

KW-DYAN: A Recurrent Dynamics-Based Network For Video Prediction

Wenqian Liu, Armand Comas, Yuexi Zhang, Octavia Camps, Mario Sznaiar
Department of Electrical & Computer Engineering, Northeastern University, Boston, USA

Motivation

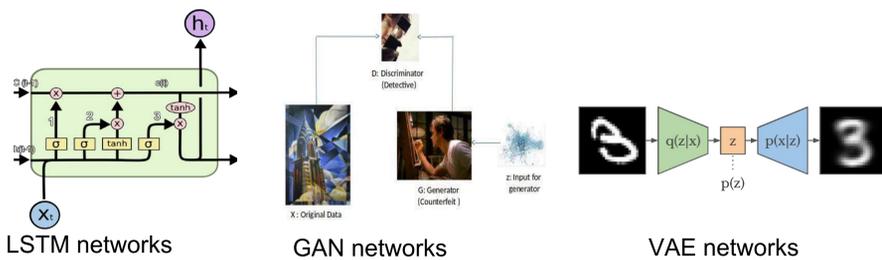
DYNAMICS:



Provide powerful cues to understand scenes and anticipate the future:



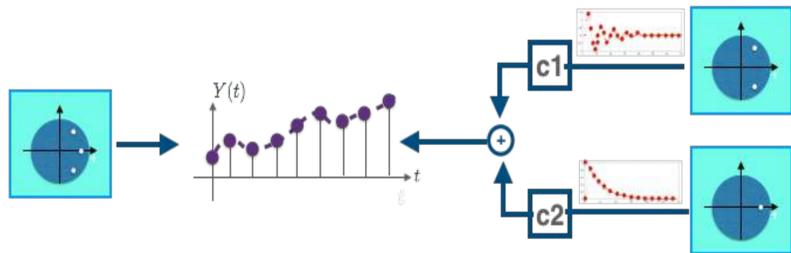
Video Frame Prediction Approaches



- Large number of parameters ✗
- Difficult to train ✗
- Blurry images ✗

Dynamics-based Representation

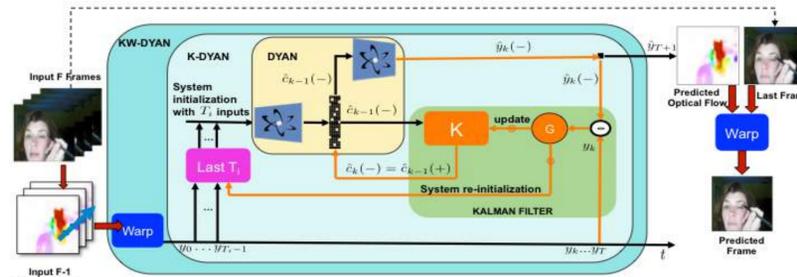
Model temporal signals as the output of Linear Time Invariant (LTI) systems:



$$Y(z) = \sum_{i=1}^n \frac{c_i z}{z - p_i} \iff y(k) = \sum_{i=1}^n c_i p_i^k$$

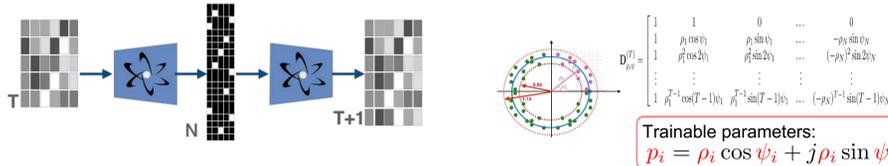
- Can be represented by n pairs (c_i, p_i) coefficient-pole
- The number of poles n is the memory of the system
- “Simpler” systems have fewer poles

KW-DYAN: A recurrent network to capture dynamics



- Small number of parameters ✓
- Easy to train ✓
- Sharp images ✓
- Accurate timing ✓
- Recurrent: can process arbitrary long sequences ✓

Uses a Dynamics-based Sparse Autoencoder



Encoder: Uses a *structured* dictionary to find the least number of poles and coefficients to **encode the optical flow inputs**:

$$c^* = \arg \min_c \frac{1}{2} \|y_{1:T} - D_{\rho, \psi}^{(T)} c\|^2 + \lambda \|c\|_1$$

Fidelity Sparsity

Decoder: Extends the dictionary to decode the sparse code and **predict future optical flows**:

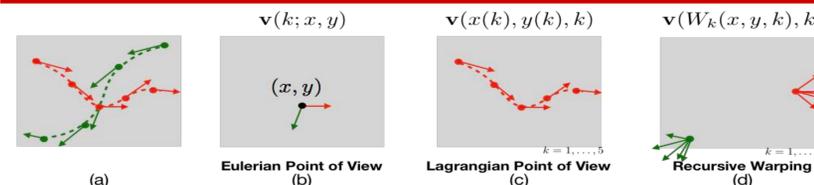
$$y_{1:T+1} = D_{\rho, \psi}^{(T+1)} c^*$$

Detects Changes in Dynamics



KW-DYAN uses a Kalman filter to process its inputs recursively, detect changes in dynamics, and **reduce prediction lagging**.

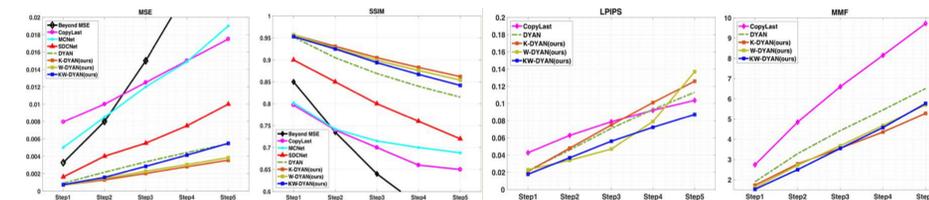
Reduces the Number of Changes in Dynamics



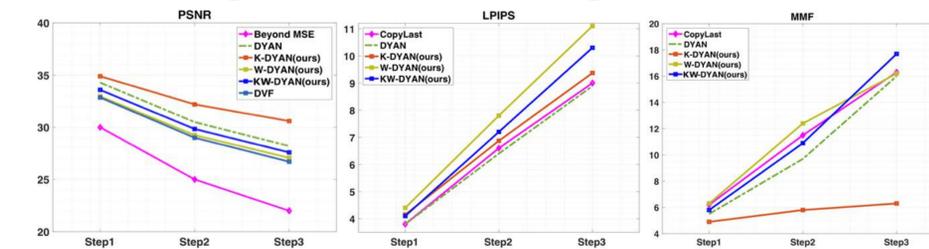
KW-DYAN applies recursive warping to the optical flow inputs to combine Lagrangian and Eulerian point of views and **reduce the number of changes in dynamics that cause lagging**.

Experiments

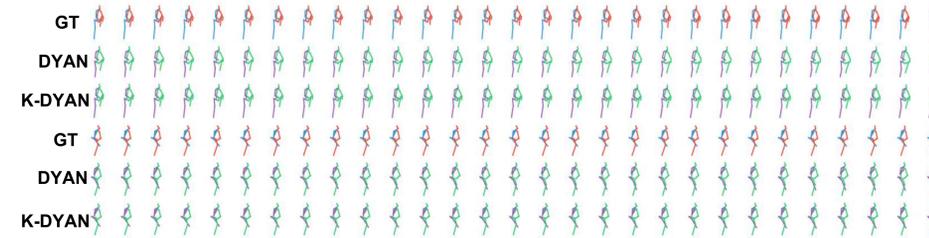
CALTECH-Car Mounted camera video dataset (10 frames input)



UCF101: Human action video dataset (4 frames input)



Human 3.6M dataset 25 steps prediction



	walking	Eating	Smoking	Discussion	Dancing	Directions	Greeting	Phoning	Posing	Purchases	Sitting	Sittingdown	Takingphoto	Waiting	Walkingdog	Walkingtogether	Average
RRNN	1.14	1.34	1.83	1.79	1.59	2.03	1.89	2.56	2.30	2.14	2.72	1.51	2.34	1.86	1.42	1.90	
CSS	0.92	0.92	1.62	1.86	1.45	1.72	1.81	2.65	2.52	1.67	2.06	2.5	1.92	1.28	1.71	1.71	
TP-RNN	0.77	1.14	1.66	1.74	1.38	1.81	1.68	2.47	2.28	1.74	1.93	1.35	2.46	1.98	1.28	1.83	
DYAN	1.21	1.34	1.86	1.72	1.44	1.84	1.83	2.40	2.32	1.95	2.30	1.44	2.35	1.84	1.56	1.83	
KDYAN	1.17	1.34	1.68	1.67	1.37	1.72	1.76	2.38	2.27	1.64	1.90	1.31	2.32	1.82	1.45	1.70	