# **TACTiS-2: Better, Faster, Simpler Attentional Copulas for Multivariate Time Series**

Arjun Ashok<sup>123</sup> \*Étienne Marcotte<sup>1 </sup> \*Valentina Zantedeschi<sup>1 </sup> †Nicolas Chapados

<sup>1</sup>ServiceNow Research <sup>2</sup>Mila-Quebec AI Institute <sup>3</sup>Université de Montréal

<sup>∗</sup>Equal Contribution †Equal Contribution

# servicenow.



## **Summary**

**Problem:** Multivariate Probabilistic Time Series Prediction, i.e., estimating the joint distribution of high-dimensional multivariate time series

**General-Purpose Models:** We seek models that support

- $\blacktriangleright$  Arbitrarily complex data distributions
- $\blacktriangleright$  Heterogeneous/irregular sampling frequencies
- $\blacktriangleright$  Hundreds of time series with missing data
- $\blacktriangleright$  Deterministic covariates for conditioning (e.g., holidays)
- $\blacktriangleright$  Tasks: forecasting, interpolation, and hybrids

## TACTiS-2 obtains better performance with lesser compute in real-world forecasting tasks, compared to TACTiS electricity 0.0 0.3 0.6 0.9 1.2 1.5 1.8 fred-md kdd-cup  $FLOPs$   $(\times 10^{16})$ **TACTIS** IACTIS-2 (ours)

 $\Pi^{(o)}$ 

## **Contributions:**

- We show that Transformer Attentional Copulas for Time Series (TACTiS) [\[1\]](#page-0-0), while flexible, are highly inefficient
- We propose a simpler and faster approach to learning valid attentional copulas and prove its correctness
- We show that this results in significantly **better** training dynamics and empirical results on real-world datasets

## **Main Takeaway**



## **Problem setting**

- $\blacktriangleright$  **Multivariate time series:** a collection of univariate time series  $\mathbf{X}\stackrel{\text{\tiny def}}{=} \{\mathbf{X}_1,\ldots,\mathbf{X}_n\}$ , where each  $\mathbf{X}_i\stackrel{\text{\tiny def}}{=} [X_{i1},\ldots,X_{i,\ell_i}]$  is a random vector representing  $\ell_i$  observations of some real-valued process in time
- $\blacktriangleright$  Additional data: for any realization  $\mathbf{x}_i \stackrel{\text{\tiny def}}{=} [x_{i1}, \ldots, x_{i,\ell_i}]$  of  $\mathbf{X}_i$ , each  $x_{ij}$  is paired with:
- a timestamp,  $t_{ij} \in \mathbb{R}$  marking its measurement time
- a vector of <u>covariates</u>  $\mathbf{c}_{ij} \in \mathbb{R}^p$  that represents arbitrary additional information available
- **Examing Tasks:** defined with the help of a mask  $m_{ij}$  ∈ {0, 1}, which determines if any  $X_{ij}$  should be considered as observed  $(m_{ij} = 1)$  or to be inferred  $(m_{ij} = 0)$
- **Coal:** estimating the joint distribution of missing values ( $m_{ij} = 0$ ), given the observed ones ( $m_{ij} = 1$ ), covariates, and timestamps:

 $\boldsymbol{\phi} = \{\phi_1, \ldots, \phi_d; \phi_c\}$  where  $\{\phi_i\}_{i=1}^d$ estimated by minimizing negative log-likelihood:

- $\pi_{c, 2} (u_{\pi_2} \mid u_{\pi_1}) \times \cdots \times c_{\phi_n}$ *π*  $u_{\pi_d}^{\pi} (u_{\pi_d} \mid u_{\pi_1}, \ldots, u_{\pi_{d-1}})$ ) (3)
	-

**Two-Stage Nonparametric Copulas (ours):** Learn marginal parameters (eq. (6)), then learn copula parameters (eq. (5)): arg min *φc* − E **x**∼**X**  $\log c_{\phi_c} \left( F_{\phi} \right)$  $s.t.$   $(\phi_1^{\star})$  $\phi_d^{\star}, \ldots, \phi_d^{\star}$ 

 $\bullet$  Proposition 2: Solving Problem [\(5\)](#page-0-3) yields a solution to Problem [\(2\)](#page-0-2) where  $c_{\phi_c}$  is a valid copula. Proof builds on Sklar's theorem [\[3\]](#page-0-1). ◆ Advantages: The model needs to fit just 1 permutation  $\rightarrow$  simpler objective with faster convergence to better solutions.



## **What is a copula?**

Informally: a mathematical construct that expresses the coupling (dependency structure) of multiple random variables, irrespective of their marginal distributions (individual structure) According to Sklar's theorem [\[3\]](#page-0-1), the **joint CDF** of any random vector  $[X_1, \ldots, X_d]$  can be expressed as combination of:

- The **marginal** CDF of each random variable  $F_i(x_i) \stackrel{\text{def}}{=} P(X_i \leq x_i)$ ,
- The copula: a distribution on the unit cube with CDF  $C$  :  $[0, 1]^d$  →  $[0, 1]$  and  $U_{[0, 1]}$  marginals

$$
P(X_1 \leq x_1, \ldots, X_d \leq x_d) = C\big(F_1(x_1), \ldots, F_d(x_d)\big)
$$

## **Improved Learning of Non-Parametric Copulas**

**Copula-Based Density Estimators [\[1\]](#page-0-0):** Joint density decompose

$$
g_{\boldsymbol{\phi}}\Big(x_1,\ldots,x_d\Big)\,\stackrel{\scriptscriptstyle\rm def}{=} c_{\phi_{\mathsf{c}}}\Big(F_{\phi_1}\!\!\left(x_1\right)\!,\ldots,F_{\phi_d}
$$



 $^d_{i=1}$  are parameters of the marginal distributions, and  $\phi_c$  are the parameters of the copula density  $c_{\phi_c}$ 

(1)

<span id="page-0-2"></span>
$$
\underset{\phi}{\arg\min} \ -\underset{\mathbf{x} \sim \mathbf{X}}{\mathbb{E}} \ ]
$$

 $\bullet$  Proposition 1: Problem [\(2\)](#page-0-2) has infinitely many invalid solutions wherein  $c_{\phi_c}$  is not the density function of a valid copula. The true marginals and copula can be entangled  $\rightarrow$  Non-trivial to learn valid non-parametric copula-based density estimators.

$$
\log g_{\phi}(x_1,\ldots,x_d) \tag{2}
$$

**Permutation-based Nonparametric Copulas (TACTiS):** Nonparametric copulas learned using a permutation-based objective. Considers an autoregressive factorization of  $c_{\phi_c}$  according to an arbitrary permutation of the variables  $\pi=[\pi_1,\ldots,\pi_d]$ :  $c_{\boldsymbol{\phi}_c}$ <sup> $\pi$ </sup>  $\frac{d}{dx}(u_1, \ldots, u_d) \stackrel{\scriptscriptstyle{\mathsf{def}}}{=} c_{\phi}$ *π*  $\pi_{c,1}^{\pi}(u_{\pi_1}) \times c_{\phi}$ *π*

where  $u_{\pi_k}=F_{\phi_{\pi_k}}\!(x_{\pi_k})$ . Optimizes a permutation-based objective over  $\Pi$ , where  $\Pi$  is the set of all  $d!$  permutations: arg min  $\phi_1$ *,...,* $\phi_d$ *,* $\phi_c^{\pi}$ − E **x**∼**X** E *π*∼Π  $\log c_{\phi_c^{\pi}}$  $\sqrt{ }$  $F_{\phi_1}\!(x_1),\ldots,F_{\phi_d}$ 

X Limitations: The model needs the capacity to fit all  $d!$  permutations  $\rightarrow$  results in slow convergence and sub-optimal solutions.

$$
\ldots, F_{\phi_d}(x_d) \big) \times f_{\phi_1}(x_1) \times \cdots \times f_{\phi_d}(x_d) \tag{4}
$$

<span id="page-0-4"></span><span id="page-0-3"></span>
$$
\phi_1^{\star}(x_1), \dots, F_{\phi_d^{\star}}(x_d) \tag{5}
$$

$$
\begin{aligned} \n\check{d}) &\in \operatorname*{arg\,min}_{\phi_1, \dots, \phi_d} -\mathop{\mathbb{E}}_{\mathbf{x} \sim \mathbf{X}} \log \prod_{i=1} f_{\phi_i}(x_i) \n\end{aligned} \tag{6}
$$



## **Putting Theory into Practice**



Architecture of TACTiS-2: The two encoders serve to parametrize the decoder in the two stages of the training curriculum (bottom right). Phase 1 solves Problem [\(6\)](#page-0-4), while Phase 2 solves Problem [\(5\)](#page-0-3).

<sup>12</sup> †Alexandre Drouin<sup>12</sup>

## **Results**

Mean CRPS-Sum values for the forecasting experiments (± standard errors). Lower is better. Best results in bold.



Mean NLL values for forecasting experiments and training FLOP counts (± standard errors). Lower is better. Best results in bold.







## **Model flexibility**





Example forecasts of TACTiS-2 on irregular and unevenly sampled data

### **References**

<span id="page-0-0"></span>[1] Alexandre Drouin, Étienne Marcotte, and Nicolas Chapados. TACTiS: Transformer-attentional copulas for time series. In *ICML*, 2022.

[2] David Salinas, Michael Bohlke-Schneider, Laurent Callot, Roberto Medico, and Jan Gasthaus. High-dimensional multivariate forecasting with low-rank Gaussian copula processes.

- 
- *NeurIPS*, 2019.
- 







<span id="page-0-1"></span>[3] Abe Sklar. Fonctions de répartition à *n* dimensions et leurs marges. *Publ. de l'Institut Statistique de l'Université de Paris*, 8:229–231, 1959.