# Deep Learning Models for Sequential Data Analysis

Mingsheng Long School of Software, Tsinghua University July 29, 2022





### Spatiotemporal Data Analysis



Radar Echo



Traffic Map



Robotics



Pedestrian Motion



Automatic stations

Temperature, Wind speed, Precipitation...

### Predictive Learning



 $\mathcal{X}_{t-2}$ 



## Predictive Learning

#### Obstacles to Progress in Al

#### Y LeCun How Much Information Does the Machine Need to Predict? Y LeCun "Pure" Reinforcement Learning (cherry) The machine predicts a scalar reward given once in a while. A few bits for some samples Supervised Learning (icing) The machine predicts a category or a few numbers for each input Predicting human-supplied data ▶ $10 \rightarrow 10,000$ bits per sample Unsupervised/Predictive Learning (cake) The machine predicts any part of its input for any observed part. Predicts future frames in videos Millions of bits per sample (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)

Physical world, digital world, people,....

Machines need to learn/understand how the world works

They need to acquire some level of common sense

They need to learn a very large amount of background knowledge

- Through observation and action
- Machines need to perceive the state of the world
- So as to make accurate predictions and planning
- Machines need to update and remember estimates of the state of the world
- Paying attention to important events. Remember relevant events

#### Machines neet to reason and plan

Predict which sequence of actions will lead to a desired state of the world

Intelligence & Common Sense =

Perception + Predictive Model + Memory + Reasoning & Planning

Credit: Yann LeCun [NIPS 2016]

### Our Research









#### Spatiotemporal Predictive Learning

Time Series Forecasting

**PredRNN-V2** for Precipitation Nowcasting

Autoformer for 2022 Beijing Olympics



A Foundation Model for General Task

## Spatiotemporal Predictive Learning



Physical world understanding [Lerer et al. ICML16; Wu et al. NeurIPS17]





#### Precipitation Nowcasting [Wang et al. NeurIPS17; CVPR19]



Action intention estimation [Wang et al. ICLR19]

Robotics control [Ha & Schmidhuber, NeurIPS18]

### Timeline



Spatiotemporal Information Modeling

## Spatiotemporal Modeling in ConvLSTM

CNN for spatial information, RNN for temporal information



## Spatiotemporal Modeling in ConvLSTM



(2) Temporal Dimension

(1) Spatial Dimension

Hierarchical

visual features





[Zeiler & Fergus. ECCV 2014]

$$g_{t} = \tanh(w_{xg} * X_{t} + w_{hg} * H_{t-1} + b_{g})$$

$$i_{t} = \sigma(w_{xi} * X_{t} + w_{hi} * H_{t-1} + w_{ci} \odot C_{t-1} + b_{i})$$

$$f_{t} = \sigma(w_{xf} * X_{t} + w_{hf} * H_{t-1} + w_{cf} \odot C_{t-1} + b_{f})$$

$$C_{t} = f_{t} \odot C_{t-1} + i_{t} \odot g_{t}$$
Capture the long- and short-term dynamics with  $C_{t}$ 

$$o_{t} = \sigma(w_{xo} * X_{t} + w_{ho} * H_{t-1} + w_{co} \odot C_{t} + b_{o})$$

$$H_{t} = o_{t} \odot \tanh(C_{t}),$$

### Timeline



PredRNN-V2

2022

- Architecture: Spatiotemporal Memory Flow to enhance state transition
- Core Module: ST-LSTM with Memory Decoupling to cover long- and short-term dynamics
- Training Strategy: Reverse Scheduled Sampling to bridge the gap between encoder and decoder

Yunbo Wang, Haixu Wu, Jianjin Zhang, Zhifeng Gao, Jianmin Wang, Philip S. Yu, Mingsheng Long. PredRNN: A Recurrent Neural Network for Spatiotemporal Predictive Learning. **TPAMI**, 2022.

## Comparison with ConvLSTM



### Comparison with ConvLSTM



### Comparison with ConvLSTM



#### PredRNN-V2



Spatiotemporal Memory M<sub>t</sub>: Transit in a (1) vertical and (2) zig-zag way

PredRNN-V2



(1) Vertical Transition

Unify the hierarchical visual features

$$g_{t} = \tanh(w_{xg} * X_{t} \mathbb{1}_{\{l=1\}} + w_{hg} * H_{t}^{l-1} + b_{g})$$
  

$$i_{t} = \sigma(w_{xi} * X_{t} \mathbb{1}_{\{l=1\}} + w_{hi} * H_{t}^{l-1} + w_{ci} * M_{t}^{l-1} + b_{i})$$
  

$$f_{t} = \sigma(w_{xf} * X_{t} \mathbb{1}_{\{l=1\}} + w_{hf} * H_{t}^{l-1} + w_{cf} * M_{t}^{l-1} + b_{f})$$
  

$$M_{t}^{l} = f_{t} \odot M_{t}^{l-1} + i_{t} \odot g_{t}$$
  

$$o_{t} = \sigma(w_{xo} * X_{t} \mathbb{1}_{\{l=1\}} + w_{ho} * H_{t}^{l-1} + w_{co} * M_{t}^{l} + b_{o})$$
  

$$H_{t}^{l} = o_{t} \odot \tanh(M_{t}^{l}),$$

PredRNN-V2



PredRNN-V2



(2) Zig-zag Transition *M* is in "Slow transition" → short-term dynamics *C* is in "Fast transition" → long-term dynamics



Achieve the long- and short-term dynamics decomposition.

PredRNN-V2



Can C and M be jointly learned as we expected?



PredRNN-V2

Red points:  $\Delta C$ ; Black Points  $\Delta M$ 



Successfully decouple C and M, release the model prediction capability.

PredRNN-V2



(a) Responses to sudden step changes

(b) Responses to long-term and short-term dynamics

Memory decoupling empowers the model with stronger modeling capability in long- and short-term dynamics



### PredRNN-V2

Encoder inputs: last frame prediction  $\hat{X}_t \rightarrow \text{ground truth } X_t$ 

 $\checkmark$  Decoder inputs: ground truth  $\chi_t \rightarrow$  last frame prediction  $\hat{\chi}_t$ 



Reverse Scheduled Sampling (RSS)

#### PredRNN-V2

Benefits of Reverse Scheduled Sampling



harder task, which can force the model to

memorize more information

(1) Optimization: Alleviates the gradient optimization problem in RNNs

## PredRNN-V2 for Traffic

- ✓ ST-LSTM can be used as a general unit and combined with U-Net
- Memory decoupling and RSS are general techniques to improve performance

		PredRNN-V2						
			$\mathbf{N}$					
Backbone	Recurrent unit	Decoupl.	RSS I	$MSE(10^{-3})$				
U-Net [88]	None	×	×	6.992				
U-Net + PredRNN	ST-LSTM ST-LSTM	× O	× O	7.035 <b>5.135</b>				
CrevNet [68]	ConvLSTM ST-LSTM ST-LSTM	× × O	× × O	6.789 6.613 6.506				



CrevNet [Yu et al. ICLR20]

### PredRNN-V2 for Robotics



Action fusion:  $\mathcal{V}_t^l = (W_{hv} * \mathcal{H}_{t-1}^l) \odot (W_{av} * \mathcal{A}_{t-1})$  $\mathcal{H}_t^l, \mathcal{C}_t^l, \mathcal{M}_t^l = \text{ST-LSTM}(\mathcal{X}_t, \mathcal{V}_t^l, \mathcal{C}_{t-1}^l, \mathcal{M}_t^{l-1}),$ 

 Action condition prediction
 Prediction model can be the world simulator for controlling





Ground truth



SV2P [Babaeizadeh et al. ICLR17]

### PredRNN-V2 for Precipitation Nowcasting



PredRNN-V2 surpasses all other five methods and achieves the best performance in long-term ( $\geq$ 60min) and high density ( $\geq$  35dbz) radar prediction

### Time Series Forecasting

#### Wide Applications



Past Observations

Future Time Series



### Transformer



Modeling temporal dependencies globally with point-wise Self-Attention

Longer time series come with high complexity for Self-Attention and complex temporal patterns to discover



### Timeline



- Decomposition architecture to split more predictable components
- Auto-Correlation to discover the temporal dependencies among sub-series with O(L log L) complexity



Haixu Wu, Jiehui Xu, Jianmin Wang, Mingsheng Long. Autoformer: Decomposition Transformers with Auto-Correlation for Long-Term Series Forecasting. NeurIPS, 2021.

## Comparison with Transformer



## Comparison with Transformer



## Comparison with Transformer

	Transformer	Autoformer			
	<ul> <li>Point-wise Self-Attention is O(L<sup>2</sup>)</li> <li>Adopt sparse version for efficiency resulting in the trade-off dilemma ×</li> </ul>	Auto-Correlation mechanism based on stochastic process theory with inherent <b>O(L log L)</b> complexity			
Computation Efficiency	Auto-Correlation From Autoformer Full Attention From Transformer LSH Attention From Reformer ProbSparse Attention From Informer 0 0 0 0 0 0 0 0 0 0 0 0 0	(b) Running Time Efficiency Analysis			



Decomposition architecture for intricate temporal patterns



Progressive decomposition capacity





Focus on seasonal part modeling, Provide cross information for decoder





Benefit from the deep decomposition, the seasonal part is highlighted with periodicity

Conduct the dependencies discovery and representation aggregation at the series level



Series-wise Auto-Correlation towards information utilization bottleneck

Discover period-based dependencies with autocorrelation in stochastic process:

$$\mathcal{R}_{\mathcal{X}\mathcal{X}}(\tau) = \lim_{L \to \infty} \frac{1}{L} \sum_{t=0}^{L-1} \mathcal{X}_t \mathcal{X}_{t-\tau}.$$

Autocorrelation reflects the time delay similarity,

and corresponds to the confidence of period estimation



Larger autocorrelation  $\mathcal{R}(\tau)$  means

- stronger time delay similarity w.r.t.  $\boldsymbol{\tau}$
- more confidence of period length as  $\boldsymbol{\tau}$



Efficient computation of autocorrelation with Wiener–Khinchin theorem by FFT

 Discover period-based dependencies with O(L log L) complexity



- 1. Select the Top-k period lengths
- 2. Softmax-normalization
- 3. Align delayed series and aggregate sub-series representations

$$\tau_{1}, \cdots, \tau_{k} = \arg \operatorname{Topk}_{\tau \in \{1, \cdots, L\}} (\mathcal{R}_{\mathcal{Q}, \mathcal{K}}(\tau))$$
$$\widehat{\mathcal{R}}_{\mathcal{Q}, \mathcal{K}}(\tau_{1}), \cdots, \widehat{\mathcal{R}}_{\mathcal{Q}, \mathcal{K}}(\tau_{k}) = \operatorname{SoftMax} (\mathcal{R}_{\mathcal{Q}, \mathcal{K}}(\tau_{1}), \cdots, \mathcal{R}_{\mathcal{Q}, \mathcal{K}}(\tau_{k}))$$
$$\operatorname{AutoCorrelation}(\mathcal{Q}, \mathcal{K}, \mathcal{V}) = \sum_{i=1}^{k} \operatorname{Roll}(\mathcal{V}, \tau_{k}) \widehat{\mathcal{R}}_{\mathcal{Q}, \mathcal{K}}(\tau_{k}),$$

Benchmark			Trar	nsforme	rs	LSTM	1s	TCN	
	Models	Autoformer	Informer[41]	LogTrans[20]	Reformer[17]	LSTNet[19]	LSTM[13]		Prediction Accuracy Relative Promotion (In MSE)
Fneray	Metric * 96 192 336 720	MSE         MAE           0.255         0.339           0.281         0.340           0.339         0.372           0.422         0.419	MSE         MAE           0.365         0.453           0.533         0.563           1.363         0.887           3.379         1.388	MSE         MAE           0.768         0.642           0.989         0.757           1.334         0.872           3.048         1.328	MSEMAE0.6580.6191.0780.8271.5490.9722.6311.242	MSEMAE3.1421.3653.1541.3693.1601.3693.1711.368	MSEMAE2.0411.0732.2491.1122.5681.2382.7201.287	B         3.041         1.330           2         3.072         1.339           3         3.105         1.348           7         3.135         1.354	Input-96-predict-336 <b>1</b> 74%
Energy	Electricity 192 336 720	0.201 0.317 0.222 0.334 0.231 0.338 0.254 0.361	0.274 0.368 0.296 0.386 0.300 0.394 0.373 0.439	0.2580.3570.2660.3680.2800.3800.2830.376	0.3120.4020.3480.4330.3500.4330.3400.420	0.680 0.645 0.725 0.676 0.828 0.727 0.957 0.811	0.375 0.437 0.442 0.473 0.439 0.473 0.980 0.814	7       0.985       0.813         8       0.996       0.821         8       1.000       0.824         4       1.438       0.784	Input-96-predict-336 <b>18%</b>
Economics	96 192 336 720	0.197 0.323 0.300 0.369 0.509 0.524 1.447 0.941	0.847 0.752 1.204 0.895 1.672 1.036 2.478 1.310	0.9680.8121.0400.8511.6591.0811.9411.127	1.0650.8291.1880.9061.3570.9761.5101.016	1.5511.0581.4771.0281.5071.0312.2851.243	1.453 1.049 1.846 1.179 2.136 1.231 2.984 1.427	3.004       1.432         3.048       1.444         3.113       1.459         3.150       1.458	Input-96-predict-336 <b>1</b> 61%
Traffic	96 192 336 720	0.613 0.388 0.616 0.382 0.622 0.337 0.660 0.408	0.719 0.391 0.696 0.379 0.777 0.420 0.864 0.472	0.6840.3840.6850.3900.7330.4080.7170.396	0.7320.4230.7330.4200.7420.4200.7550.423	1.107 0.685 1.157 0.706 1.216 0.730 1.481 0.805	0.843 0.453 0.847 0.453 0.853 0.455 1.500 0.805	81.4380.78481.4630.79451.4790.79951.4990.804	Input-96-predict-336 <b>1</b> 5%
Weather	Meather 192 336 720	0.266 0.336 0.307 0.367 0.359 0.395 0.419 0.428	0.300 0.384 0.598 0.544 0.578 0.523 1.059 0.741	0.4580.4900.6580.5890.7970.6520.8690.675	0.6890.5960.7520.6380.6390.5961.1300.792	0.594 0.587 0.560 0.565 0.597 0.587 0.618 0.599	0.369 0.406 0.416 0.435 0.455 0.454 0.535 0.520	5       0.615       0.589         5       0.629       0.600         6       0.639       0.608         0       0.639       0.610	Input-96-predict-336 1 21%
Disease	$\begin{array}{c c}1&24\\36\\48\\60\end{array}$	3.4831.2873.1031.1482.6691.0852.7701.125	5.764 1.677 4.755 1.467 4.763 1.469 5.264 1.564	4.4801.4444.7991.4674.8001.4685.2781.560	4.4001.3824.7831.4484.8321.4654.8821.483	6.0261.7705.3401.6686.0801.7875.5481.720	5.914 1.734 6.631 1.845 6.736 1.857 6.870 1.879	6.624 1.830 6.858 1.879 6.968 1.892 7.127 1.918	Input-24-predict-48 1 43%
* <i>ETT</i> means the ETTm2. See supplementary materials for the <b>full benchmark</b> of ETTh1, ETTh2, ETTm1.									

#### Ablation of Decomposition

Input-96	Trans	sforme	r[35]	Info	ormer[4	41]	Log	Trans[	17]	Ref	ormer[	20]	Prom	otion
Predict-O	Origin	Sep	Ours	Origin	Sep	Ours	Origin	Sep	Ours	Origin	Sep	Ours	Sep	Ours
96	0.604	0.311	0.204	0.365	0.490	0.354	0.768	0.862	0.231	0.658	0.445	0.218	0.069	0.347
192	1.060	0.760	0.266	0.533	0.658	0.432	0.989	0.533	0.378	1.078	0.510	0.336	0.300	0.562
336	1.413	0.665	0.375	1.363	1.469	0.481	1.334	0.762	0.362	1.549	1.028	0.366	0.434	1.019
720	2.672	3.200	0.537	3.379	2.766	0.822	3.048	2.601	0.539	2.631	2.845	0.502	0.079	2.332

#### Ablation of Auto-Correlation

Input Leng	th I		96			192			336	
Prediction Ler	ngth O	336	720	1440	336	720	1440	336	720	1440
Auto-	MSE	0.339	0.422	0.555	0.355	0.429	0.503	0.361	0.425	0.574
Correlation	MAE	0.372	0.419	0.496	0.392	0.430	0.484	0.406	0.440	0.534
Full	MSE	0.375	0.537	0.667	0.450	0.554	-	0.501	0.647	-
Attention[35]	MAE	0.425	0.502	0.589	0.470	0.533		0.485	0.491	-
LogSparse	MSE	0.362	0.539	0.582	0.420	0.552	0.958	0.474	0.601	-
Attention[20]	MAE	0.413	0.522	0.529	0.450	0.513	0.736	0.474	0.524	
LSH	MSE	0.366	0.502	0.663	0.407	0.636	1.069	0.442	0.615	-
Attention[17]	MAE	0.404	0.475	0.567	0.421	0.571	0.756	0.476	0.532	
ProbSparse	MSE	0.481	0.822	0.715	0.404	1.148	0.732	0.417	0.631	1.133
Attention[41]	MAE	0.472	0.559	0.586	0.425	0.654	0.602	0.434	0.528	0.691

- Decomposition outperforms separately forecasting convention especially in the long-term setting
- Decomposition architecture can be generalized to other Transformers

 Under various input-predict settings, Auto-Correlation outperforms Self-Attention and its variants

## Autoformer for 2022 Beijing Olympics









Indoors: Temperature

Outdoors: Wind speed

Provides online forecast service of temperature and wind speed for the 2022 Beijing Winter Olympics, assists athletes preparation and schedule planning, works as a solid support for the competition.

Achieves 10-minute real-time temperature and wind speed forecast based on meteorological observation, and achieves 23% lower forecast error than the mainstream numerical prediction methods.

#### Foundation Models



[Data Universal] Learn from various modalities [Task Universal] Adapt to a wide range of downstream tasks

Bommasani et al. On the Opportunities and Risks of Foundation Models. Arxiv 2021.

#### General Relation Modeling of Transformers



#### Quadratic Complexity in Self-Attention



#### Quadratic Complexity in Self-Attention



Pair-wise Relation Modeling: Attention $(Q, K, V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$ 

 $\mathcal{O}(n^2 d)$ Can we remove Softmax function?  $(QK^T)V = Q(K^TV) \implies \mathcal{O}(n^2 d) \rightarrow \mathcal{O}(nd^2)$  Recap: Softmax function



Bridle et al. Training stochastic model recognition algorithms as networks can lead to maximum mutual information estimation of parameters. NeurIPS 1989.

Recap: Softmax function



Bridle et al. Training stochastic model recognition algorithms as networks can lead to maximum mutual information estimation of parameters. NeurIPS 1989.

#### Flow Network Theory

![](_page_53_Figure_1.jpeg)

#### Attention: A Flow Network View

![](_page_54_Figure_1.jpeg)

![](_page_54_Figure_2.jpeg)

(a) Inner View

#### Attention: A Flow Network View

![](_page_55_Figure_1.jpeg)

#### Flowformer: Conservation in Attention

![](_page_56_Figure_1.jpeg)

[Incoming Flow Conservation]: Competition among Source tokens [Outgoing Flow Conservation]: Competition among Sink tokens

![](_page_57_Figure_1.jpeg)

![](_page_58_Figure_1.jpeg)

![](_page_58_Figure_2.jpeg)

Incoming flow:  $I_i = \phi(Q_i) \sum_j \phi(K_j)^T$ 

Incoming flow conservation:  $\frac{\phi(Q)}{I}$ 

Incoming flow:  $\frac{\phi(Q_i)}{I_i} \sum_j \phi(K_j)^T = \frac{I_i}{I_i} = 1$ 

![](_page_59_Figure_1.jpeg)

![](_page_59_Figure_2.jpeg)

Incoming flow:  $I_i = \phi(Q_i) \sum_j \phi(K_j)^T$ 

Incoming flow conservation:  $\frac{\phi(Q)}{I}$ 

Conserved outgoing flow:  $\hat{\boldsymbol{o}} = \phi(\boldsymbol{K}) \sum_{i} \frac{\phi(Q_i)^T}{I_i}$ 

![](_page_60_Figure_1.jpeg)

![](_page_60_Figure_2.jpeg)

Outgoing flow:  $O_i = \phi(K_i) \sum_j \phi(Q_j)^T$ 

Outgoing flow conservation:  $\frac{\phi(K)}{O}$ 

Outgoing flow: 
$$\frac{\phi(K_i)}{O_i} \sum_j \phi(Q_j)^T = \frac{O_i}{O_i} = 1$$

![](_page_61_Figure_1.jpeg)

![](_page_61_Figure_2.jpeg)

Outgoing flow:  $O_i = \phi(K_i) \sum_j \phi(Q_j)^T$ 

Outgoing flow conservation:  $\frac{\phi(K)}{o}$ 

Conserved incoming flow:  $\hat{I} = \phi(Q) \sum_{j} \frac{\phi(K_{j})^{T}}{o_{j}}$ 

![](_page_62_Figure_1.jpeg)

Competition: 
$$\widehat{\mathbf{V}} = \operatorname{Softmax}(\widehat{\mathbf{O}}) \odot \mathbf{V}$$
  
Aggregation:  $\mathbf{A} = \frac{\phi(\mathbf{Q})}{\mathbf{I}} (\phi(\mathbf{K})^{\mathsf{T}} \widehat{\mathbf{V}})$   
Allocation:  $\mathbf{R} = \operatorname{Sigmoid}(\widehat{\mathbf{I}}) \odot \mathbf{A}$ ,

Successfully bring the <u>Competition Mechanism</u> Into Attention design to avoid trivial attention

#### Flowformer: Efficiency and Universality

![](_page_63_Figure_1.jpeg)

[Efficiency]: All the calculations are in linear complexity.

[Universality]: The whole design is based on flow network without specific inductive biases.

#### Flowformer Experiments

![](_page_64_Picture_1.jpeg)

BENCHMARKS	TASK	VERSION	Length
LRA (2020C)	SEQUENCE	NORMAL	1000~4000
WIKITEXT (2017)	LANGUAGE	CAUSAL	512
IMAGENET (2009)	VISION	NORMAL	49~3136
UEA (2018)	TIME SERIES	NORMAL	29~1751
D4RL (2020)	OFFLINE RL	CAUSAL	60

- Extensive tasks (covering 5 mainstream tasks)
- Normal and causal versions
- Various sequence lengths (29-4000)
- Extensive baselines (20+)

#### Flowformer Experiments

Task	Metrics	Flowformer	Performer	Reformer	Vanilla Transformer
Long Sequence Modeling (LRA)	Avg Acc (%) ↑	56.48	51.41	50.67	OOM
Vision Recognization (ImageNet-1K)	Top-1 Acc (%) ↑	80.6	78.1	79.6	78.7
Language Modeling (WikiText-103)	Perplexity $\downarrow$	30.8	37.5	33.6	33.0
Time series classification (UEA)	Avg Acc (%) ↑	73.0	71.5	71.9	71.9
Offline RL (D4RL)	Avg Reward $\uparrow$ Avg Deviation $\downarrow$	<b>73.5</b> ± 2.9	$63.8\pm7.6$	$63.9\pm2.9$	72.2 ± <b>2.6</b>
			•		

Strong performance on all five mainstream tasks within the linear complexity

#### Flowformer

![](_page_66_Figure_1.jpeg)

Quadratic Complexity

![](_page_66_Figure_3.jpeg)

Flowformer

#### Linear complexity w.r.t. sequence length Based on flow network & without specific inductive biases Strong performance in Long Sequence, CV, NLP, Time Series, RL

Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, Mingsheng Long. Flowformer: Linearizing Transformers with Conservation Flows. ICML, 2022.

### Thank You!

![](_page_67_Picture_1.jpeg)

![](_page_67_Picture_2.jpeg)

![](_page_67_Picture_3.jpeg)

![](_page_67_Picture_4.jpeg)

吴海旭

![](_page_67_Picture_6.jpeg)

![](_page_67_Picture_7.jpeg)

![](_page_67_Picture_8.jpeg)

#### 大数据系统软件国家工程研究中心

![](_page_67_Picture_10.jpeg)

![](_page_67_Picture_11.jpeg)