

# Representation Learning for Time Series

---



***Sogang University***

*Vision & Display Systems Lab, Dept. of Electronic Engineering*



***Presented By***

*Yein Park*

# Outline

- Abstract
- Background
- Problem Formulation
- Introduction
- Proposed Method
- Experimental Results
- Conclusion

"SimTS: Rethinking Contrastive Representation Learning for Time Series Forecasting." *arXiv*, 2023.

# Abstract

- Contrastive learning
  - Meaningful representations for image or time series classification
  - Less effectiveness for time series forecasting
    - Task of predicting future state from history context
    - Instance discrimination optimization that cannot be applied directly
  - Generalization limits for different types of time series data
    - High dependence on specific time series characteristics of the construction of positive and negative pairs
- SimTS, a simple representation learning approach for improving time series forecasting
  - Learning to predict the future from the past in the latent space
  - No reliance on negative pairs or specific assumptions

# Background

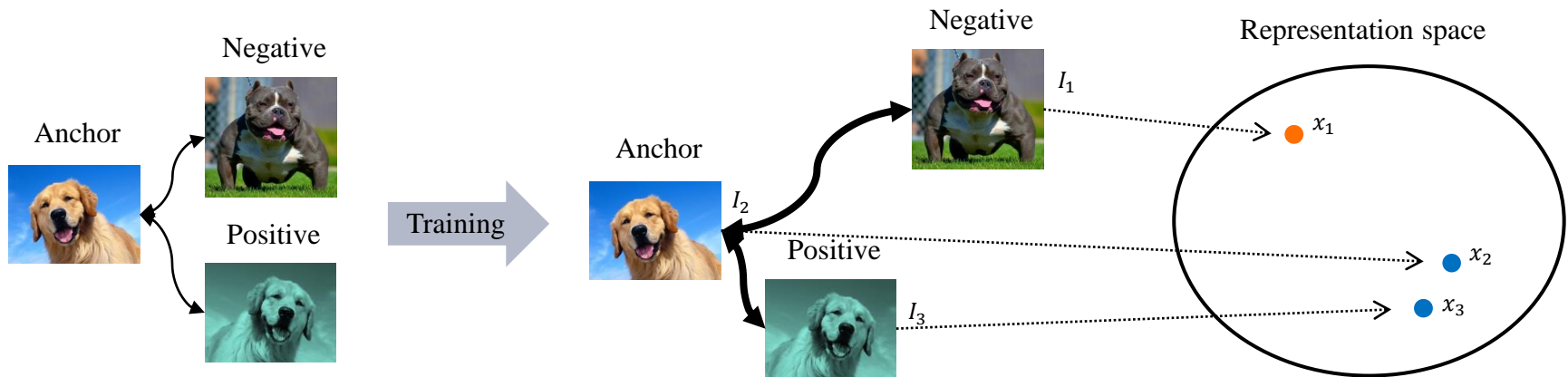
- Contrastive learning

- A method of representation learning

- Representation learning through comparison between input samples

- Learning the representation space so that "similar" data is close and "different" data is far

- For multiple input pairs, the similarity is learned by label



# Background

- Contrastive learning

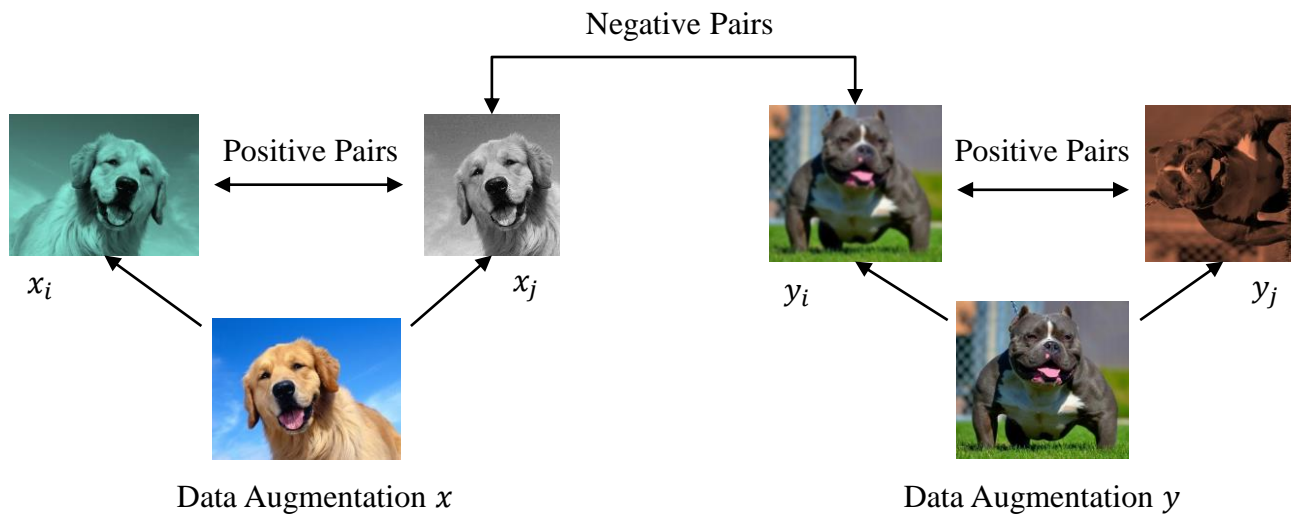
- Generation of input pairs through data augmentation

- Data augmentation in image domain

- ⌘ Positive pair → data augmented on the same image

- ⌘ Negative pair → data augmented on the different images

- ⌘ Random crop, rotation, resizing, shifting, noising, blur, color distortion, perspective distortion



# Background

- Contrastive learning

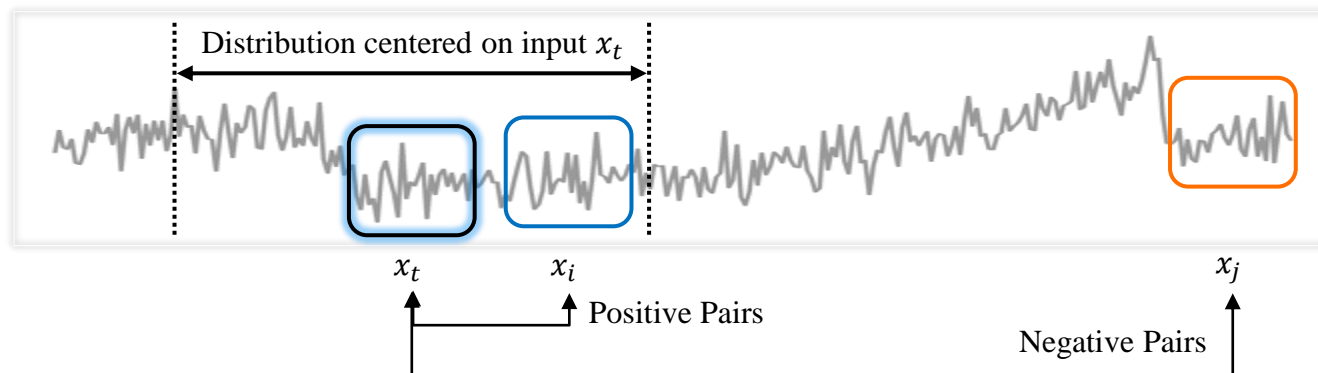
- Generation of input pairs through data augmentation

- Data augmentation in time series domain

- ✧ Input-centered distribution

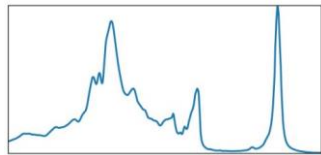
- ✓ Positive pair → belonging to input-centered distribution, neighborhood

- ✓ Negative pair → not belonging to input-centered distribution, non-neighborhood

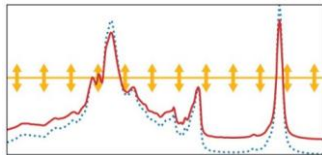


# Background

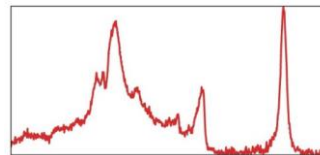
- Contrastive learning
  - Generation of input pairs through data augmentation
    - Data augmentation in time series domain
      - ⊛ Scaling, jittering, window slicing, time warping
      - ⊛ Random masking
      - ⊛ Random warping
      - ⊛ Random reordering



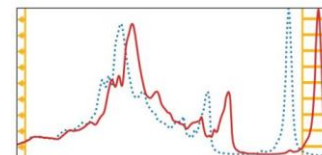
Original



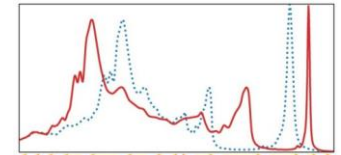
Scaling



Jittering



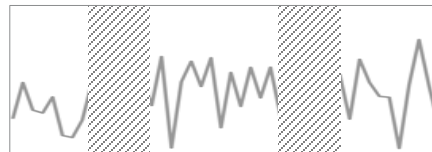
Window slicing



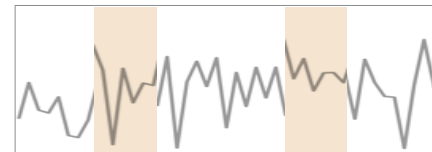
Time warping



Original



Random masking



Random warping

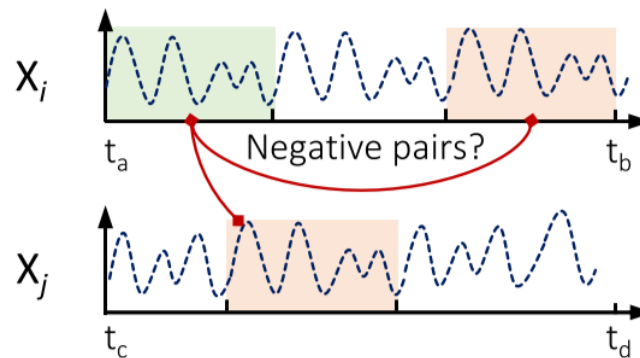


Random reordering



# Problem Formulation

- Contrastive learning
  - Learning representations close to positive samples and far from negative samples
    - Well suited for classification tasks
      - ⊛ The resulting representation contains information to better distinguish different instances of the time series
    - Ineffective for forecasting tasks
      - ⊛ Forecasting the future based on past data
  - Negative sample methods of existing methods that are difficult to trust
    - Time stamp difference, different time stamp



< Existing method of selecting negative pairs >

# Introduction

- Key questions
  - What is important for time series forecasting with contrastive learning?
  - How can we adapt contrastive ideas more effectively to time series forecasting tasks?
  - Are the existing assumptions and techniques for constructing positive and negative pairs reasonable?

# Proposed Method

- Motivation

- What is important for time series forecasting with contrastive learning?
  - Existing methods
    - ✧ They ignore the possibility that a repeating patterns exist within a time series
    - ✧ They disregard the possibility that distinct time series contain similar patterns
- In forecasting task, a good representation should effectively capture the temporal dependencies between past segments and future predictions.
- We emphasize that the temporal dependency has greater significance than the similarity between positive and negative pairs.
- Negative pairs inducing the issue of false repulsion
  - Patterns that repeat across different samples
- → We train an encoder to learn time series representations by predicting its future from historical segments in the latent space.

# Proposed Method

- SimTS: Simple Representation Learning for Time Series Forecasting

- Siamese neural network architecture

- Two identical networks sharing parameters
- *History* encoding path, *future* encoding path

- Multi-scale encoder network  $F_\theta$

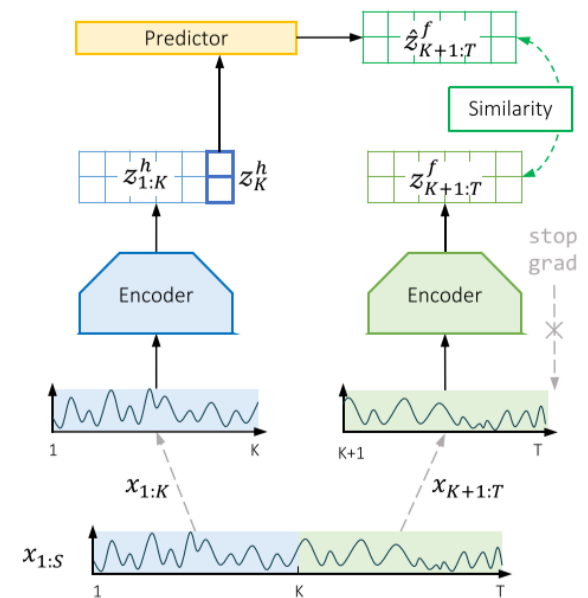
- Projecting raw features into a high dimensional space
- Multiple CNN blocks with different kernel sizes

- Prediction network  $G_\phi$

- Input: last column of encoded history view
- Predicting future in latent space

- Cosine similarity loss

- Considering only positive samples



# Proposed Method

- SimTS

- Objective

- Learning latent representation of *history* segment  
 $X^h = [x_1, x_2, \dots, x_K], \quad 0 < K(201) < T(402)$

- Predicting *future* segment  $X^f = [x_{K+1}, x_{K+2}, \dots, x_T]$

- Multi-scale Encoder network  $F_\theta$

- Inputs: *history* segment  $X^h$ , *future* segment  $X^f$

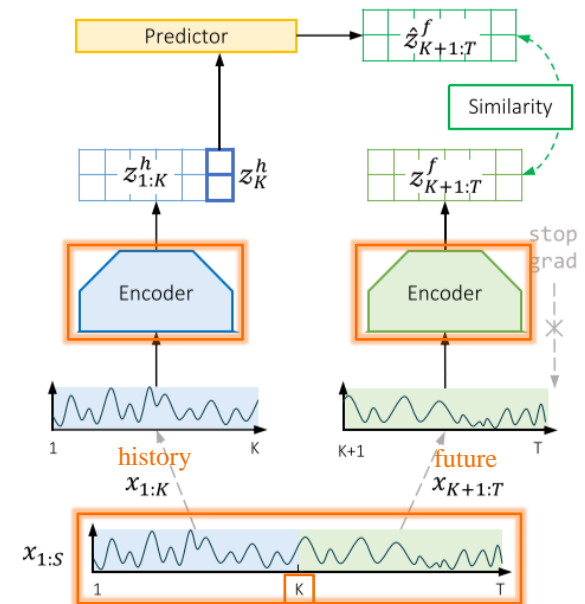
- ⊛ Learning to map them to their latent representations  
 $Z^h, Z^f$

- *History* encoding path

- ⊛  $Z^h = F_\theta(X^h) = [z_1^h, z_2^h, \dots, z_K^h] \in R^{C' \times K}, (C':320)$

- *Future* encoding path

- ⊛  $Z^f = F_\theta(X^f) = [z_{K+1}^f, z_{K+2}^f, \dots, z_T^f] \in R^{C' \times (T-K)}$



# Proposed Method

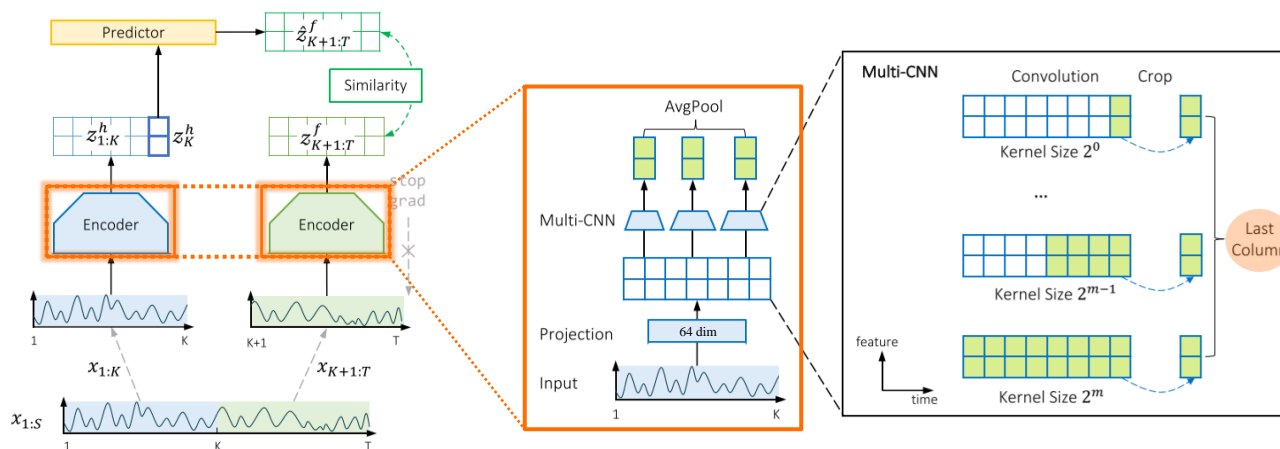
- SimTS

- Multi-scale Encoder network  $F_\theta$

- Convolutional network with multiple filters that have various kernel sizes
    - Extracting both local/global patterns
    - For a time series  $X$  with length  $K$ ,  $m \equiv \lfloor \log_2 K \rfloor + 1$  parallel convolution layers, kernel size  $2^{i-1}$  of  $i$ th convolution

$\because K = 8, 16, 32, 64, \dots \rightarrow m = 4, 5, 6, 7, \dots,$   
 $\rightarrow \text{kernel size} = (1, 2, 4, 8), (1, 2, 4, 8, 16), \dots,$

- Average pooling of the last column of each CNN layer



# Proposed Method

- SimTS

- Prediction network  $G_\phi$

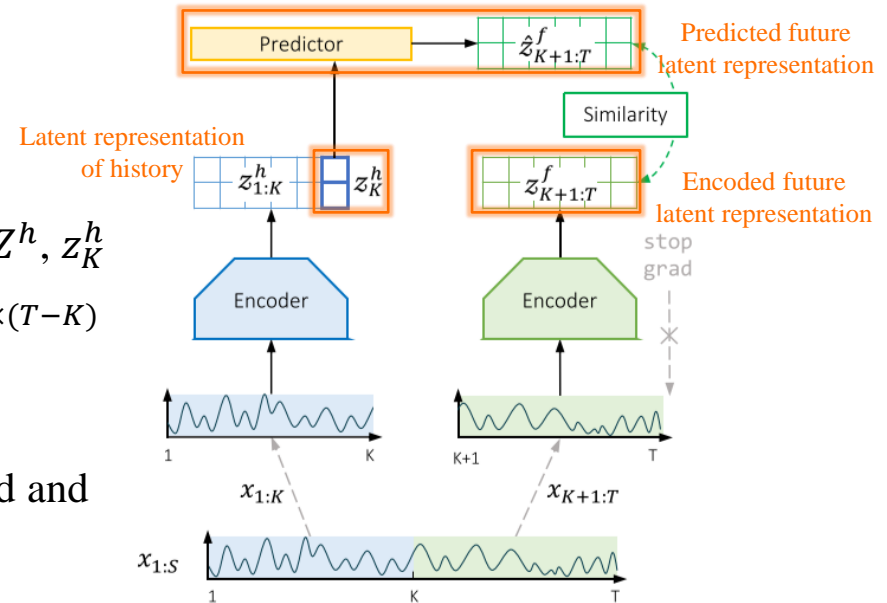
- Predicting *future* latent representations
    - Input: last column of encoded *history* view  $Z^h, z_K^h$
    - $\hat{Z}^f = G_\phi(z_K^h) = [\hat{z}_{K+1}^f, \hat{z}_{K+2}^f, \dots, \hat{z}_T^f] \in R^{C' \times (T-K)}$

- Similarity

- Maximizing the similarity between predicted and encoded latent features

- ⊛ Forcing them to get closer
      - ⊛ Learning historical representation that is informative for the future
      - ⊛ Regarding the predicted  $\hat{Z}^f$  and the encoded  $Z^f$  as the **positive pair**
      - ⊛ Cosine similarity

$$\surd Sim(\hat{Z}^f, Z^f) = -\frac{1}{T-K} \sum_{i=K+1}^T \frac{\hat{z}_i^f}{\|\hat{z}_i^f\|_2} \cdot \frac{z_i^f}{\|z_i^f\|_2}, \quad \|\cdot\|_2: l_2\text{-norm}$$



# Proposed Method

- SimTS

- Stop-gradient Operation

- Using same encoder for both *history* and *future* encoding path

- ⚡ Problem of optimizing the encoder by pushing encoded *future*  $Z^f$  towards predicted *future*  $\hat{Z}^f$

- Applying it to *future* encoding path

- ⚡  $\hat{Z}^f$  can only move towards  $Z^f$  in the latent space

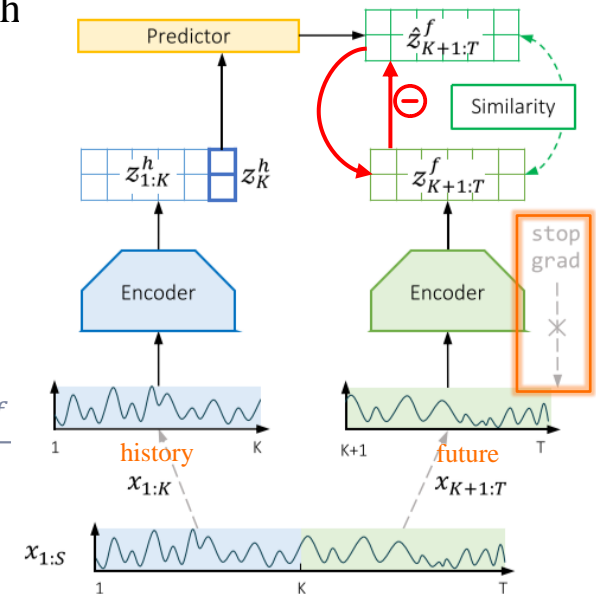
- Encoder network

- ⚡ Unable to receive updates from future representations  $Z^f$

- ⚡ Constrained to only optimize the history representation  $Z^h$  and its prediction  $\hat{Z}^f$

- Loss function

$$\text{⚡ } L_{\theta, \phi}(X^h, X^f) = \text{Sim}\left(G_{\theta}\left(F_{\theta}(X^h)\right), F_{sg(\theta)}(X^f)\right) = \text{Sim}(\hat{Z}^f, sg(Z^f))$$





# Experimental Results

- Experimental setup
  - Input size
    - Length  $T = 402$
    - First 201 corresponding to history view, subsequent 201 corresponding to future view
    - Projected to a 64-dim latent space in projection
    - Projected to a 320-dim latent space in multi-scale encoder
  - Datasets
    - Electricity Transformer Temperature
      - ⌘ ETTh1, ETTh2, ETTm1, ETTm2
    - Exchange-Rate
    - Weather
  - Training, validation, test sets in the ratio of 6:2:2
  - 500 epoch, learning rate 0.001, batch size 8

# Experimental Results

using negative pairs

Methods		Random masking						Scaling, shifting, jittering					
		Unsupervised Representation Learning						End-to-end Forecasting					
		Ours		TS2Vec		TNC		CoST		Informer		TCN	
L		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	24	<b>0.377</b>	<b>0.422</b>	0.590	0.531	0.708	0.592	<u>0.386</u>	<u>0.429</u>	0.577	0.549	0.583	0.547
	48	<b>0.427</b>	<b>0.454</b>	0.624	0.555	0.749	0.619	<u>0.437</u>	<u>0.464</u>	0.685	0.625	0.670	0.606
	168	<b>0.638</b>	<b>0.577</b>	0.762	0.639	0.884	0.699	<u>0.643</u>	<u>0.582</u>	0.931	0.752	0.811	0.680
	336	<u>0.815</u>	<b>0.685</b>	0.931	0.728	1.020	0.768	<b>0.812</b>	<u>0.679</u>	1.128	0.873	1.132	0.815
	720	<b>0.956</b>	<b>0.771</b>	1.063	0.799	1.157	0.830	<u>0.970</u>	<u>0.771</u>	1.215	1.869	1.165	0.813
ETTh2	24	<b>0.336</b>	<b>0.434</b>	<u>0.424</u>	<u>0.489</u>	0.612	0.595	0.447	0.502	0.720	0.665	0.935	0.754
	48	<b>0.564</b>	<b>0.571</b>	<u>0.619</u>	<u>0.605</u>	0.840	0.716	0.699	0.637	1.457	1.001	1.300	0.911
	168	<b>1.407</b>	<b>0.926</b>	1.845	1.074	2.359	1.213	<u>1.549</u>	<u>0.982</u>	3.489	1.515	4.017	1.579
	336	<b>1.640</b>	<b>0.996</b>	2.194	1.197	2.782	1.349	<u>1.749</u>	<u>1.042</u>	2.723	1.340	3.460	1.456
	720	<b>1.878</b>	<b>1.065</b>	2.636	1.370	2.753	1.394	<u>1.971</u>	<u>1.092</u>	3.467	1.473	3.106	1.381
ETTm1	24	<b>0.232</b>	<b>0.314</b>	0.453	0.444	0.522	0.472	<u>0.246</u>	<u>0.329</u>	0.323	0.369	0.522	0.472
	48	<b>0.311</b>	<b>0.368</b>	0.592	0.521	0.695	0.567	<u>0.381</u>	<u>0.386</u>	0.494	0.503	0.542	0.508
	96	<b>0.360</b>	<b>0.402</b>	0.635	0.554	0.731	0.595	<u>0.378</u>	<u>0.419</u>	0.678	0.614	0.666	0.578
	288	<b>0.450</b>	<b>0.467</b>	0.693	0.597	0.818	0.649	<u>0.472</u>	<u>0.486</u>	1.056	0.786	0.991	0.735
	672	<b>0.612</b>	<b>0.563</b>	0.782	0.653	0.932	0.712	<u>0.620</u>	<u>0.574</u>	1.192	0.926	1.032	0.756
ETTm2	24	<b>0.108</b>	<b>0.223</b>	0.180	0.293	0.185	0.297	<u>0.122</u>	<u>0.244</u>	0.173	0.301	0.180	0.324
	48	<b>0.164</b>	<b>0.285</b>	0.244	0.350	0.264	0.360	<u>0.183</u>	<u>0.305</u>	0.303	0.409	0.204	0.327
	96	<b>0.271</b>	<b>0.376</b>	0.360	0.427	0.389	0.458	<u>0.294</u>	<u>0.394</u>	0.365	0.453	3.041	1.330
	288	<b>0.716</b>	<b>0.646</b>	<u>0.723</u>	<u>0.639</u>	0.920	0.788	<u>0.723</u>	0.652	1.047	0.804	3.162	1.337
	672	<b>1.600</b>	<b>0.979</b>	<u>1.753</u>	<u>1.007</u>	2.164	1.135	1.899	1.073	3.126	1.302	3.624	1.484
Exchange	24	<b>0.059</b>	<b>0.172</b>	0.108	0.252	<u>0.105</u>	<u>0.236</u>	0.136	0.291	0.611	0.626	2.483	1.327
	48	<b>0.135</b>	<b>0.265</b>	0.200	0.341	<u>0.162</u>	<u>0.270</u>	0.250	0.387	0.680	0.644	2.328	1.256
	168	0.713	0.635	<u>0.412</u>	<u>0.492</u>	<b>0.397</b>	<b>0.480</b>	0.924	0.762	1.097	0.825	2.372	1.279
	336	1.409	0.938	<u>1.339</u>	<u>0.901</u>	<b>1.008</b>	<b>0.866</b>	1.774	1.063	1.672	1.036	3.113	1.459
	720	<b>1.628</b>	<b>1.056</b>	2.114	1.125	<u>1.989</u>	<u>1.063</u>	2.160	1.209	2.478	1.310	3.150	1.458
Weather	24	<b>0.298</b>	<b>0.359</b>	<u>0.308</u>	0.364	0.320	0.373	<b>0.298</b>	<u>0.360</u>	0.335	0.381	0.321	0.367
	48	<b>0.359</b>	<b>0.410</b>	<u>0.375</u>	0.417	0.380	0.421	<b>0.359</b>	<u>0.411</u>	0.395	0.459	0.386	0.423
	168	<b>0.426</b>	<b>0.461</b>	0.496	0.506	0.479	0.495	<u>0.464</u>	<u>0.491</u>	0.608	0.567	0.491	0.501
	336	0.504	0.520	0.532	0.533	0.505	0.514	<b>0.497</b>	<b>0.517</b>	0.702	0.620	<u>0.502</u>	<u>0.507</u>
	720	0.535	<u>0.542</u>	0.567	0.558	0.543	0.547	<u>0.533</u>	<u>0.542</u>	0.831	0.731	<b>0.498</b>	<b>0.508</b>
Avg.		<b>0.664</b>	<b>0.562</b>	0.818	0.632	0.909	0.669	<u>0.746</u>	<u>0.603</u>	1.151	0.812	1.560	0.884

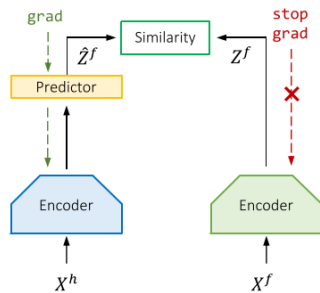
# Experimental Results

- Negative samples
  - Required careful construction
  - TS2Vec and CoST using negative pairs for contrastive learning with poor performance
    - Selection of negative pairs may be inaccurate

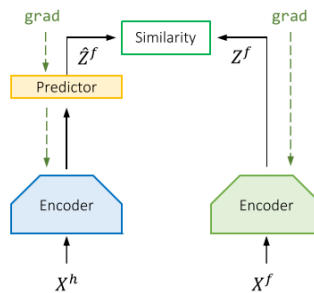
Datasets	ETTh1		ETTh2		ETTm1		ETTm2		Exchange		Weather	
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
SimTS	0.642	0.582	1.165	0.798	0.393	0.423	0.572	0.502	0.789	0.613	0.424	0.458
SimTS w/ neg <sup>†</sup>	0.685	0.632	1.544	0.938	0.392	0.441	0.747	0.572	1.405	0.769	0.434	0.468

# Experimental Results

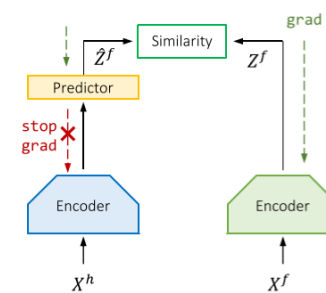
- Stop-Gradient Operation
  - SimTS: applying to *future* encoding path
  - SimTS w/o SG: not applying to any path
  - RevSimTS: applying to *history* encoding path



(a) SimTS



(b) SimTS w/o SG



(c) RevSimTS

Model	SimTS		SimTS w/o SG <sup>†</sup>		RevSimTS	
	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	<b>0.642</b>	<b>0.582</b>	0.783	0.663	0.762	0.634
ETTh2	<b>1.165</b>	<b>0.798</b>	2.940	1.490	3.128	1.449
ETTm1	<b>0.393</b>	<b>0.432</b>	0.681	0.609	0.551	0.525
ETTm2	<b>0.572</b>	<b>0.502</b>	1.315	0.863	1.186	0.796
Exchange	<b>0.789</b>	<b>0.613</b>	1.808	1.062	1.398	0.900
Weather	<b>0.424</b>	<b>0.458</b>	0.605	0.592	0.485	0.512

# Conclusion

- SimTS for time series forecasting
  - Using a simple encoder to learn representations in latent space without negative pairs
- The effectiveness and generalizability of SimTS
- Our goal to challenge the assumptions and components that are widely used
- Current representation learning methods that cannot be universally applicable to different types of time series data

Thanks!