# SRVP:
# Strong Recollection Video Prediction Model Using Attention-Based Spatiotemporal Correlation Fusion

Yuseon Kim     Kyongseok Park
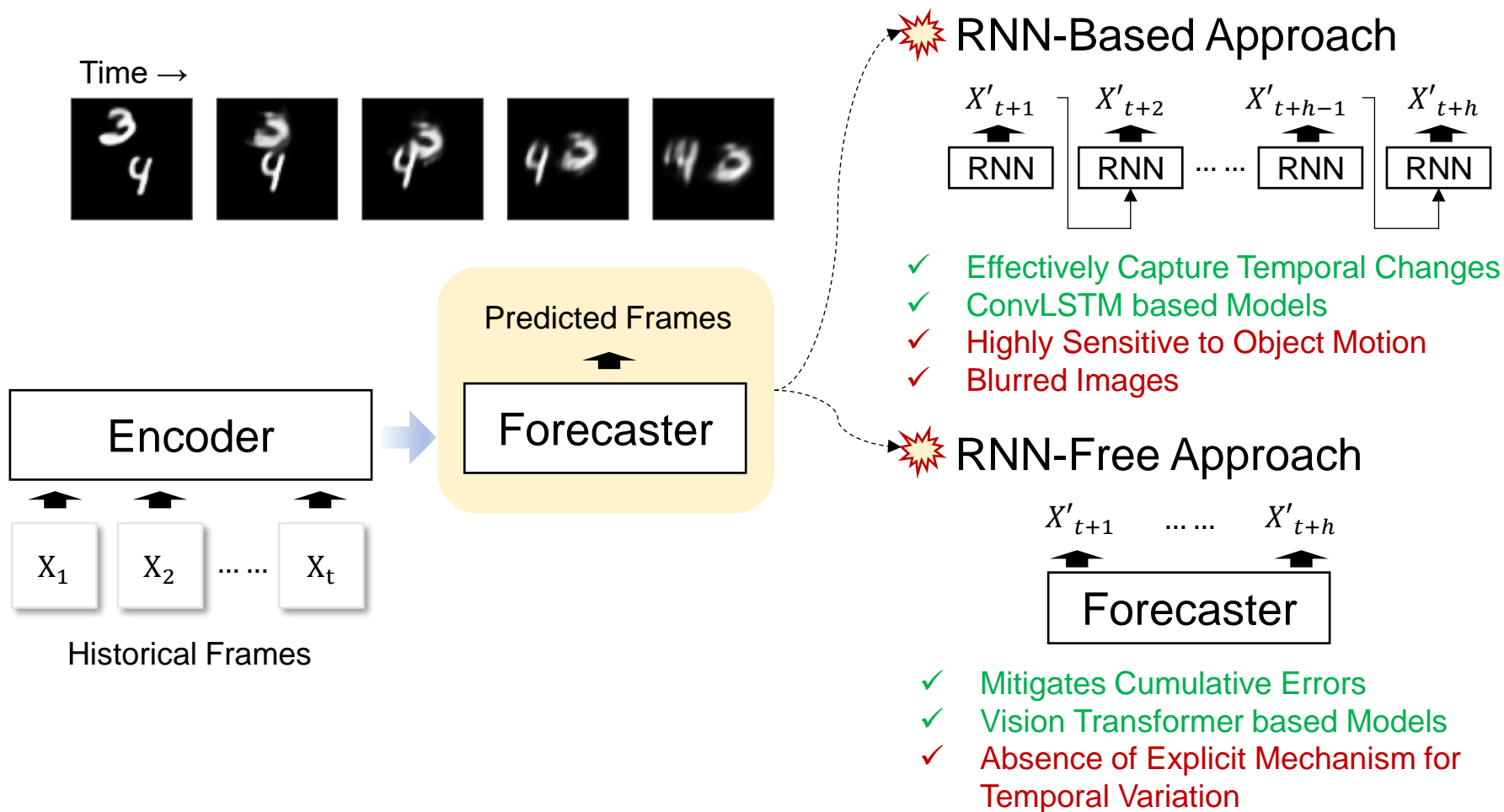
KiSTi
www.kisti.re.kr
**Korea Institute of Science and Technology Information**

UST
**UNIVERSITY OF SCIENCE & TECHNOLOGY**

# Video Prediction

Time →



**Encoder** → **Forecaster** (Predicted Frames)

$X_1$ $X_2$ ...... $X_t$

Historical Frames

## 💥 RNN-Based Approach

$$X'_{t+1} \quad X'_{t+2} \quad X'_{t+h-1} \quad X'_{t+h}$$

RNN — RNN ····· RNN — RNN

- ✓ **Effectively Capture Temporal Changes**
- ✓ **ConvLSTM based Models**
- ✓ **Highly Sensitive to Object Motion**
- ✓ **Blurred Images**

## 💥 RNN-Free Approach

$$X'_{t+1} \quad ...... \quad X'_{t+h}$$

**Forecaster**

- ✓ **Mitigates Cumulative Errors**
- ✓ **Vision Transformer based Models**
- ✓ **Absence of Explicit Mechanism for Temporal Variation**
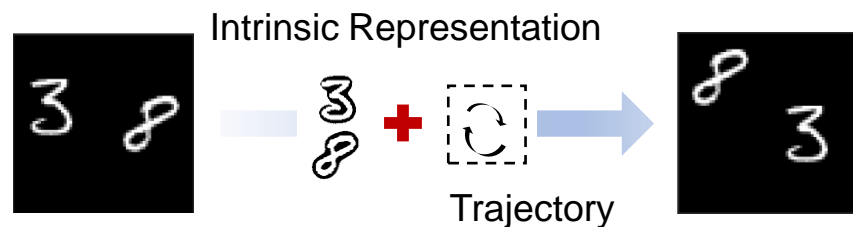
# Strong Recollection Video Prediction



✓ **Baseline Architecture**
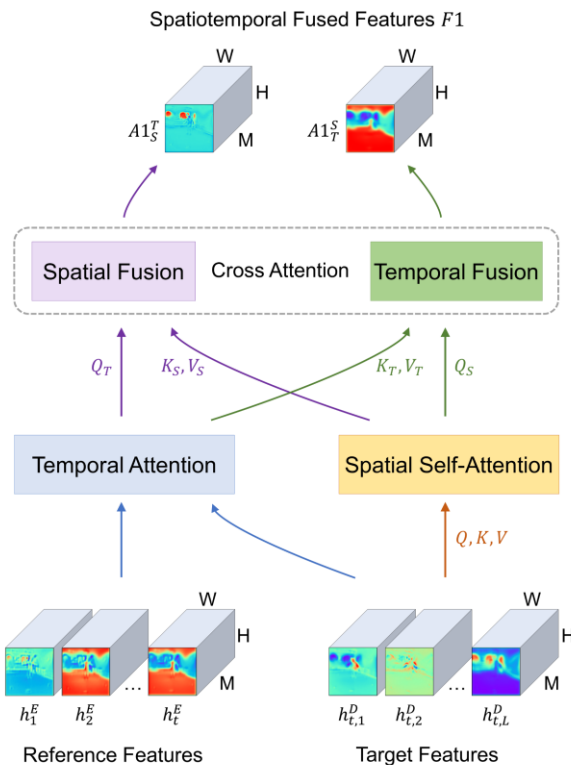   - ConvGRU based Encoder-Forecaster

✓ **Temporal Dynamics**

✓ **Spatiotemporal Fused Features**

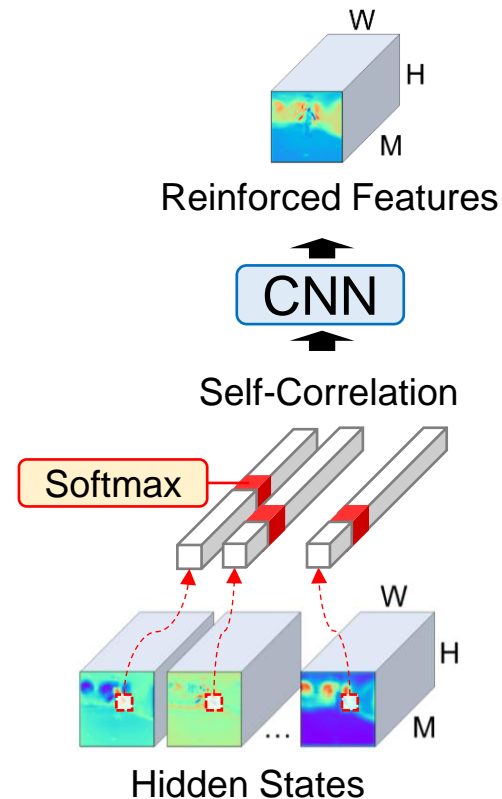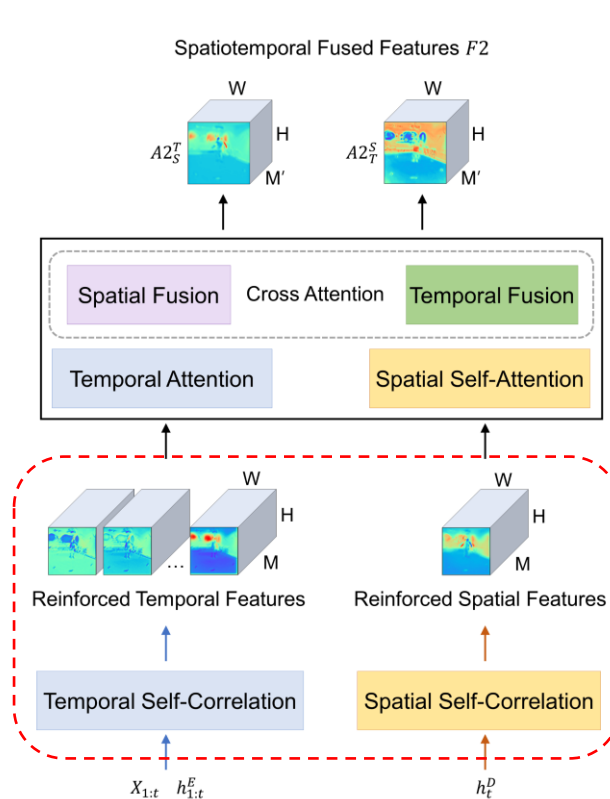Intrinsic Representation

Trajectory

✓ **Scaled-Dot Product Attention**
   - Vaswani, Ashish, et al. "Attention is all you need." In NeurIPS (2017).

# Attention-based Spatiotemporal Fusion



## Standard Attention

Spatiotemporal Fused Features $F1$

Spatial Fusion — Cross Attention — Temporal Fusion

Temporal Attention    Spatial Self-Attention

Reference Features    Target Features

## Reinforced Feature Attention

Spatiotemporal Fused Features $F2$

Spatial Fusion — Cross Attention — Temporal Fusion

Temporal Attention    Spatial Self-Attention

Reinforced Temporal Features    Reinforced Spatial Features

Temporal Self-Correlation    Spatial Self-Correlation

Reinforced Features

CNN

Self-Correlation

Softmax

Hidden States

# Experiment Materials

✓ Benchmarks



**Moving MNIST**
- 64x64 Gray Image
- 10 Frames → 10 Frames



**KTH Action**
- 64x64 Gray Image
- 10 Frames → 10 Frames



**Human3.6M**
- 100x100 Color Image
- 4 Frames → 4 Frames

✓ Evaluation Metrics
- Mean Squared Error (MSE) ↓
- Peak Signal-to-Noise Ratio (PSNR) ↑
- Structural Similarity Index Measure (SSIM) ↑

✓ Reference Models

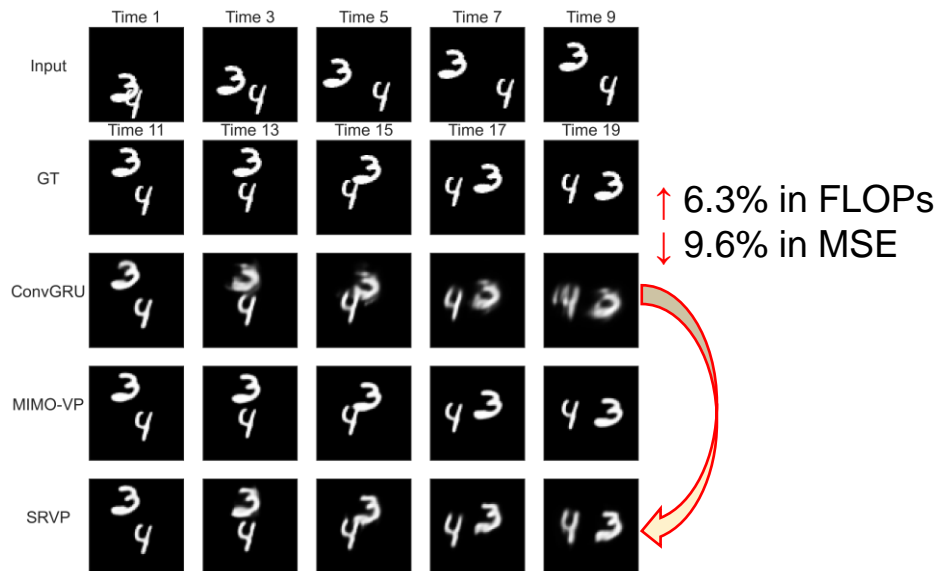| Approach | Model | Conference |
|---|---|---|
| RNN-based | ConvLSTM | NIPS 2015 |
| | ST-LSTM | NIPS 2017 |
| | Causal LSTM | ICML 2018 |
| | MIM | CVPR 2019 |
| | ConvGRU | - |
| RNN-free | MIMO-VP | AAAI 2023 |
| | SimVP | CVPR 2022 |

- MIMO-VP → Transformer-based VP
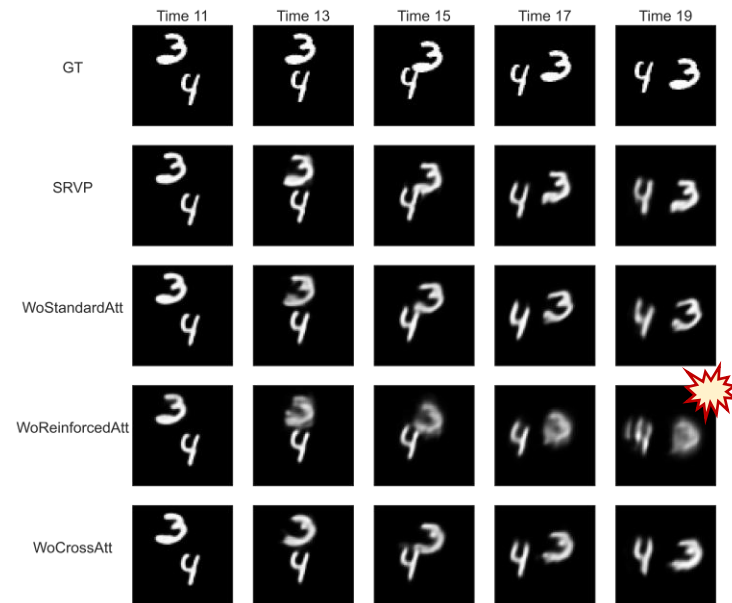- SimVP → CNN-based VP

# Notable Findings on SRVP

## ✓ Outperform than RNN-based Models

| Method | 10 → 10 | | 10 → 30 | |
|---|---|---|---|---|
| | MSE ↓ | SSIM ↑ | MSE ↓ | SSIM ↑ |
| ConvLSTM [19] | 1419.41 | 0.5661 | 1272.60 | 0.4963 |
| ST-LSTM [25] | 1533.24 | 0.5481 | 1324.94 | 0.4636 |
| Causal LSTM [26] | 1162.65 | 0.6693 | 1203.48 | 0.4886 |
| MIM [28] | 1385.35 | 0.5950 | 1305. | |
| ConvGRU (Sec. 3.2) | 1033.93 | 0.7082 | 1257.3 | |
| MIMO-VP [16] | **401.59** | **0.9049** | **1021.5** | |
| SimVP [8] | 673.01 | 0.8283 | 1280. | |
| SRVP | 935.02 | 0.7474 | 73.18 | 0.5533 |

↓ 39.02% in MSE
↑ 36.36% in SSIM



↑ 6.3% in FLOPs
↓ 9.6% in MSE

## ✓ Impact of Feature Reinforcement



| Method | MSE ↓ | PSNR ↑ | SSIM ↑ |
|---|---|---|---|
| SRVP | 935.02 | 18.6160 | 0.7474 |
| Without-SA | 968.64 | 18.4719 | 0.7365 |
| Without-RFA | **1023.84** | **18.2027** | **0.6983** |
| Without-CrossAtt | 973.16 | 18.4690 | 0.7314 |

# Outperforming on Complex Motion

✓ **SRVP Achieved Superior Performance**

| Method | MSE ↓ | PSNR ↑ | SSIM ↑ |
|---|---|---|---|
| ConvLSTM [19] | 936.69 | 17.1198 | 0.5322 |
| ST-LSTM [25] | 386.21 | 21.9665 | 0.7294 |
| Causal LSTM [26] | 408.14 | 21.5750 | 0.7158 |
| MIM [28] | 445.27 | 21.3214 | 0.6947 |
| ConvGRU (Sec. 3.2) | 418.09 | 21.5636 | 0.7115 |
| MIMO-VP [16] | 502.15 | 21.0167 | 0.7276 |
| SimVP [8] | 452.65 | 21.4365 | 0.7179 |
| SRVP | **383.16** | **22.0322** | **0.7360** |

↓ 59.09% in MSE
↑ 29.69% in PSNR
↑ 38.29% in SSIM

✓ **Outperform than RNN-free Models**

| Method | Horizon | | | |
|---|---|---|---|---|
| | 1 | 2 | 3 | 4 |
| ConvLSTM [19] | 0.8960 | 0.8799 | 0.8617 | 0.8440 |
| ST-LSTM [25] | 0.9656 | **0.9519** | **0.9397** | **0.9285** |
| Causal LSTM [26] | 0.9639 | 0.9502 | 0.9375 | 0.9262 |
| MIM [28] | 0.9599 | 0.9446 | 0.9313 | 0.9185 |
| ConvGRU (Sec. 3.2) | 0.9643 | 0.9489 | 0.9354 | 0.9235 |
| MIMO-VP [16] | 0.9303 | 0.9154 | 0.9034 | 0.8924 |
| SimVP [8] | 0.9415 | 0.9286 | 0.9194 | 0.9122 |
| SRVP | **0.9659** | 0.9512 | 0.9386 | 0.9272 |

1.64 to 3.9% in SSIM

**Stronger Robustness and Consistency**



SRVP

(GT-P)

VS

MIMO-VP

(GT-P)

**Producing Errors even in Static Regions**

# Conclusion

**Ground Truth**

**RNN-based**

➢ Error Accumulation
➢ Blurring

**RNN-free**

➢ Temporal Inconsistency (Inaccurate Position)
➢ Sensitive to image size changes

**SRVP**

➢ Improving Blurring Problems
➢ Consistent Performance
➢ Future Work
   → Enhancing Recollection Ability

Thank you for your attention!

https://github.com/yuseonk/SRVP