

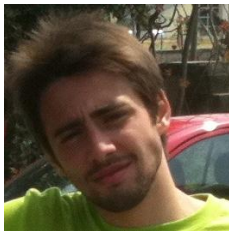
# Whittle Networks: A Deep Likelihood Model for Time Series



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT



Zhongjie Yu<sup>1</sup>



Fabrizio Ventola<sup>1</sup>



Kristian Kersting<sup>1,2</sup>

<sup>1</sup>Department of Computer Science, TU Darmstadt, Germany

<sup>2</sup>Centre for Cognitive Science, TU Darmstadt, and Hessian Center for AI

{yu, ventola, kersting}@cs.tu-darmstadt.de

MADESI

KompA+KI

AIPHES  
ADAPTIVE PROCESSING OF INFORMATION FROM HETEROGENEOUS SOURCES

LOEWE  
Exzellente Forschung für  
Hessens Zukunft

Centre for  
Cognitive  
Science

hessian.AI

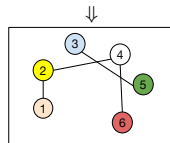
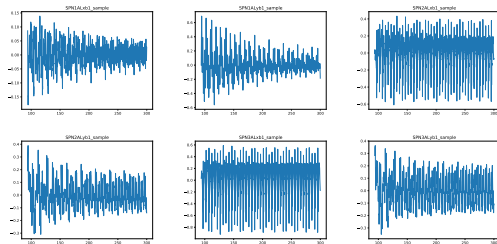


# A Real World Motivation

## Multivariate Time Series Analysis on the Edge



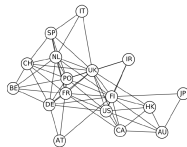
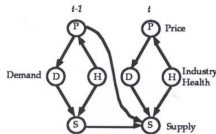
sensor data from wind turbine



Limitations of time series graphical models (TGMs)<sup>[1,2,3]</sup>:

- Inference is exponential in the worst case
- Sample size required for accurate learning is also exponential
- Learning requires inference as subroutine i.e. can take exponential time

⇒ In other words, **TGMs are intractable!**



[1] Tank, A., Foti, N. J., and Fox, E. B. Bayesian structure learning for stationary time series. In UAI, 2015.

[2] Dahlhaus, R. Graphical interaction models for multivariate time series. *Metrika*, 2000.

[3] Bach, F. R. and Jordan, M. I. Learning graphical models for stationary time series. *IEEE Transactions on Signal Processing*, 2004.



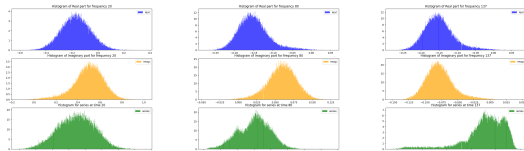
We propose:

- The first tractable probabilistic circuit for modeling the joint distribution of multivariate time series, called Whittle sum-product networks (WSPNs), by introducing complex-valued SPNs.
- Deep likelihood functions for training deep neural networks for time series in an end-to-end fashion, called Whittle Networks.



# Why Whittle Likelihood?

## Time series statistics



The Fourier coefficients (real and imaginary parts) follow Gaussian distribution while the time series in the time domain at each step could follow an arbitrary distribution

Whittle approximation<sup>[4]</sup> – the Fourier coefficients from discrete Fourier transform are independent complex normal distributed:

$$d_{n,k} \sim \mathcal{N}(0, S_k), \quad k = 0, \dots, T-1$$

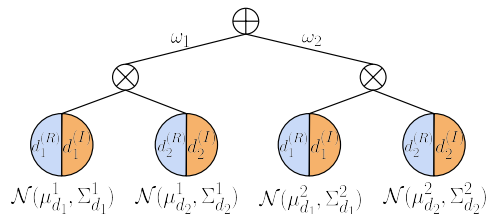
where  $S_k$  is the spectral density matrix:

$$S_k = \sum_{h=-\infty}^{\infty} \Gamma(h) e^{-i\lambda_k h}$$

Thus, the Whittle likelihood given  $N$  independent realizations is defined as:

$$\prod_{n=1}^N \prod_{k=0}^{T-1} \frac{1}{\pi^p |S_k|} e^{-d_{n,k}^* S_k^{-1} d_{n,k}} \quad (1)$$

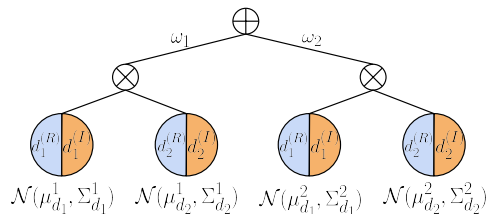
[4] Whittle, P. The analysis of multiple stationary time series. Journal of the Royal Statistical Society: Series B (Methodological), 1953.



CoSPN:

- using pairwise Gaussian leaf nodes to model complex random variables
- using an adapted non-parametric independence test for structure learning<sup>[5]</sup>

[5] Gens, R. and Domingos, P. Learning the Structure of Sum- Product Networks. ICML, 2013.



CoSPN:

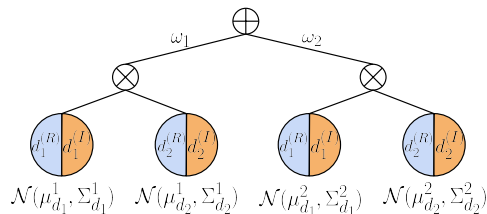
- using pairwise Gaussian leaf nodes to model complex random variables
- using an adapted non-parametric independence test for structure learning<sup>[5]</sup>

The Whittle Approximation applies for stationary time series. What about non-stationary time series?

We propose the following relaxations for general time series

- The mean of each frequency need not be 0
- The Fourier coefficients of different frequencies need not be independent

[5] Gens, R. and Domingos, P. Learning the Structure of Sum- Product Networks. ICML, 2013.



## CoSPN:

- using pairwise Gaussian leaf nodes to model complex random variables
- using an adapted non-parametric independence test for structure learning<sup>[5]</sup>

The Whittle Approximation applies for stationary time series. What about non-stationary time series?

We propose the following relaxations for general time series

- The mean of each frequency need not be 0
- The Fourier coefficients of different frequencies need not be independent

WSPN: CoSPN which jointly models the Fourier coefficients of time series

[5] Gens, R. and Domingos, P. Learning the Structure of Sum- Product Networks. ICML, 2013.



# Experiments

## Probabilistic Modeling of Time Series



		LearnSPN	WSPN-Pair	WSPN-2d	ResSPN	ResWSPN-Pair	ResWSPN-2d	MADE
<i>Sine</i>	train $\uparrow$	-0.47	2.65	<b>6.67</b>	<b>-60.62</b>	-148.48	-135.94	-105.91
	test $\uparrow$	-0.75	1.85	<b>5.75</b>	<b>-63.13</b>	-150.90	-138.86	-108.64
	ood $\downarrow$	$-\infty$	$-\infty$	$-\infty$	<b>-5880.85</b>	-4010.04	-4227.18	-11646865.93
<i>MNIST</i>	train $\uparrow$	256.11	272.84	<b>277.50</b>	249.47	<b>254.46</b>	<b>254.30</b>	336.03
	test $\uparrow$	254.99	270.40	<b>274.42</b>	245.67	251.74	<b>252.54</b>	327.22
	ood $\downarrow$	<b>125.19</b>	160.29	155.76	<b>204.93</b>	218.25	216.01	136.98
<i>Billiards</i>	train $\uparrow$	54.73	63.75	<b>65.01</b>	-367.83	-318.10	<b>-213.13</b>	-204.23
	test $\uparrow$	52.80	<b>54.14</b>	<b>54.12</b>	-377.38	-324.78	<b>-219.04</b>	-252.51
	ood $\downarrow$	-1984.38	-2348.57	<b>-2435.70</b>	-1003.49	-1052.21	<b>-2113.68</b>	-89521.82
<i>S&amp;P</i>	train $\uparrow$	-191.64	113.06	<b>174.45</b>	308.22	194.57	<b>1831.91</b>	359.52
<i>Stock</i>	train $\uparrow$	-615.76	328.90	<b>417.81</b>	257.03	496.07	<b>1172.85</b>	639.10

WSPNs **capture densities** over time series **better** than baselines!

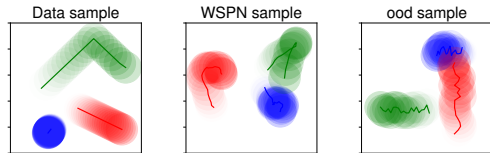


# Experiments

## Probabilistic Modeling of Time Series

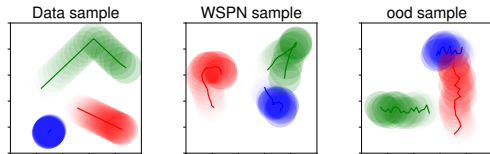


TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

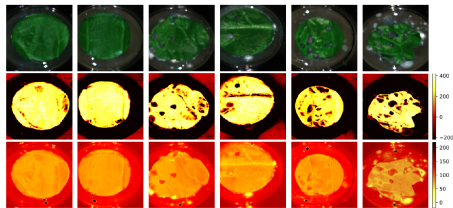


# Experiments

## Probabilistic Modeling of Time Series

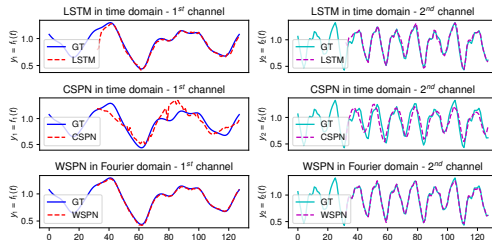
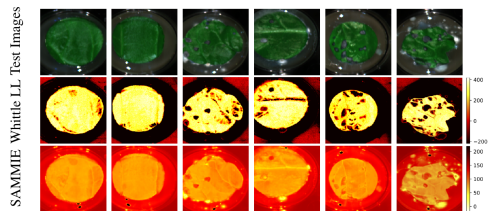
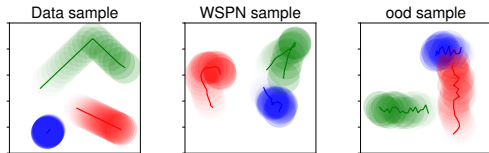


SAMMIE Whittle LL Test Images



# Experiments

## Probabilistic Modeling of Time Series



WSPNs are great for time series modeling and forecasting!

For directed acyclic graphs (DAGs):

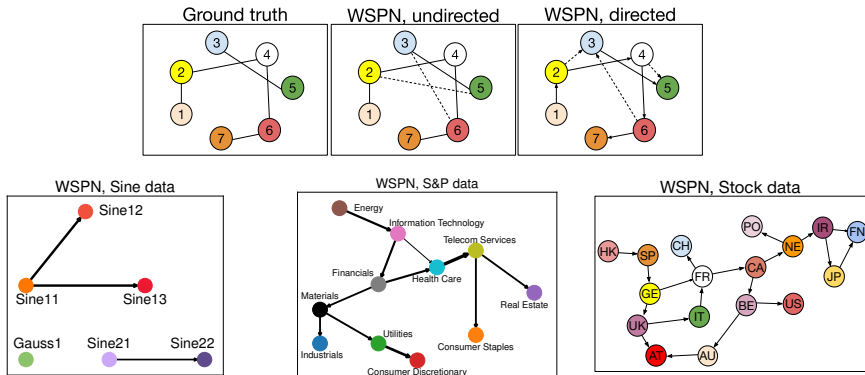
$$p(X_{1:N} | G, Co) \approx \prod_{n=1}^N \frac{\prod_{v_i \in V} p(d_n^{\{v_i \cup Pa_G(v_i)\}} | Co)}{\prod_{v_i \in V} p(d_n^{\{Pa_G(v_i)\}} | Co)}$$

For undirected graphs:

$$p(X_{1:N} | G, Co) \approx \prod_{n=1}^N \frac{\prod_{c_i \in C} p(d_n^{\{c_i\}} | Co)}{\prod_{s_i \in S} p(d_n^{\{s_i\}} | Co)}$$

# Experiments

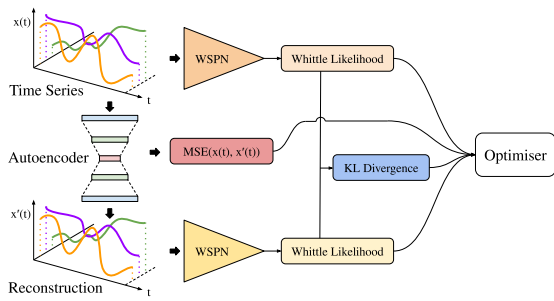
## Conditional Independence Structure



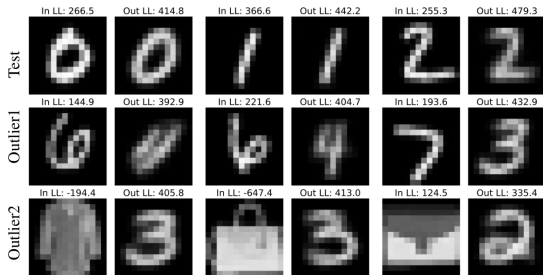
WSPNs can successfully extract the conditional independence of time series components!

# Whittle Networks

Providing Meaningful Probabilities to Deep Neural Networks

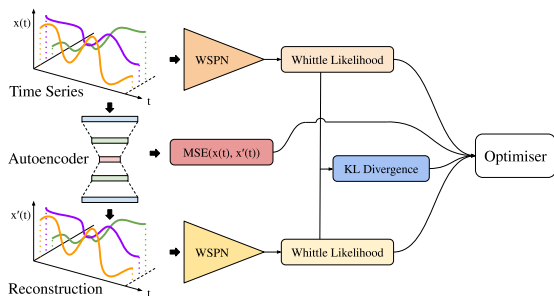


Whittle Autoencoder



# Whittle Networks

Providing Meaningful Probabilities to Deep Neural Networks



Whittle Autoencoder

	WSPN Input	WSPN Output
<b>train</b>	295.49	411.32
<b>test</b>	295.22	411.10
<b>outlier1</b>	<b>239.78</b>	401.54
<b>outlier2</b>	<b>48.84</b>	397.58

Whittle Networks can provide meaningful probabilities for deep neural networks!





We introduce:

- CoSPN – The first complex-valued SPN for modeling complex random variables
- WSPN – The first tractable probabilistic circuit for modeling the joint distribution of multivariate time series, exploiting the Whittle approximation
- Conditional independence structure can be extracted efficiently from WSPNs
- Whittle Networks – meaningful probabilities for deep neural networks



# Thanks

Zhongjie Yu

yu@cs.tu-darmstadt.de

Fabrizio Ventola

ventola@cs.tu-darmstadt.de

Kristian Kersting

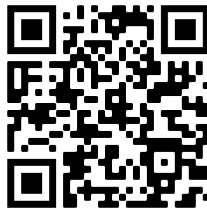
kersting@cs.tu-darmstadt.de

AIML lab



<https://www.aiml.informatik.tu-darmstadt.de>

Code



<https://github.com/ml-research/WhittleNetworks>