

# A Fully Probabilistic Tensor Network for Regularized Volterra System Identification

Afra Kilic\* Kim Batselier\*

\* Delft University of Technology, Delft, Netherlands  
(e-mail: {h.a.kilic, k.batselier}@tudelft.nl).

**Abstract:** Modeling nonlinear systems with Volterra series is challenging since the number of kernel coefficients grows exponentially with the model order. This work introduces Bayesian Tensor Network Volterra kernel machines (BTN-V), extending the Bayesian Tensor Network (BTN) framework to Volterra system identification. BTN-V represents Volterra kernels via canonical polyadic decomposition, reducing model complexity from  $\mathcal{O}(I^D)$  to  $\mathcal{O}(DIR)$ . By treating all tensor components and hyperparameters as random variables, BTN-V provides predictive uncertainty estimation at no extra computational cost. Sparsity-inducing hierarchical priors enable automatic rank determination and learning of fading-memory behavior directly from data, improving interpretability and avoiding overfitting. Empirical results demonstrate competitive accuracy, enhanced uncertainty quantification, and reduced computational cost.

**Keywords:** nonlinear system identification, Volterra series, tensor network kernel machines, variational inference, Bayesian methods

## 1. INTRODUCTION

Modeling nonlinear systems remains a key challenge in system identification. The Volterra series, a nonlinear extension of the finite impulse response model, is commonly used when the exact system behavior is unknown. However, its application is typically limited to weakly nonlinear systems, where linear dynamics dominate. This limitation arises not from the Volterra framework itself but from the exponential growth of kernel coefficients with increasing model order.

Several methods exist to reduce the complexity of Volterra kernels. One approach incorporates prior knowledge of the kernel structure within a Bayesian framework (Chen et al., 2011), encoding the expected smooth decay of Volterra kernels via Bayesian priors (Pillonetto and De Nicolao, 2010). This was extended to parametric Volterra models with decaying covariance matrices (Birpoutsoukis et al., 2018), though limited to third-order kernels, while higher-order non-parametric models using similar structures remain restricted to small datasets (Libera et al., 2021).

An alternative approach imposes sparsity by assuming most Volterra coefficients are negligible, implemented through sparse Bayesian learning (Miao et al., 2019). However, this strong sparsity assumption applies only to the specific system considered in that study. A further alternative is kernel compression using low-rank tensor structures, which represent all Volterra coefficients through a small set of parameters from which the full set can be reconstructed. Favier et al. (2012) apply both low-rank Tucker and canonical polyadic decompositions. Batselier et al. (2017) has proposed a multiple-input-multiple-output (MIMO) Volterra tensor network (TN) that uses tensor-train decomposition to represent all Volterra ker-

nels at once. By exploiting multilinearity and applying alternating least squares (ALS), the MIMO Volterra TN avoids the exponential growth of coefficients and enables efficient identification of high-order (up to 10th-order) Volterra systems within seconds on standard hardware.

Although the MIMO Volterra TN offers low computational complexity, it does not account for parameter uncertainty and therefore cannot provide uncertainty quantification for predictions. It also requires the tensor rank to be specified a priori and tends to overfit when higher ranks are used. Rank selection is often performed via trial and error, which is both costly and imprecise, while more systematic approaches, such as maximum likelihood estimation, may still lead to overfitting, and cross-validation becomes computationally expensive when multiple hyperparameters are involved. The first probabilistic treatment of the Volterra TN with Tikhonov regularization has been proposed to avoid overfitting and to quantify uncertainty in the predictions (Memmel et al., 2023); however, it models only a single TN core as a random variable, treating the remaining cores as deterministic, and is thus not fully probabilistic. Moreover, despite the regularization, this approach still requires manual tuning of the tensor rank. Recently, Kilic and Batselier (2025) presented the Bayesian Tensor Network (BTN) kernel machines framework, a fully probabilistic approach that employs sparsity-inducing hierarchical priors on the TN components to automatically infer model complexity. Specifically, these sparsity-inducing priors allow the model, during identification, to automatically determine the effective tensor rank and identify the most relevant features for prediction, thereby enhancing interpretability.

**In this paper,** we propose the Bayesian Tensor Network Volterra kernel machines (BTN-V) by extending the BTN kernel machines framework (Kilic and Batselier, 2025) to Volterra system identification. BTN-V treats all tensor components and model hyperparameters as random

\* This publication is part of the project *Sustainable Learning for AI from Noisy Large-Scale Data* (project number VI.Vidi.213.017), which is financed by the Dutch Research Council (NWO).

variables, allowing prediction uncertainty to be estimated without additional computational cost. Through the use of sparsity-inducing priors, BTN-V offers two main benefits: it automatically determines the tensor rank during the identification process and learns the fading-memory behavior of Volterra systems directly from the data, helping to prevent overfitting by reducing the influence of distant past inputs.

## 2. TENSOR BASICS AND NOTATION

The order of a tensor is the number of its dimensions, also referred to as ways or modes. Scalars are denoted by lowercase letters, e.g.,  $a$ ; vectors (first-order tensors) by bold lowercase letters, e.g.,  $\mathbf{a}$ ; matrices (second-order tensors) by bold uppercase letters, e.g.,  $\mathbf{A}$ ; and higher-order tensors (order  $\geq 3$ ) by bold calligraphic letters, e.g.,  $\mathcal{A}$ .

A  $D$ th-order tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_D}$  has entries  $a_{i_1 i_2 \dots i_D}$ , with indices  $i_d = 1, \dots, I_d$  for  $d \in [1, D]$ . Often, it is computationally convenient to vectorize tensors. The vectorization  $\text{vec}(\mathcal{A}) \in \mathbb{R}^{I_1 I_2 \dots I_D}$  maps each entry as

$$\text{vec}(\mathcal{A})_i = a_{i_1 i_2 \dots i_D}, \quad i = i_1 + \sum_{d=2}^D (i_d - 1) \prod_{k=1}^{d-1} I_k. \quad (1)$$

When applied to a matrix, the operator  $\text{vec}(\cdot)$  performs column-wise vectorization. The operator  $\text{diag}(\cdot)$  returns a diagonal matrix from a vector. The Kronecker product of  $\mathbf{A} \in \mathbb{R}^{I \times J}$  and  $\mathbf{B} \in \mathbb{R}^{K \times L}$  is  $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{KI \times LJ}$ . The Khatri-Rao product  $\mathbf{A} \odot \mathbf{B}$  of  $\mathbf{A} \in \mathbb{R}^{I \times J}$  and  $\mathbf{B} \in \mathbb{R}^{K \times J}$  is an  $IK \times J$  matrix formed by column-wise Kronecker products. The Hadamard (element-wise) product of  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{I \times J}$  is  $\mathbf{A} \oslash \mathbf{B} \in \mathbb{R}^{I \times J}$ . The identity matrix is denoted by  $\mathbf{I}$ , with size either inferred or specified.

**Definition 1.** A rank- $R$  Canonical Polyadic Decomposition (CPD) of  $\mathbf{w} = \text{vec}(\mathbf{W}) \in \mathbb{R}^{I^D}$  consists of  $D$  factor matrices  $\mathbf{W}^{(d)} \in \mathbb{R}^{I \times R}$ , such that

$$\mathbf{w} = (\mathbf{W}^{(1)} \odot \mathbf{W}^{(2)} \odot \dots \odot \mathbf{W}^{(D)}) \mathbf{1}_R = \sum_{r=1}^R \mathbf{w}_r^{(1)} \otimes \mathbf{w}_r^{(2)} \otimes \dots \otimes \mathbf{w}_r^{(D)}, \quad (2)$$

where  $\mathbf{w}_r^{(d)} \in \mathbb{R}^I$  denotes the  $r$ th column of the matrix  $\mathbf{W}^{(d)}$ . The CPD is unique under mild conditions (Kruskal, 1977), and its storage complexity scales as  $\mathcal{O}(\text{DIR})$ .

## 3. TENSOR NETWORK SISO VOLTERRA SYSTEM

In this section, we briefly introduce the TN formulation of the single-input single-output (SISO) Volterra system (Batselier et al., 2017). We focus on the SISO case for simplicity; the extension to MIMO is trivial but not discussed for brevity.

Consider a discrete time truncated  $D$ th order Volterra system with memory length  $M$ :

$$y(n) = \sum_{d=0}^D \sum_{m_1, \dots, m_d=0}^{M-1} \mathbf{w}_d(m_1, \dots, m_d) \prod_{j=1}^d u(n - m_j) + e(n), \quad (3)$$

where  $u(n)$  is the input,  $y(n)$  is the output, and  $e(n)$  represents additive noise, for  $n = 1, \dots, N$ . The

set  $\{\mathbf{W}_d\}_{d=0}^D$  represents the Volterra kernels of different orders. These kernels grow exponentially with the system order  $D$ . This curse of dimensionality can be mitigated by representing all Volterra kernels simultaneously with a low-rank TN (Batselier et al., 2017), which transforms the problem into a set of smaller linear systems, each solved iteratively during the identification process. To obtain a low-rank TN representation of the Volterra kernels, the first step is to rewrite the inputs  $u(n)$  following Batselier et al. (2017) as

$$\mathbf{u}_n = \begin{pmatrix} 1 & u(n) & \dots & u(n - M + 1) \end{pmatrix}^\top \in \mathbb{R}^I, \quad (4)$$

where  $I := M + 1$ , with the corresponding system output is  $y(n) \in \mathbb{R}$ . Including the constant term as the first term in  $\mathbf{u}_n$  allows us to define a feature map  $\mathbf{u}_n^D$  as the  $D$ -times Kronecker product of  $\mathbf{u}_n$ , such that

$$\mathbf{u}_n^D := \underbrace{\mathbf{u}_n \otimes \mathbf{u}_n \otimes \dots \otimes \mathbf{u}_n}_{D \text{ times}} \in \mathbb{R}^{I^D}, \quad (5)$$

which contains all monomials of the input from degree 0 up to  $D$ . With this definition, the output at time  $n$  can be expressed as

$$y(n) = (\mathbf{u}_n^D)^\top \mathbf{w} + e_n, \quad (6)$$

where  $\mathbf{w} \in \mathbb{R}^{I^D}$  vector containing all coefficients from all Volterra kernels. Since the Volterra feature map  $\mathbf{u}_n^D$  in (5) has a Kronecker structure, this property can be exploited to express the Volterra coefficients  $\mathbf{w}$  using the CPD model in Definition 1, which reduces the number of learnable parameters from  $\mathcal{O}(I^D)$  to  $\mathcal{O}(\text{DIR})$  (Harshman, 1970). In this model,  $\mathbf{w}$  is represented by  $D$  factor matrices of size  $\mathbb{R}^{I \times R}$ , which are estimated iteratively using the ALS algorithm without explicitly constructing  $\mathbf{w}$ . When the rank  $R$  matches the true CP rank, the exact solution of (6) is recovered. However, ALS assumes deterministic factor matrices and yields only point estimates. Following Kilic and Batselier (2025), we instead treat the factor matrices as random variables with sparsity-inducing priors. This probabilistic formulation allows uncertainty quantification in the predictions and the model to automatically infer both the effective tensor rank  $R$  and the fading-memory behavior of the Volterra series during identification.

## 4. PROBABILISTIC TENSOR NETWORK SISO VOLTERRA SYSTEMS

In this section, we present the main contribution of this paper: BTN-V, the BTN kernel machines framework extended to SISO Volterra systems. For a detailed discussion of the priors and the mean-field variational inference-based probabilistic identification procedure, the reader is referred to Kilic and Batselier (2025).

### 4.1 Probabilistic Model and Priors

Expressing equation (6) for  $n = 0, 1, \dots, N$  leads to the following matrix equation:

$$\mathbf{y} = \mathbf{U}^{D^T} \mathbf{w} + \mathbf{e}, \quad (7)$$

where  $\mathbf{y} \in \mathbb{R}^N$  is the vector of measured outputs and  $\mathbf{e} \in \mathbb{R}^N$  denotes Gaussian white noise with  $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \tau^{-1} \mathbf{I}_N)$ . Let  $\mathbf{U} \in \mathbb{R}^{I \times N}$ , where the  $n$ th column contains  $\mathbf{u}_n$ . The matrix  $\mathbf{U}^D = \mathbf{U} \odot \mathbf{U} \odot \dots \odot \mathbf{U} \in \mathbb{R}^{I^D \times N}$  represents the  $D$ -fold row-wise Khatri-Rao product of  $\mathbf{U}$  with itself,

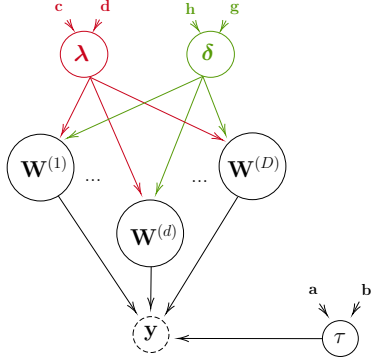


Fig. 1. Probabilistic graphical model of the Volterra TN, where the CPD-decomposed coefficients  $\mathbf{w}$  are represented by factor matrices  $\{\mathbf{W}^{(d)}\}_{d=1}^D$ . Dashed, solid, and unbounded nodes denote observed data  $\mathbf{y}$ , random variables, and Gamma hyperparameters, respectively.

thus  $\mathbf{U}^D$  the input matrix whose  $n$ th column corresponds to the Kronecker vector  $\mathbf{u}_n^D$ . Assuming a CP-decomposed-Volterra-coefficients vector  $\mathbf{w}$  together with the Gaussian noise model, the likelihood of the observed outputs is given by

$$p(\mathbf{y} | \{\mathbf{W}^{(d)}\}_{d=1}^D, \tau) = \prod_{n=1}^N \mathcal{N}(y(n) | (\mathbf{u}_n^D)^\top \mathbf{w}, \tau^{-1}), \quad (8)$$

where  $\tau$  denotes the noise precision and  $\mathbf{W}^{(d)} \in \mathbb{R}^{I \times R}$  are the factor matrices of the CP-decomposed model weights (Definition 1), whose  $r$ th column is  $\mathbf{w}_r^{(d)}$  and  $i$ th row is  $\mathbf{w}_i^{(d)}$ , with  $r = 1, \dots, R$  and  $i = 1, \dots, I$ . Determining an appropriate tensor rank  $R$  and feature dimension  $I$  is generally a nontrivial and computationally intensive task.

To address this, Kilic and Batselier (2025) use a hierarchical sparsity-inducing prior to automatically infer  $R$  and  $I$ , avoiding overfitting. Building on this methodology, we define two sets of sparsity parameters:  $\boldsymbol{\lambda} := [\lambda_1, \dots, \lambda_R]$ , with each  $\lambda_r$  controlling the regularization of the  $r$ th column  $\mathbf{w}_r^{(d)}$ , and  $\boldsymbol{\delta} := [\delta_1, \dots, \delta_I]$ , with each  $\delta_i$  regulating the  $i$ th row  $\mathbf{w}_i^{(d)}$  of  $\mathbf{W}^{(d)}$ , for all  $d \in [1, D]$ . Together,  $\boldsymbol{\lambda}$  and  $\boldsymbol{\delta}$  govern sparsity across the columns and rows of the factor matrices, respectively. The prior distribution for the vectorized factor matrices is a zero mean Gaussian prior

$$p(\text{vec}(\mathbf{W}^{(d)}) | \boldsymbol{\lambda}, \boldsymbol{\delta}) = \mathcal{N}(\mathbf{0}, \boldsymbol{\Lambda}^{-1} \otimes \boldsymbol{\Delta}^{-1}) \quad (9)$$

where  $\boldsymbol{\Lambda} \otimes \boldsymbol{\Delta} = \text{diag}(\boldsymbol{\lambda}) \otimes \text{diag}(\boldsymbol{\delta})$  represents the inverse covariance matrix (precision matrix). Under the zero-mean prior assumption, higher values in  $\boldsymbol{\lambda}$  or  $\boldsymbol{\delta}$  force the corresponding columns or rows of  $\mathbf{W}^{(d)}$  toward zero, thereby regularizing the associated components or input dimensions. This follows the automatic relevance determination (ARD) principle (Neal, 1996) and its tensor-based extensions (Zhao et al., 2015).

For a fully probabilistic formulation, we further specify Gamma hyperpriors over the sparsity parameters as  $p(\boldsymbol{\lambda}) = \prod_{r=1}^R \text{Ga}(\lambda_r | c_0, d_0)$  and  $p(\boldsymbol{\delta}) = \prod_{i=1}^I \text{Ga}(\delta_i | g_0, h_0)$ , where  $\text{Ga}(x | a, b)$  denotes the Gamma distribution. Similarly, the noise precision  $\tau$  is assigned a Gamma prior  $p(\tau) = \text{Ga}(\tau | a_0, b_0)$ .

**Remark.** In the original BTN-Kernel machines formulation, a separate  $\delta_d$  was defined for each factor matrix

$\mathbf{W}^{(d)}$ , since the feature map was formed by the Kronecker product of  $D$  different feature vectors  $\mathbf{u}_n^{(d)}$ . In contrast, the Volterra feature map in (5) uses the same input vector  $\mathbf{u}_n$  repeated  $D$  times in the Kronecker product. Thus, we use a single vector  $\boldsymbol{\delta}$  that applies row-wise regularization across all factor matrices. Since each row corresponds to a specific memory lag,  $\boldsymbol{\delta}$  is directly associated with the temporal structure of the model (i.e., the memory lags). By modeling  $\boldsymbol{\delta}$  as a random variable, the model learns which lags are more important during identification, thereby regularizing the temporal structure by reducing the influence of less informative lags and highlighting those that are more relevant.

Denoting all latent variables and hyperparameters by  $\Theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(D)}, \boldsymbol{\delta}, \boldsymbol{\lambda}, \tau\}$ , the joint distribution  $p(\mathbf{y}, \Theta)$  can be expressed as

$$p(\mathbf{y} | \{\mathbf{W}^{(d)}\}_{d=1}^D, \tau) \prod_{d=1}^D p(\mathbf{W}^{(d)} | \boldsymbol{\lambda}, \boldsymbol{\delta}) p(\boldsymbol{\delta}) p(\boldsymbol{\lambda}) p(\tau). \quad (10)$$

Figure 1 illustrates the probabilistic graphical model corresponding to the joint distribution in (10). The input matrices  $\mathbf{U}$  are omitted from the figure since they are not random variables. The factor matrices  $\{\mathbf{W}^{(d)}\}_{d=1}^D$ , and the shape and scale hyperparameters of the Gamma distributions, represented by unbounded nodes in the figure, must be initialized before identification. Then the full posterior distribution of all variables in  $\Theta$  given the observed data is

$$p(\Theta | \mathbf{y}) = \frac{p(\mathbf{y}, \Theta)}{\int p(\mathbf{y}, \Theta) d\Theta}. \quad (11)$$

Based on the posterior distribution of  $\Theta$ , the predictive distribution over unseen data points, denoted  $\tilde{y}_i$  can be inferred by

$$p(\tilde{y}_i | \mathbf{y}) = \int p(\tilde{y}_i | \Theta) p(\Theta | \mathbf{y}) d\Theta. \quad (12)$$

#### 4.2 Identification Process

An exact Bayesian inference of (11) and (12) requires integrating over all latent variables and hyperparameters, making it analytically intractable. Therefore, as presented in (Kilic and Batselier, 2025) we perform a Bayesian identification using mean-field variational inference (Winn and Bishop, 2005). The key idea is to approximate the true posterior  $p(\Theta | \mathbf{y})$  with a variational distribution  $q(\Theta)$  by minimizing the Kullback-Leibler (KL) divergence. This is equivalent to maximizing a lower bound  $\mathcal{L}(q)$  on the model evidence  $\ln p(\mathbf{y})$ , defined as  $\mathcal{L}(q) = \mathbb{E}_{q(\Theta)}[\ln p(\mathbf{y}, \Theta)]$ . The maximum of the lower bound occurs when the KL divergence vanishes, implying  $q(\Theta) = p(\Theta | \mathbf{y})$ . The mean-field assumption factorizes the variational distribution over each variable  $\theta_j \in \Theta$  as

$$q(\Theta) = \prod_{d=1}^D q_{\mathbf{W}^{(d)}}(\mathbf{W}^{(d)}) q_{\boldsymbol{\delta}}(\boldsymbol{\delta}) q_{\boldsymbol{\lambda}}(\boldsymbol{\lambda}) q_{\tau}(\tau). \quad (13)$$

In other words, this factorization assumes that all variables  $\theta_j \in \Theta$  are independent. The functional form of each factor can then be derived analytically in turn. Specifically, the optimal  $\theta_j$  is obtained by maximizing  $\mathcal{L}(q)$ , with the maximum occurring when

$$\ln q_j(\theta_j) = \mathbb{E}_{q(\Theta \setminus \theta_j)}[\ln p(\mathbf{y}, \Theta)] + \text{const}, \quad (14)$$

---

**Algorithm 1**


---

**Require:** Input data  $\mathbf{x} = \{x_n\}_{n=1}^N$ , output data  $\mathbf{y} = \{y_n\}_{n=1}^N$

- 1: **Initialize:**  $R, I, \mathbf{W}^{(d)}, \Sigma^{(d)}, \forall d \in [1, D], a_0, b_0, \mathbf{c}_0, \mathbf{d}_0, \mathbf{g}_0, \mathbf{h}_0$  and set  $\tau = a_0/b_0, \lambda_r = c_0^r/d_0^r, \forall r \in [1, R], \lambda_I = g_0^I/h_0^I, \forall i \in [1, I]$ .
- 2: **repeat**
- 3:   **for**  $d = 1$  to  $D$  **do**
- 4:     Update posterior  $q_d(\text{vec}(\mathbf{W}^{(d)}))$  using Eq. (16)
- 5:   **end for**
- 6:   Update posterior  $q(\delta)$  using Eq. (17)
- 7:   Update posterior  $q(\lambda)$  using Eq. (18)
- 8:   Evaluate variational lower bound
- 9:   **if** truncation criterion satisfied **then**
- 10:     Reduce rank  $R$  by removing zero-columns in  $\mathbf{W}^{(d)}, \forall d \in [1, D]$
- 11:   **end if**
- 12: **until** convergence
- 13: Update posterior  $q(\tau)$  using Eq. (19)
- 14: Compute predictive distribution using Eq. (20)

---

where  $\mathbb{E}_{q(\Theta \setminus \theta_j)}[\cdot]$  denotes the expectation taken with respect to the variational distributions of all variables except  $\theta_j$ . Since all parameter distributions are in the exponential family and conjugate to their priors, this yields closed-form posterior updates for each factor. Identification proceeds by initializing each  $q_j(\theta_j)$  and iteratively updating them using the update rules until convergence. In the following, we present the posterior update rules for each  $q_j(\theta_j)$ , for all  $\theta_j \in \Theta$ . For full derivations, proofs, and explicit formulas, see Kilic and Batselier (2025).

*Posterior Distribution of Factor Matrices* Because of the multilinear nature of the CPD, a tensor represented in CPD form can be expressed as a function that is linear with respect to any one of its factor matrices. Thus, data-fitting term  $\mathbf{U}^{D^T} \mathbf{w}$  in (7) can be rewritten linearly in terms of the unknown  $d$ th factor matrix  $\mathbf{W}^{(d)}$

$$\begin{aligned} \mathbf{U}^{D^T} \mathbf{w} &= \text{vec}(\mathbf{W}^{(d)})^T \mathbf{G}^{(d)}, \\ \mathbf{G}^{(d)} &= \mathbf{U} \odot \left( \bigotimes_{k \neq d} \mathbf{W}^{(k)^T} \mathbf{U} \right). \end{aligned} \quad (15)$$

where  $\mathbf{G}^{(d)} \in \mathbb{R}^{IR \times N}$  is the design matrix. By applying (14), the posterior mean  $\text{vec}(\tilde{\mathbf{W}}^{(d)})$  and covariance  $\Sigma^{(d)}$  of

$$q_{\mathbf{W}^{(d)}}(\text{vec}(\mathbf{W}^{(d)})) = \mathcal{N}\left(\text{vec}(\mathbf{W}^{(d)}) \mid \text{vec}(\tilde{\mathbf{W}}^{(d)}), \Sigma^{(d)}\right),$$

are updated by

$$\begin{aligned} \text{vec}(\tilde{\mathbf{W}}^{(d)}) &= \mathbb{E}_q[\tau] \Sigma^{(d)} \mathbb{E}_q[\mathbf{G}^{(d)}] \mathbf{y}, \\ \Sigma^{(d)} &= \left[ \mathbb{E}_q[\tau] \mathbb{E}_q[\mathbf{G}^{(d)} \mathbf{G}^{(d)^T}] + \mathbb{E}_q[\Lambda] \otimes \mathbb{E}_q[\Delta] \right]^{-1}. \end{aligned} \quad (16)$$

*Posterior distributions of  $\lambda$  and  $\delta$*  The posterior of the row precision  $\delta$  is an independent Gamma distribution over rows  $q_\delta(\delta) = \prod_{i=1}^I \text{Ga}(\delta_i \mid g_N^i, h_N^i)$  where the posterior parameters are updated by

$$g_N^i = g_0^i + \frac{DR}{2}, \quad h_N^i = h_0^i + \sum_{d=1}^D \mathbb{E}_q \left[ \mathbf{w}_i^{(d)^T} \Lambda \mathbf{w}_i^{(d)} \right]. \quad (17)$$

The posterior expectation of  $\delta$  can be obtained by  $\mathbb{E}_q[\delta] = [g_N^1/h_N^1, \dots, g_N^I/h_N^I]^T$ , and thus  $\mathbb{E}_q[\Delta] = \text{diag}(\mathbb{E}_q[\delta])$ .

The posterior of  $\lambda$  is an independent Gamma distribution over columns  $q_\lambda(\lambda) = \prod_{r=1}^R \text{Ga}(\lambda_r \mid c_N^r, d_N^r)$  where the posterior parameters are updated by

$$c_N^r = c_0^r + \frac{DI}{2}, \quad d_N^r = d_0^r + \frac{1}{2} \sum_{d=1}^D \mathbb{E}_q \left[ \mathbf{w}_r^{(d)^T} \Delta \mathbf{w}_r^{(d)} \right]. \quad (18)$$

The posterior expectation of  $\lambda$  can be obtained by  $\mathbb{E}_q[\lambda] = [c_N^1/d_N^1, \dots, c_N^R/d_N^R]^T$ , and thus  $\mathbb{E}_q[\Lambda_R] = \text{diag}(\mathbb{E}_q[\lambda_R])$ .

*Posterior distribution of noise precision  $\tau$*  The noise precision  $\tau$  is inferred by combining information from the observed data, all factor matrices, and its hyperprior. Its variational posterior is a Gamma distribution  $q_\tau(\tau) = \text{Ga}(\tau \mid a_N, b_N)$ , where the posterior parameters are updated by

$$a_N = a_0 + \frac{N}{2}, \quad b_N = b_0 + \frac{1}{2} \mathbb{E}_q[\|\mathbf{y} - \mathbf{U}^{D^T} \mathbf{w}\|_F^2]. \quad (19)$$

The posterior expectation of  $\tau$  is then  $\mathbb{E}_q[\tau] = a_N/b_N$ .

*Lower Bound of Model Evidence* The variational inference framework maximizes the evidence lower bound (ELBO)  $\mathcal{L}(q)$ , which is guaranteed not to decrease across iterations and can therefore serve as a convergence criterion. It is defined as  $\mathcal{L}(q) = \mathbb{E}_{q(\Theta)}[\ln p(\mathbf{y}, \Theta)] + H(q(\Theta))$ , where the first term is the posterior expectation of the joint distribution and the second term is the entropy of the variational posterior distributions.

#### 4.3 Predictive Distribution

The predictive distribution for unseen data, given training data, is approximated using the variational posterior as

$$p(\tilde{y}_i \mid \mathbf{y}) \simeq \int \int p(\tilde{y}_i \mid \{\mathbf{W}^{(d)}\}, \tau^{-1}) q(\{\mathbf{W}^{(d)}\}) q(\tau) d\{\mathbf{W}^{(d)}\} d\tau. \quad (20)$$

This yields a Student's t-distribution  $\tilde{y}_i \mid \mathbf{y} \sim \mathcal{T}(\tilde{y}_i, \mathcal{S}_i, \nu_y)$ , with parameters

$$\begin{aligned} \tilde{y}_i &= \bigotimes_{d=1}^D \tilde{\mathbf{W}}^{(d)} \varphi_i^{(d)}, \quad \nu_y = 2a_N, \\ \mathcal{S}_i &= \left\{ \frac{b_N}{a_N} + \sum_d \mathbf{g}^{(d)}(x_n)^T \Sigma^{(d)} \mathbf{g}^{(d)}(x_n) \right\}^{-1}. \end{aligned} \quad (21)$$

The predictive variance can be obtained from the Student's t-distribution as  $\text{Var}(y_i) = \frac{\nu_y}{\nu_y - 2} \mathcal{S}_i^{-1}$ . **Algorithm 1** summarizes the BTN-V presented above.

## 5. NUMERICAL EXPERIMENTS

In this section, we evaluate the performance of BTN-V in terms of predictive accuracy, uncertainty quantification, and computational efficiency, measured by the Root Mean Squared Error (RMSE), Negative Log-Likelihood (NLL), and training runtime, respectively. The NLL is defined as  $\text{NLL} = -\frac{1}{N} \sum_{n=1}^N \log p(y_n \mid \theta)$ , and it quantifies how well the predicted probability distribution fits the observed data, penalizing both inaccurate and overconfident predictions.

In the Volterra model, increasing the memory length ( $M$ ) and polynomial degree ( $D$ ) increases the number of parameters to be learned. Especially when the dataset is short, there are not enough samples to estimate these parameters effectively, which can lead to poor predictive performance. Using the short Cascaded Tanks Benchmark

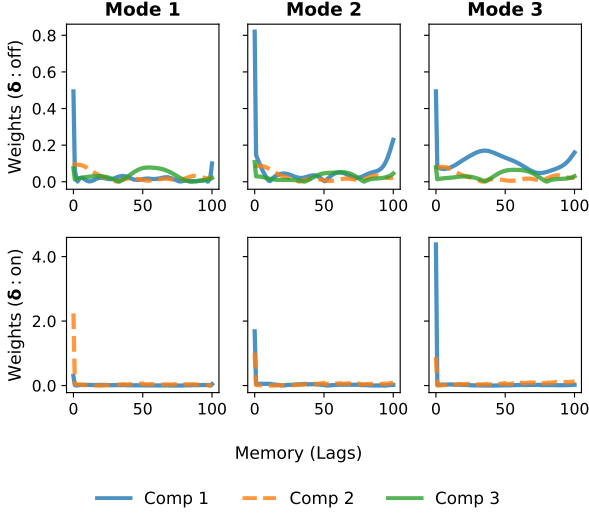


Fig. 2. Absolute values of the column weights of the factor matrices without (top) and with (bottom)  $\delta$  regularization. The regularization of  $\delta$  enforces decay across lags, consistent with the fading memory property of Volterra kernels.

dataset (Schoukens et al., 2016), we first demonstrate that regularization via  $\delta$  and  $\lambda$  yields a more parsimonious model with fewer effective parameters and improved predictive accuracy. We then compare the performance of BTN-V with three state-of-the-art methods for Volterra system identification: BMVALS (Mommel et al., 2023), RVS (Birpoutsoukis et al., 2018) and SED-MPK (Libera et al., 2021)

For all models, the Volterra order is set to  $D = 3$  and the memory length to  $M = 100$ . For BMVALS, which is a tensor-based Volterra kernel method, we use a rank of  $R = 48$ , as specified in the original paper. For the BTN-V models, the initial CPD rank is set to  $R = 20$ , and the hyperparameters are initialized following the procedure described in Kilic and Batselier (2025). To ensure numerical stability during identification, all input data in the BTN-V models are normalized to the range  $[0, 1]$ , and the output data are standardized to have zero mean and unit variance. The normalization on the targets is removed for prediction, and the validation performance of each method is reported. All computations are performed on an Apple MacBook Pro with an Apple M2 Pro chip and 16 GB of RAM running macOS 15.7.1. The Python code enabling the reproduction of all experiments in this section is available at [github.com/afraakilic/BTN\\_Volterra\\_Sys\\_ID](https://github.com/afraakilic/BTN_Volterra_Sys_ID).

### 5.1 Benchmark Description

The Cascaded Tanks Benchmark is a nonlinear system with a short dataset ( $N = 1024$ ). The setup includes two water tanks: water is pumped from a reservoir into the upper tank, flows into the lower tank, and then returns to the reservoir. The input signal controls the pump, while the output signal measures the water level in the lower tank. When too much water is pumped from the reservoir, the tanks may overflow. Thus, this dataset exhibits both soft and hard nonlinearities, where the soft nonlinearities arise from the smooth hydraulic flow between the tanks and the hard nonlinearities are caused by tank overflow.

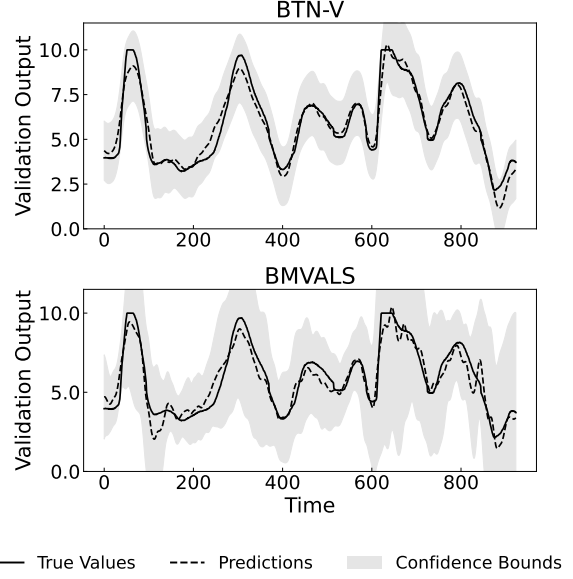


Fig. 3. Predicted and actual validation outputs for BTN-V and BMVALS. The solid line shows the true observations, the dashed line shows the predictive mean, and the shaded area represents  $\pm 3$  standard deviations from the mean.

In the following, we model the system’s input–output relationship using a Volterra series.

### 5.2 Regularization Through $\delta$ and $\lambda$

First to examine the BTN-V’s ability to perform automatic rank inference, we analyze whether column-wise regularization of the factor matrices via  $\lambda$  promotes convergence to low-rank solutions. Since the factor matrices are randomly initialized, we evaluated the model under 10 different random initializations to account for this randomness. With the initial rank fixed at  $R = 20$ , BTN-V yields an average final rank of  $2.5 \pm 0.5$ , demonstrating convergence to low-rank solutions.

Next, we present the effect of row-wise regularization through  $\delta$ , which is associated with memory (lags). In the Volterra framework, past observations are expected to have a diminishing influence on the current state. To examine the effect  $\delta$  on the factor matrices, we train the models both when the row-wise penalization term  $\delta$  is enabled and disabled. Specifically, setting  $\delta = \mathbf{I}$  and keeping it fixed throughout the identification process results in no penalization on the rows; we refer to this setting as  $\delta : \text{off}$ . Conversely, when  $\delta$  is updated during identification, we denote it as  $\delta : \text{on}$ . We present weights in the resulting factor matrices columnwise in Figure 2. In Figure 2, we plot the factor matrix weights columnwise. For  $D = 3$  modes, there are three factor matrices of size  $I \times R$ , where  $I = 101$  for a memory length of  $M = 100$ . Each column represents a CP component of the Volterra kernel coefficients, and each row corresponds to a specific memory element. The top row of Figure 2 shows the column weights without  $\delta$  regularization, where no clear decay with increasing lag is observed. This is most evident in the first component (blue line), which even increases at higher lags, while other components exhibit similarly irregular patterns. In contrast, with  $\delta$  regularization (bottom row),

Method	RMSE	NLL	Time (s)
BTN-V ( $\delta$ : on)	$0.51 \pm 0.02$	<b><math>0.77 \pm 0.05</math></b>	$13.68 \pm 1.73$
BTN-V ( $\delta$ : off)	$0.69 \pm 0.04$	$1.12 \pm 0.10$	<b><math>12.52 \pm 2.12</math></b>
BMVALS	$0.66 \pm 0.00$	$1.23 \pm 0.00$	$38.99 \pm 1.46$
RVS*	0.54	NA	$\approx 23400$
SED-MPK**	<b>0.48</b>	NA	$\approx 120$

Table 1. Performance comparison of BTN-V and existing methods in terms of RMSE, NLL, and computation time. **Bold** indicates the best result. \*RVS and \*\*SED-MPK values are taken from (Birpoutsoukis et al., 2018; Libera et al., 2021).

the estimated rank decreases to 2, and the column weights show a clear decaying trend across lags, as expected for Volterra kernels with fading memory. Unlike RVS and SED-MPK, which rely on fixed exponentially decaying priors, the proposed BTN-V model learns fading memory behavior directly from the data through  $\delta$  regularization. As shown in Table 1, applying row-wise regularization  $\delta$  also results in lower RMSE and NLL values, indicating improved predictive performance.

### 5.3 Predictive Performance

Table 1 presents a comparison between the proposed BTN-V and the three state-of-the-art methods on validation data: BMVALS, RVS and SED-MPK in terms of RMSE, NLL, and computation time. Since the tensor components of BTN-V and BMVALS are randomly initialized, results are reported as the mean  $\pm$  standard deviation over 10 runs to account for this randomness.

As mentioned before, RVS and SED-MPK use fixed exponentially decaying priors to model the fading-memory property of Volterra kernels, which improves their predictive accuracy and generally leads to lower RMSE. As shown in Table 1, BTN-V achieves a similar RMSE ( $0.51 \pm 0.02$ ) to SED-MPK (0.48) and better than RVS (0.54), while being much faster, requiring only 13.68 s for identification time compared to 120 s for SED-MPK and 23,400 s for RVS. Furthermore, we compare BTN-V with BMVALS, a probabilistic tensor-based Volterra kernel method. BTN-V has an average effective rank of  $2.5 \pm 0.5$  over 10 runs, while BMVALS uses a fixed rank of  $R = 48$ , making BTN-V a much simpler model. Despite this, BTN-V achieves a lower RMSE ( $0.51 \pm 0.02$  vs. 0.66) and yields smoother predictive curves, as shown in Figure 3. For uncertainty quantification, BTN-V also reports a lower NLL ( $0.77 \pm 0.05$  vs.  $1.23 \pm 0.00$ ). As illustrated in Figure 3, BTN-V provides tighter, more consistent uncertainty bounds, and is about three times faster than BMVALS (13.68 s vs. 38.99 s).

## 6. CONCLUSION

This paper presents a probabilistic method for nonlinear system identification using the TN SISO Volterra series. Using sparsity-inducing hierarchical priors, BTN-V automatically determines the effective tensor rank and fading-memory behavior from data. Experiments on the Cascaded Tanks Benchmark show that BTN-V offers competitive accuracy, improved uncertainty quantification, and reduced computational cost. Future work will focus on

extending BTN-V to MIMO systems. In this study, we initialized  $\Delta = \mathbf{I}$  as an uninformative prior and showed that  $\Delta$  learned the fading-memory behavior directly from data. We also plan to explore alternative prior structures, such as decaying priors, and investigate representing the Volterra coefficients using other tensor-network formats, such as tensor-train decomposition, instead of CPD.

## REFERENCES

- Batselier, K., Chen, Z., and Wong, N. (2017). Tensor network alternating linear scheme for mimo volterra system identification. *Automatica*, 84, 26–35.
- Birpoutsoukis, G., Csurscia, P.Z., and Schoukens, J. (2018). Efficient multidimensional regularization for volterra series estimation. *Mechanical Systems and Signal Processing*, 104, 896–914.
- Chen, T., Ohlsson, H., and Ljung, L. (2011). On the estimation of transfer functions, regularizations and gaussian processes – revisited. *IFAC Proceedings Volumes*, 44(1).
- Favier, G., Kibangou, A.Y., and Bouillot, T. (2012). Non-linear system modeling and identification using volterra-parafac models. *International Journal of Adaptive Control and Signal Processing*, 26(1), 30–53.
- Harshman, R.A. (1970). Foundations of the parafac procedure: Model and conditions for an “Explanatory” multi-mode factor analysis. Technical Report 16, UCLA Working Papers in Phonetics.
- Kilic, A. and Batselier, K. (2025). Interpretable bayesian tensor network kernel machines with automatic rank and feature selection. ArXiv:2507.11136.
- Kruskal, J.B. (1977). Three-way arrays: Rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics. *Linear Algebra and Its Applications*, 18(2), 95–138.
- Libera, A.D., Carli, R., and Pillionetto, G. (2021). Kernel-based methods for volterra series identification. *Automatica*, 129, 109686.
- Mommel, E., Menzen, C., and Batselier, K. (2023). Bayesian framework for a mimo volterra tensor network. *IFAC-PapersOnLine*, 56(2), 7294–7299.
- Miao, P., Qi, C., Jin, Y., Song, K., and Yu, T. (2019). Kernels pruning for volterra digital predistortion using sparse bayesian learning. In *2019 11th International Conference on Wireless Communications and Signal Processing (WCSP)*, 1–6.
- Neal, R.M. (1996). *Bayesian learning for neural networks*, volume 118 of *Lecture Notes in Statistics*. Springer, New York, NY.
- Pillionetto, G. and De Nicolao, G. (2010). A new kernel-based approach for linear system identification. *Automatica*, 46(1), 81–93.
- Schoukens, M., Mattsson, P., Wigren, T., and Noel, J.P. (2016). Cascaded tanks benchmark combining soft and hard nonlinearities. In *Proceedings of the Workshop on Nonlinear System Identification Benchmarks*.
- Winn, J.M. and Bishop, C.M. (2005). Variational message passing. *Journal of Machine Learning Research*, 6, 661–694.
- Zhao, Q., Zhang, L., and Cichocki, A. (2015). Bayesian cp factorization of incomplete tensors with automatic rank determination. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9), 1751–1763.