

# Spatio-temporal and Video Prediction with State-based Neural Models February 1st, 2022

Jean-Yves Franceschi,<sup>1</sup> Edouard Delasalles,<sup>1</sup> Jérémie Donà,<sup>1</sup> Mickael Chen,<sup>1</sup> Sylvain Lamprier,<sup>1</sup> Patrick Gallinari<sup>1,2</sup>

<sup>1</sup>Sorbonne Université, CNRS, ISIR, Paris, France <sup>2</sup>Criteo AI Lab, Paris, France





Machine Learning & Deep Learning for nformation Access





# Spatio-temporal and Video Prediction



Given conditioning frames, predict the distribution of future frames.





### Given conditioning frames, predict the distribution of future frames.

### Applications

- ▶ Visual Forecast, Analysis (e.g., climate, health, crowds, etc.)
- ▶ Model-based Reinforcement Learning (Gregor et al., 2019)

# Challenges

- Capture visual and dynamic representations of the world
- Generation of realistic images
- ► Long-term prediction (not only at t+1)
- Account for uncertainty in the future
- Disentangling static/dynamic objects



### Given conditioning frames, predict the distribution of future frames.

### Applications

- ▶ Visual Forecast, Analysis (e.g., climate, health, crowds, etc.)
- ▶ Model-based Reinforcement Learning (Gregor et al., 2019)

# Challenges

- Capture visual and dynamic representations of the world
- Generation of realistic images
- ► Long-term prediction (not only at t+1)
- Account for uncertainty in the future
- Disentangling static/dynamic objects



### Given conditioning frames, predict the distribution of future frames.

### Applications

- ▶ Visual Forecast, Analysis (e.g., climate, health, crowds, etc.)
- ▶ Model-based Reinforcement Learning (Gregor et al., 2019)

# Challenges

- Capture visual and dynamic representations of the world
- Generation of realistic images
- ► Long-term prediction (not only at t+1)
- Account for uncertainty in the future
- Disentangling static/dynamic objects

Stochastic Video Prediction

Jean-Yves Franceschi, Edouard Delasalles, Mickaël Chen, Sylvain Lamprier, Patrick Gallinari: Stochastic Latent Residual Video Prediction. ICML 2020: 3233-3246



#### Autoregressive Models

- + Easy to learn, powerful
- Temporal model tied to generation
- Ex.: Denton and Fergus (2018)



#### State-Space Models

- + Decoupled dynamics and prediction, interpretable
- Open-loop: Harder to train

Ex.: Fraccaro et al. (2017)



### Our Approach



#### State-space Model with Residual Updates

- High-order Dynamics in state  $y_t$  (memory)
- From 1st order ODE inspiration :  $\frac{\mathrm{d}\boldsymbol{y}}{\mathrm{d}t} = f_{\boldsymbol{z}_{\lfloor t \rfloor + 1}}(\boldsymbol{y})$
- State-dependent noise  $z_t$  (stochasticity not only on  $y_1$ )

#### Model



Update Rule:

$$egin{pmatrix} oldsymbol{z}_{t+1} &\sim \mathcal{N}ig(\mu_{ heta}(oldsymbol{y}_t), \sigma_{ heta}(oldsymbol{y}_t)Iig) \ oldsymbol{y}_{t+1} &= oldsymbol{y}_t + f_{ heta}(oldsymbol{y}_t, oldsymbol{z}_{t+1}) \end{split}$$



### State-space Model with Residual Updates

### Variable w to Separate stochasticity from dynamics

Better gradient backprop



 $\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{y}_t + f_{\theta}(\mathbf{y}_t, \mathbf{z}_{t+1}) \\ \mathbf{z}_{t+1} &\sim \mathcal{N}\big(\mu_{\theta}(\mathbf{y}_t), \sigma_{\theta}(\mathbf{y}_t)I\big) \end{aligned}$ 

### Variational Autoencoders (VAEs)



• Generative model  $p_{\theta}(\boldsymbol{x} \mid \boldsymbol{z})$  for  $\boldsymbol{x}$  from a latent  $\boldsymbol{z} \sim p_{\theta}(\boldsymbol{z})$ :

$$p_{\theta}(\boldsymbol{x}) = \int_{\boldsymbol{z}} p_{\theta}(\boldsymbol{x} \mid \boldsymbol{z}) p(\boldsymbol{z}) \, \mathrm{d}\boldsymbol{z}$$

▶ Intractable objective: sampling from p(z) very inefficient



• Introducing  $q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x})$  to approximate  $p_{\theta}(\boldsymbol{z} \mid \boldsymbol{x}) \propto p_{\theta}(\boldsymbol{z}, \boldsymbol{x})$ :

$$D_{\mathrm{KL}}(q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x}) \parallel p_{\theta}(\boldsymbol{z} \mid \boldsymbol{x})) = p_{\theta}(\boldsymbol{x}) - \int q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x}) \log \frac{p_{\theta}(\boldsymbol{z}, \boldsymbol{x})}{q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x})} dz$$

$$\Rightarrow \arg\min_{\phi} D_{\mathrm{KL}} (q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x}) \mid p_{\theta}(\boldsymbol{z} \mid \boldsymbol{x})) = \arg\max_{\phi} \int q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x}) \log \frac{p_{\theta}(\boldsymbol{z}, \boldsymbol{x})}{q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x})} dz$$



### Evidence Lower Bound (ELBO)

#### Optimized objective:

$$\underset{\phi,\theta}{\arg\max} \underbrace{\mathbb{E}_{\boldsymbol{z} \sim q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x})} \left[ \log p_{\theta}(\boldsymbol{x} \mid \boldsymbol{z}) \right]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}} \left( q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x}) \parallel p(\boldsymbol{z}) \right)}_{\text{KL term}}.$$





#### State-space Model with Residual Updates

► Initial Condition  $y_1$  unknown, Stochastic Noise z unknown ⇒  $p(y_1, z_{2:T} | x_{1:T})$  to be inferred via VAE

#### Inference



#### Variational Distribution:

$$q(\boldsymbol{y}_{1}, \boldsymbol{z}_{2:T} \mid \boldsymbol{x}_{1:T}) =$$

$$q_{\phi}^{y}(\boldsymbol{y}_{1} \mid \boldsymbol{x}_{1:k}) \prod_{t=2}^{T} q_{\phi}^{z}(\boldsymbol{z}_{t} \mid \boldsymbol{x}_{1:t})$$

$$\begin{pmatrix} \boldsymbol{y}_{1} \sim \mathcal{N} \left( \mu_{\phi}^{y}(\boldsymbol{x}_{1:k}), \sigma_{\phi}^{y}(\boldsymbol{x}_{1:k})I \right) \\ \boldsymbol{z}_{t} \sim \mathcal{N} \left( \mu_{\phi}^{z}(\boldsymbol{x}_{1:t}), \sigma_{\phi}^{z}(\boldsymbol{x}_{1:t})I \right) \end{pmatrix}$$

### Generative Model





$$\begin{cases} \boldsymbol{y}_1 \sim \mathcal{N}(\boldsymbol{0}, I) \\ \boldsymbol{z}_{t+1} \sim \mathcal{N}(\boldsymbol{\mu}_{\theta}(\boldsymbol{y}_t), \sigma_{\theta}(\boldsymbol{y}_t)I) \\ \boldsymbol{y}_{t+1} = \boldsymbol{y}_t + f_{\theta}(\boldsymbol{y}_t, \boldsymbol{z}_{t+1}) \\ \boldsymbol{x}_t \sim \mathcal{N}(g_{\theta}(\boldsymbol{y}_t), \nu I) \end{cases}$$

(initial condition)
(random variable prediction)
(latent state prediction)
(decoding)

Inference





$$q(\boldsymbol{z}_{2:T}, \boldsymbol{y}_1 \mid \boldsymbol{x}_{1:T}) = \underbrace{q(\boldsymbol{y}_1 \mid \boldsymbol{x}_{1:k})}_{\text{Init. Cond.}} \prod_{t=2}^T \underbrace{q(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t})}_{\text{LSTM}}$$

Inference





$$q(\boldsymbol{z}_{2:T}, \boldsymbol{y}_1 \mid \boldsymbol{x}_{1:T}) = \underbrace{q(\boldsymbol{y}_1 \mid \boldsymbol{x}_{1:k})}_{\text{Init. Cond.}} \prod_{t=2} \underbrace{q(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t})}_{\text{LSTM}}$$

Training done using variational inference within an ELBO objective.

#### Content Variable w





- Store static information (e.g., background and object shapes)
- $\blacktriangleright$  Computed from randomly sampled frames  $\rightarrow$  temporal invariance
- Skip connections between encoder and decoder

#### Evaluation





- Conditioning frames are used to infer dynamic variables
- Prediction follows using the forward model



Madala	Stoc	hastic	Deterr	Deterministic		
Models	PSNR	PSNR SSIM		SSIM		
SVG	14.50	0.7090	12.85	0.6185		
Ours	16.93	0.7799	18.25	0.8300		
Ours - GRU	15.80	0.7464	13.17	0.6237		
Ours - MLP	16.55	0.7694	16.70	0.7876		
Ours - w/o <i>z</i>			14.99	0.4757		

#### SVG: Denton and Fergus (2018)

	t = 1	<i>t</i> =	= 3	t = 5	t = 6	t = 8	t = 10	t=12	t=14	t=16	t=18	t=20	t=22	t=24
Ground Truth	<sup>2</sup> 5	२	5	35	35	35	æ	3	cy,	8	S	(FP	~~)	~5
				SVG	35	ž	6	a	61	g	2	る	36	36
				Ours	35	35	ŵ	ţ,n	ŝ	67	3	3	~~)	~5



Metric	Dataset	SV2P	SAVP	SVG	SVRNN	Ours
	KTH	636	374	377		222
FVD (↓)	H3.6M				556	416
	BAIR	965	152	255		163
LPIPS (↓)	KTH	0.2049	0.1120	0.0923	_	0.0736
	H3.6M	—		—	0.0557	0.0509
	BAIR	0.0912	0.0634	0.0609		0.0574
PSNR (↑)	KTH	28.19	26.51	28.06		29.69
	H3.6M	—	—	—	24.46	25.30
	BAIR	20.39	18.44	18.95		19.59

SV2P: Finn et al. (2016), SAVP: Babaeizadeh et al. (2018), SVRNN: Minderer et al. (2019)





Using the Euler approximation scheme:

$$\frac{\mathrm{d}\boldsymbol{y}}{\mathrm{d}t} = f_{\boldsymbol{z}_{\lfloor t \rfloor + 1}}(\boldsymbol{y}) \quad \Rightarrow \quad \boldsymbol{y}_{t + \Delta t} = \boldsymbol{y}_t + \Delta t \cdot f_{\theta} \Big( \boldsymbol{y}_t, \boldsymbol{z}_{\lfloor t \rfloor + 1} \Big)$$

### Experimental Results: Latent Space Properties







- The separation of dynamic variables y and the content variable w can be seen as spatiotemporal disentanglement.
- Spatiotemporal disentanglement provides interpretability.
- It can help to improve prediction performance.

#### Question

What is spatiotemporal disentanglement?

- Often complex and seldom analyzed, achieved through:
  - KL-based separation in VAEs: ours, Hsieh et al. (2018) and Yingzhen et al. (2018);
  - adversarial losses: Villegas et al. (2017) and Denton and Birodkar (2017).
- We aim at grounding spatiotemporal disentanglement on stronger foundations, with fewer implicit assumptions

# Spatiotemporal Disentanglement

Jérémie Donà, Jean-Yves Franceschi, Sylvain Lamprier, and Patrick Gallinari: PDE-Driven Spatiotemporal Disentanglement. ICLR 2021.

### Deterministic Setting





#### Hidden State u

Function of continuous coordinates following a PDE:

 $u{:}\,(x,t)\mapsto u(x,t).$ 

Ex.: physical state of an action, ocean temperature.

### Observations $\boldsymbol{v}$

Vectorial  $v_{t_0}, \ldots, v_{t_1}$ , spatial measurements of u:

$$\boldsymbol{v}_t = \boldsymbol{\zeta} \circ \boldsymbol{u}(\cdot, t).$$

Ex.: pixel values, punctual surface temperatures.



- We propose a novel interpretation of spatiotemporal disentanglement based on the separation of variables in PDEs
- Example with the heat equation:

$$\frac{\partial u}{\partial t} = c^2 \frac{\partial^2 u}{\partial x^2}, \quad u(0,t) = u(L,t) = 0, \quad u(x,0) = f(x),$$

with separable solutions:

$$u(x,t) = \underbrace{\mu \sin\left(\frac{n\pi}{L}x\right)}_{\phi(x)} \times \underbrace{\exp\left(-\left(\frac{cn\pi}{L}\right)^2 t\right)}_{\psi(t)} = \xi(\phi(x),\psi(t))$$

•  $\phi$  and  $\psi$  can be found by solving an ODE on t and a PDE on x

# Spatiotemporal Disentanglement as Separation of Variables



For spatiotemporal modeling from observations v:



But directly learning U from partial observations can be too difficult:

- $\blacktriangleright$  No always access to coordinates x
- Many unobserved variables

Thus, focus on learning  $\phi$ ,  $\psi$  and decoder  $D(\phi, \psi(t)) = (\zeta \circ U \circ \xi) (\phi(\cdot), \psi(t))$ :







• We learn representations S,  $T_t$  such that:

$$\phi \equiv S \in \mathbb{R}^d, \qquad \psi \equiv T : t \mapsto T_t \in \mathbb{R}^p$$

• S and  $T_{t_0}$  are inferred with encoders  $E_S$  and  $E_T$  from conditioning frames:

$$V_{\tau}(t_0) = \left(v_{t_0}, \dots, v_{t_0+\tau}\right)$$

### Forecasting





•  $T \equiv \psi$  is driven by an ODE:

$$\frac{\partial T_t}{\partial t} = f(T_t) \qquad \Leftrightarrow \qquad T_t = T_{t_0} + \int_{t_0}^t f(T_{t'}) \,\mathrm{d}t'$$

► Forecasting loss:

$$\widehat{v}_t = D\left(S, T_{t_0} + \int_{t_0}^t f(T_{t'}) \,\mathrm{d}t'\right), \quad \mathcal{L}_{\text{pred}} = \frac{1}{\nu + 1} \sum_{i=0}^{\nu} \frac{1}{m} \|\widehat{v}_{t_0 + i\Delta t} - v_{t_0 + i\Delta t}\|_2^2$$

### Forecasting





•  $T \equiv \psi$  is driven by an ODE:

$$\frac{\partial T_t}{\partial t} = f(T_t) \qquad \Leftrightarrow \qquad T_t = T_{t_0} + \int_{t_0}^t f(T_{t'}) \, \mathrm{d}t'$$

► Alignment loss:

$$\mathcal{L}_{AE} = \frac{1}{m} \left\| D\left(S, E_T\left(V_{\tau}(t_0 + i\Delta t)\right)\right) - v_{t_0 + i\Delta t} \right\|_2^2, \text{ with } i \sim \mathcal{U}\left(\llbracket 0, \nu - \tau \rrbracket\right).$$

### Invariance of $\boldsymbol{S}$ and Spatiotemporal Disentanglement





► From a strict *S* invariance constraint to a weaker one to take into account variations of observable content:

$$\frac{\partial E_S(V_\tau(t))}{\partial t} = 0 \Rightarrow \mathcal{L}_{reg}^S = \left\| E_S(V_\tau(t_0)) - E_S(V_\tau(t_1 - \tau)) \right\|_2^2$$

Disentanglement loss:

$$\mathcal{L}_{\text{reg}}^{T} = \left\| T_{t_0} \right\|_{2}^{2} = \left\| E_{T} \left( V_{\tau}(t_0) \right) \right\|_{2}^{2}$$

### Experimental Results: Multiview & Prediction





### Experimental Results: Multiview & Prediction



	t=0	t=4	t=6	t=10	t=14	t=24	t=54	t=99
Input frames	<b>s</b> 5	<b>s</b> 5	5	5 <b>5</b>	5 <b>5</b>	5	<b>s</b> 5	5
		MTM	5	$5_{s}$	58	F	100	
		town town	5	uļn	9 5	նդեր	9 5	5 <sup>5</sup>
			ঠ	5 <b>5</b>	5 <b>5</b>	5 s	5 8	5 8
			\$	58	55	(49)	1	
			5	5 <b>5</b>	5 <sub>8</sub>	5	<b>б</b> 5	<b>8</b> 5
Content input	z	<b>~</b> ~	Z	87	67	5	7	76

### Experimental Results: Multiview & Prediction







Table 1: Forecasting performance on WaveEq-100, WaveEq and SST of compared models with respect to indicated prediction horizons. Bold scores indicate the best performing method.

	WaveEq-100	WaveEq	SST			
Models	MS	N	ISE	SSIM		
	t + 40	t + 40	t+6	t + 10	t+6	t + 10
PKnl	_		1.28	2.03	0.6686	0.5844
PhyDNet	_	_	1.27	1.91	0.5782	0.4645
SVG	_	_	1.51	2.06	0.6259	0.5595
MIM	—	—	0.91	1.45	0.7406	0.6525
Ours	$\textbf{4.33}\times\textbf{10^{-5}}$	$1.44  imes 10^{-4}$	0.86	1.43	0.7466	0.6577
Ours (without $S$ )	$1.33  imes 10^{-4}$	$5.09  imes 10^{-4}$	0.95	1.50	0.7204	0.6446



Table 2: Prediction and content swap scores of all compared models on Moving MNIST. Bold scores indicate the best performing method.

Models	Pred. (	Pred. $(t + 10)$		Pred. $(t + 95)$		Swap $(t+10)$		(t + 95)
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
SVG	18.18	0.8329	12.85	0.6185	_	_	_	_
MIM	24.16	0.9113	16.50	0.6529		_		_
DrNet	14.94	0.6596	12.91	0.5379	14.12	0.6206	12.80	0.5306
DDPAE	21.17	0.8814	13.56	0.6446	18.44	0.8256	13.25	0.6378
PhyDNet	23.12	0.9128	16.46	0.3878	12.04	0.5572	13.49	0.2839
Ours	21.70	0.9088	17.50	0.7990	18.42	0.8368	16.50	0.7713



- We study the relevance of state-space models for spatiotemporal data, with a focus on:
  - stochastic generation for videos;
  - variable separation for physical phenomena.
- Such models allow us to simulatenously perform accurate prediction and learn meaningful representations.
- Perspectives
  - Stochastic Differential Equations (e.g. Brownian Motion)
  - Neural ODE Formalism (invertible transformations?)
  - Adaptative Steps of Noise Introduction (i.e., asynchronously with frame rate)

#### Further Resources: code, models, animated samples

- https://sites.google.com/view/srvp/
- https://github.com/JeremDona/spatiotemporal\_variable\_separation

### References



- Babaeizadeh, Mohammad et al. (2018). "Stochastic Variational Video Prediction". In: International Conference on Learning Representations.
- Denton, Emily and Vighnesh Birodkar (2017). "Unsupervised Learning of Disentangled Representations from Video". In: Advances in Neural Information Processing Systems. Ed. by Isabelle Guyon et al. Vol. 30. Curran Associates, Inc., pp. 4417–4426.
- Denton, Emily and Rob Fergus (July 2018). "Stochastic Video Generation with a Learned Prior". In: Proceedings of the 35th International Conference on Machine Learning. Ed. by Jennifer Dy and Andreas Krause. Vol. 80. Proceedings of Machine Learning Research. Stockholmsmässan, Stockholm, Sweden: PMLR, pp. 1174–1183.
- Finn, Chelsea, Ian Goodfellow, and Sergey Levine (2016). "Unsupervised Learning for Physical Interaction through Video Prediction". In: Advances in Neural Information Processing Systems. Ed. by Daniel D. Lee et al. Vol. 29. Curran Associates, Inc., pp. 64–72.
- Fraccaro, Marco et al. (2017). "A Disentangled Recognition and Nonlinear Dynamics Model for Unsupervised Learning". In: Advances in Neural Information Processing Systems 30. Ed. by Isabelle Guyon et al. Vol. 30. Curran Associates, Inc., pp. 3604–3613.
- Gregor, Karol et al. (2019). "Temporal Difference Variational Auto-Encoder". In: International Conference on Learning Representations.
- Hsieh, Jun-Ting et al. (2018). "Learning to Decompose and Disentangle Representations for Video Prediction". In: Advances in Neural Information Processing Systems. Ed. by Samy Bengio et al. Vol. 31. Curran Associates, Inc., pp. 515–524.
- Minderer, Matthias et al. (2019). "Unsupervised learning of object structure and dynamics from videos". In: Advances in Neural Information Processing Systems. Ed. by Hanna Wallach et al. Vol. 32. Curran Associates, Inc., pp. 92–102.
- Villegas, Ruben et al. (2017). "Decomposing motion and content for natural video sequence prediction". In: International Conference on Learning Representations.
- Yingzhen, Li and Stephan Mandt (July 2018). "Disentangled Sequential Autoencoder". In: Proceedings of the 35th International Conference on Machine Learning. Ed. by Jennifer Dy and Andreas Krause. Vol. 80. Proceedings of Machine Learning Research. Stockholmsmässan, Stockholm, Sweden: PMLR, pp. 5670–5679.