mathbf{D}_1 \rightarrow \mathbf{D}_2$	$\mathbf{D}_1 \to \mathbf{D}_3$	$\mathbf{D}_2  ightarrow \mathbf{D}_1$	$\mathbf{D}_2 \to \mathbf{D}_3$	$  \mathbf{D}_3 \rightarrow \mathbf{D}_1$	$\mathbf{D}_3 \to \mathbf{D}_2$	Avg.
	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	18.61	26.62	22.92	34.00	30.29	33.79	30.49	28.51	34.27	28.83
$\checkmark$		$\checkmark$	$\checkmark$	$\checkmark$	83.06	93.52	48.07	40.93	43.33	43.91	51.13	41.84	52.67	55.38
$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$	85.83	93.52	65.12	46.67	48.56	49.43	55.34	45.52	54.53	60.60
$\checkmark$	$\checkmark$	$\checkmark$		$\checkmark$	83.89	95.80	64.89	48.53	48.25	48.96	54.21	45.52	55.73	60.64
$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		87.22	94.40	64.64	49.87	48.25	49.66	56.47	47.59	55.07	61.46
$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	87.78	98.95	66.10	50.53	50.31	50.34	58.62	48.04	58.00	63.19

The loss integration study and t-SNE plots verify that all constraints serve to achieve good disentanglement effect with domain divergence minimization.

	(b) 🍂 🜼	(c) 🦛	(d)	1	(e) 🗼	#	Leva	$\mathcal{L}_{cle}$	$\mathcal{L}_{adv}$	$\mathcal{L}_{mi}$	$\mathcal{L}_{ctc}$	PL	Accuracy (%)
	Store C		100		. 1 .	$\mathcal{S}_{only}$	Jova		auv	III	cit		80.27
*		× 👌 💊	18	AN		(a)	$\checkmark$	$\checkmark$					82.50 (+2.23)
See.	W.	1 No. 1		100 M	10 B	(b)	~	$\checkmark$	$\checkmark$				84.44 (+4.17)
\$	1 - 3			1	*	(c)	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$			85. 56 (+5. 29)
Ś	AND THE REAL OF		4.11	Martin Con	- 1 -	(d)	√	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		<b>87.22</b> (+6.95)
*				1		(e)	√	√	$\checkmark$	$\checkmark$	√	$\checkmark$	87.78 (+7.61)
See.	*	The state				$\mathcal{T}_{sup}$							95.00
*	(b) **	(c) 🌸	(d)		(e)	#	$\mathcal{L}_{svac}$	$_{e} \mathcal{L}_{cls}$	$\mathcal{L}_{adv}$	$\mathcal{L}_{\mathrm{mi}}$	$\mathcal{L}_{ctc}$	PL	Accuracy (%)
*	(b) *	(c) *	(d)		(e)	# S <sub>only</sub>	L <sub>sva</sub>	$_{\rm e}$ $\mathcal{L}_{\rm cls}$	$\mathcal{L}_{adv}$	$\mathcal{L}_{mi}$	$\mathcal{L}_{ ext{ctc}}$	PL	Accuracy (%) 88.79
	(b)	(c)	(d)		(e)	# S <sub>only</sub> (a)	L <sub>svae</sub>	e L <sub>cls</sub>	$\mathcal{L}_{adv}$	$\mathcal{L}_{\mathrm{mi}}$	$\mathcal{L}_{ ext{ctc}}$	PL	Accuracy (%) 88.79 92.29 (+3.50)
	(b) **	(c)	(d)		(e)	# S <sub>only</sub> (a) (b)	L <sub>sva</sub> √	e L <sub>cls</sub>	L <sub>adv</sub>	£ <sub>mi</sub>	$\mathcal{L}_{ ext{ctc}}$	PL	Accuracy (%) 88.79 92.29 (+3.50) 93.52 (+4.73)
	(b)	(c)	(d)		(e)	# <i>S</i> <sub>only</sub> (a) (b) (c)	L <sub>sva</sub> √ √	e L <sub>cls</sub>	L <sub>adv</sub>	L <sub>mi</sub>	L <sub>ctc</sub>	PL	Accuracy (%) 88.79 92.29 (+3.50) 93.52 (+4.73) 93.70 (+4.91)
	(b) *		(d)		(e)	# <i>S</i> <sub>only</sub> (a) (b) (c) (d)	∠svat       ✓       ✓       ✓       ✓       ✓       ✓	e L <sub>cls</sub> ✓ ✓ ✓ ✓ ✓ ✓	L <sub>adv</sub> √ √	L <sub>mi</sub>	L <sub>ctc</sub>	PL	Accuracy (%) 88.79 92.29 (+3.50) 93.52 (+4.73) 93.70 (+4.91) 94.40 (+5.61)
	(b) **		(d)		(e)	# <i>S</i> <sub>only</sub> (a) (b) (c) (d) (e)	∠svaa       √       √       √       √       √       √       √	e £ <sub>cls</sub> ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓	L <sub>adv</sub> √ √ √ √	L <sub>mi</sub> ✓ ✓ ✓	L <sub>ctc</sub>	PL	Accuracy (%) 88.79 92.29 (+3.50) 93.52 (+4.73) 93.70 (+4.91) 94.40 (+5.61) 98.95 (+10.16)
	(b) **		(d)		(e)	# $ $	∠svaa           √           √           √           √           √           √           √           √           √           √		L <sub>adv</sub> √ √ √	L <sub>mi</sub> ✓ ✓ ✓	L <sub>ctc</sub>	PL	Accuracy (%) 88.79 92.29 (+3.50) 93.52 (+4.73) 93.70 (+4.91) 94.40 (+5.61) 98.95 (+10.16) 96.85

