

# Learning Disentangled Representations of Video with Missing Data

Armand Comas Massagué<sup>1†</sup>, Chi Zhang<sup>2‡</sup>, Zlatan Feric<sup>3‡</sup>, Octavia Camps<sup>4†</sup>, Rose Yu<sup>5\*</sup>

<sup>†</sup> College of Electrical and Computer Engineering, <sup>‡</sup> Khoury College of Computer Sciences, Northeastern University, MA, USA, \* Computer Science & Engineering, University of California San Diego, CA, USA.  
<sup>1</sup> comasmassague.a@northeastern.edu, <sup>2</sup> zhang.chi13@northeastern.edu, <sup>3</sup> feric.z@northeastern.edu <sup>4</sup> camps@coe.neu.edu <sup>5</sup> roseyu@ucsd.edu

UC San Diego



## Problem

- How to learn representations of video sequences in the presence of missing data?

## Our Results

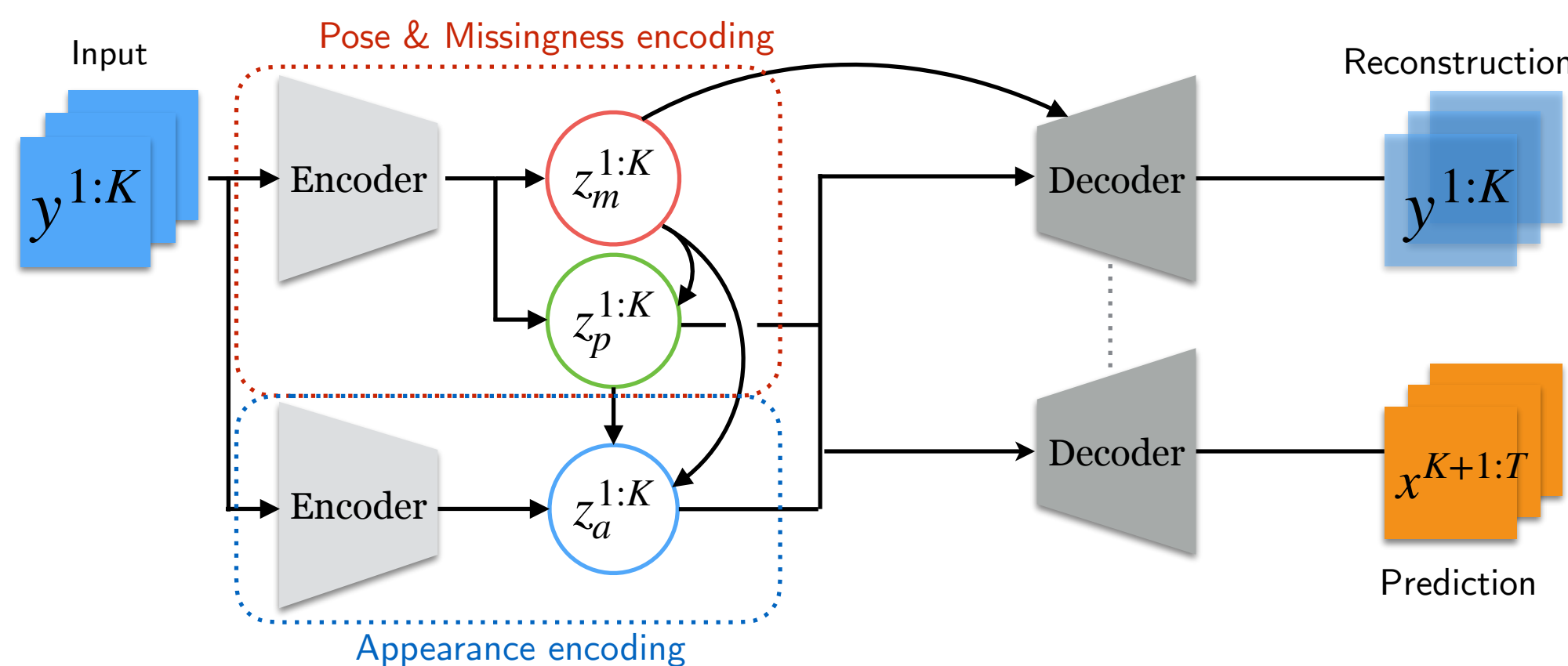
We propose a deep generative model: *Disentangled-Imputed-Video-autoEncoder (DIVE)*.

- Learns representations factorized into **appearance**, **pose** and **missingness** latent variables;
- Imputes** missing data by sampling from the learned latent variables;
- Performs unsupervised stochastic **video prediction** using the imputed hidden representation;
- Robustly generates objects even when their appearances are changing by modeling the **static** and **dynamic appearances** separately.
- outperforms the state of the art by a substantial margin on a moving MNIST dataset with various missing scenarios, and on a real-world MOTChallenge pedestrian dataset.

Code: <https://github.com/Rose-STL-Lab/DIVE>

## Disentangled-Imputed-Video-autoEncoder (DIVE)

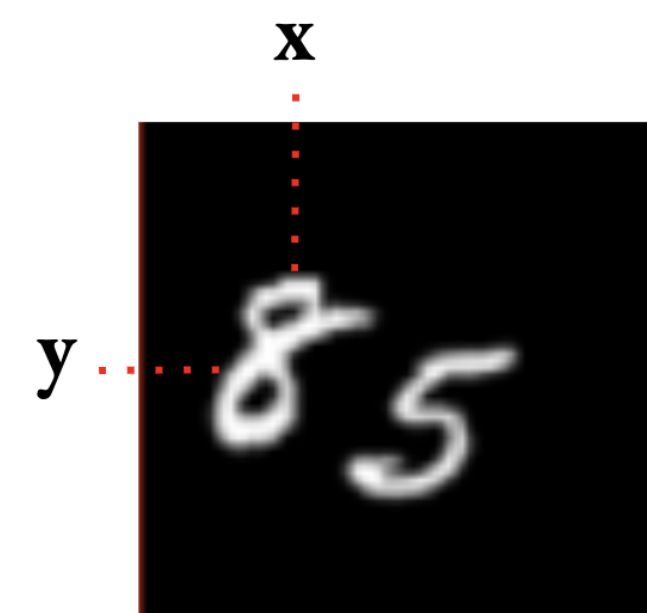
**Deep Generative Model** Denote a video sequence with missing data as  $(y^1, \dots, y^t)$  where each  $y^t \in \mathbb{R}^d$  is a frame.  $x^t \in \mathbb{R}^d$  are the complete frames in the future.



**Figure:** Overall architecture of DIVE, which takes the input video with missing data, infers the missingness (red), pose (green) and appearance (blue) latent variables. Two separate decoders reconstruct and predict the future sequences.

## Disentangled Representation

$$z_i^t = [z_{i,a}^t, z_{i,p}^t, z_{i,m}^t], \quad z_{i,a}^t \in \mathbb{R}^h, z_{i,p}^t \in \mathbb{R}^3, z_{i,m}^t \in \mathbb{Z} \quad (1)$$



**miss = 0**

**miss = 1**

- $z_{i,a}^t$  Appearance vector.
- $z_{i,p}^t$  Pose vector (x, y, size).
- $z_{i,m}^t$  Missingness label.

## Generative Model and Learning

- The **generative distribution** is given by:

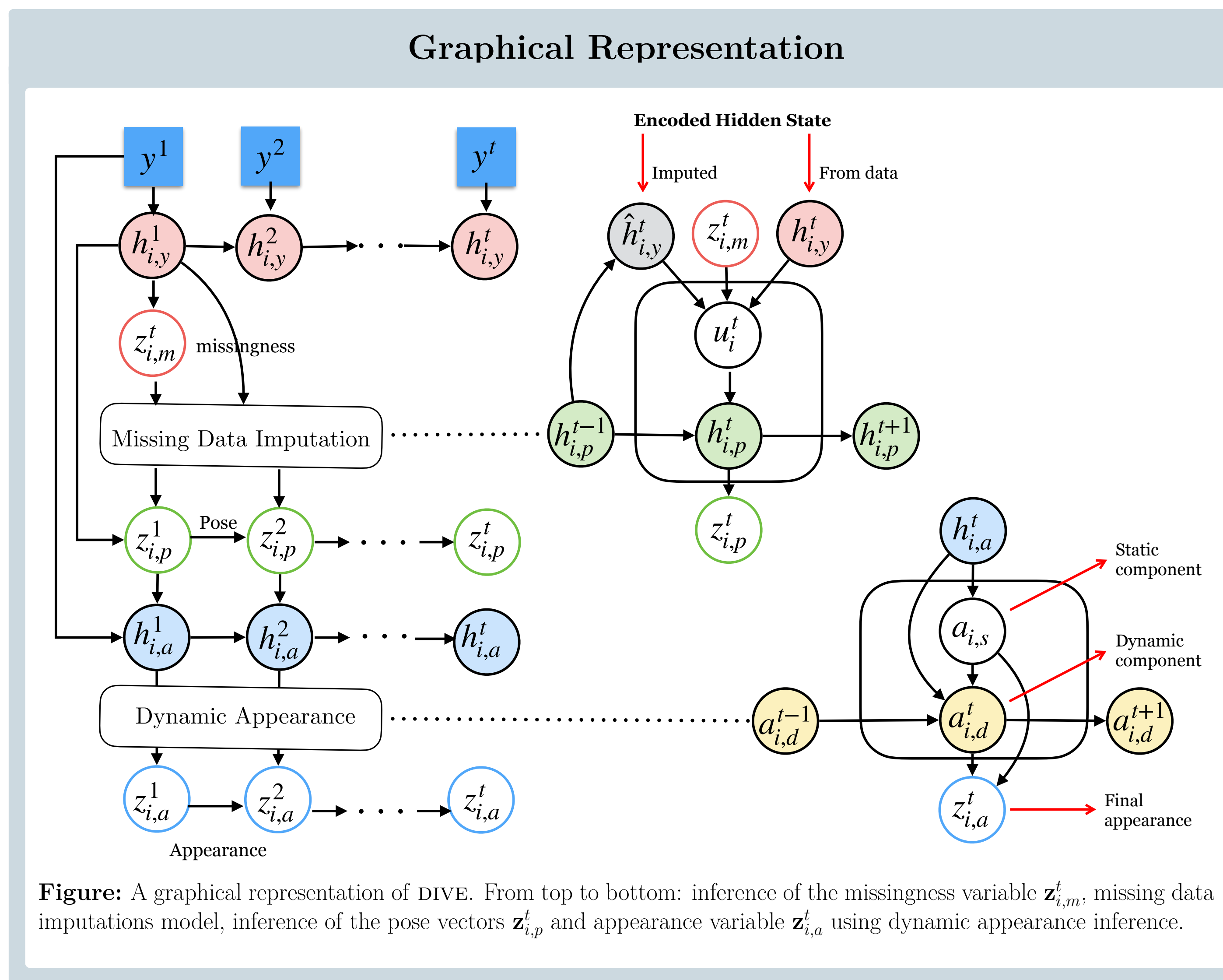
$$p(y^{1:K}, x^{K+1:T} | z^{1:T}) = \prod_{i=1}^N p(y_i^{1:K} | z_i^{1:K}) p(x_i^{K+1:T} | z_i^{K+1:T}) \quad (2)$$

- We **locate** each decoded object with its estimated pose; **mask** it with the missingness label; and **sum** over all objects.

$$p(y_i^t | z_{i,a}^t) = \mathcal{T}(f_{\text{dec}}(z_{i,a}^t; z_{i,p}^t) \circ (1 - z_{i,m}^t)) \quad (3)$$

Following the VAE framework, we train the model by maximizing the evidence lower bound (ELBO).

## Imputation and Inference Model



**Figure:** A graphical representation of DIVE. From top to bottom: inference of the missingness variable  $z_{i,m}^t$ , missing data imputations model, inference of the pose vectors  $z_{i,p}^t$  and appearance variable  $z_{i,a}^t$  using dynamic appearance inference.

## Variational Inference

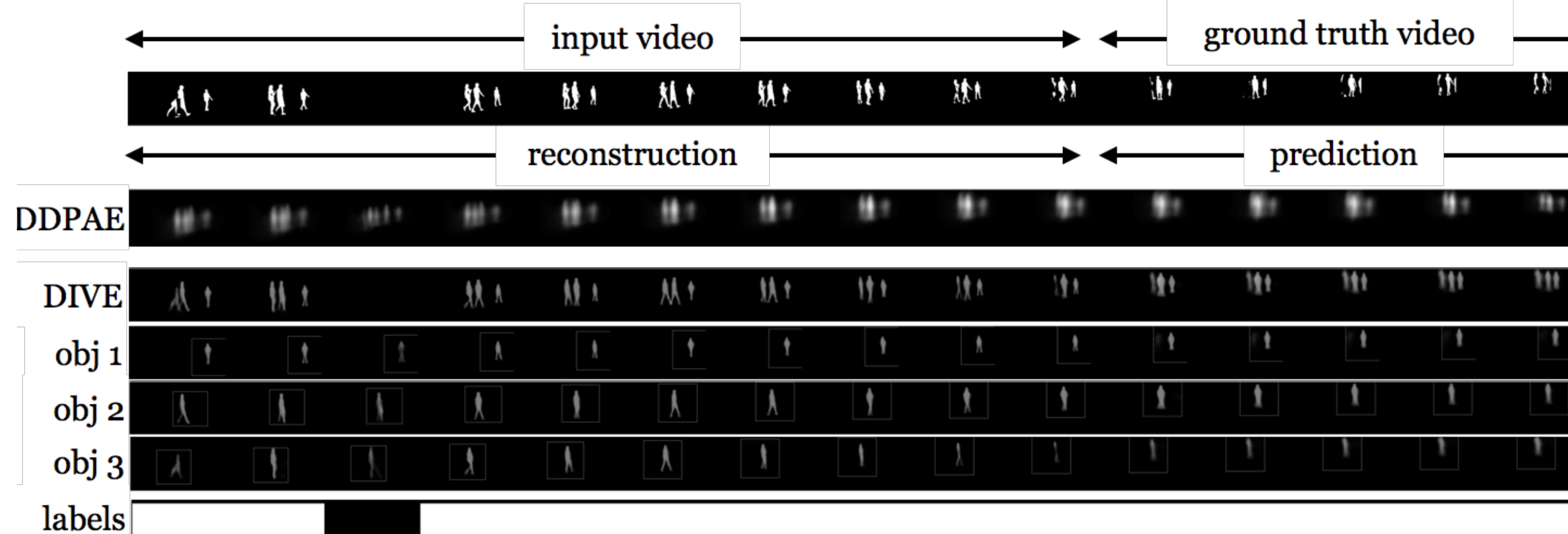
**Missingness:**  $z_{i,m}^t$ , we leverage the input encoding and binarize it with a Heavyside step function.  $z_{i,m}^t = \text{Heavyside}(\text{FC}(h_{i,y}^t))$ .  $h_{i,y}^t$  is a hidden representation of the input data.

**Pose:** We obtain  $z_{i,p}^{1:K}$  from the pose hidden representation as  $h_{i,p}^t = \text{LSTM}(h_{i,p}^{t-1}, u_i^t)$ .  $u_i^t = h_{i,y}^t$  if  $z_{i,m}^t = 0$ . In case of missing data, we **impute**  $h_{i,y}^t$  instead of relying on the input data.

**Appearance:** We decompose it into a **Static component**  $a_{i,s}$ : It captures the inherent **semantics**; and a **Dynamic component**  $a_{i,d}$ : that models the **nuanced variations** in shape. For the **final appearance** vector, we concatenate and mix the dynamic and static components.

## Experiments

### MOTS Pedestrian Dataset



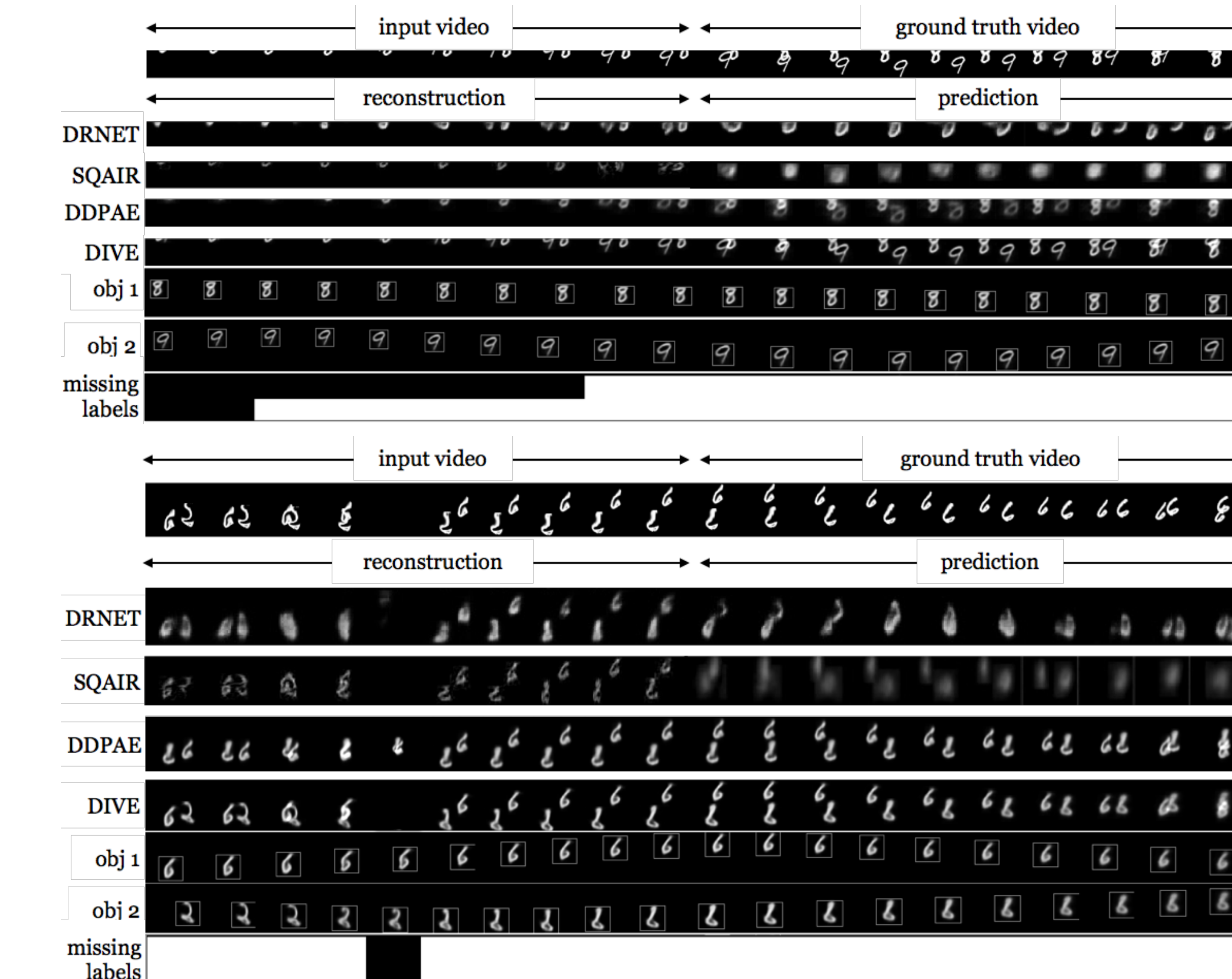
**Figure:** MOTs dataset qualitative results. Note that our method successfully identifies the missing time step, decomposes the objects and keeps track of the missing pedestrians.

Model	BCE ↓	MSE ↓	PSNR ↑	SSIM ↑	NELBO ↓
DDPAE	2495.08	560.37	22.22	0.90	-0.24
DIVE	<b>1355.89</b>	<b>328.96</b>	<b>24.82</b>	<b>0.96</b>	<b>-0.26</b>

**Table:** Quantitative comparison on MOTs pedestrian dataset for DDPAE and DIVE.

### Moving MNIST

**Scenario 1: Partial occlusion:** We can only see 50% of the frame.  
**Scenario 2: Out of scene:** Objects disappear independently for 2 consecutive time-steps.  
**Scenario 3: Varying:** Out of scene 1 time-step + appearance variation in time.



**Figure:** Qualitative results. *Obj 1* and *Obj 2* show DIVEs individual object generations and *missing labels* indicate whether each object is estimated completely missing in the scene. Note that objects are well decomposed, sharply generated and the labels properly predicted. **Scenario 1** and **Scenario 3** respectively.

Scenario 1	BCE ↓		MSE ↓		PSNR ↑		SSIM ↑		NELBO ↓
Model	Rec	Pred	Rec	Pred	Rec	Pred	Rec	Pred	
DRNET[1]	482.07	852.59	72.21	96.36	7.99	6.89	0.76	0.72	/
SQAIR[3]	178.71	967.20	21.84	84.73	13.19	9.96	<b>0.90</b>	0.73	-0.16
DDPAE[2]	182.66	<b>417.00</b>	39.09	67.41	17.56	15.49	0.77	0.72	-0.09
DIVE	<b>119.25</b>	459.10	<b>19.73</b>	<b>64.49</b>	<b>20.64</b>	<b>15.85</b>	<b>0.90</b>	<b>0.78</b>	<b>-0.18</b>
Scenario 2									
DRNET	392.33	1402.45	90.64	187.72	9.59	9.88	0.80	0.67	/
SQAIR	468.22	927.09	73.13	137.04	10.33	8.21	0.84	0.69	-0.17
DDPAE	266.03	409.26	58.37	89.57	18.64	16.94	0.87	0.77	-0.17
DIVE	<b>165.42</b>	<b>321.29</b>	<b>27.03</b>	<b>64.17</b>	<b>22.15</b>	<b>18.56</b>	<b>0.93</b>	<b>0.83</b>	<b>-0.21</b>
Scenario 3									
DRNET	421.72	1304.53	90.46	176.28	9.91	7.33	0.75	0.70	/
SQAIR	560.51	1518.61	74.30	163.25	10.80	7.64	0.83	0.62	-0.16
DDPAE	322.23	403.48	63.63	82.71	18.29	17.22	0.81	<b>0.78</b>	-0.18
DIVE	<b>272.74</b>	<b>374.59</b>	<b>42.81</b>	<b>74.87</b>	<b>20.08</b>	<b>17.61</b>	<b>0.87</b>	<b>0.78</b>	<b>-0.19</b>

**Table:** Quantitative comparison

**Acknowledgments** This work was supported in part by NSF under Grants IIS#1850349, IIS#1814631, ECCS#1808381 and CMMI#1638234, the U. S. Army Research Office under Grant W911NF-20-1-0334 and the Alert DHS Center of Excellence under Award Number 2013-ST-061-ED0001.

## References

- [1] E. L. Denton et al. Unsupervised learning of disentangled representations from video. In *Advances in neural information processing systems*, pages 4414–4423, 2017.
- [2] J.-T. Hsieh, B. Liu, D.-A. Huang, L. F. Fei-Fei, and J. C. Nibbles. Learning to decompose and disentangle representations for video prediction. In *Advances in Neural Information Processing Systems*, pages 517–526, 2018.
- [3] A. Kosiorek, H. Kim, Y. W. Teh, and I. Posner. Sequential attend, infer, repeat: Generative modelling of moving objects. In *Advances in Neural Information Processing Systems*, pages 8606–8616, 2018.